

C4.5

▶Decision Tree

Cannot-Link Constraint

A pairwise constraint between two items indicating that they should be placed into different clusters in the final partition.

Candidate-Elimination Algorithm

Mitchell's, (1982, 1997) candidate-elimination algorithm performs a bidirectional search in the ▶hypothesis space. It maintains a set, *S*, of most specific hypotheses that are consistent with the training data and a set, *G*, of most general hypotheses consistent with the training data. These two sets form two boundaries on the version space. See ▶Learning as Search.

Recommended Reading

Mitchell, T. M. (1982). Generalization as search. Artificial Intelligence, 18(2), 203–226.

Mitchell, T. M. (1997). Machine learning. New York: McGraw-Hill.

Cascade-Correlation

THOMAS R. SHULTZ¹, SCOTT E. FAHLMAN²

¹McGill University, Montréal, QC, Canada

²Carnegie Mellon University, Pittsburgh, PA, USA

Synonyms

Cascor: CC

Definition

Cascade-Correlation (often abbreviated as "Cascor" or "CC") is a ▶supervised learning algorithm for ▶artificial neural networks. It is related to the ▶back-propagation algorithm ("backprop"). CC differs from backprop in that a CC network begins with no hidden units, and then adds units one-by-one, as needed during learning.

Each new hidden unit is trained to correlate with residual error in the network built so far. When it is added to the network, the new unit is frozen, in the sense that its input weights are fixed. The hidden units form a *cascade*: each new unit receives weighted input from all the original network inputs and from the output of every previously created hidden unit. This cascading creates a network that is as deep as the number of hidden units. Stated another way, the CC algorithm is capable of efficiently creating complex, higher-order nonlinear basis functions – the hidden units – which are then combined to form the desired outputs.

The result is an algorithm that learns complex input/output mappings very fast compared to backprop, and that builds a multi-layer network structure that is customized for the problem at hand.

Motivation and Background

Cascade-Correlation was designed (Fahlman & Lebiere, 1990) to address two well-known problems with the popular back-propagation algorithm ("backprop"). First, a backprop user has to guess what network structure – the number of hidden layers and the number of units in each layer – would be best for a given learning problem. If the network is too small or too shallow, it won't solve the problem; if it is too large or too deep, training is very slow, and the network is prone to overfitting the training data. Because there is no reliable way to choose a good structure before training begins, most backprop users have to train many different structures before finding one that is well-matched to the task.

Second, even if a backprop user manages to choose a good network structure, training is generally very slow. That is particularly true in networks with many hidden units or with more than one hidden layer. One cause of slow learning in backprop is the use of a uniform learning-rate parameter for updating network weights. This problem was addressed with the Quickprop algorithm (Fahlman, 1988), an approximation to Newton's method that adapts the learning rate for each weight parameter depending on the first two derivatives of the local error surface. Quickprop improved learning speed, sometimes dramatically, but learning was still too slow in large or deep networks.

Another cause of slow learning in backprop is the "herd effect" (Fahlman & Lebiere, 1990). If the solution to a network problem requires, say, 30 hidden units, each of these units must be trained to do a different job - that is, to compute a different nonlinear basis function. Each hidden unit starts with a different and randomly chosen set of input weights; but if the units are all trained at once, they all see the same error signal. There is no central authority telling each unit to do a separate job, so they tend to drift toward the same part of parameter space, forming a herd that moves around together. Eventually, the units may drift apart and begin to differentiate, but there is nothing to compel this, so the process is slow and unreliable. Usually, in selecting an initial topology for a backprop net, it is necessary to include many extra hidden units to increase the odds that each job will be done by some unit.

CC addresses this problem by introducing and training hidden units one by one. Each hidden unit sees a strong, clear error gradient, not confused by the simultaneous movement of other hidden units. A new hidden unit can thus move quickly and decisively to a position in parameter space where it can perform a useful function, reducing the residual error. One by one, cascor-hidden units take up distinct jobs, instead of milling about together competing to do the same job.

Structure of Learning System

The Algorithm

The CC architecture is illustrated in Fig. 1. It begins with some inputs and one or more output units, but no hidden units. The numbers of inputs and outputs are dictated by the problem. As in backprop, the output

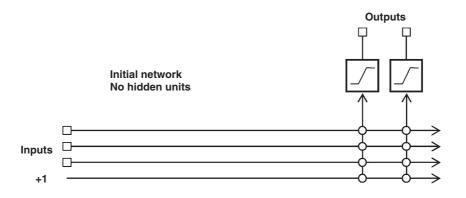
units generally have a sigmoid activation function, but could alternatively have a linear activation function. Every input is connected to every output unit by a connection with an adjustable weight. There is also a *bias* input, permanently set to +1.

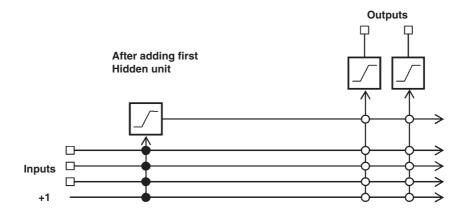
Hidden units are added to the network one by one. Each new hidden unit receives a weighted connection from each of the network's original inputs and also from every existing hidden unit. Each new unit therefore adds a new single-unit layer to the network. This makes it possible to create high-order nonlinear feature detectors, customized for the problem at hand.

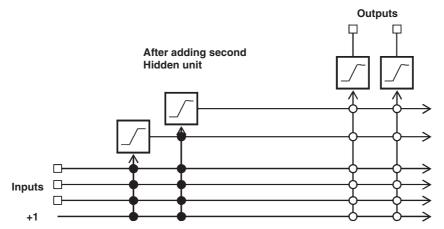
As noted, learning begins without hidden units. The direct input-output connections are trained as well as possible over the entire set of training examples, using Quickprop. At some point, this training approaches an asymptote. When no significant error reduction has occurred after a certain number of training cycles, this output phase is terminated and there is a shift to input phase to recruit a new hidden unit, using the unit-creation algorithm to be described. The new unit is added to the net, its input weights are frozen, and all the output weights are once again trained using Quickprop. This cycle repeats until the error is acceptably small, in the sense that all network outputs for all training patterns are within a specified threshold of their target values.

To create a new hidden unit, input phase begins with several *candidate units* that receive trainable input connections from all of the network inputs and from all existing hidden units. The outputs of these candidates are not yet connected to the network. There are a number of passes over the examples of the training set, adjusting the candidate unit's input weights after each pass. The goal of these adjustments, using Quickprop, is to maximize the correlation between each candidate's output and the residual error.

When these correlation measures show no further significant improvement, input phase stops, the best-correlating candidate's input weights are frozen, and that unit is installed in the network. The remaining candidates are discarded and the algorithm then retrains the output weights, making use of this new feature as well as all the old ones. As the new unit's output correlates well with some component of the residual error, its output weights can be quickly adjusted to reduce that







Cascade-Correlation. Figure 1. The Cascade-Correlation (CC) architecture, as new hidden units are added. Black circles are frozen connection weights, white circles are weights trained during output-training phase. The vertical lines sum all incoming activation

component. So after adding each new hidden unit, the network's residual error should be smaller than before.

Using several candidates, each with differently-initialized input weights, greatly reduces the chances of installing a bad hidden unit that gets the network stuck in a local optimum far from the global optimum value. All candidates receive the same input signals and see the same residual error for each training pattern. Because they do not interact with one another or affect the network during training, these candidates can be trained in parallel. In a pool of four to eight candidates, there are almost always several high-quality candidates with nearly equal correlation values.

Hidden units continue to be recruited until network error reaches an acceptable level, or until crossvalidation signals a stop. Because only a single layer of weights is adjusted at a time, rather than backpropagating an error signal through several layers of shifting units, CC training proceeds very quickly.

Performance

CC is designed to produce a network just large enough to solve the problem, and to do so much faster than backprop and related algorithms. In many reported cases that require hidden units, CC learns the desired behavior 10-100 times faster than standard backprop (Fahlman & Lebiere, 1990). One striking example of this is the two-spirals problem, an artificial benchmark designed to be very difficult for neural networks with sigmoid units. At the time CC was developed, the best known backprop solutions for two-spirals required a network with three hidden layers of five units each. CC typically solves this problem with 12 hidden units, and has found solutions with as few as nine hidden units. In terms of runtime, CC training was about 50 times faster than standard backprop and 23 times faster than Quickprop used within a static network.

Variants of Cascade-Correlation

Flat Cascade-Correlation In standard CC, each new hidden unit receives inputs from every existing unit, so the net becomes one level deeper every time a unit is added. This is a powerful mechanism, creating increasingly complex feature detectors as the network learns. But sometimes this added depth is not required for the problem, creating a very deep network that performs no better than a shallow one. The resulting network might

have more weights than are required for the problem, raising concern about over-fitting. Another concern was that the cascaded non-linearity of CC might also compromise generalization. To address these concerns, a flat variant of cascor adds new recruited units onto a single layer (i.e., cascaded connections are eliminated), limiting the depth of the network and eliminating all cascaded weights between hidden units.

Comparison of flat to standard CC on generalization in particular learning problems yielded inconsistent results, but a more problem-neutral, student-teacher approach found no generalization differences between flat and standard versions of CC (Dandurand, Berthiaume, & Shultz, 2007). Here, flat and standard student CC networks learned the inputoutput mappings of other, randomly initialized flat and standard CC teacher networks, where task complexity was systematically manipulated. Both standard and flat CC student networks learned and generalized well on problems of varying complexity. In low-complexity tasks, there were no significant performance differences between flat and standard CC student networks. For high-complexity tasks, flat CC student networks required fewer connection weights and learned with less computational cost than did standard CC student networks.

Sibling-Descendant Cascade-Correlation (SDCC) SDCC (Baluja & Fahlman, 1994) provides a more general solution to the problem of network depth. In the candidate pool there are two kinds of candidate units: descendant units that receive inputs from all existing hidden units, and sibling units that receive the same inputs as the deepest hidden units in the current net. When the time comes to choose a winning candidate, the candidate with the best correlation wins, but there is a slight preference for sibling units. So unless a descendant unit is clearly superior, a sibling unit will be recruited, making the active network larger, but not deeper. In problems where standard CC produces a network with 15 or 20 hidden units and an equal number of layers, SDCC often produces a network with only two or three hidden layers.

Recurrent Cascade-Correlation (RCC) Standard CC produces a network that maps its *current* inputs to outputs. The network has no memory of recent inputs, so this

architecture is not able to learn to recognize a sequence of inputs. In the RCC algorithm, each candidate and hidden unit takes the same inputs as in standard CC, but it also takes an additional input: the unit's own previous output, delayed by one time interval (Fahlman, 1991). The weight on this time-delayed input is trained by the same algorithm as all the other inputs.

This delayed loop gives RCC networks a way of remembering past inputs and internal states, so they can learn to recognize sequences of input patterns. In effect, the architecture builds a finite-state machine tailored specifically to recognize the pattern sequences in the training set. For example, an RCC net learned to recognize characters in Morse code.

Knowledge-Based Cascade-Correlation (KBCC) KBCC is a variant that can recruit previously-learned networks or indeed any differentiable function, in competition with single hidden units (Shultz & Rivest, 2001; Shultz, Rivest, Egri, Thivierge, & Dandurand, 2007). The recruit is the candidate whose output correlates best with residual network error, just as in ordinary CC. The candidate pool usually has a number of randomly initialized sigmoid units and a number of candidate source networks, i.e., networks previously trained on other tasks. The input weights to multiple copies of the source networks are usually randomly initialized to improve optimization. Of these copies, one is typically connected with an identity matrix with off-diagonal zeros, to enable quick learning of the target task when exact knowledge is available. A hypothetical KBCC network is shown in Fig. 2.

Software Most CC algorithms are available in a variety of formats and languages, including:

CASCOR: Lisp and C implementations of Cascadecorrelation

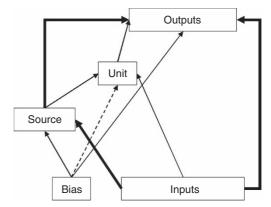
http://www.cs.cmu.edu/afs/cs/project/ai-repository/ai/areas/neural/systems/cascor/0.html

Free Lisp and C implementations of cascade-correlation.

Cascade Neural Network Simulator

http://www.cs.cmu.edu/~sef/sefSoft.htm

A public domain C program that implements various psychologicascade-correlation and recurrent cascade-correlation, particular domain.



Cascade-Correlation. Figure 2. Hypothetical knowledgebased cascade-correlation (KBCC) network that has recruited a source network and then a sigmoid unit, each installed on a separate layer. The dashed line represents a single connection weight, thin solid lines represent weight vectors, and thick solid lines represent weight matrices

plus experimental versions of cascade 2 and recurrent cascade 2.

LNSC Cascade-correlation Simulator Applet http://www.psych.mcgill.ca/perpg/fac/shultz/cdp/lnsc_applet.htm

A Java applet allowing direct comparisons of cascade-correlation and back-propagation algorithms on some benchmark problems, also permitting entry of text-edited custom training and test patterns.

LNSC Java Code Library

http://www.lnsclab.org/

Free compiled Java versions of BP, CC, SDCC, and KBCC neural-network software, along with a tutorial

Applications

C

Partly because of its ability to grow its own networks and build new learning on top of existing knowledge, CC has been used to simulate many phenomena in cognitive development. These characteristics embody the constructivism that developmental psychologists often discussed but did not formulate precisely. Simulations are typically evaluated by how well they capture the various psychological phenomena that characterize a particular domain.

The balance-scale task involves presenting a child with a rigid beam balanced on a fulcrum with pegs spaced at equal intervals to the left and right of the fulcrum. A number of identical weights are placed on a peg on the left side and a peg on the right side, and the child is asked to predict which side will descend when the beam is released from its moorings. CC networks passed through the stages observed with children and captured the so-called torque-difference effect, the tendency to do better on problems with large absolute torque differences than on problems with small torque differences (Shultz, Mareschal, & Schmidt, 1994; Shultz and Takane, 2007).

The conservation task presents a child with two quantities of objects that the child judges to be equal and then transforms one set in a way that either changes that relationship or conserves it. CC networks captured four important conservation regularities (Shultz, 1998):

- A shift from nonconservation to conservation beliefs
- 2. A sudden spurt in performance during acquisition
- 3. Emergence of correct conservation judgments for small quantities before larger quantities
- 4. Young children's choice of the longer row as having more items than the shorter row

Analysis of network solutions at various points in development revealed a gradual shift from perceptual (how the sets of items look) to cognitive (whether or not the transformation changed a quantity) solutions, similar to what had been found with children.

The seriation task requires a child to order a disordered collection of sticks of different lengths. CC networks passed through the four stages seen in children (total failure, partial sort, trial-and-error sort, and systematic sort) and captured the tendency for sets with smaller differences to be more difficult to sort (Mareschal & Shultz, 1999). Analysis of network solutions revealed early success at the short end of the series that was gradually extended to the longer end, as in children.

The transitivity problem typically also employs sticks of different length. Here the child is trained on all pairs of sticks that are adjacent in length and then is asked to infer the relative length of untrained pairs. Five psychological regularities were captured when CC networks were trained to compare the relative sizes of adjacent pairs (Shultz & Vogel, 2004):

- Learning short or long adjacent pairs before adjacent pairs of medium length.
- 2. Faster inferences with pairs farther apart in length than with pairs close together in length, an effect that diminished with age. A constraint-satisfaction network module simulated reaction times by inputting the output of a CC network and settling over time cycles into a low-energy solution that satisfied the constraints supplied by connection weights and inputs, effectively cleaning up the output of the CC network.
- Faster inferences with pairs containing the shortest or longest stick.
- 4. Faster inferences when the expression used in the question (e.g., shorter) is compatible with an end stick (e.g., the shortest stick) in the compared pair than when the question term (e.g., shorter) is incompatible with an end stick (e.g., the longest stick) in the compared pair.
- Older children learned adjacent pairs faster and made inference comparisons faster and more accurately than did young children.

The computational bases for these effects were revealed by examining the pattern of connection weights within the CC network module. The pattern of these weights formed a cubic shape, symmetrical for the two sticks being compared, in which discrimination was better at the ends of the array than in the middle and became sharper with deeper learning.

Another task calls for integration of cues for moving objects, governed by the equation velocity = distance/time. Children were presented with information on two of those quantities and asked to infer the third. Three stages involved first using the quantity that varied positively with the quantity to be inferred, second adding or subtracting the known quantities, and finally multiplying or dividing the known quantities. Already documented stages were captured and others were correctly predicted by CC networks (Buckingham & Shultz, 2000).

Semantic rules for deictic personal pronouns specify that *me* refers to the person using the pronoun and *you* refers to the person who is being addressed. Although

most children acquire these pronouns without notable errors, a few reverse these pronouns, persistently calling themselves *you* and the mother *me*. Such reversals in children are produced by lack of opportunity to overhear these pronouns used by other people, where the shifting reference can be observed. CC networks covered these phenomena and generated predictions for effective therapy to correct reversal errors (Oshima-Takane, Takane, & Shultz, 1999).

Discrimination shift learning tasks repeatedly present pairs of stimuli with mutually exclusive attributes on several binary dimensions, such as color, shape, and position, and a child learns to select the correct stimulus in each pair, e.g., square. Feedback is given and learning continues until the child reaches a success criterion, e.g., 8/10 correct. Then reward contingencies shift, usually without warning. A reversal shift stays within the initially relevant dimension, e.g., from square to circle. A nonreversal shift is to another dimension, such as from square to blue. There are related tasks that use new stimulus values in the shift phase. These are called intradimensional shifts if the shift remains within the initial dimension, e.g., square to triangle, or extradimensional if there is a change to another dimension, e.g., from square to yellow. The optional shift task presents only two stimulus pairs in the shift phase, making it ambiguous whether the shift is a reversal or nonreversal shift. The pattern of subsequent choices allows determination of whether the child interprets this as a reversal or a nonreversal shift.

Age differences in the large literature on these shifts indicate that older children learn a reversal shift faster than a nonreversal shift, learn an intradimensional shift faster than an extradimensional shift, make a reversal shift in the optional task, and are initially impaired on unchanged pairs during a nonreversal shift. Younger children learn reversal and nonreversal shifts equally fast, learn an intra-dimensional shift faster than an extra-dimensional shift, make a nonreversal shift in the optional task, and are unimpaired on unchanged pairs during a nonreversal shift. These findings were simulated by CC networks (Sirois & Shultz, 1998), which also generated predictions that were later confirmed.

When infants repeatedly experience stimuli from a particular class, their attention decreases, but it recovers to stimuli from a different class. This familiarize-andtest paradigm is responsible for most of the discoveries of infant psychological abilities. CC networks simulated findings on infant attention to syntactic patterns in an artificial language (Shultz & Bale, 2006) and age differences in infant categorization of visual stimuli (Shultz & Cohen, 2004), and generated several predictions, some of which were tested and confirmed.

SDCC

Because of SDCC's ability to create a variety of network topologies, it is beginning to be used in psychology simulations: infant learning of word-stress patterns in artificial languages (Shultz & Bale, 2006), syllable boundaries (Shultz & Bale, 2006), visual concepts (Shultz, 2006), and false-belief tasks; learning the structure of mathematical groups (Schlimm & Shultz, 2009); replication of the results of the CC simulation of conservation acquisition (Shultz, 2006); and concept acquisition.

CC and SDCC networks capture developmental stages by growing in computational power and by being sensitive to statistical patterns in the training environment (Shultz, 2003). The importance of growth was demonstrated by comparisons with static backprop networks, designed with the same final topology as successful CC networks, that learn only by adjusting connection weights (Shultz, 2006). Coupled with the variety of successful SDCC topologies, this suggests that the constructive process is more important than precise network topologies. Capturing stages is challenging because the system has to not only succeed on the task but also make the same mistakes on the road to success that children do. CC and SDCC arguably produced the best data coverage of any models applied to the foregoing phenomena. Both static and constructive networks capture various perceptual effects by virtue of their sensitivity to quantitative variation in stimulus inputs (Shultz, 2003).

Comparison of the two algorithms in psychological modeling indicates that SDCC provides the same functionality as CC but with fewer connection weights and shallower and more variable network topologies (Shultz, 2006).

KBCC

KBCC also has potential for simulating psychological development, but it has so far been applied mainly to toy and engineering problems. Exploration of a C

variety of toy problems was important in understanding the behavior of this complex algorithm. Some toy problems involved learning about two-dimensional geometric shapes under various transformations such as translation, rotation, and size changes, as well as compositions of complex shapes from simpler shapes (Shultz & Rivest, 2001). Networks had to learn to distinguish points within a target shape from points outside the shape. Learning time without relevant knowledge was up to 16 times longer than with relevant knowledge on these problems. There was a strong tendency to recruit relevant knowledge whenever it was available. Direct comparison revealed that KBCC learned spatial translation problems faster than Multitask Learning networks did.

Parity problems require a network to activate an output unit only when an odd number of binary inputs are activated. When parity-4 networks were included in the candidate source pool, KBCC learned parity-8 problems (with eight binary inputs) faster and with fewer recruits than did CC networks. Parity-4 networks were recruited by these KBCC target networks whenever available.

KBCC also learned complex chessboard shapes from knowledge of simpler chessboards. As with parity, networks used simpler previous knowledge to compose a solution to a more complex problem and learning was speeded accordingly.

In a more realistic vein, KBCC networks recruiting knowledge of vowels from one sort of speaker (e.g., adult females) learned to recognize vowels spoken by other sets of speakers (e.g., children and adult males) faster than did knowledge-free networks.

KBCC learned an efficient algorithm for detecting prime numbers by recruiting previously-learned knowledge of divisibility (Shultz et al., 2007). This well-known detection algorithm tests the primality of an integer n by checking if n is divisible by any integers between 2 and the integer part of \sqrt{n} . Starting with small primes is efficient because the smaller the prime divisor, the more composites are detected in a fixed range of integers. The candidate pool contained networks that had learned whether an integer could be divided by each of a range of integers, e.g., a divide-by-2 network, a divide-by-3 network, etc., up to a divisor of 20. KBCC target networks trained on 306 randomly-selected integers from 21 to 360 recruited only source

networks involving prime divisors below the square root of 360, in order from small to large divisors. KBCC avoided recruiting single hidden units, source networks with composite divisors, any divisors greater than the square root of 360 even if prime, and divisor networks with randomized connection weights. KBCC never recruited a divide-by-2 source network because it instead learned to check the last binary digit of n to determine if *n* was odd or even, an effective shortcut to dividing by 2. Such KBCC networks learned the training patterns in about one third the time required by knowledge-free networks, with fewer recruits on fewer network layers, and they generalized almost perfectly to novel test integers. In contrast, even after mastering the training patterns, CC networks generalized less well than automatic guessing that the integer was composite, which was true for 81% of integers in this range. As predicted by the simulation, adults testing primality also used mainly prime divisors below \sqrt{n} and ordered divisors from small to large.

This work underscores the possibility of neuralnetwork approaches to compositionality because KBCC effectively composed a solution to prime-number detection by recruiting and organizing previously learned parts of the problem, in the form of divisibility networks.

Future Directions

One new trend is to inject symbolic rules or functions into KBCC source networks. This is similar to KBANN, but more flexible because a KBCC target network decides whether and how to recruit these functions. This provides one method of integrating symbolic and neural computation and allows for simulation of the effects of direct instruction.

Cross References

- ► Artificial Neural Networks
- **▶**Backpropagation

Recommended Reading

Baluja, S., & Fahlman, S. E. (1994). Reducing network depth in the cascade-correlation learning architecture. Pittsburgh, PA: School of Computer Science, Carnegie Mellon University.

Buckingham, D., & Shultz, T. R. (2000). The developmental course of distance, time, and velocity concepts: A generative connectionist model. *Journal of Cognition and Development*, 1, 305-345.

Dandurand, F., Berthiaume, V., & Shultz, T. R. (2007). A systematic comparison of flat and standard cascade-correlation using a student-teacher network approximation task. *Connection Sci*ence, 19, 223–244.

Fahlman, S. E. (1988). Faster-learning variations on back-propagation: An empirical study. In D. S. Touretzky, G. E. Hinton, & T. J. Sejnowski (Eds.), Proceedings of the 1988 connectionist models summer school (pp. 38-51). Los Altos, CA: Morgan Kaufmann.

- Fahlman, S. E. (1991). The recurrent cascade-correlation architecture. In D. S. Touretzky (Ed.), *Advances in neural information processing systems*. (Vol. 3) Los Altos CA: Morgan Kaufmann.
- Fahlman, S. E., & Lebiere, C. (1990). The cascade-correlation learning architecture. In D. S. Touretzky (Ed.), Advances in neural information processing systems (Vol. 2, pp. 524–532). Los Altos, CA: Morgan Kaufmann.
- Mareschal, D., & Shultz, T. R. (1999). Development of children's seriation: A connectionist approach. Connection Science, 11, 149–186
- Oshima-Takane, Y., Takane, Y., & Shultz, T. R. (1999). The learning of first and second pronouns in English: Network models and analysis. *Journal of Child Language*, 26, 545-575.
- Schlimm, D., & Shultz, T. R. (2009). Learning the structure of abstract groups. In N. A. Taatgen & H. V. Rijn (Eds.), Proceedings of the 31st annual conference of the cognitive science society (pp. 2950-2955). Austin, TX: Cognitive Science Society.
- Shultz, T. R. (1998). A computational analysis of conservation. Developmental Science, 1, 103–126.
- Shultz, T. R. (2003). Computational developmental psychology. Cambridge, MA: MIT Press.
- Shultz, T. R. (2006). Constructive learning in the modeling of psychological development. In Y. Munakata & M. H. Johnson (Eds.), Processes of change in brain and cognitive development: Attention and performance XXI (pp. 61–86). Oxford, UK: Oxford University Press.
- Shultz, T. R., & Bale, A. C. (2006). Neural networks discover a nearidentity relation to distinguish simple syntactic forms. *Minds* and Machines, 16, 107–139.
- Shultz, T. R., & Cohen, L. B. (2004). Modeling age differences in infant category learning. *Infancy*, 5, 153-171.
- Shultz, T. R., Mareschal, D., & Schmidt, W. C. (1994). Modeling cognitive development on balance scale phenomena. *Machine Learning*, 16, 57–86.
- Shultz, T. R., & Rivest, F. (2001). Knowledge-based cascade-correlation: Using knowledge to speed learning. Connection Science, 13, 1-30.
- Shultz, T. R., Rivest, F., Egri, L., Thivierge, J.-P., & Dandurand, F. (2007). Could knowledge-based neural learning be useful in developmental robotics? The case of KBCC. *International Journal of Humanoid Robotics*, 4, 245–279.
- Shultz, T. R., & Takane, Y. (2007). Rule following and rule use in simulations of the balance-scale task. *Cognition*, 103, 460-472.
- Shultz, T. R., & Vogel, A. (2004). A connectionist model of the development of transitivity. In Proceedings of the twenty-sixth annual conference of the cognitive science society (pp. 1243–1248). Mahwah, NJ: Erlbaum.
- Sirois, S., & Shultz, T. R. (1998). Neural network modeling of developmental effects in discrimination shifts. *Journal of Experimental Child Psychology*, 71, 235–274.

CART

▶Decision Tree

Cascor

▶ Cascade-Correlation

Case

▶Instance

Case-Based Learning

►Instance-Based Learning

Case-Based Reasoning

Susan Craw

The Robert Gordon University, Scotland, UK

Synonyms

CBR; Experience-based reasoning; Lessons-learned systems; Memory-based learning

Definition

Case-based reasoning solves problems by retrieving similar, previously solved problems and reusing their solutions. Experiences are memorized as cases in a case base. Each experience is learned as a problem or situation together with its corresponding solution or action. The experience need not record *how* the solution was reached, simply that the solution was used for the problem. The case base acts as a memory, and remembering is achieved using similarity-based retrieval and reuse of the retrieved solutions. Newly solved problems may be retained in the case base and so the memory is able to grow as problem-solving occurs.

CII

Motivation and Background

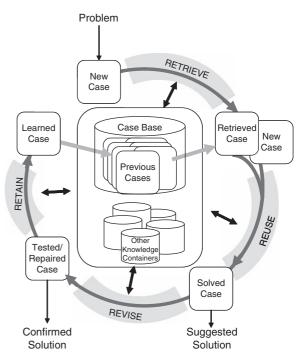
Case-based reasoning (CBR) is inspired by memorybased human problem-solving in which instances of earlier problem-solving are remembered and applied to solve new problems. For example, in Case Law, the decisions in trials are based on legal precedents from previous trials. In this way specific experiences are memorized, and remembered and reused when appropriate (Leake, 1996). This contrasts with rule-based or theory-based problem-solving in which knowledge of how to solve a problem is applied. A doctor diagnosing a patient's symptoms may apply knowledge about how diseases manifest themselves, or she may remember a previous patient who demonstrated similar symptoms. Schank's ▶dynamic memory model was highly influential in early CBR systems (Kolodner, 1993; Riesbeck & Schank, 1989). Its emphasis on the use of specific experiences to underpin problem-solving and enable learning is replicated in CBR.

The fundamental assumption of CBR is that Similar problems have similar solutions. For example, a patient with similar symptoms will have the same diagnosis, the price of a house with similar accommodation and location will be similar, the design for a kitchen with a similar shape and size can be reused, a journey plan for a nearby destination is similar to the earlier trip. A related assumption is that the world is a regular place, and what holds true today will probably hold true tomorrow. A further assumption relevant to memory is that situations repeat, because if they do not, there is no point in remembering them!

CBR is an example of Lazy Learning because there is no learned model to apply to solve new problems. Instead, the generalization needed to solve unseen problems happens when a new problem is presented and the similarity-based retrieval identifies relevant previous experiences. The lack of a learned model and the reliance on stored experiences mean that CBR is particularly relevant in domains which are ill-defined, not well understood, or where no underlying theory is available.

Structure of the Learning System

Figure 1 shows the structure of a CBR system (Aamodt & Plaza, 1994). A case base of Previous Cases is the primary knowledge source in a CBR system, with



Case-Based Reasoning. Figure 1. CBR system (adapted from Aamodt and Plaza (1994))

additional knowledge being used to identify similar cases in the RETRIEVE stage, and to REUSE and REVISE the Retrieved Case. A CBR system learns as it solves new problems when a Learned Case is created from the New Case and its Confirmed Solution, and RETAINed as a new case in the case base.

Knowledge Containers

Case knowledge is the primary source of knowledge in a CBR system. However, case knowledge is only one of four knowledge containers identified by Richter and Aamodt (2005).

- Vocabulary: The representation language used to describe the cases captures the concepts involved in the problem-solving.
- Similarity Knowledge: The similarity measure defines how the distances between cases are computed so that the nearest neighbors are identified for retrieval.
- Adaptation Knowledge: Reusing solutions from retrieved cases may require some adaptation to enable them to fit the new problem.
- Case Base: The stored cases capture the previous problem-solving experiences.

The content of each knowledge container is not fixed and knowledge in one container can compensate for a lack of knowledge in another. It is easy to see that a more sophisticated knowledge representation could be less demanding on the content of the case base. Similarly, vocabulary can make similarity assessment during retrieval easier, or a more complete case base could reduce the demands on adaptation during reuse. Further knowledge containers (e.g., maintenance) are proposed by others.

Cases may be represented as simple feature vectors containing nominal or numeric values. A case capturing a whisky tasting experience might contain features such as sweetness, color, nose, palate, and the ▶classification as the distillery where it was made.

Sweetness	Peatiness	Color	Nose	Palate	Distillery
6	5	amber	full	medium-dry	Dalmore

More structured representations can use frame-based or object-oriented cases. The choice of representation depends on the complexity of the experiences being remembered and is influenced by how similarity should be determined. Hierarchical case representations allow cases to be remembered at different levels of abstraction, and retrieval and reuse may occur at these different levels (Bergmann & Wilke, 1996).

In ▶ classification, the case base can be considered to contain exemplars of problem-solving. Aha et al.'s (1991) family of instance-based learning algorithms IB1, IB2 and IB3 apply increasingly informed selection methods to determine whether a classification experience should become part of the case base. IB1 simply remembers all experiences, IB2 stores an experience only if it would be wrongly solved by the existing stored experiences, and IB3 keeps a score for the reuse of each experience and discards those whose classification success is poor. This notion of exemplar confirms a CBR case base as a source of knowledge; it contains only those experiences that are believed to be useful for problem-solving. A similar view is taken for non-classification domains; for example, the case base contains useful designs that can be used for re-design, and plans for re-planning.

One of the advantages of CBR is that a case base is composed of independent cases that each captures some local problem-solving knowledge that has been experienced. Therefore, the "knowledge acquisition bottleneck" of many rule-based and model-based systems is reduced for CBR. However, the Other Knowledge Containers provide additional knowledge acquisition demands that may lessen the advantage of CBR for some domains.

CBR Cycle

Aamodt and Plaza (1994) propose a four-stage CBR cycle for problem-solving and learning (Fig. 1). It is commonly referred to as the "Four REs" or "R⁴" cycle to recognize the following stages.

- RETRIEVE: The cases that are most similar to the New Case defined by the description of the new problem are identified and retrieved from the case base. The RETRIEVE stage uses the similarity knowledge container in addition to the case base.
- REUSE: The solutions in the Retrieved (most similar) Cases are reused to build a Suggested Solution to create the Solved Case from the New Case. The REUSE stage may use the adaptation knowledge container to refine the retrieved solutions.
- REVISE: The Suggested Solution in the Solved Case is evaluated for correctness and is repaired if necessary to provide the Confirmed Solution in the Tested/Repaired Case. The REVISE stage may be achieved manually or may use adaptation knowledge, but it should be noted that a revision to a Suggested Solution is likely to be a less demanding task than solving the problem from scratch.
- RETAIN: The Repaired Case may be retained in the
 case base as a newly Learned Case if it is likely to be
 useful for future problem-solving. Thus the primary
 knowledge source for CBR may be added to during problem-solving and is an evolving, self-adaptive
 collection of problem-solving experiences.

Retrieval

CBR retrieval compares the problem part of the new case with each of the cases in the case base to establish the distance between the new case and the stored cases. The cases closest to the new case are retrieved for reuse. Retrieval is a major focus of López de Mántaras et al.'s (2005) review of research contributions related to the CBR cycle.

Similarity- and distance-based neighborhoods are commonly used interchangeably when discussing CBR

retrieval. Similarity and distance are inverses: the similarity is highest when the distance is close to 0, and the similarity is 0 when the distance is large. Several functions may be applied to define a suitable relationship between a distance *d* and a similarity *s*, including the following simple versions:

Linear:
$$s = 1 - d$$
 (for normalized d)
Inverse: $s = \frac{1}{d}$ (for $d \neq 0$).

It is common to establish the distance between each pair of feature values and then to use a distance metric, often Euclidean or \blacktriangleright Manhattan distance (see also \blacktriangleright Similarity Measure), to calculate the distance between the feature vectors for the New and Retrieved Cases. The distance between two numeric feature values v and w for a feature F is normally taken to be the distance between the normalized values:

$$d(v,w) = \frac{|v-w|}{F_{max} - F_{min}}$$

where F_{max}/F_{min} are the maximum/minimum values of the feature F.

For nominal values *v* and *w* the simplest approach is to apply a binary distance function:

$$d(v, w) = \begin{cases} 0 & \text{if } v = w \\ 1 & \text{otherwise} \end{cases}$$

For ordered nominal values a more fine-grained distance may be appropriate. The distance between the *i*th value v_i and the *j*th value v_j in the ordered values v_1, v_2, \ldots, v_n may use the separation in the ordering to define the distance:

$$d(v_i,v_j)=\frac{\mid i-j\mid}{n-1}.$$

Extending this to arbitrary nominal values, a distance matrix D may define the distance between each pair of nominal values by assigning the distance $d(v_i, v_j)$ to d_{ii} .

Returning to the whisky-tasting example, suppose Sweetness and Peatiness score values 0–10, Color takes ordered values {pale, straw, gold, honey, amber}, Palate uses binary distance, and Nose is defined by the following distance matrix.

Nose Distance Matrix					
Distances	peat	fresh soft		full	
peat	0	0.3	1	0.5	
fresh	0.3	0	0.5	0.7	
soft	1	0.5	0	0.3	
full	0.5	0.7	0.3	0	

Dalmore whisky above can be compared with Laphroaig and The Macallan as follows:

Sweetness	Peatiness	Color	Nose	Palate	Distillery
2	10	amber	peat	medium-dry	Laphroaig
7	4	gold	full	big-body	The Macallan

The Manhattan distances are:

$$d(Dalmore,Laphroaig) = 0.4 + 0.5 + 0 + 0.5 + 0 = 1.4;$$

$$d(Dalmore, The Macallan) = 0.1 + 0.1 + 0.5 + 0 + 1 = 1.7.$$

Taking all the whisky features with equal importance, Dalmore is more similar to Laphroaig than to The Macallan.

In situations where the relative importance of features should be taken into account, a weighted version of the distance function should be used; for example, the weighted Manhattan distance between two normalized vectors $\mathbf{x} = (x_1, x_2, \dots x_n)$ and $\mathbf{y} = (y_1, y_2, \dots y_n)$ with weight w_i for the ith feature is

$$d(\mathbf{x},\mathbf{y}) = \frac{\sum_{i=1}^{n} w_i \mid x_i - y_i \mid}{\sum_{i=1}^{n} w_i}$$

In the example above if Peatiness has weight 4 and the other features have weight 1 then the weighted Manhattan distances are:

$$d(Dalmore,Laphroaig) = (0.4 + 4 \times 0.5 + 0 + 0.5 + 0)/8 = 0.36;$$

 $d(Dalmore,The Macallan) = (0.1 + 4 \times 0.1 + 0.5 + 0 + 1)/8 = 0.25.$

Therefore, emphasizing the distinctive Peatiness feature, Dalmore is more similar to The Macallan than to Laphroaig.

The similarity knowledge container contains knowledge to calculate similarities. For simple feature vectors a weighted sum of distances is often sufficient, and the weights are similarity knowledge. However, even our whisky tasting domain had additional similarity knowledge containing the distance matrix for the Nose feature. Structured cases require methods to calculate the similarity of two cases from the similarities of components. CBR may use very knowledge-intensive methods to decide similarity for the retrieval stage. Ease of reuse or revision may even be incorporated as part of the assessment of similarity. Similarity knowledge may also define how ▶missing values are handled: the feature may be ignored, the similarity may be maximally pessimistic, or a default or average value may be used to calculate the distance.

A CBR case base may be indexed to avoid similarity matching being applied to all the cases in the case base. One approach uses kd trees to partition the case base according to hyper-planes. Decision Tree algorithms may be used to build the kd tree by using the cases as training data, partitioning the cases according to the chosen decision nodes, and storing the cases in the appropriate leaf nodes. Retrieval first traverses the decision tree to select the cases in a leaf node, and similarity matching is applied to only this partition. Case Retrieval Nets are designed to speed up retrieval by applying spreading activation to select relevant cases. In Case Retrieval Nets the feature value nodes are linked via similarity to each other and to cases. Indexes can speed up retrieval but they also pre-select cases according to some criteria that may differ from similarity.

Reuse and Revision

Reuse may be as simple as copying the solution from the Retrieved Case. If *k* nearest neighbors are retrieved then a vote of the classes predicted in the retrieved cases may be used for **>** classification, or the average of retrieved values for **>** regression. A weighted vote or weighted average of the retrieved solutions can take account of the nearness of the retrieved cases in the calculation. For more complex solutions, such as designs or plans, the amalgamation of the solutions from the Retrieved Cases may be more knowledge intensive.

If the New Case and the Retrieved Case are different in a significant way then it may be that the solution from the Retrieved Case should be adapted before

being proposed as a Suggested Solution. Adaptation is designed to recognize significant differences between the New and Retrieved Cases and to take account of these by adapting the solution in the Retrieved Case.

In classification domains, it is likely that all classes are represented in the case base. However, different problem features may alter the classification and so adaptation may correct for a lack of cases. In constructive problem-solving like design and planning, however, it is unlikely that all solutions (designs, plans, etc.) will be represented in the case base. Therefore, a retrieved case suggests an initial design or plan, and adaptation alters it to reflect novel feature values.

There are three main types of adaptation that may be used as part of the reuse step, to refine the solution in the Retrieved Case to match better the new problem, or as part of the revise stage to repair the Suggested Solution in the Solved Case.

- Substitution: Replace parts of the retrieved solution. In Hammond's (1990) CHEF system to plan Szechuan recipes, the substitution of ingredients enables the requirements of the new menu to be achieved. For example, the beef and broccoli in a retrieved recipe is substituted with chicken and snowpeas.
- Transformation: Add, change, or remove parts of the retrieved solution. CHEF adds a skinning step to the retrieved recipe that is needed for chicken but not for beef.
- Generative Adaptation: Replay the method used to derive the retrieved solution. Thus the retrieved solution is not adapted but a new solution is generated from reusing the retrieved method for the new circumstances. This approach is similar to reasoning by analogy.

CHEF also had a clear REVISE stage where the Suggested Solution recipe was tested in simulation and any faults were identified, explained, and repaired. In one recipe a strawberry soufflé was too liquid. CHEF has a set of Thematic Organization Packets (TOPs) that are templates for repairs for different types of explained failures. (TOPs continue the experience template theme of bdynamic memory model MOPs.) One repair for the soufflé is to drain the strawberry pulp and this

transformation adaptation is one REVISE operation that could be applied.

The adaptation knowledge container is an important source of knowledge for some CBR systems, particularly for design and planning, where refining an initial design or plan is expected. Acquiring adaptation knowledge can be onerous. The CHEF examples above indicate that the knowledge must store refinements to the solutions initially proposed from the retrieved cases. Learning adaptation knowledge from the implicit refinement information captured in the case base has been effective for substitution adaptation in component-based design (Craw, Wiratunga, & Rowe, 2006).

Retention and Maintenance

Retention of new cases during problem-solving is an important advantage of CBR systems. However, it is not always advantageous to retain all new cases. The ► Utility Problem – that the computational benefit from additional knowledge must not outweigh the cost of applying it - in CBR refers to cases and the added cost of retrieval. The case base must be kept "lean and mean," and so new cases are not retained automatically, and cases that are no longer useful are removed. New cases should be retained if they add to the competence of the CBR system by providing problem-solving capability in an area of the problem space that is currently sparse. Conversely, existing cases should be reviewed for the role they play and forgetting cases is an important maintenance task. Existing cases may contain outdated experiences and so should be removed, or they may be superseded by new cases.

Case base maintenance manages the contents of the case base to achieve high competence. Competence depends on the domain and may involve

- quality of solution;
- user confidence in solution; or
- efficiency of solution prediction (e.g., speed-up learning).

Case base maintenance systems commonly assume that the case base contains a representative sample of the problem-solving experiences. They exploit this by using a leave-one-out approach where repeatedly for each case in the case base, the one extracted case is used as a new case to be solved, and the remaining cases become

the case base. This enables the problem-solving competence of the cases in the case base to be estimated using the extracted cases as representative new cases to be solved. Smyth & McKenna's (2001) competence model uses this approach to identify competence groups of cases with similar problem-solving ability. This model is used to underpin maintenance algorithms to prioritize cases for deletion and to identify areas where new cases might be added. There are several trade-offs to be managed by case base maintenance algorithms: larger case bases contain more experiences but take longer for retrieval; smaller case bases are likely to lack some key problem-solving ability; cases whose solution is markedly different from their nearest neighbors may be noisy or may be an important outlier.

CBR Tools

myCBR (mycbr-project.net) and jCOLIBRI (www.sourceforge.net) are open source CBR tools. Both provide state-of-the-art CBR functionality, and jColibri also incorporates a range of facilities for textual CBR. Several commercial CBR tools are available including Empolis: Information Access Suite (www.attensity.com), Kaidara's Advisor (www.servigistics.com), and ISoft's ReCall (www.alice-soft.com).

Applications

Several successful deployed applications of CBR are described in Cheetham and Watson (2005), including Lockheed's CLAVIER for designing layouts for autoclave ovens, Compaq's SMART help-desk system, Boeing's CASSIOPÉE for trouble-shooting aircraft engines and General Electric's FormTool for plastic color matching (Cheetham, 2005). The development of the INRECA methodology for engineering CBR systems was based on a range of industrial applications (Bergmann et al., 2003).

The wide range of CBR applications is demonstrated by the following list of application types.

 Classification – Medical diagnosis systems use patient records as a source of reusable experiences. Examples include SHRINK for psychiatry, CASEY for cardiac disease, ICONS for antibiotic therapy for intensive care, and BOLERO for pneumonia. Other diagnostic systems include failure prediction of rails

for Dutch railways from ultrasonic NDT, and failure analysis of semiconductors at National Semiconductor. Classification systems include PROTOS for audiologic disorders.

- Design Architectural design was a popular early domain: CYCLOPS, ARCHIE, FABEL, and CADsyn are all early case-based design systems. Other design applications include CADET and KRITIK for engineering design, JULIA for recipes, chemical formulation for tyre rubber and pharmaceutical products, Déjà Vu for plant control software.
- Planning PRODIGY is a general purpose planner that uses analogical reasoning to adapt retrieved plans. Other planning applications include PARIS (Bergmann & Wilke, 1996) for manufacturing planning in mechanical engineering, HICAP for evacuation planning, planning for forest fire management, mission planning for US navy, and route planning for DaimlerChrysler cars.
- Conversational CBR Conversational systems extract the problem specification from the user through an interactive case-based dialogue and suggest solutions that match the partial specification. Examples include help-desk support, CaseAdvisor and CBR Strategist for fault diagnosis, and Wasabi, Sermo and ShowMe product recommender systems.
- Personalization Personalized compilations of news, stories, music tracks, TV listings reuse previous experiences of the individual or others who have similar tastes. Other forms of personalized systems using CBR include route and travel planning, SPAM filtering and email management, and ClixSmart Navigator for mobile devices.
- Textual CBR Legal decision support systems were an important early application domain for textual CBR. Examples include HYPO (Ashley & Rissland, 1988), CATO, GREBE, and SMILE. Question answering was another fruitful text-based domain: FAQ-Finder and FA11Q. More recently, textual CBR is used for decision support systems based on textual reports; for example, SOPHIA.

Future Directions

The drivers for Ubiquitous Computing – wireless communication and small devices – also affect future developments in CBR. The local, independent knowledge of

case bases make them ideal to collect experiences, and to deliver experience-based knowledge for reuse.

Textual CBR systems are becoming increasingly important for extracting and representing knowledge captured in textual documents. This is particularly influenced by the availability of electronic documents and the Web as an information source.

Cross References

- ► Explanation-Based Learning
- ►Instance-Based Learning
- ► Lazy Learning
- ► Nearest Neighbor
- ► Similarity Metrics

Recommended Reading

- Aamodt, A., & Plaza, E. (1994). Case-based reasoning: Foundational issues, methodological variations, and system approaches. *AI Communications*, 7, 39–59. citeseerx. ist.psu.edu/viewdoc/summary?doi=10.1.1.39.1670.
- Aha, D. W., Kibler, D., & Albert, M. K. (1991). Instance-based learning algorithms. *Machine Learning*, 6(1), 37–66. doi: 10.1023/A:1022689900470.
- Ashley, K., & Rissland, E. L. (1988). A case-based approach to modeling legal expertise. *IEEE Expert*, 3(3), 70-77. doi: 10.1109/64.21892.
- Bergmann, R., Althoff, K.-D., Breen, S., Göker, M., Manago, M., & Traphöner, R. (Eds.). (2003). Developing industrial case-based reasoning applications. LNCS (Vol. 1612). Berlin/Heidelberg: Springer. doi:10.1007/b94998.
- Bergmann, R., & Wilke, W. (1996). On the role of abstraction in case-based reasoning. In I. Smith & B. Faltings (Eds.), Proceedings of the third European workshop on case-based reasoning, Lausanne, Switzerland (pp. 28–43), LNCS (Vol. 1168). Berlin/Heidelberg: Springer.
- Cheetham, W. (2005). Tenth anniversary of the plastics color formulation tool. *AI Magazine*, 26(3), 51–61. www.aaai.org/Papers/Magazine/Vol26/26-03/AIMag26-03-007.pdf.
- Cheetham, W., & Watson, I. (2005). Fielded applications of case-based reasoning. *Knowledge Engineering Review*, 20(3), 321–323. doi:10.1017/S0269888906000580.
- Craw, S., Wiratunga, N., & Rowe, R. C. (2006). Learning adaptation knowledge to improve case-based reasoning. *Artificial Intelligence*, 170(16–17), 1175–1192. doi: 10.1016/j.artint.2006.09.001.
- Hammond, K. J. (1990). Explaining and repairing plans that fail. Artificial Intelligence, 45(1-2), 173-228.
- Kolodner, J. L. (1993). Case-based reasoning. San Mateo, CA: Morgan Kaufmann.
- Leake, D. (1996). CBR in context: The present and future. In D. Leake (Ed.), Case-based reasoning: Experiences, lessons, and future directions (pp. 3-30). Menlo Park, CA: AAAI Press. citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.35.9114.
- López de Mántaras, R., McSherry, D., Bridge, D., Leake, D., Smyth, B., Craw, S., Faltings, B., Maher, M.L., Cox, M.T., Forbus, K., Aamodt, A., & Watson, I. (2005). Retrieval, reuse, revision,

C

154 Categorical Attribute

and retention in case-based reasoning. Knowledge Engineering Review, 20(3), 215–240. doi: 10.1017/S0269888906000646.

Richter, M. M., & Aamodt, A. (2005). Case-based reasoning foundations. Knowledge Engineering Review, 20(3), 203–207. doi:10.1017/S0269888906000695.

Riesbeck, C. K., & Schank, R. C. (1989). *Inside case-based reasoning*. Hillsdale, NJ: Lawrence Erlbaum.

Smyth, B., & McKenna, E. (2001). Competence models and the maintenance problem. *Computational Intelligence*, 17(2), 235–249. doi:10.1111/0824-7935.00142.

Categorical Attribute

Synonyms

Qualitative attribute

Categorical attributes are attributes whose values can be placed into distinct categories. See ▶Attribute and ▶Measurement Scales.

Categorical Data Clustering

Periklis Andritsos¹, Panayiotis Tsaparas²
¹Toronto, ON, Canada
²Mountain View, CA, USA

Synonyms

Clustering of nonnumerical data; Grouping

Definition

Data clustering is informally defined as the problem of partitioning a set of objects into groups, such that the objects in the same group are similar, while the objects in different groups are dissimilar. Categorical data clustering refers to the case where the data objects are defined over **categorical attributes**. A categorical attribute is an attribute whose domain is a set of discrete values that are not inherently comparable. That is, there is no single ordering or inherent distance function for the categorical values, and there is no mapping from categorical to numerical values that is semantically meaningful.

Motivation and Background

Clustering is a problem of great practical importance that has been the focus of substantial research in several domains for decades. As storage capacities grow, we have at hand larger amounts of data available for analysis and mining. Clustering plays an instrumental role in this process. This trend has created a surge of research activity in devising new clustering algorithms that can handle large amounts of data and produce results of high quality.

In data clustering, we want to partition objects into groups such that similar objects are grouped together while dissimilar objects are grouped separately. This objective assumes that there is some well-defined notion of similarity, or distance, between data objects, and a way to decide if a group of objects is a homogeneous cluster. Most of the clustering algorithms in the literature focus on data sets where objects are defined over numerical attributes. In such cases, the similarity (or dissimilarity) of objects and the quality of a cluster can be defined using well-studied measures that are derived from geometric analogies. These definitions rely on the semantics of the data values themselves. For example, the values \$100,000 and \$110,000 are more similar than \$100,000 and \$50,000, and intuitively more similar than \$10,000 and \$1. The existence of a distance measure allows us to define a quality measure for a clustering such as the mean square distance between each point and the representative of its cluster. Clustering then becomes the problem of grouping together points such that the quality measure is optimized.

However, there are many data sets where the data objects are defined over attributes, which are neither numerical nor inherently comparable in any way. We term such data sets *categorical*, since they represent values of certain categories. As a concrete example, consider the toy data set in Table 1 that stores information about movies. For the purpose of exposition, a movie is characterized by the attributes "director," "actor/actress," and "genre." In this setting, it is not immediately obvious what the distance, or similarity, is between the values "Coppola" and "Scorsese," or the movies "Vertigo" and "Harvey."

There are plenty of examples of categorical data: product data, where products are defined over attributes such as brand, model, or color; census data, where information about individuals includes attributes such as marital status, education, and occupation; ecological data where plants and animals can be described with attributes such as shape of petals or type of habitat.

C

Categorical Data Clustering. Table 1 An instance of a movie database

	Director	Actor	Genre
t ₁ (Godfather II)	Coppola	De Niro	Crime
t ₂ (Good fellas)	Scorsese	De Niro	Crime
t ₃ (Vertigo)	Hitchcock	Stewart	Thriller
t ₄ (N by NW)	Hitchcock	Grant	Thriller
t ₅ (Bishop's wife)	Koster	Grant	Comedy
t ₆ (Harvey)	Koster	Stewart	Comedy

There is a plethora of such data sets, and there is always a need for clustering and analyzing them.

The lack of an inherent distance or similarity measure between categorical data objects, makes categorical data clustering a challenging problem. The challenge lies in defining a quality measure for categorical data clustering that captures the human intuition of what it means for categorical data objects to be similar. In the next sections, we present an overview of the different efforts at addressing this problem and the resulting clustering algorithms.

Structure of the Learning System

Generic Data Clustering System

We first describe the outline for a generic data clustering system, not necessarily of categorical data. In the next section we focus on categorical data specific challenges.

Data clustering is not a one-step process. In one of the seminal texts on Cluster Analysis, Jain and Dubes divide the clustering process into the following stages (Jain & Dubes, 1988):

Data collection: Includes the careful extraction of relevant data objects from the underlying data sources. In our context, data objects are distinguished by their individual values for a set of attributes.

Initial screening: Refers to the massaging of data after its extraction from the source or sources. This stage is closely related to the process of data cleaning in databases (Jarke, Lenzerini, Vassiliou, & Vassiliadis, 1999).

Representation: Includes the proper preparation of the data in order to become suitable for the clustering algorithm. In this step, the similarity measure is chosen, and the characteristics and dimensionality of the data are examined.

Clustering tendency: Checks whether the data at hand has a natural tendency to cluster or not. This stage is often ignored, especially for large data sets.

Clustering strategy: Involves the careful choice of clustering algorithm and initial parameters, if any.

Validation: Validation is often based on manual examination and visual techniques. However, as the amount of data and its dimensionality grow, we may have no means to compare the results with preconceived ideas or other clusterings.

Interpretation: This stage includes the combination of clustering results with other analyses of the data (e.g., classification), in order to draw conclusions and suggest further analysis.

In this chapter, we are interested in problems relating to Representation and Clustering Strategy. These lie in the heart of the data clustering problem, and there has been considerable research effort in these areas within the Data Mining and Machine Learning communities. More specifically, we consider the following two subproblems.

Formal formulation of the clustering problem: In order to devise algorithms for clustering, we need to mathematically formulate the intuition captured in the informal definition of the clustering problem that similar objects should be grouped together and dissimilar objects should be grouped separately. The problem formulation typically requires at least one of the following:

- A measure of similarity or distance between two data objects.
- A measure of similarity or distance between a data object and a cluster of objects. This is often defined by defining a representative for a cluster as a (new) data object and comparing the data object with the representative.
- A measure of the quality of a cluster of data objects.

The result of the formulation step is to define a clustering optimization criterion that guides the grouping of the objects into clusters.

When the data is defined over numerical attributes, these measures are defined using geometric analogies. For example, if each object is a point in the Euclidean space, then the distance between two points can be

defined as the Euclidean distance, and the representative of a cluster as the mean Euclidean vector. The quality of a cluster can be defined with respect to the "variance" of the cluster, that is, the sum of squares of the distances between each object and the mean of the cluster. The optimization criterion then becomes to minimize the variance over all clusters of the clustering.

The clustering algorithm: Once we have a mathematical formulation of the clustering problem, we need an algorithm that will find the optimal clustering in an efficient manner. In most cases, finding the optimal solution is an NP-hard problem efficient heuristics or approximation algorithms are considered. There is a large literature on this subject that approaches the problem from different angles.

There exist a large number of different clustering techniques and algorithms. We now selectively describe some broad classes of clustering algorithms and problems. A thorough categorization of clustering techniques can be found in Han and Kamber (2001), where different clustering problems, paradigms, and techniques are discussed.

Hierarchical clustering algorithms: This is a popular clustering technique, since it is easy to implement, and it lends itself well to visualization and interpretation. Hierarchical algorithms create a hierarchical decomposition of the objects. They are either agglomerative (bottom-up) or divisive (top-down). Agglomerative algorithms start with each object being a separate cluster itself, and successively merge groups according to a distance measure. Divisive algorithms follow the opposite strategy. They start with one group of all objects and successively split groups into smaller ones, until each object falls into one cluster, or as desired. The hierarchical dendrogram produced is often in itself the output of the algorithm, since it can be used for visualizing the data. Most of the times, both approaches suffer from the fact that once a merge or a split is committed, it cannot be undone or refined.

Partitional clustering algorithms: Partitional clustering algorithms define a clustering optimization criterion and then seek the partition that optimizes this criterion. Exhaustive search over all partitions is infeasible, since even for few data objects the number of possible partitions is huge. Partitional clustering algorithms often start with an initial, usually random, partition and proceed with its refinement by locally

improving the optimization criterion. The majority of such algorithms could be considered as greedy-like algorithms. They suffer from the fact that they can easily get stuck to local optima.

Spectral clustering algorithms: Spectral algorithms view the data set to be clustered as a two dimensional matrix of data objects and attributes. The entries in the matrix may be the raw values or some normalized form of these values. The principal eigenvectors of the matrix have been shown to capture the main clusters in the data. There is a rich literature on different types of spectral algorithms.

Graph clustering: Graph clustering defines a range of clustering problems, where the distinctive characteristic is that the input data is represented as a graph. The nodes of the graph are the data objects, and the (possibly weighted) edges capture the similarity or distance between the data objects. The data may come naturally in the form of a graph (e.g., a social network), or the graph may be derived in some way from the data (e.g., link two products if they appear together in a transaction). Some of the techniques we described above are directly applicable to graph data. We can also use techniques from graph theory for finding a good clustering.

Categorical Data Clustering System

In the clustering paradigm we outlined, a step of fundamental importance is to formally formulate the clustering problem, by defining a clustering optimization criterion. As we detailed above, for this step we need a measure of distance or similarity between the data objects, or a measure of cluster quality for a group of data objects. For categorical data there exists no inherent ordering or distance measure, and no natural geometric analogies we can explore, causing the clustering paradigm to break down. Research efforts on categorical data clustering have focused on addressing this problem by imposing distance measures on the categorical data and defining clustering quality criteria. We now outline some of these approaches.

Overlap-Based Similarity Measures A simple and intuitive method for comparing two categorical data objects is based on counting the overlap between the categorical values of the objects. The higher the overlap, the more similar the two objects are. This intuition leads to the

use of well-known measures such as the (*generalized*) *Hamming distance* (Jain & Dubes, 1988), which measures the number of attributes that take different values between two tuples, or the *Jaccard* similarity measure, which is defined as the intersection over the union of the values in the two tuples. In the example of Table 1, tuples \mathbf{t}_1 (Godfather II) and \mathbf{t}_2 (Good fellas) have Hamming distance 1 and Jaccard coefficient 1/2.

Two algorithms that make use of overlap-based measures are *k-modes* (Huang, 1998), and *RObust Clustering using linKs* (*ROCK*) (Guha, Rastogi, & Shim, 1999). The *k*-modes algorithm is a partitional algorithm inspired by the *k-means* algorithm, a well-known clustering algorithm for numerical data. The representative of a categorical data cluster is defined to be a data object where each attribute takes the *mode* emphasize the mode value of an attribute is the most frequent value for that attribute in the cluster.

The ROCK algorithm makes use of the Jaccard coefficient to define *links* between data objects. The data is then represented in the form of a graph, and the problem becomes essentially a graph clustering problem. Given two clusters of categorical data, ROCK measures the similarity of two clusters by comparing their *aggregate interconnectivity* against a user-specified model, thus avoiding the problem of defining a cluster representative.

Context-Based Similarity Measures One way to define relationships between categorical values is by comparing the context in which they appear. For two categorical attribute values we define the context as the values of the remaining attributes with which they co-occur in the data set. The more similar these two contexts are, the more similar the attribute values are. For example, in Table 1, Scorsese and Coppola are close since they appear in exactly the same context ("De Niro", "Crime"), while Scorsese and Hitchcock are far since their contexts are completely disjoint. Defining a distance between value contexts can be done using overlap similarity measures (Das & Mannila, 2000) or by using information-theoretic measures, i.e., comparing the distributions defined by the two contexts (Andritsos, Tsaparas, Miller, Kenneth, & Sevcik, 2004). Once we have the relationships between the values we can use standard clustering techniques for solving the clustering problem.

There are various algorithms that make use of the idea that similar values should appear in similar contexts in order to cluster categorical data. The *Clustering cAteCorigal daTa Using Summaries* (*CACTUS*) algorithm (Ganti, Gehrke, & Ramakrishnan, 1999) creates groups of attribute values based on the similarity of their context. It then uses a hierarchical greedy algorithm for grouping tuples and attributes.

In a slightly different fashion, the STIRR algorithm (Sieving Through Iterated Relational Reinforcement) [GKR98] uses the idea that similar tuples should contain co-occurring values and similar values should appear in tuples with high overlap. This idea is implemented via a dynamical system, inspired by Information Retrieval techniques (Kleinberg Jon, 1999). When the dynamical system is linear, the algorithm is similar to spectral clustering algorithms.

CLICKS (Zaki, Peters, Assent, & Seidl, 2005) is an algorithm that is similar to STIRR. Rather than a measure of similarity/distance, it uses a graph-theoretic approach to find k disjoint sets of vertices in a graph constructed for a particular data set. One special characteristic of this algorithm is that it discovers clusters in a subset of the underlying set of attributes.

Information-Theoretic Clustering Criteria The information content in a data set can be quantified through the well-studied notions of entropy and mutual information (Cover & Thomas, 1991). Entropy measures the uncertainty in predicting the values of the data when drawn from the data distribution. If we view each tuple, or cluster of tuples, as a distribution over the categorical values, then we can define the conditional entropy of the attribute values given a set of tuples, as the uncertainty of predicting the values in this set of tuples. If we have a single tuple, then the entropy is zero, since we can accurately predict the values. For tuple t_1 we know the director, the actor, and the genre with full certainty. As we group tuples together the uncertainty (and entropy) increases. Grouping together tuples \mathbf{t}_1 and \mathbf{t}_2 creates uncertainty about the director attribute, while grouping \mathbf{t}_1 and \mathbf{t}_3 creates uncertainty about all attributes. Hence the latter grouping has higher entropy than the former. Information-theoretic criteria for clustering aim at generating clusters with low entropy, since this would imply that the clusters are homogeneous, and there is little

C

information loss as a result of the clustering. This formulation allows for defining the distance between sets of tuples, using entropy-based distance measures such as the *Jensen–Shannon* divergence (Cover & Thomas, 1991). The Jensen–Shannon divergence captures the informational distances in categorical data, in a similar way that Euclidean distance captures geometric distances inherent in numerical data.

Two algorithms that make use of this idea are COOLCAT (Barbarà, Couto, & Li, 2002) and LIMBO ($scaLable\ InforMation\ Bottleneck$) [ATMR04]. COOLCAT is a partitional algorithm that performs a local search for finding the partition with k clusters with the lowest entropy. LIMBO works by constructing a summary of the data set that preserves as much information about the data as possible and then produces a hierarchical clustering of the summary. It is a scalable algorithm that can be used in both static and streaming environments.

A related approach is adopted by the COBWEB algorithm (Fisher, 1987; Gluck & Corter, 1985), a divisive hierarchical algorithm that optimizes the *category utility* measure, which measures how well particular values can be predicted given the clustering as opposed to having them in the original data set unclustered.

Categorical Clustering as Clustering Aggregation A different approach to the categorical data clustering problem is to view it as a clustering aggregation problem. Given a collection of clusterings of the data objects, the clustering aggregation problem looks for the single clustering that agrees as much as possible with the input clusterings. The problem of clustering aggregation has been shown to be equivalent to categorical data clustering (Gionis, Mannila, & Tsaparas, 2007), where each categorical attribute defines a clustering of the data objects, grouping all the objects with the same value together. For example, in Table 1, the attribute "genre" defines three clusters: the Crime cluster, the Thriller cluster, and the Comedy cluster. Similarly, the attribute "actor" defines three clusters, and the attribute "director" defines four clusters.

Various definitions have been considered in the literature for the notion of agreement between the output clustering and the input clusterings. One definition looks at all pairs of objects, and defines a

disagreement between two clusterings if one clustering places the two objects in the same cluster, while the other places them in different clusters; an agreement is defined otherwise. The clustering criterion is then to minimize the number of disagreements (or maximize the number of agreements). Other definitions are also possible, which make use of information-theoretic measures, or mappings between the clusters of the two clusterings. There is a variety of algorithms for finding the best aggregate cluster, many of which have also been studied theoretically.

Cross References

- **►**Clustering
- **▶**Data Mining
- ▶ Graph clustering
- ►Instance-Based Learning
- ▶ Partitional clustering

Recommended Reading

Andritsos, P., Tsaparas, P., Miller, R. J., Kenneth, C., & Sevcik, K. C. (2004). LIMBO: Scalable clustering of categorical data. In Proceedings of the 9th international conference on extending database technology (EDBT) (pp. 123-146). Heraklion, Greece.

Barbarà, D., Couto, J., & Li, Y. (2002). COOLCAT: An entropybased algorithm for categorical clustering. In *Proceedings of* the 11th international conference on information and knowledge management (CIKM) (pp. 582-589). McLean, VA.

Cover, T. M., & Thomas, J. A. (1991). Elements of information theory. New York: Wiley.

Das, G., & Mannila, H. (2000). Context-based similarity measures for categorical databases. In Proceedings of the 4th European conference on principles of data mining and knowledge discovery (PKDD) (pp. 201–210). Lyon, France.

Fisher, D. H. (1987). Knowledge acquisition via incremental conceptual clustering. *Machine Learning*, 2, 139–172.

Ganti, V., Gehrke, J., & Ramakrishnan, R. (1999). CACTUS: Clustering categorical data using summaries. In Proceedings of the 5th international conference on knowledge discovery and data mining (KDD) (pp. 73-83). San Diego, CA.

Gionis, A., Mannila, H., & Tsaparas, P. (2007). Clustering aggregation. ACM Transactions on Knowledge Discovery from Data, 1(1), Article No 4.

Gluck, M., & Corter, J. (1985). Information, uncertainty, and the utility of categories. In Proceedings of the 7th annual conference of the cognitive science society (COGSCI) (pp. 283–287). Irvine, CA.

Guha, S., Rastogi, R., & Shim, K. (1999). ROCK: A robust clustering algorithm for categorical attributes. In *Proceedings of the* 15th international conference on data engineering (pp. 512–521). Sydney, Australia.

Han, J., & Kamber, M. (2001). Data mining: Concepts and techniques. San Francisco: Morgan Kaufmann.

Huang, Z. (1998). Extensions to the k-means algorithm for clustering large data sets with categorical values. Data Mining and Knowledge Discovery, 2(3), 283-304.

Jain, A. K., & Dubes, R. C. (1988). Algorithms for clustering data. Englewood Cliffs, NJ: Prentice-Hall.

Jarke, M., Lenzerini, M., Vassiliou, Y., & Vassiliadis, P. (1999).
Fundamentals of data warehouses. Berlin: Springer.

Kleinberg, Jon (1999). Authoritative sources in a hyperlinked environment". *Journal of the ACM 46*(5): 604–632.

Zaki, M. J., Peters, M., Assent, I., & Seidl, T. (2005). CLICKS: An effective algorithm for mining subspace clusters in categorical datasets. In Proceeding of the 11th international conference on knowledge discovery and data mining (KDD) (pp. 736-742). Chicago, IL.

Categorization

- **►**Classification
- ► Concept Learning

Category

►Class

Causal Discovery

► Learning Graphical Models

Causality

RICARDO SILVA University College London, London, UK

Definition

The main task in causal inference is predicting the outcome of an intervention. For example, a treatment assigned by a doctor that will change the patient's heart condition is an intervention. Predicting the change in the patient's condition is a causal inference task. In general, an intervention is an action taken by an external agent that changes the original values, or the probability distributions, of some of the variables in the system. Besides predicting outcomes of actions, causal inference

is also concerned with explanation: identifying which were the causes of a particular event that happened in the past.

Motivation and Background

Many problems in machine learning are prediction problems. Given a feature vector \mathbf{X} , the task is to provide an estimate of some output vector \mathbf{Y} , or its conditional probability distribution $P(\mathbf{Y}|\mathbf{X})$. This typically assumes that the distribution of \mathbf{Y} given \mathbf{X} during learning is the same distribution at prediction time. There are many scenarios where this is not the case.

Epidemiology and several medical sciences provide counterexamples. Consider two seemingly straightforward learning problems. In the first example, one is given epidemiological data where smokers are clearly more propense than nonsmokers to develop lung cancer. Can I use this data to learn that smoking causes cancer? In the second example, consider a group of patients suffering from a type of artery disease. In this group, those that receive a bypass surgery are likely to survive longer than those that receive a particular set of drugs with no surgery.

There is no fundamental problem on using such datasets to predict the probability of a smoker developing lung cancer, or the life expectancy of someone who went through surgery. Yet, the data does not necessarily tell you if smoking is a cause of lung cancer, or that nationwide the government should promote surgery as the treatment of choice for that particular heart disease. What is going on?

There are reasons to be initially suspicious of such claims. This is well-known in statistics as the expression "association is not causation" (Wasserman, 2004, p. 253). The data-generating mechanism for our outcome **Y** ("developing lung cancer," "getting cured from artery disease") given the relevant inputs **X** ("smoking habit," "having a surgery") might change under an *intervention* for reasons such as follows.

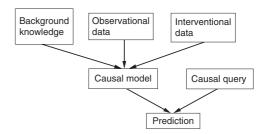
In the smoking example, the reality might be that there are several *hidden common causes* that are responsible for the observed association. A genetic factor, for instance: the possibility that there is a class of genotypes on which people are more likely to pick up smoking *and* develop lung cancer, without any direct causal connection between the two variables. In the artery disease

example, surgery might not be the best choice to be made by a doctor. It might have been the case that so far patients in better shape were more daring in choosing, by themselves, the surgery treatment. This *selection bias* will favor surgery over drug treatment, since from the outset the patients that are most likely to improve take that treatment.

When treatment is enforced by an external agent (the doctor), such selection bias disappears, and the resulting P(Y|X) will not be the same. One way of learning this relationship is through randomized trials (Rosenbaum, 2002). The simplest case consists on flipping a coin for each patient on the training set. Each face of the coin corresponds to a possible treatment, and assignment is done accordingly. Since assignment does not depend on any hidden common cause or selection bias, this provides a basis for learning causal effects. Machine learning and statistical techniques can be applied directly in this case (e.g., \blacktriangleright logistic regression). Data analysis performed with a randomized trial is sometimes called an interventional study.

The smoking case is more complicated: a direct intervention is not possible, since it is not acceptable to force someone to smoke or not to smoke. The inquiry asks only for a *hypothetical intervention*, that is, if someone is forced to smoke, will his or her chances of developing lung cancer increase? Such an intervention will not take place, but this still has obvious implications in public policy. This is the heart of the matter in issues such as deciding on raising tobacco taxes, or forbidding smoking in public places. However, data that measures this interventional data-generation mechanism will never be available for ethical reasons. The question has to be addressed through an *observational study*, that is, a study for causal predictions without interventional data.

Observational studies arise not only under the impossibility of performing interventions, but also in the case where performing interventions is too expensive or time consuming. In this case, observational studies, or a combination of observational and interventional studies, can provide extra information to guide an experimental analysis (Cooper & Yoo, 1999; Eaton & Murphy, 2007; Eberhardt, Glymour, & Scheines, 2005; Sachs, Prez, Pe'er, Lauffenburger, & Nolan, 2005). The use of observational data, or the combination of several interventional datasets, is where the greatest contributions of machine learning to causal inference rest.



Structure of the Learning System

Structure of Causal Inference

In order to use observational data, a causal inference system needs a way of linking the state of the world under an intervention to the *natural* state of the world. The natural state is defined as the one to which no external intervention is applied. In the most general formulation, this link between the natural state and the manipulated world is defined for interventions in any subset of variables in the system.

A common language for expressing the relationship between the different states of the world is a *causal graph*, as explained in more detail in the next section. A causal model is composed of the graph and a probability distribution that factorizes according to the graph, as in a standard praphical model. The only difference between a standard graphical model and a causal graphical model is that in the latter extra assumptions are made. The graphical model can be seen as a way of encoding such assumptions.

The combination of assumptions, observational, and interventional data generates such a causal graphical model. In the related problem of reinforcement learning, the agent has to maximize a specific utility function and typically has full control on which interventions (actions) can be performed. Here we will focus on the unsupervised problem of learning a causal model for a fixed input of observational and interventional data.

Because only some (or no) interventional data might be available, the learning system might not be able to answer some causal queries. That is, the system will not provide an answer for some prediction tasks.

Languages and Assumptions for Causal Inference Directed acyclic graphs (DAGs) are a popular language in machine learning to encode qualitative statements about causal relationships. A DAG is composed of a set of vertices and a set of directed edges. The notions

of parents, children, ancestors, and descendants are the usual ones found in graphical modeling literature.

In terms of causal statements, a directed edge $A \rightarrow B$ states that A is a *direct* cause of B: that is, different interventions on A will result in different distributions for B, even if we intervene on all other variables. The assumption that A is a cause of B is not used in noncausal paraphical models.

A causal DAG *G* satisfies the *causal Markov condition* if and only if a vertex is independent of all of its nondescendants given its direct causes (parents). In Fig. 1(a), *A* is independent of *D*, *E*, and *F* given its parents, *B* and *C*. It may or may not be independent of *G* given *B* and *C*.

The causal Markov condition implies several other conditional independence statements. For instance, in Fig. 1(a) we have that *H* is independent of *F* given *A*. Yet, this is not a statement about the parents of any vertex. Pearl's d-separation criterion (Pearl, 2000) is a sound and complete way of reading off independencies, out of a DAG, which are entailed by the causal Markov condition. We assume that the joint probability distribution over the vertice variables is *Markov* with respect to the graph, that is, any independence statement that is encoded by the graph should imply the corresponding independence in the distribution.

Representing Interventions

The local modularity given by the causal Markov condition leads to a natural notion of intervention. Random variable V, represented by a particular vertex in the graph, is following a *local mechanism*: its direct causes determine the distribution of V before its direct effects are generated. The role of an intervention is to *override* the natural local mechanism. An external agent substitutes the natural P(V|Parents(V)) by a new distribution $P_{Man}(V|Parents(V))$ while keeping the rest of the model unchanged ("Man" here stands for a particular

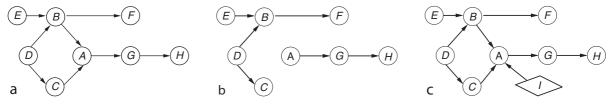
manipulation). The notion of intervening by changing a single local mechanism is sometimes known as an *ideal intervention*. Other general types of interventions can be defined (Eaton & Murphy, 2007), but the most common frameworks for calculating causal effects rely on this notion.

A common type of intervention is the point mass intervention, which happens when V is set to some constant v. This can be represented graphically by "wiping out" all edges into V. Figure 1(b) represents the resulting graph in (a) under a point manipulation of A. Notice that A is now d-separated from its direct causes under this regime. It is also probabilistically independent, since A is now constant. This allows for a graphical machinery that can read off independencies out of a manipulated graph (i.e., the one with removed edges). It is the idea of representing the natural state of the world with a single causal graph, and allowing for modifications in this graph according to the intervention of choice, that links the different regimes obtained under different interventions.

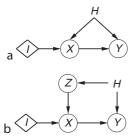
For the general case where a particular variable V is set to a new distribution, a manipulation node is added as an extra parent of V: this represents that an external agent is acting over that particular variable (Dawid, 2002; Pearl, 2000; Spirtes, Glymour, & Scheines, 2000), as illustrated in Fig. 1(c). P(V|Parents(V)) under intervention I is some new distribution $P_{Man}(V|Parents(V), I)$.

Calculating Distributions under Interventions

The notion of independence is a key aspect of probabilistic graphical models, where it allows for efficient computation of marginal probabilities. In causal graphical models, it also fulfills another important role: independencies indicate that the effect of some interventions can be estimated using observational data.



Causality. Figure 1. (a) A causal DAG. (b) Structure of the causal graph under some intervention that sets the value of A to a constant. (c) Structure of the causal graph under some intervention that changes the distribution of A



Causality. Figure 2. (a) X and Y have a hidden common cause H. (b) Y is dependent on the intervention node I given X, but conditioning on Z and marginalizing it out will allow us to eliminate the "back-door" path that links X and Y through the hidden common cause H

We will illustrate this concept with a simple example. One of the key difficulties in calculating a causal effect is *unmeasured confounding*, that is, hidden common causes. Consider Fig. 2(a), where X is a direct cause of Y, H is a hidden common cause of both and I is an intervention vertex. Without extra assumptions, there is no way of estimating the effect of X on Y using a training set that is sampled from the observed marginal P(X, Y). This is more easily seen in the case where the model is multivariate Gaussian with zero mean. In this case, each variable is a linear combination of its parents with standard Gaussian additive noise

$$X = aH + \epsilon_X$$
$$Y = bX + cH + \epsilon_Y$$

where H is also a standard normal random variable. The manipulated distribution $P_{Man}(Y|X,I)$, where I is a point intervention setting X=x, is a Gaussian distribution with mean $b \cdot x$. Value x is given by construction, but one needs to learn the unknown value b.

One can verify that the covariance of X and Y in the natural state is given by a + bc. Observational data, that is, data sampled from P(X, Y), can be used to estimate the covariance of X and Y in the natural state, but from that it is not possible to infer the value of b: there are too many degrees of freedom.

However, there *are* several cases where the probability of **Y** given some intervention on **X** can be estimated with observational data and a given causal graph. Consider the graph in Fig. 2(b). The problem again is to learn the distribution of Y given X under regime I, that is, P(Y|X,I). It can be read from the graph that

I and Y are not independent given X, which means $P(Y|X,I) \neq P(Y|X)$. How can someone then estimate P(Y|X,I) if no data for this process has been collected? The answer lies on *reducing the "causal query" to a "probabilistic query"* where the dependence on I disappears (and, hence, the necessity of having data measured under the I intervention). This is done by relying on the assumptions encoded by the graph:

$$P(Y|X,I) = \sum_{z} P(Y|X,I,z)P(Z=z|X,I)$$

$$(Z \text{ is discrete in this example})$$

$$= \sum_{z} P(Y|X,z)P(Z=z|X,I)$$

$$(Y \text{ and } I \text{ are independent given } Z)$$

$$\propto \sum_{z} P(Y|X,z)P(X|z,I)P(Z=z|I)$$

$$(\text{By Bayes' rule})$$

$$= \sum_{z} P(Y|X,z)P(X|z,I)P(Z=z)$$

$$(Z \text{ and } I \text{ are marginally independent})$$

In the last line, we have P(Y|X,Z) and P(Z), which can be estimated with observational data, since no intervention variable I appears on the expression. P(X|Z,I) is set by the external agent: its value is known by construction. This means that the causal distribution P(Y|X,I) can be learned even in this case where X and Y share a hidden common cause H.

There are several notations for denoting an interventional distribution such as P(Y|X,I). One of the earliest was that of Spirtes et al. (2000), who used the notation

$$P(Y|set X = x) \tag{1}$$

to represent the distribution under an intervention I that fixed the value of X to some constant x. Pearl (2000) uses the operator do with a similar meaning.

$$P(Y|do(X=x)) (2)$$

Pearl's *do*-calculus is essentially a set of operations for reducing a probability distribution that is a function of some intervention to a probability distribution that does not refer to any intervention. All reductions are conditioned on the independencies encoded in a given causal graph. This is in the same spirit of the example presented above.

C

Causality 163

The advantage of such notations is that, for point interventions, they lead to simple yet effective transformations (or to a report that no transformation is possible). Spirtes et al. (2000) and Pearl (2000) provide a detailed account of such prediction tools. By making a clear distinction between P(Y|X) (X under the natural state) and P(Y|do(X)) (X under some intervention), much of the confusion that conflates causal and noncausal predictions disappears.

It is important to stress that methods such as the *do*-calculus are nonparametric, in the sense that they rely on conditional independence constraints only. More powerful reductions are possible if one is willing to provide extra information, such as assuming linearity of causal effects. For such cases, other parametric constraints can be exploited (Pearl, 2000; Spirtes et al., 2000).

Learning Causal Structure

In all of the previous section, we assumed that a causal graph was available. Since background knowledge is often limited, learning such graph structures becomes an important task. However, this cannot be accomplished without extra assumptions. To see why, consider again the example in Fig. 2(a): if a + bc = 0, it follows that the X and Y are independent in the natural state. However, Y is *not* causally independent of X (if $b \neq 0$): $P(Y|do(X = x_1))$ and $P(Y|do(X = x_2))$ will be two different Gaussians with means $b \cdot x_1$ and $b \cdot x_2$, respectively.

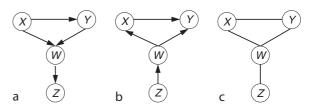
This example demonstrates that an independence constraint that is testable by observational data does not warrant causal independence, at least based on the causal Markov condition only. However, an independence constraint that arises from particular identities such as a + bc = 0 is not *stable*, in the sense that it does not follow from the qualitative causal relations entailed by the Markov condition: a change in any of the parameter values will destroy such a constraint.

The artificiality of unstable independencies motivates an extra assumption: the *faithfulness* condition (Spirtes et al., 2000), also known as the *stability* condition (Pearl, 2000). We say that a distribution *P* is faithful to a causal graph *G* if *P* is Markov with respect to *G*, *and* if each conditional independence in *P* corresponds to some d-separation in *G*. That is, on top of the causal Markov condition we assume that all independencies in *P* are entailed by the causal graph *G*.

The faithfulness condition allows us to reconstruct classes of causal graphs from observational data. In the simplest case, observing that X and Y are independent entails that there is no causal connection between X and Y. Consequently, P(Y|do(X)) = P(Y|X) = P(Y). No interventional data was necessary to arrive at this conclusion, given the faithfulness condition.

In general, the solution is undetermined: more than one causal graph will be compatible with a set of observable independence constraints. Consider a simple example, where data is generated by a causal model with a causal graph given as in Fig. 3(a). This graph entails some independencies: for instance, that X and Zare independent given W, or that X and Y are not independent given any subset of $\{W, Z\}$. However, several other graphs entail the same conditional independencies. The graph in Fig. 3(b) is one example. The learning task is then discovering an equivalence class of graphs, not necessarily a particular graph. This is in contrast with the problem of learning the structure of noncausal graphical models: the fact that there are other structures compatible with the data is not important in this case, since we will not use such graphical models to predict the effect of some hypothetical intervention. An equivalence class might not be enough information to reduce a desired causal query to a probabilistic query, but it might require much less prior knowledge than specifying a full causal graph.

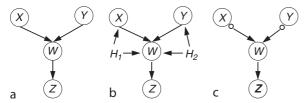
Assume for now that no hidden common causes exist in this domain. In particular, the graphical object in Fig. 3(c) is a representation of the equivalence class of graphs that are compatible with the independencies encoded in Fig. 3(a) (Pearl, 2000; Spirtes et al., 2000). All members of the equivalence class will have the same



Causality. Figure 3. (a) A particular causal graph which entails a few independence constraints, such as X and Z being independent given W. (b) A different causal graph that entails exactly the same independence constraints as (a). (c) A representation for all graphs that entail the same conditional independencies as (a) and (b)

skeleton of this representation, that is, the same adjacencies. An undirected edge indicates that there are two members in the equivalence class where directionality of this particular edge goes in opposite directions. Some different directions are illustrated in Fig. 3(b). One can verify from the properties of d-separation that if an expert or an experiment indicates that X - W should be directed as $X \to W$, then the edge W - Z is compelled to be directed as $W \to Z$: the direction $W \leftarrow Z$ is incompatible with the simultaneous findings that X and Z are independent given W, and that X causes W.

More can be discovered if more independence constraints exist. In Fig. 4(a), X is not a cause of Y. If we assume no hidden common causes exist in this domain, then no other causal graph is compatible with the independence constraints of Fig. 4(a): the equivalence class is this graph only. However, the assumption of no hidden common causes is strong and undesirable. For instance, the graph in Fig. 4(b), where H_1 and H_2 are hidden, is in the same equivalence class of (a). Yet, the graph in Fig. 4(a) indicates that P(W|do(X)) = P(W|X), which can be arbitrarily different from the real P(W|do(X)) if Fig. 4(b) is the real graph. Some equivalence class representations, such as the Partial Ancestral Graph representation (Spirtes et al., 2000), are robust to hidden common causes: in Fig. 4(c), an edge that has a circle as endpoint indicates that is not known if there is a causal path into both, for example, X and W (which would be the case for a hidden common cause of *X* and *W*). The arrow into *W* does indicate that W cannot be a cause of X. A fully directed edge such as $W \rightarrow Z$ indicates total information: W is a cause of Z, Z is



Causality. Figure 4. (a) A particular causal graph with no other member on its equivalence class (assuming there are no hidden common causes). (b) Graph under the presence of two hidden common causes H_1 and H_2 . (c) A representation for all graphs that entail the same conditional independencies as (a), without assuming the nonexistence of hidden common causes

not a cause of W, and W and Z have no hidden common causes.

Given equivalence class representations and background knowledge, different types of algorithms explore independence constraints to learn an equivalence class. It is typically assumed that the true graph is acyclic. The basic structure is to evaluate how well a set of conditional independence hypotheses are supported by the data. Depending on which constraints are judged to hold in the population, we keep, delete, or orient edges accordingly. Some algorithms, such as the PC algorithm (Spirtes et al., 2000), test a single independence hypothesis at a time, and assemble the individual outcomes in the end into an equivalence class representation. Other algorithms such as the GES algorithm (Chickering, 2002; Meek, 1997) start from a prior distribution for graphs and parameters, and proceed to compare the marginal likelihood of members of different equivalence classes (which can be seen as a Bayesian joint test of independence hypotheses). In the end, this reduces to a search for the maximum a posteriori equivalence class estimator. Both PC and GES have consistency properties: in the limit of infinite data, they return the right equivalence class under the faithfulness assumption. However, both PC and GES, and most causal discovery algorithms, assume that there are no hidden common causes in the domain. The fast causal inference (FCI) algorithm of Spirtes et al. (2000) is able to generate equivalence class representations as in Fig. 4(c). As in the PC algorithm, this is done by testing a single independence hypothesis at a time, and therefore is not very robust given small samples. A GES-like algorithm with the consistency properties of FCI is not currently known. An algorithm that allows for cyclic networks is discussed by Richardson (1996). More details of the fundamentals of structure learning algorithms are given by Scheines (1997).

Our examples relied on conditional independence constraints. In this case, the equivalence class is known as the *Markov equivalence class*. Markov equivalence classes are "nonparametric," in the sense that they do not refer to any particular probability family. In practice, this advantage is limited by our ability on evaluating independence hypotheses within flexible probability families. Another shortcoming of Markov equivalence classes is that they might be poorly informative if few independence constraints exist in the population. This

will happen, for instance, if a single hidden variable is a common cause of all observed variables. If one is willing to incorporate further assumptions, such as linearity of causal relationships, parametric constraints can be used to define other types of equivalence classes that are more discriminative than the Markov equivalence class. Silva, Scheines, Glymour, & Spirtes (2006) describe how some rank constraints in the covariance matrix of the observed variables can be used to learn the structure of linear models, even if no independence constraints are observable.

Evaluating the success of a structure learning algorithm is difficult, since ultimately it depends on interventional data. A promising area of application is molecular biology, where the large number of variables makes the use of graphical models a promising venue for decomposing complex biological systems, and for combining multiple sources of observational and interventional data. Sachs et al. (2005) describe a potential application, with further analysis discussed by Eaton and Murphy (2007). Other applications are discussed in the volume edited by Glymour and Cooper (1999).

Confidence Intervals Several causal learning algorithms such as the PC and FCI algorithms (Spirtes et al., 2000) are consistent, in the sense that they can recover the correct equivalence class given the faithfulness assumption and an infinite amount of data. Although point estimates of causal effects are important, it is also important to provide confidence intervals. Bayesian confidence intervals are readily available by having priors over parameters and graphs. ▶Markov chain Monte Carlo algorithms, however, might be problematic due to the high-dimensional and discrete graph space. A practical algorithm that relies on a prior over *orderings* of variables (such that for a given ordering, a graph is not allowed to have vertex *X* as an ancestor of *Y* if *Y* antecedes *X* in the ordering) is given by Friedman and Koller (2003).

Such methods do not necessarily guarantee good frequentist properties. As a matter of fact, it has been shown that no such method can exist given the faithfulness assumption only (Robins, Scheines, Spirtes, & Wasserman, 2003). An intuitive explanation is as follows: consider the model such as the one in Fig. 2(a). For any given sample size, there is at least one model such that the associations due to the paths $X \leftarrow H \rightarrow Y$

and $X \rightarrow Y$ nearly cancel each other (faithfulness is still preserved), making the covariance of X and Y statistically undistinguishable from zero. In order to achieve uniform consistency, causal inference algorithms will need assumptions stronger than faithfulness. Zhang and Spirtes (2003) provide some directions.

Other Languages and Tasks in Causal Learning

A closely related language for representing causal models is the *potential outcomes* framework popularized by Rubin (2004). In this case, random variables for a same variable *Y* are defined for each possible state of the intervened variable *X*. Notice that, by definition, only one of the possible *Y* outcomes can be observed for any specific data point. This model is popular in statistics literature as a type of missing data model. The relation between potential outcomes and graphical models is discussed by Pearl (2000).

A case where potential outcomes become more clearly motivated is in causal explanation. In this setup, the model is asked for the probability that a particular event in time was the cause of a particular outcome. This is often cast as a counterfactual question: had A been false, would B still have happened? Questions in History and law are of this type: the legal responsibility of an airplane manufacturer in an accident depends on technical malfunction being an actual cause of the accident. Ultimately, such issues of causal explanation, actual causation and other counterfactual answers, are untestable. Although machine learning can be a useful tool to derive the consequences of assumptions combined with data about other events of the same type, in general the answers will not be robust to changes in the assumptions, and the proper assumptions ultimately cannot be selected with the available data. Some advances in generating explanations with causal models are described by Halpern and Pearl (2005).

Cross References

- ► Graphical Models
- ► Learning Graphical Models

Recommended Reading

Chickering, D. (2002). Optimal structure identification with greedy search. *Journal of Machine Learning Research*, 3, 507–554.

Cooper, G., & Yoo, C. (1999). Causal discovery from a mixture of experimental and observational data. In *Uncertainty in Artifi*cial Intelligence (UAI). 166 CBR

- Dawid, A. P. (2002). Influence diagrams for causal modelling and inference. *International Statistical Review, 70*, 161–189.
- Eaton, D., & Murphy, K. (2007). Exact Bayesian structure learning from uncertain interventions. In Artificial Intelligence and Statistics (AISTATS).
- Eberhardt, F., Glymour, C., & Scheines, R. (2005). On the number of experiments sufficient and in the worst case necessary to identify all causal relations among N variables. In *Uncertainty in Artificial Intelligence (UAI)* (pp. 178–184).
- Friedman, N., & Koller, D. (2003). Being Bayesian about network structure: A Bayesian approach to structure discovery in Bayesian networks. *Machine Learning*, 50, 95–126.
- Glymour, C., & Cooper, G. (1999). Computation, causation and discovery. Cambridge, MA: MIT Press.
- Halpern, J., & Pearl, J. (2005). Causes and explanations: A structural-model approach. Part II: Explanations. British Journal for the Philosophy of Science, 56, 889-911.
- Meek, C. (1997). Graphical models: Selecting causal and statistical models. PhD thesis, Carnegie Mellon University.
- Pearl, J. (2000). Causality: Models, reasoning and inference. Cambridge: Cambridge University Press.
- Richardson, T. (1996). A discovery algorithm for directed cyclic graphs. In *Proceedings of 12th conference on Uncertainty in Artificial Intelligence*.
- Robins, J., Scheines, R., Spirtes, P., & Wasserman, L. (2003). Uniform consistency in causal inference. *Biometrika*, 90, 491–515.
- Rosenbaum, P. (2002). Observational studies. New York: Springer. Rubin, D. (2004). Direct and indirect causal effects via potential outcomes. Scandinavian Journal of Statistics, 31, 161-170.
- Sachs, K., Prez, O., Pe'er, D., Lauffenburger, D., & Nolan, G. (2005).
 Causal protein-signaling networks derived from multiparameter single-cell data. Science, 308, 523–529.
- Scheines, R. (1997). An introduction to causal inference. In V. McKim & S. Turner (Eds.), *Causality in Crisis?* (pp. 185-200).
- Silva, R., Scheines, R., Glymour, C., & Spirtes, P. (2006). Learning the structure of linear latent variable models. *Journal of Machine Learning Research*, 7, 191–246.
- Spirtes, P., Glymour, C., & Scheines, R. (2000). Causation, prediction and search. Cambridge, MA: Cambridge University Press.
- Wasserman, L. (2004). All of statistics. New York: Springer.
- Zhang, J., & Spirtes, P. (2003). Strong faithfulness and uniform consistency in causal inference. In *Uncertainty in Artificial Intelligence*.

CBR

► Case-Based Reasoning

CC

▶ Cascade-Correlation

Certainty Equivalence Principle

▶Internal Model Control

Characteristic

▶ Attribute

City Block Distance

► Manhattan Distance

Class

CHRIS DRUMMOND
National Research Council of Canada

Synonyms

Category; Class; Collection; Kind; Set; Sort; Type

Definition

A class is a collection of things that might reasonably be grouped together. Classes that we commonly encounter have simple names so, as humans, we can easily refer to them. The class of dogs, for example, allows me to say "my dog ate my newspaper" without having to describe a particular dog, or indeed, a particular newspaper. In machine learning, the name of the class is called the class label. Exactly what it means to belong to a class, or category, is a complex philosophical question but often we think of a class in terms of the common properties of its members. We think particularly of those properties which seperate them from other things which are in many ways similar, e.g., cats mieow and dogs bow-wow. We would be unlikely to form a class from a random collection of things, as they would share no common properties. Knowing something belonged to such a collection would be of no particular benefit. Although

Class Imbalance Problem 167

many every day classes will have simple names, we may construct them however we like, e.g., "The things I like to eat for breakfast on a Saturday morning." As there is no simple name for such a collection, in machine learning we would typically refer to it as the positive class, and all occurences of it are positive examples; the negative class would be everything else.

Motivation and Background

The idea of a class is important in learning. If we discover something belongs to a class, we suddenly know quite a lot about it even if we have not encountered that particular example before. In machine learning, our use of the term accords closely with the mathematical definition of a class, as a collection of sets unambiguously defined by a property that all its members share. It also accords with the idea of equivalence classes, which group similar things. Sets have an intension, the description of what it means to be a member, and an extension, things that belong to the set, useful properties of a class in machine learning. Class is also a term used extensively in knowledge bases to denote an important relationship between groups, of sub-class and super class. Learning is often viewed as a way of solving the knowledge acquisition bottleneck (Buchanan et al., 1983) in knowledge bases and the use of the term class in machine learning highlights this connection.

Recommended Reading

Buchanan, B., Barstow, D., Bechtel, R., Bennett, J., Clancey, W., Kulikowski, C., et al. (1983) Constructing an expert system. In F. Hayes-Roth, D.A. Waterman, & D.B. Lenat (Eds.), Building expert systems (pp. 127-167). Reading, MA: Addison-Wesley.

Class Imbalance Problem

CHARLES X. LING, VICTOR S. SHENG The University of Western Ontario Canada

Definition

Data are said to suffer the *Class Imbalance Problem* when the class distributions are highly imbalanced. In this context, many **b**classification learning algorithms

have low predictive accuracy for the infrequent class. Cost-sensitive learning is a common approach to solve this problem.

Motivation and Background

Class imbalanced datasets occur in many real-world applications where the class distributions of data are highly imbalanced. For the two-class case, without loss of generality, one assumes that the minority or rare class is the positive class, and the majority class is the negative class. Often the minority class is very infrequent, such as 1% of the dataset. If one applies most traditional (costinsensitive) classifiers on the dataset, they are likely to predict everything as negative (the majority class). This was often regarded as a problem in learning from highly imbalanced datasets.

However, Provost (2000) describes two fundamental assumptions that are often made by traditional costinsensitive classifiers. The first is that the goal of the classifiers is to maximize the accuracy (or minimize the error rate); the second is that the class distribution of the training and test datasets is the same. Under these two assumptions, predicting everything as negative for a highly imbalanced dataset *is often the right thing to do.* Drummond and Holte (2000) show that it is usually very difficult to outperform this simple classifier in this situation.

Thus, the imbalanced class problem becomes meaningful only if one or both of the two assumptions above are not true; that is, if the cost of different types of error (false positive and false negative in the binary classification) is not the same, or if the class distribution in the test data is different from that of the training data. The first case can be dealt with effectively using methods in cost-sensitive meta-learning (see Costsensitive learning).

In the case when the misclassification cost is not equal, it is usually more expensive to misclassify a minority (positive) example into the majority (negative) class, than a majority example into the minority class (otherwise it is more plausible to predict everything as negative). That is, *FNcost* > *FPcost*. Thus, given the values of *FNcost* and *FPcost*, a variety of cost-sensitive meta-learning methods can be, and have been, used to solve the class imbalance problem (Japkowicz & Stephen, 2002; Ling & Li, 1998). If the values of

168 Classification

FNcost and *FPcost* are not unknown explicitly, *FNcost* and *FPcost* can be assigned to be proportional to the number of positive and negative training cases (Japkowicz & Stephen, 2002).

In case the class distributions of training and test datasets are different (e.g., if the training data is highly imbalanced but the test data is more balanced), an obvious approach is to sample the training data such that its class distribution is the same as the test data. This can be achieved by oversampling (creating multiple copies of examples of) the minority class and/or undersampling (selecting a subset of) the majority class (Provost, 2000).

Note that sometimes the number of examples of the minority class is too small for classifiers to learn adequately. This is the problem of insufficient (small) training data and different from that of imbalanced datasets.

Recommended Reading

Drummond, C., & Holte, R. (2000). Exploiting the cost (in)sensitivity of decision tree splitting criteria. In Proceedings of the seventeenth international conference on machine learning (pp. 239-246).

Drummond, C., & Holte, R. (2005). Severe class imbalance: Why better algorithms aren't the answer. In *Proceedings of the sixteenth European conference of machine learning, LNAI* (Vol. 3720, pp. 539-546).

Japkowicz, N., & Stephen, S. (2002). The class imbalance problem: A systematic study. *Intelligent Data Analysis*, 6(5), 429-450.

Ling, C. X., & Li, C. (1998). Data mining for direct marketing – Specific problems and solutions. In Proceedings of fourth international conference on Knowledge Discovery and Data Mining (KDD-98) (pp. 73-79).

Provost, F. (2000). Machine learning from imbalanced data sets 101. In Proceedings of the AAAI'2000 workshop on imbalanced data.

Classification

CHRIS DRUMMOND
National Research Council of Canada

Synonyms

Categorization; Generalization; Identification; Induction; Recognition

Definition

In common usage, the word classification means to put things into categories, group them together in some useful way. If we are screening for a disease, we would group people into those with the disease and those without. We, as humans, usually do this because things in a group, called a >class in machine learning, share common characteristics. If we know the class of something, we know a lot about it. In machine learning, the term classification is most commonly associated with a particular type of learning where examples of one or more >classes, labeled with the name of the class, are given to the learning algorithm. The algorithm produces a classifier which maps the properties of these examples, normally expressed as >attribute-value pairs, to the class labels. A new example whose class is unknown is classified when it is given a class label by the classifier based on its properties. In machine learning, we use the word classification because we call the grouping of things a class. We should note, however, that other fields use different terms. In philosophy and statistics, the term categorization is more commonly used. In many areas, in fact, classification often refers to what is called **\rightarrow** clustering in machines learning.

Motivation and Background

Classification is a common, and important, human activity. Knowing something's class allows us to predict many of its properties and so act appropriately. Telling other people its class allows them to do the same, making for efficient communication. This emphasizes two commonly held views of the objectives of learning. First, it is a means of peneralization, to predict accurately the values for previously unseen examples. Second, it is a means of compression, to make transmission or communication more efficient. Classification is certainly not a new idea and has been studied for some considerable time. From the days of the early Greek philosophers such as Socrates, we had the idea of categorization. There are essential properties of things that make them what they are. It embodies the idea that there are natural kinds, ways of grouping things, that are inherent in the world. A major goal of learning, therefore, is recognizing natural kinds, establishing the necessary and sufficient conditions for belonging to a category. This "classical" view of categorization, most

C

Classification 169

often attributed to Aristotle, is now strongly disputed. The main competitor is prototype theory; things are categorized by their similarity to a prototypical example (Lakoff, 1987), either real or imagined. There is also much debate in psychology (Ashby & Maddox, 2005), where many argue that there is no single method of categorization used by humans.

As much of the inspiration for machine learning originated in how humans learn, it is unsurprising that our algorithms reflect these distinctions. Nearest neighbor algorithms would seem to have much in common with prototype theory. These have been part of pattern recognition for some time (Cover & Hart, 1967) and have become popular in machine learning, more recently, as ▶instance-based learners (Aha, Kiber, & Albert, 1991). In machine learning, we measure the distance to one or more members of a concept rather a specially constructed prototype. So, this type of learning is perhaps more a case of the exemplar learning found in the psychological literature, where multiple examples represent a category. The closest we have to prototype learning occurs in clustering, a type of ▶unsupervised learning, rather than classification. For example, in >k-means clustering group membership is determined by closeness to a central value.

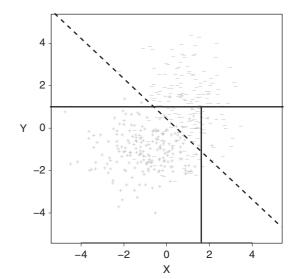
In the early days of machine learning, our algorithms (Mitchell, 1977; Winston, 1975) had much in common with the classical theory of categorization in philosophy and psychology. It was assumed that the data were consistent, there were no examples with the same attribute values but belonging to different classes. It was quickly realized that, even if the properties where necessary and sufficient to capture the class, there was often noise in the attribute and perhaps the class values. So, complete consistency was seldom attainable in practice. New ▶ classification algorithms were designed, which could tolerate some noise, such as ▶decision trees (Breiman, Friedman, Olshen, & Stone, 1984; Quinlan, 1986, 1993) and rule-based learners (see ▶Rule Learning) (Clark & Niblett, 1989; Holte, 1993; Michalski, 1983).

Structure of the Learning System

Whether one uses instance-based learning, rule-based learning, decision trees, or indeed any other classification algorithm, the end result is the division of the input space into regions belonging to a single class. The input space is defined by the Cartesian product of the attributes, all possible combinations of possible values.

As a simple example, Fig. 1 shows two classes + and -, each a random sample of a normal distribution. The attributes are *X* and *Y* of real type. The values for each attribute range from $\pm \infty$. The figure shows a couple of alternative ways that the space may be divided into regions. The bold dark lines, construct regions using lines that are parallel to the axes. New examples that have Y less than 1 and X less than 1.5 with be classified as +, all others classified as -. Decision trees and rules form this type of boundary. A ▶linear discriminant function, such as the bold dashed line, would divide the space into half-spaces, with new examples below the line being classified as + and those above as -. Instance-based learning will also divide the space into regions but the boundary is implicit. Classification occurs by choosing the class of the majority of the nearest neighbors to a new example. To make the boundary explicit, we could mark the regions where an example would be classified as + and those classified as -. We would end up with regions bounded by polygons.

What differs among the algorithms is the shape of the regions, and how and when they are chosen. Sometimes the regions are implicit as in lazy learners (see Lazy Learning) (Aha, 1997), where the boundaries are not decided until a new example is being classified.



Classification. Figure 1. Dividing the input space

170 Classification

Sometimes the regions are determined by decision theory as in generative classifiers (see ▶Generative Learners) (Rubinstein & Hastie, 1997), which model the full joint distribution of the classes. For all classifiers though, the input space is effectively partitioned into regions representing a single class.

Applications

One of the reasons that classification is an important part of machine learning is that it has proved to be a very useful technique for solving practical problems. Classification has been used to help scientists in the exploration, and comprehension, of their particular domains of interest. It has also been used to help solve significant industrial problems. Over the years a number of authors have stressed the importance of applications to machine learning and listed many successful examples (Brachman, Khabaza, Kloesgen, Piatetsky-Shapiro, & Simoudis, 1996; Langley & Simon, 1995; Michie, 1982). There have also been workshops on applications (Aha & Riddle, 1995; Engels, Evans, Herrmann, & Verdenius, 1997; Kodratoff, 1994) at major machine learning conferences and a special issue of Machine Learning (Kohavi & Provost, 1998), one of the main journals in the field. There are now conferences that are highly focused on applications. Collocated with major artificial intelligence conferences is the Innovative Applications of Artificial Intelligence conference. Since 1989, this conference has highlighted practical applications of machine learning, including classification (Schorr & Rappaport, 1989). In addition, there are now at least two major knowledge discovery and ▶data mining conferences (Fayyad & Uthurusamy, 1995; Komorowski & Zytkow, 1997) with a strong focus on applications.

Future Directions

In machine learning, there are already a large number of different classification algorithms, yet new ones still appear. It seems unlikely that there is an end in sight. The "no free lunch theory" (Wolpert & Macready, 1997) indicates that there will never be a single best algorithm, better than all others in terms of predictive power. However, apart from their predictive performance, each classifier has its own attractive properties which are important to different groups of people. So,

new algorithms are still of value. Further, even if we are solely concerned about performance, it may be useful to have many different algorithms, all with their own biases (see ►Inductive Bias). They may be combined together to form an ensemble classifier (Caruana, Niculescu-Mizil, Crew, & Ksikes, 2004), which outperforms single classifiers of one type (see ►Ensemble Learning).

Limitations

Classification has been critical part of machine research for some time. There is a concern that the emphasis on classification, and more generally on supervised learning, is too strong. Certainly much of human learning does not use, or require, labels supplied by an expert. Arguably, unsupervised learning should play a more central role in machine learning research. Although classification does require a label, it does necessarily need an expert to provide labeled examples. Many successful applications rely on finding some, easily identifiable, property which stands in for the class.

Recommended Reading

- Aha, D. W. (1997). Editorial. Artificial Intelligence Review, 11(1–5), 1–6.
- Aha, D. W., Kibler, D., & Albert, M. K. (1991). Instance-based learning algorithms. *Machine Learning*, 6(1), 37-66.
- Aha, D. W., & Riddle, P. J. (Eds.). (1995). Workshop on applying machine learning in practice. In *Proceedings of the 12th international conference on machine learning*.
- Ashby, F. G., & Maddox, W. T. (2005). Human category learning. Annual Review of Psychology, 56, 149–178.
- Bishop, C. M. (2007). Pattern recognition and machine learning. New York: Springer.
- Brachman, R. J., Khabaza, T., Kloesgen, W., Piatetsky-Shapiro, G., & Simoudis, E. (1996). Mining business databases. *Communications of the ACM*, 39(11), 42-48.
- Breiman, L., Friedman, J. H., Olshen, R. A., & Stone, C. J. (1984). Classification and regression trees. Belmont, CA: Wadsworth.
- Caruana, R., Niculescu-Mizil, A., Crew, G., & Ksikes, A. (2004). Ensemble selection from libraries of models. In Proceedings of the 21st international conference on machine learning (pp. 137-144).
- Clark, P., & Niblett, T. (1989). The CN2 induction algorithm. Machine Learning, 3, 261–284.
- Cover, T., & Hart, P. (1967). Nearest neighbor pattern classification.

 IEEE Transactions on Information Theory, 13, 21–27.
- Dietterich, T., & Shavlik, J. (Eds.). Readings in machine learning. San Mateo, CA: Morgan Kaufmann.
- Engels, R., Evans, B., Herrmann, J., & Verdenius, F. (Eds.). (1997). Workshop on machine learning applications in the real world;

Classification Tree 171

C

- methodological aspects and implications. In *Proceedings of the 14th international conference on machine learning.*
- Fayyad, U. M., & Uthurusamy, R. (Eds.). (1995). Proceedings of the first international conference on knowledge discovery and data mining.
- Holte, R. C. (1993). Very simple classification rules perform well on most commonly used datasets. *Machine Learning*, 11(1), 63–91.
- Kodratoff, Y. (Ed.). (1994). Proceedings of MLNet workshop on industrial application of machine learning.
- Kodratoff, Y., & Michalski, R. S. (1990). Machine learning: An artificial intelligence approach, (Vol. 3). San Mateo, CA: Morgan Kaufmann.
- Kohavi, R., & Provost, F. (1998). Glossary of terms. Editorial for the special issue on applications of machine learning and the knowledge discovery process. *Machine Learning*, 30(2/3).
- Komorowski, H. J., & Zytkow, J. M. (Eds.). (1997). Proceedings of the first European conference on principles of data mining and knowledge discovery.
- Lakoff, G. (1987). Women, fire and dangerous things. Chicago, IL: University of Chicago Press.
- Langley, P., & Simon, H. A. (1995). Applications of machine learning and rule induction. Communications of the ACM, 38(11), 54-64.
- Michalski, R. S. (1983). A theory and methodology of inductive learning. In R. S. Michalski, T. J. Carbonell, & T. M. Mitchell (Eds.), Machine learning: An artificial intelligence approach (pp. 83–134). Palo Alto, CA: TIOGA Publishing.
- Michalski, R. S., Carbonell, J. G., & Mitchell, T. M. (Eds.). (1983).
 Machine learning: An artificial intelligence approach. Palo Alto,
 CA: Tioga Publishing Company.
- Michalski, R. S., Carbonell, J. G., & Mitchell, T. M. (Eds.). (1986).

 Machine learning: An artificial intelligence approach, (Vol. 2).

 San Mateo, CA: Morgan Kaufmann.
- Michie, D. (1982). Machine intelligence and related topics. New York: Gordon and Breach Science Publishers.
- Mitchell, T. M. (1977). Version spaces: A candidate elimination approach to rule learning. In *Proceedings of the fifth international joint conferences on artificial intelligence* (pp. 305-310).
- Mitchell, T. M. (1997). *Machine learning*. Boston, MA: McGraw-Hill. Quinlan, J. R. (1986). Induction of decision trees. *Machine Learning*, *1*. 81–106
- Quinlan, J. R. (1993). C4.5 programs for machine learning. San Mateo, CA: Morgan Kaufmann.
- Rubinstein, Y. D., & Hastie, T. (1997). Discriminative vs informative learning. In *Proceedings of the third international conference on knowledge discovery and data mining* (pp. 49–53).
- Russell, S., & Norvig, P. (2003). Artificial intelligence: A modern approach. Upper Saddle River, NJ: Prentice-Hall.
- Schorr, H., & Rappaport, A. (Eds.). (1989). Proceedings of the first conference on innovative applications of artificial intelligence.
- Winston, P. H. (1975). Learning structural descriptions from examples. In P. H. Winston (Ed.), *The psychology of computer vision* (pp. 157–209). New York: McGraw-Hill.
- Witten, I. H., & Frank, E. (2005). Data mining: Practical machine learning tools and techniques. San Fransisco: Morgan Kaufmann.
- Wolpert, D. H., & Macready, W. G. (1997). No free lunch theorems for optimization. *IEEE Transactions on Evolutionary Computation*, *1*(1), 67–82.

Classification Algorithms

There is a very large number of classification algorithms, including ▶decision trees, ▶instance-based learners, ▶support vector machines, ▶rule-based learners, ▶neural networks, ▶Bayesian networks. There also ways of combining them into ▶ensemble classifiers such as ▶boosting, ▶bagging, ▶stacking, and ▶forests of trees.

To delve deeper into classifiers and their role in machine learning, a number of books are recommended covering machine learning (Bishop, 2007; Mitchell, 1997; Witten & Frank, 2005) and artificial intelligence (Russell & Norvig, 2003) in general. Seminal papers on classifiers can be found in collections of papers on machine learning (Dietterich & Shavlik, 1990; Kodratoff & Michalski, 1990; Michalski, Carbonell, & Mitchell, 1983, 1986).

Recommended Reading

- Bishop, C. M. (2007). Pattern recognition and machine learning. New York: Springer.
- Dietterich, T., & Shavlik, J. (Eds.). *Readings in machine learning*. San Mateo, CA: Morgan Kaufmann.
- Kodratoff, Y., & Michalski, R. S. (1990). Machine learning: An artificial intelligence approach, (Vol. 3). San Mateo, CA: Morgan Kaufmann.
- Michalski, R. S., Carbonell, J. G., & Mitchell, T. M. (Eds.). (1983).
 Machine learning: An artificial intelligence approach. Palo Alto,
 CA: Tioga Publishing Company.
- Michalski, R. S., Carbonell, J. G., & Mitchell, T. M. (Eds.). (1986). Machine learning: An artificial intelligence approach, (Vol. 2). San Mateo, CA: Morgan Kaufmann.
- Mitchell, T. M. (1997). Machine learning. Boston, MA: McGraw-Hill. Russell, S., & Norvig, P. (2003). Artificial intelligence: A modern approach. Upper Saddle River, NJ: Prentice-Hall.
- Witten, I. H., & Frank, E. (2005). Data mining: Practical machine learning tools and techniques. San Fransisco: Morgan Kaufmann.

Classification Learning

► Concept Learning

Classification Tree

▶Decision Tree

172 Classifier Systems

Classifier Systems

PIER LUCA LANZI Politecnico di Milano, Milano, Italy

Synonyms

Genetics-based machine learning; Learning classifier systems

Definition

Classifier systems are rule-based systems that combine ▶temporal difference learning or ▶supervised learning with a genetic algorithm to solve classification and ▶reinforcement learning problems. Classifier systems come in two flavors: Michigan classifier systems, which are designed for online learning, but can also tackle offline problems; and Pittsburgh classifier systems, which can only be applied to offline learning.

In Michigan classifier systems (Holland, 1976), learning is viewed as an online adaptation process to an unknown environment that represents the problem and provides feedback in terms of a numerical reward. Michigan classifier systems maintain a single candidate solution consisting of a set of rules, or a population of classifiers. Michigan systems apply (1) temporal difference learning to distribute the incoming reward to the classifiers that are accountable for it; and (2) a genetic algorithm to select, recombine, and mutate individual classifiers so as to improve their contribution to the current solution.

In contrast, in Pittsburgh classifier systems (Smith, 1980), learning is viewed as an offline optimization process in which a genetic algorithm alone is applied to search for the best solution to a given problem. In addition, Pittsburgh classifier systems maintain not one, but a set of candidate solutions. While in the Michigan classifier system each individual classifier represents a part of the overall solution, in the Pittsburgh system each individual is a complete candidate solution (itself consisting of a set of classifiers). The fitness of each Pittsburgh individual is computed offline by testing it on a representative sample of problem instances. The individuals compete among themselves through selection, while crossover and mutation recombine solutions to search for better solutions.

Motivation and Background

Machine learning is usually viewed as a search process in which a solution space is explored until an appropriate solution to the target problem is found (Mitchell, 1982) (see ▶Learning as Search). Machine learning methods are characterized by the way they represent solutions (e.g., using ▶decision trees, rules), by the way they evaluate solutions (e.g., classification accuracy, information gain) and by the way they explore the solution space (e.g., using a ▶general-to-specific strategy or a ▶specific-to-general strategy).

Classifier systems are methods of genetics-based machine learning introduced by Holland, the father of ▶ genetic algorithms. They made their first appearance in Holland (1976) where the first diagram of a classifier system, labeled "cognitive system," was shown. Subsequently, they were described in detail in the paper "Cognitive Systems based on Adaptive Algorithms" (Holland and Reitman, 1978). Classifier systems are characterized by a rule-based representation of solutions and a genetics-based exploration of the solution space. While other ▶rule learning methods, such as CN2 (Clark & Niblett, 1989) and FOIL (Quinlan & Cameron-Jones, 1995), generate one rule at a time following a sequential covering strategy (see ▶Covering Algorithm), classifier systems work on one or more solutions at once, and they explore the solution space by applying the principles of natural selection and genetics.

In classifier systems (Holland, 1976; Holland and Reitman, 1978; Wilson, 1995), machine learning is modeled as an online adaptation process to an unknown environment, which provides feedback in terms of a numerical reward. A classifier system perceives the environment through its detectors and, based on its sensations, it selects an action to be performed in the environment through its effectors. Depending on the efficacy of its actions, the environment may eventually reward the system. A classifier system learns by trying to maximize the amount of reward it receives from the environment. To pursue such a goal, it maintains a set (a population) of condition-action-prediction rules, called classifiers, which represents the current solution. Each classifier's condition identifies some part of the problem domain; the classifier's action represents a decision on the subproblem identified by its condition; and the classifier's prediction, or strength, estimates the value of the action in terms of future

Classifier Systems 173

rewards on that subproblem. Two separate components, credit assignment and rule discovery, act on the population with different goals. Credit assignment, implemented either by methods of temporal difference or supervised learning, exploits the incoming reward to estimate the action values in each subproblem so as to identify the best classifiers in the population. At the same time, rule discovery, usually implemented by a genetic algorithm, selects, recombines, and mutates the classifiers in the population to improve the current solution.

Classifier systems were initially conceived as modeling tools. Given a real system with unknown underlying dynamics, for instance a financial market, a classifier system would be used to generate a behavior that matched the real system. The evolved rules would provide a plausible, human readable model of the unknown system - a way to look inside the box. Subsequently, with the developments in the area of machine learning and the rise of reinforcement learning, classifier systems have been more and more often studied and presented as alternatives to other machine learning methods. Wilson's XCS (1995), the most successful classifier system to date, has proven to be both a valid alternative to other reinforcement learning approaches and an effective approach to classification and data mining (Bull, 2004; Bull & Kovacs, 2005; Lanzi, Stolzmann, & Wilson, 2000).

Kenneth de Jong and his students (de Jong, 1988; Smith, 1980, 1983) took a different perspective on genetics-based machine learning and modeled learning as an optimization process rather than an adaptation process as done in Holland (1976). In this case, the solution space is explored by applying a genetic algorithm to a population of individuals each representing a complete candidate solution - that is, a set of rules (or a production system, de Jong, 1988; Smith, 1980). At each cycle, a critic is applied to each individual (to each set of rules) to obtain a performance measure that is then used by the genetic algorithm to guide the exploration of the solution space. The individuals in the population compete among themselves through selection, while crossover and mutation recombine solutions to search for better ones.

The approaches of Holland (Holland, 1976; Holland and Reitman, 1978) and de Jong (de Jong, 1988; Smith, 1980, 1983) have been extended and improved

in several ways (see Lanzi et al. (2000) for a review). The models of classifier systems that are inspired by the work of Holland (1976) at the University of Michigan are usually called Michigan classifier systems; the ones that are inspired by Smith (1980, 1983) and de Jong (1988) at the University of Pittsburgh are usually termed Pittsburgh classifier systems – or briefly, Pitt classifier systems.

Pittsburgh classifier systems separate the evaluation of candidate solutions, performed by an external critic, from the genetic search. As they evaluate candidate solutions as a whole, Pittsburgh classifier systems can easily identify and emphasize sequentially cooperating classifiers, which is particularly helpful in problems involving partial observability. In contrast, in Michigan classifier systems the credit assignment is focused, due to identification of the actual classifiers that produce the reward, so learning is much faster but sequentially cooperating classifiers are more difficult to spot. As Pittsburgh classifier systems apply the genetic algorithm to a set of solutions, they only work offline, whereas Michigan classifier systems work online, although they can also tackle offline problems. Finally, the design of Pittsburgh classifier systems involves decisions as to how an entire solution should be represented and how solutions should be recombined – a task which can be daunting. In contrast, the design of Michigan classifier systems involves simpler decisions about how a rule should be represented and how two rules should be recombined. Accordingly, while the representation of solutions and its related issues play a key role in Pittsburgh models, Michigan models easily work with several types of representations (Lanzi, 2001; Lanzi & Perrucci, 1999; Mellor, 2005).

Structure of the Learning System

Michigan and Pittsburgh classifier systems were both inspired by the work of Holland on the broadcast language (Holland, 1975). However, their structures reflect two different ways to model machine learning: as an adaptation process in the case of Michigan classifier systems; and as an optimization problem, in the case of Pittsburgh classifier systems. Thus, the two models, originating from the same idea (Holland's broadcast language), have radically different structures.

C

174 Classifier Systems

Michigan Classifier Systems

Holland's classifier systems define a general paradigm for genetics-based machine learning. The description in Holland and Reitman (1978) provides a list of principles for online learning through adaptation. Over the years, such principles have guided researchers who developed several models of Michigan classifier systems (Butz, 2002; Wilson, 1994, 1995, 2002) and applied them to a large variety of domains (Bull, 2004; Lanzi & Riolo, 2003; Lanzi et al., 2000). These models extended and improved Holland's original ideas, but kept all the ingredients of the original recipe: a population of classifiers, which represents the current system knowledge; a performance component, which is responsible for the short-term behavior of the system; a credit assignment (or reinforcement) component, which distributes the incoming reward among the classifiers; and a rule discovery component, which applies a genetic algorithm to the classifiers to improve the current knowledge.

Knowledge Representation

In Michigan classifier systems, knowledge is represented by a population of classifiers. Each classifier is usually defined by four main parameters: the *condition*, which identifies some part of the problem domain; the *action*, which represents a decision on the subproblem identified by its condition; the *prediction* or strength, which estimates the amount of reward that the system will receive if its action is performed; and finally, the *fitness*, which estimates how good the classifier is in terms of problem solution.

The knowledge representation of Michigan classifier systems is extremely flexible. Each one of the four classifier components can be tailored to fit the need of a particular application, without modifying the main structure of the system. In problems involving binary inputs, classifier conditions can be simply represented using strings defined over the alphabet {0, 1, #}, as done in Holland and Reitman (1978), Goldberg (1989), and Wilson (1995). In problems involving real inputs, conditions can be represented as disjunctions of intervals, similar to the ones produced by other rule learning methods (Clark & Niblett, 1989) Conditions can also be represented as general-purpose symbolic expressions

(Lanzi, 2001; Lanzi & Perrucci, 1999) or first-order logic expressions (Mellor, 2005). Classifier actions are typically encoded by a set of symbols (either binary strings or simple labels), but continuous real-valued actions are also available (Wilson, 2007). Classifier prediction (or strength) is usually encoded by a parameter (Goldberg, 1989; Holland & Reitman, 1978; Wilson, 1995). However, classifier prediction can also be computed using a parameterized function (Wilson, 2002), which results in solutions represented as an ensemble of local approximators – similar to the ones produced in generalized reinforcement learning (Sutton & Barto, 1998).

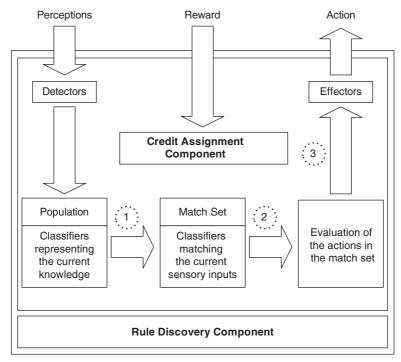
Performance Component

A simplified structure of Michigan classifier systems is shown in Fig. 1. We refer the reader to Goldberg (1989) and Holland and Reitman (1978) for a detailed description of the original model and to Butz (2002) and Wilson (1994, 1995, 2001) for descriptions of recent classifier system models.

A classifier system learns through trial and error interactions with an unknown environment. The system and the environment interact continually. At each time step, the classifier system perceives the environment through its detectors; it builds a *match set* containing all the classifiers in the population whose condition matches the current sensory input. The match set typically contains classifiers that advocate contrasting actions; accordingly, the classifier system evaluates each action in the match set, and selects an action to be performed balancing exploration and exploitation. The selected action is sent to the effectors to be executed in the environment; depending on the effect that the action has in the environment, the system receives a scalar reward.

Credit Assignment

The credit assignment component (also called reinforcement component, Wilson 1995) distributes the incoming reward to the classifiers that are accountable for it. In Holland and Reitman (1978), credit assignment is implemented by Holland's bucket brigade algorithm (Holland, 1986), which was partially inspired by the credit allocation mechanism used by Samuel in his



Classifier Systems. Figure 1. Simplified structure of a Michigan classifier system. The system perceives the environment through its detectors and (1) it builds the match set containing the classifiers in the population that match the current sensory inputs; then (2) all the actions in the match set are evaluated, and (3) an action is selected to be performed in the environment through the effectors

pioneering work on learning checkers-playing programs (Samuel, 1959).

In the early years, classifier systems and the bucket brigade algorithm were confined to the evolutionary computation community. The rise of reinforcement learning increased the connection between classifier systems and temporal difference learning (Sutton, 1988; Sutton & Barto, 1998): in particular, Sutton (1988) showed that the bucket brigade algorithm is a kind of temporal difference learning, and similar connections were also made in Watkins (1989) and Dorigo and Bersini (1994). Later, the connection between classifier systems and reinforcement learning became tighter with the introduction of Wilson's XCS (1995), in which credit assignment is implemented by a modification of Watkins Q-learning (Watkins, 1989). As a consequence, in recent years, classifier systems are often presented as methods of reinforcement learning with genetics-based generalization (Bull & Kovacs, 2005).

Rule Discovery Component

The *rule discovery component* is usually implemented by a genetic algorithm that selects classifiers in the population with probability proportional to their fitness; it copies the selected classifiers and applies genetic operators (usually crossover and mutation) to the offspring classifiers; the new classifiers are inserted in the population, while other classifiers are deleted to keep the population size constant.

Classifiers selection plays a central role in rule discovery. Classifier selection depends on the definition of classifier fitness and on the subset of classifiers considered during the selection process. In Holland and Reitman (1978), classifier fitness coincides with classifier prediction, while selection is applied to all the classifiers in the population. This approach results in a pressure toward classifiers predicting high returns, but typically tends to produce overly general solutions. To avoid such solutions, Wilson (1995) introduced the XCS classifier system in which accuracy-based fitness is

176 Classifier Systems

coupled with a niched genetic algorithm. This approach results in a pressure toward accurate maximally general classifiers, and has made XCS the most successful classifier system to date.

Pittsburgh Classifier Systems

The idea underlying the development of Pittsburgh classifier systems was to show that interesting behaviors could be evolved using a simpler model than the one proposed by Holland with Michigan classifier systems (Holland, 1976; Holland & Reitman, 1978).

In Pittsburgh classifier systems, each individual is a set of rules that encodes an entire candidate solution; each rule has a fixed length, but each rule set (each individual) usually contains a variable number of rules. The genetic operators, crossover and mutation, are tailored to the rule-based, variable-length representation. The individuals in the population compete among themselves, following the selection-recombination-mutation cycle that is typical of genetic algorithms (Goldberg, 1989; Holland, 1975). While in Michigan classifier systems individuals in the population (the single rules) cooperate, in Pittsburgh classifier systems there is no cooperation among individuals (the rule sets), so that the genetic algorithm operation is simpler for Pittsburgh models. However, as Pittsburgh classifier systems explore a much larger search space, they usually require more computational resources than Michigan classifier systems.

The pseudo-code of a Pittsburgh classifier system is shown in Fig. 2. At first, the individuals in the population are randomly initialized (line 2). At time t, the

individuals are evaluated by an external critic, which returns a performance measure that the genetic algorithm exploits to compute the fitness of individuals (lines 3 and 10). Following this, selection (line 6), recombination, and mutation (line 7) are applied to the individuals in the population – as done in a typical genetic algorithm. The process stops when a termination criterion is met (line 4), usually when an appropriate solution is found.

The design of Pittsburgh classifier systems follows the typical steps of genetic algorithm design, which means deciding how a rule set should be represented, what genetic operators should be applied, and how the fitness of a set of rules should be calculated. In addition, Pittsburgh classifier systems need to address the bloat phenomenon (Tackett, 1994) that arises with any variable-sized representation, like the rule sets evolved by Pittsburgh classifier systems. Bloat can be defined as the growth of individuals without an actual fitness improvement. In Pittsburgh classifier systems, bloat increases the size of candidate solutions by adding useless rules to individuals, and it is typically limited by introducing a parsimony pressure that discourages large rule sets (Bassett & de Jong, 2000). Alternatively, Pittsburgh classifier systems can be combined with multi-objective optimization, so as to separate the maximization of the rule set performance and the minimization of the rule set size.

Examples of Pittsburgh classifier systems include SAMUEL (Grefenstette, Ramsey, & Schultz, 1990), the Genetic Algorithm Batch-Incremental Concept Learner (GABIL) (de Jong & Spears, 1991), GIL (Janikow, 1993), GALE (Llorà, 2002), and GAssist (Bacardit, 2004).

```
1.
     t := 0
2.
     Initialize the population P(t)
     Evaluate the rules sets in P(t)
3.
     While the termination condition is not satisfied
4.
5.
     Begin
6.
        Select the rule sets in P(t) and generate Ps(t)
7.
        Recombine and mutate the rule sets in Ps(t)
8.
        P(t+1) := Ps(t)
9.
        t := t+1
10.
        Evaluate the rules sets in P(t)
11.
```

Classifier Systems 177

C

Applications

Classifier systems have been applied to a large variety of domains, including computational economics (e.g., Arthur, Holland, LeBaron, Palmer, & Talyer, 1996), autonomous robotics (e.g., Dorigo & Colombetti, 1998), classification (e.g., Barry, Holmes, & Llora, 2004), fighter aircraft maneuvering (Bull, 2004; Smith, Dike, Mehra, Ravichandran, & El-Fallah, 2000), and many others. Reviews of classifier system applications are available in Lanzi et al. (2000), Lanzi and Riolo (2003), and Bull (2004).

Programs and Data

The major sources of information about classifier systems are the LCSWeb maintained by Alwyn Barry, which can be reached through, and www.learning-classifier-systems.org_maintained by Xavier Llorà.

Several implementations of classifier systems are freely available online. The first standard implementation of Holland's classifier system in Pascal was described in Goldberg (1989), and it is available at http://www.illigal.org/; a C version of the same implementation, developed by Robert E. Smith, is available at http://www.etsimo.uniovi.es/ftp/pub/EC/CFS/src/. Another implementation of an extension of Holland's classifier system in C by Rick L. Riolo is available at http://www.cscs.umich.edu/Software/Contents. html. Implementations of Wilson's XCS (1995) are distributed by Alwyn Barry at the LCSWeb, by Martin V. Butz (at www.illigal.org), and by Pier Luca Lanzi (at xcslib.sf.net). Among the implementations of Pittsburgh classifier systems, the Samuel system is available from Alan C. Schultz at http://www.nrl.navy.mil/; Xavier Llorà distributes GALE (Genetic and Artificial Life Environment) a fine-grained parallel genetic algorithm for data mining at www.illigal.org/xllora.

Cross References

- ► Credit Assignment
- ► Genetic Algorithms
- ▶ Reinforcement Learning
- ►Rule Learning

Recommended Reading

Arthur, B. W., Holland, J. H., LeBaron, B., Palmer, R., & Talyer, P. (1996). Asset pricing under endogenous expectations in an artificial stock market. Technical Report, Santa Fe Institute.

- Bacardit i Peñarroya, J. (2004). Pittsburgh genetic-based machine learning in the data mining era: Representations, generalization, and run-time. PhD thesis, Computer Science Department, Enginyeria i Arquitectura La Salle Universitat Ramon Llull, Barcelona.
- Barry, A. M., Holmes, J., & Llora, X. (2004). Data mining using learning classifier systems. In L. Bull (Ed.), Applications of learning classifier systems, studies in fuzziness and soft computing (Vol. 150, pp. 15–67). Pagg: Springer.
- Bassett, J. K., & de Jong, K. A. (2000). Evolving behaviors for cooperating agents. In Proceedings of the twelfth international symposium on methodologies for intelligent systems, LNAI (Vol. 1932). Berlin: Springer.
- Booker, L. B. (1989). Triggered rule discovery in classifier systems. In J. D. Schaffer (Ed.), Proceedings of the 3rd international conference on genetic algorithms (ICGA89). San Francisco: Morgan Kaufmann.
- Bull, L. (Ed.). (2004). Applications of learning classifier systems, studies in fuzziness and soft computing (Vol. 150). Berlin: Springer, ISBN 978-3-540-21109-9.
- Bull, L., & Kovacs, T. (Eds.). (2005). Foundations of learning classifier systems, studies in fuzziness and soft computing (Vol. 183). Berlin: Springer, ISBN 978-3-540-25073-9.
- Butz, M. V. (2002). Anticipatory learning classifier systems. Genetic algorithms and evolutionary computation. Boston, MA: Kluwer Academic Publishers.
- Clark, P., & Niblett, T. (1989). The CN2 induction algorithm. Machine Learning, 3(4), 261–283.
- de Jong, K. (1988). Learning with genetic algorithms: An overview. Machine Learning, 3(2-3), 121-138.
- de Jong, K. A., & Spears, W. M. (1991). Learning concept classification rules using genetic algorithms. In *Proceedings of the international joint conference on artificial intelligence* (pp. 651–656). San Francisco: Morgan Kaufmann.
- Dorigo, M., & Bersini, H. (1994). A comparison of Q-learning and classifier systems. In D. Cliff, P. Husbands, J.-A. Meyer, & S. W. Wilson (Eds.), From animals to animats 3: Proceedings of the third international conference on simulation of adaptive behavior (pp. 248-255). Cambridge, MA: MIT Press.
- Dorigo, M., & Colombetti, M. (1998). Robot shaping: An experiment in behavior engineering. Cambridge, MA: MIT Press/Bradford Books.
- Goldberg, D. E. (1989). Genetic algorithms in search, optimization, and machine learning. Reading, MA: Addison-Wesley.
- Grefenstette, J. J., Ramsey, C. L., & Schultz, A. (1990) Learning sequential decision rules using simulation models and competition. *Machine Learning*, 5(4), 355–381.
- Holland, J. (1986) Escaping brittleness: The possibilities of general-purpose learning algorithms applied to parallel rule-based systems. In R. S. Michalski, J. G. Carbonell, & T. M. Mitchell (Eds.), Machine learning, an artificial intelligence approach (Vol. II, Chap. 20) (pp. 593–623). San Francisco: Morgan Kaufmann.
- Holland, J. H. (1975). Adaptation in natural and artificial systems. Ann Arbor, MI: University of Michigan Press (Reprinted by the MIT Press in 1992).
- Holland, J. H. (1976). Adaptation. Progress in Theoretical Biology, 4, 263–293.
- Holland, J. H., & Reitman, J. S. (1978). Cognitive systems based on adaptive algorithms. In D. A. Waterman & F. Hayes-Roth (Eds.), Pattern-directed inference systems. New York: Academic Press.

178 Clause

- (Reprinted from Evolutionary computation. The fossil record. D. B. Fogel (Ed.), IEEE Press (1998)).
- Janikow, C. Z. (1993). A knowledge-intensive genetic algorithm for supervised learning. *Machine Learning*, 13(2-3), 189-228.
- Lanzi, P. L. (2001). Mining interesting knowledge from data with the XCS classifier system. In L. Spector, E. D. Goodman, A. Wu, W. B. Langdon, H.-M. Voigt, M. Gen, et al. (Eds.), Proceedings of the genetic and evolutionary computation conference (GECCO-2001) (pp. 958-965). San Francisco: Morgan Kaufmann.
- Lanzi, P. L. (2005). Learning classifier systems: A reinforcement learning perspective. In L. Bull & T. Kovacs (Eds.), Foundations of learning classifier systems, studies in fuzziness and soft computing (pp. 267-284). Berlin: Springer.
- Lanzi, P. L., & Perrucci, A. (1999). Extending the representation of classifier conditions part II: From messy coding to Sexpressions. In W. Banzhaf, J. Daida, A. E. Eiben, M. H. Garzon, V. Honavar, M. Jakiela, & R. E. Smith (Eds.), Proceedings of the genetic and evolutionary computation conference (GECCO 99) (pp. 345-352). Orlando, FL: Morgan Kaufmann.
- Lanzi, P. L., & Riolo, R. L. (2003). Recent trends in learning classifier systems research. In A. Ghosh & S. Tsutsui (Eds.), Advances in evolutionary computing: Theory and applications (pp. 955–988). Berlin: Springer.
- Lanzi, P. L., Stolzmann, W., & Wilson, S. W. (Eds.). (2000). Learning classifier systems: From foundations to applications. Lecture notes in computer science (Vol. 1813). Berlin: Springer.
- Llorá, X. (2002). Genetics-based machine learning using fine-grained parallelism for data mining. PhD thesis, Enginyeria i Arquitectura La Salle, Ramon Llull University, Barcelona.
- Mellor, D. (2005). A first order logic classifier system. In H. Beyer (Ed.), Proceedings of the 2005 conference on genetic and evolutionary computation (GECCO '05), (pp. 1819–1826). New York: ACM Press.
- Quinlan, J. R., & Cameron-Jones, R. M. (1995). Induction of logic programs: FOIL and related systems. New Generation Computing, 13(3&4), 287-312.
- Samuel, A. L. (1959). Some studies in machine learning using the game of checkers. In E. A. Feigenbaum & J. Feldman (Eds.), Computers and thought. New York: McGraw-Hill.
- Smith, R. E., Dike, B. A., Niehra, R. K., Ravichandran, B., & El-Fallah, A. (2000). Classifier systems in combat: Two-sided learning of maneuvers for advanced fighter aircraft. Computer Methods in Applied Mechanics and Engineering, 186(2-4), 421-437
- Smith, S. F. (1980) A learning system based on genetic adaptive algorithms. Doctoral dissertation, Department of Computer Science, University of Pittsburgh.
- Smith, S. F. (1983). Flexible learning of problem solving heuristics through adaptive search. In Proceedings of the eighth international joint conference on artificial intelligence (pp. 421–425). Los Altos, CA: Morgan Kaufmann.
- Sutton, R. S. (1988). Learning to predict by the methods of temporal differences. *Machine Learning*, 3, 9–44.
- Sutton, R. S., & Barto, A. G. (1998). Reinforcement learning: An introduction. Cambridge, MA: MIT Press.
- Tackett, W. A. (1994). Recombination, selection, and the genetic construction of computer programs. Unpublished doctoral dissertation, University of Southern California.
- Watkins, C. (1989). Learning from delayed rewards. PhD thesis, King's College.

- Wilson, S. W. (1995). Classifier fitness based on accuracy. *Evolutionary Computation*, 3(2), 149–175.
- Wilson, S. W. (2002). Classifiers that approximate functions. *Natural Computing*, 1(2-3), 211-234.
- Wilson, S. W. (2007). "Three architectures for continuous action" learning classifier systems. International workshops, IWLCS 2003-2005, revised selected papers. In T. Kovacs, X. Llorà, K. Takadama, P. L. Lanzi, W. Stolzmann, & S. W. Wilson (Eds.), Lecture notes in artificial intelligence 4399 Vol. (pp. 239-257). Berlin: Springer.

Clause

A *clause* is a logical rule in a ▶logic program. Formally, a clause is a disjunction of (possibly negated) literals, such as

```
grandfather(x, y) \lor \neg father(x, z) \lor \neg parent(z, y).
```

In the logic programming language ▶Prolog this clause is written as

```
grandfather(X,Y) :- father(X,Z),
parent(Z,Y).
```

The part to the left of : - ("if") is the *head* of the clause, and the right part is its *body*. Informally, the clause asserts the truth of the head given the truth of the body. A clause with exactly one literal in the head is called a *Horn clause* or *definite clause*; logic programs mostly consist of definite clauses. A clause without a body is also called a *fact*; a clause without a head is also called a *denial*, or a *query* in a proof by refutation. The clause without head or body is called the *empty clause*: it signifies inconsistency or falsehood and is denoted □. Given a set of clauses, the *resolution* inference rule can be used to deduce logical consequences and answer queries (see ►First-Order Logic).

In machine learning, clauses can be used to express classification rules for structured individuals. For example, the following definite clause classifies a molecular compound as carcinogenic if it contains a hydrogen atom with charge above a certain threshold.

Cluster Optimization 179

C

Cross References

- ▶First-Order Logic
- ►Inductive Logic Programming
- ► Learning from Structured Data
- ►Logic Program
- **▶**Prolog

Clause Learning

In ▶speedup learning, clause learning is a ▶deductive learning technique used for the purpose of ▶intelligent backtracking in satisfiability solvers. The approach analyzes failures at backtracking points and derives clauses that must be satisfied by the solution. The clauses are added to the set of clauses from the original satisfiability problem and serve to prune new search nodes that violate them.

Click-Through Rate (CTR)

CTR measures the success of a ranking of search results, or advertisement placing. Given the number of *impressions*, the number of times a web result or ad has been displayed, and the number of *clicks*, the number of users who clicked on the result/advertisement, CTR is the number of clicks divided by the number of impressions.

Clonal Selection

The clonal selection theory (CST) is the theory used to explain the basic response of the adaptive immune system to an antigenic stimulus. It establishes the idea that only those cells capable of recognizing an antigenic stimulus will proliferate, thus being selected against those that do not. Clonal selection operates on both T-cells and B-cells. When antibodies on a B-cell bind with an antigen, the B-cell becomes activated and begins to proliferate. New B-cell clones are produced that are an exact copy of the parent B-cell, but then they undergo somatic hypermutation and produce antibodies that are specific to the invading antigen. The B-cells, in addition to proliferating or differentiating into plasma cells, can differentiate into long-lived B memory cells. Plasma cells produce large amounts of antibody which will attach

themselves to the antigen and act as a type of *tag* for T-cells to pick up on and remove from the system. This whole process is known as *affinity maturation*. This process forms the basis of many artificial immune system algorithms such as AIRS and aiNET.

Closest Point

► Nearest Neighbor

Cluster Editing

The Cluster Editing problem is almost equivalent to Correlation Clustering on complete instances. The idea is to obtain a graph that consists only of cliques. Although Cluster Deletion requires us to delete the smallest number of edges to obtain such a graph, in Cluster Editing we are permitted to add as well as remove edges. The final variant is Cluster Completion in which edges can only be added: each of these problems can be restricted to building a specified number of cliques.

Cluster Ensembles

Cluster ensembles are an unsupervised Pensemble learning method. The principle is to create multiple different clusterings of a dataset, possibly using different algorithms, then aggregate the opinions of the different clusterings into an ensemble result. The final ensemble clustering should be in theory more reliable than the individual clusterings.

Cluster Optimization

► Evolutionary Clustering

180 Clustering

Clustering

Clustering is a type of variety of variety of variety called clusters. Intuitively, the examples within a cluster are more similar to each other than to examples from other clusters. In order to measure the similarity between examples, clustering algorithms use various distortion or distance measures. There are two major types clustering approaches: generative and discriminative. The former assumes a parametric form of the data and tries to find the model parameters that maximize the probability that the data was generated by the chosen model. The latter represents graph-theoretic approaches that compute a similarity matrix defined over the input data.

Cross References

- ► Categorical Data Clustering
- **►**Cluster Editing
- ►Cluster Ensembles
- ► Clustering from Data Streams
- **▶**Constrained Clustering
- ► Consensus Clustering
- ► Correlation Clustering
- ► Cross-Language Document Clustering
- ▶ Density-Based Clustering
- **▶**Dirichlet Process
- ► Document Clustering
- ► Evolutionary Clustering
- ► Graph Clustering
- ▶ *k*-Means Clustering
- ▶ *k*-Mediods Clustering
- ► Model-Based Clustering
- ▶ Partitional Clustering
- ▶ Projective Clustering
- ► Sublinear Clustering

Clustering Aggregation

► Consensus Clustering

Clustering Ensembles

► Consensus Clustering

Clustering from Data Streams

João Gama University of Porto, Porto, Portugal

Definition

▶ Clustering is the process of grouping objects into different groups, such that the common properties of data in each subset is high, and between different subsets is low. The data stream clustering problem is defined as to maintain a consistent good clustering of the sequence observed so far, using a small amount of memory and time. The issues are imposed by the continuous arriving data points, and the need to analyze them in real time. These characteristics require incremental clustering, maintaining cluster structures that evolve over time. Moreover, the data stream may evolve over time and new clusters might appear, others disappear reflecting the dynamics of the stream.

Main Techniques

Major clustering approaches in data stream cluster analysis include:

- *Partitioning* algorithms: construct a partition of a set of objects into *k* clusters, that minimize some objective function (e.g., the sum of squares distances to the centroid representative). Examples include *k*-means (Farnstrom, Lewis, & Elkan, 2000), and *k*-medoids (Guha, Meyerson, Mishra, Motwani, & O'Callaghan, 2003)
- Microclustering algorithms: divide the clustering process into two phases, where the first phase is online and summarizes the data stream in local models (microclusters) and the second phase generates a global cluster model from the microclusters. Examples of these algorithms include BIRCH (Zhang, Ramakrishnan, & Livny, 1996) and CluStream (Aggarwal, Han, Wang, & Yu, 2003)

Basic Concepts

A powerful idea in clustering from data streams is the concept of *cluster feature*, CF. A cluster feature, or *micro-cluster*, is a compact representation of a set of points. A CF structure is a triple (N, LS, SS), used to store the sufficient statistics of a set of points:

181

- *N* is the number of data points
- *LS* is a vector, of the same dimension of data points, that store the linear sum of the *N* points
- *SS* is a vector, of the same dimension of data points, that store the square sum of the *N* points

The properties of cluster features are:

• Incrementality

If a point *x* is added to the cluster, the sufficient statistics are updated as follows:

$$LS_A \leftarrow LS_A + x$$

$$SS_A \leftarrow SS_A + x^2$$

$$N_A \leftarrow N_A + 1$$
.

Additivity

If A_1 and A_2 are disjoint sets, merging them is equal to the sum of their parts. The additive property allows us to merge subclusters incrementally.

$$LS_C \leftarrow LS_A + LS_B,$$

$$SS_C \leftarrow SS_A + SS_B$$

$$N_C \leftarrow N_A + N_B$$
.

A CF entry has sufficient information to calculate the norms

$$L_1 = \sum_{i=1}^{n} |x_{a_i} - x_{b_i}|,$$

$$L_2 = \sqrt{\sum_{i=1}^{n} (x_{a_i} - x_{b_i})^2}$$

and basic measures to characterize a cluster.

• *Centroid*, defined as the gravity center of the cluster:

$$\vec{X}0 = \frac{LS}{N}$$
.

 Radius, defined as the average distance from member points to the centroid:

$$R = \sqrt{\frac{\sum_{1}^{N} (\vec{x}_{i} - \vec{X}0)^{2}}{N}}.$$

Partitioning Clustering

k-means is the most widely used clustering algorithm. It constructs a partition of a set of objects into k clusters that minimize some objective function, usually a squared error function, which imply round-shape clusters. The input parameter k is fixed and must be given in advance that limits its real applicability to streaming and evolving data.

Farnstrom et al. (2000) proposed a single pass k-means algorithm. The main idea is to use a buffer where points of the dataset are kept compressed. The data stream is processed in blocks. All available space on the buffer is filled with points from the stream. Using these points, find k centers such that the sum of distances from data points to their closest center is minimized. Only the *k* centroids (representing the clustering results) are retained, with the corresponding k cluster features. In the following iterations, the buffer is initialized with the k-centroids, found in previous iteration, weighted by the k cluster features, and incoming data points from the stream. The Very Fast k-means (VFKM) algorithm (Domingos & Hulten, 2001) uses the Hoeffding bound to determine the number of examples needed in each step of a k-means algorithm. VFKM runs as a sequence of k-means runs, with increasing number of examples until the Hoeffding bound is satisfied.

Guha et al. (2003) present an analytical study on k-median clustering data streams. The proposed algorithm makes a single pass over the data stream and uses small space. It requires O(nk) time and $O(n\epsilon)$ space where k is the number of centers, n is the number of points, and $\epsilon < 1$. They have proved that any k-median algorithm that achieves a constant factor approximation cannot achieve a better run time than O(nk).

Micro Clustering

The idea of dividing the clustering process into two layers, where the first layer generates local models (microclusters) and the second layer generates global models from the local ones, is a powerful idea that has been used elsewhere.

The BIRCH system (Zhang et al., 1996) builds a hierarchical structure of data, the CF-tree, where each node contains a set of cluster features. These CF's contain the sufficient statistics describing a set of points in the data set, and all information of the cluster features below in

182 Clustering from Data Streams

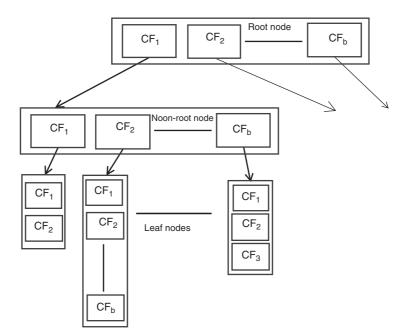
the tree. The system requires two user defined parameters: B the branch factor or the maximum number of entries in each non-leaf node; and T the maximum diameter (or radius) of any CF in a leaf node. The maximum diameter T defines the examples that can be absorbed by a CF. Increasing T, more examples can be absorbed by a micro-cluster and smaller CF-Trees are generated (Fig. 1).

When an example is available, it traverses down the current tree from the root it finds the appropriate leaf. At each non-leaf node, the example follow the *closest*-CF path, with respect to norms L_1 or L_2 . If the closest-CF in the leaf cannot absorb the example, make a new CF entry. If there is no room for new leaf, split the parent node. A leaf node might be expanded due to the constraints imposed by B and T. The process consists of taking the two farthest CFs and creates two new leaf nodes. When traversing backup the CFs are updated.

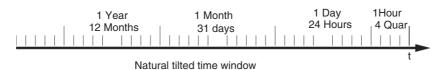
Monitoring the Evolution of the Cluster Structure

The *CluStream* Algorithm (Aggarwal et al., 2003) is an extension of the BIRCH system designed for data streams. Here, the CFs include temporal information: the time-stamp of an example is treated as a feature. CFs are initialized offline, using a standard *k*-means, with a large value for *k*. For each incoming data point, the distance to the centroids of existing CFs are computed. The data point is absorbed by an existing CF if the distance to the centroid falls within the *maximum boundary* of the CF. The *maximum boundary* is defined as a factor *t* of the *radius* deviation of the CF; otherwise, the data point starts a new micro-cluster.

CluStream can generate approximate clusters for any user defined time granularity. This is achieved by storing the CFT at regular time intervals, referred to as snapshots. Suppose the user wants to find clusters in the stream based on a history of length h, the off-line



Clustering from Data Streams. Figure 1. The clustering feature tree in BIRCH. **B** is the maximum number of CFs in a level of the tree



Clustering from Data Streams. Figure 2. The figure presents a *natural tilted time window*. The most recent data is stored with high-detail, older data is stored in a compressed way. The degree of detail decreases with time

Coevolution 183

C

component can analyze the snapshots stored at the snapshots t, the current time, and (t - h) by using the addictive property of CFT. An important problem is when to store the snapshots of the current set of microclusters. For example, the natural time frame (Fig. 2) stores snapshots each quarter, four quarters are aggregated in hours, 24 h are aggregated in days, etc. The aggregation level is domain-dependent and explores the addictive property of CFT.

Tracking the Evolution of the Cluster Structure

Promising research lines are tracking change in clusters. Spiliopoulou, Ntoutsi, Theodoridis, and Schult (2006) present system MONIC, for detecting and tracking change in clusters. MONIC assumes that a cluster is an object in a geometric space. It encompasses changes that involve more than one cluster, allowing for insights on cluster change in the whole clustering. The transition tracking mechanism is based on the degree of overlapping between the two clusters. The concept of overlap between two clusters, X and Y, is defined as the normed number of common records weighted with the age of the records. Assume that cluster X was obtained at time t_1 and cluster Y at time t_2 . The degree of overlapping between the two clusters is given by: overlap $(X, Y) = \sum_{a \in X \cap Y} age(a, t_2) / \sum_{x \in X} age(x, t_2)$. The degree of overlapping allows inferring properties of the underlying data stream. Cluster transition at a given time point is a change in a cluster discovered at an earlier timepoint. MONIC considers transitions as Internal and external transitions, that reflect the dynamics of the stream. Examples of cluster transitions include: the cluster survives, the cluster is absorbed; a cluster disappears; a new cluster emerges (Fig. 2).

Recommended Reading

Aggarwal, C., Han, J., Wang, J., & Yu, P. (2003). A framework for clustering evolving data streams. In *Proceedings of the 29th international conference on very large data bases* (pp. 81–92). San Mateo, MA: Morgan Kaufmann.

Domingos, P., & Hulten, G. (2001). A general method for scaling up machine learning algorithms and its application to clustering. In *Proceedings of international conference on machine learning* (pp. 106–113). San Mateo, MA: Morgan Kaufmann.

Farnstrom, F., Lewis, J., & Elkan, C. (2000). Scalability for clustering algorithms revisited. SIGKDD Explorations, 2(1), 51–57.

Guha, S., Meyerson, A., Mishra, N., Motwani, R., & O'Callaghan, L. (2003). Clustering data streams: Theory and practice. IEEE Transactions on Knowledge and Data Engineering, 15(3), 515-528. Spiliopoulou, M., Ntoutsi, I., Theodoridis, Y., & Schult, R. (2006). Monic: Modeling and monitoring cluster transitions. In Proceedings of ACM SIGKDD international conference on knowledge discovery and data mining (pp. 706-711). New York: ACM Press.

Zhang, T., Ramakrishnan, R., & Livny, M. (1996). Birch: An efficient data clustering method for very large databases. In *Proceedings* of ACM SIGMOD international conference on management of data (pp. 103-114). New York: ACM Press.

Clustering of Nonnumerical Data

► Categorical Data Clustering

Clustering with Advice

► Correlation Clustering

Clustering with Constraints

► Correlation Clustering

Clustering with Qualitative Information

► Correlation Clustering

Clustering with Side Information

► Correlation Clustering

CN₂

► Rule Learning

^l Co-Training

► Semi-Supervised Learning

Coevolution

► Coevolutionary Learning

184 Coevolutionary Computation

Coevolutionary Computation

► Coevolutionary Learning

Coevolutionary Learning

R. PAUL WIEGAND University of Central Florida, Orlando, FL, USA

Synonyms

Coevolution; Coevolutionary computation

Definition

Coevolutionary learning is a form of evolutionary learning (see Fevolutionary Algorithms) in which the fitness evaluation is based on interactions between individuals. Since the evaluation of an individual is dependent on interactions with other evolving entities, changes in the set of entities used for evaluation can affect an individual's ranking in a population. In this sense, coevolutionary fitness is *subjective*, while fitness in traditional evolutionary learning systems typically uses an *objective* performance measure.

Motivation and Background

Ideally, coevolutionary learning systems focus on relevant areas of a search space by making adaptive changes between interacting, concurrently evolving parts. This can be particularly helpful when problem spaces are very large - infinite search spaces in particular. Additionally, coevolution is useful when applied to problems when no intrinsic objective measure exists. The interactive nature of evaluation makes them natural methods to consider for problems such as the search for gameplaying strategies (Fogel, 2001). Finally, some coevolutionary systems appear natural for search spaces which contain certain kinds of complex structures (Potter, 1997; Stanley, 2004), since search on smaller components in a larger structure can be emphasized. In fact, there is reason to believe that coevolutionary systems may be well suited for uncovering complex structures within a problem (Bucci & Pollack, 2002).

Still, the dynamics of coevolutionary learning can be quite complex, and a number of pathologies often plague naïve users. Indeed, because of the subjective nature of coevolution, it can be easy to apply a particular coevolutionary learning system without a clear understanding of what kind of solution one expects a coevolutionary algorithm to produce. Recent theoretical analysis suggests that a clear concept of solution and a careful implementation of an evaluation process consistent with this concept can produce a coevolutionary system capable of addressing many problems (de Jong & Pollack, 2004; Ficici, 2004; Panait, 2006; Wiegand, 2003). Accordingly, a great deal of research in this area focuses on evaluation and progress measurement.

Structure of Learning System

Coevolutionary learning systems work in much the same way that an evolutionary learning system works: individuals encode some aspect of potential solutions to a problem, those representatives are altered during search using genetic-like operators such as mutation and crossover, and the search is directed by selecting better individuals as determined by some kind of fitness assessment. These heuristic methods gradually refine solutions by repeatedly cycling through such steps, using the ideas of heredity and survival of the fittest to produce new generations of individuals, with increased quality of solution. Just as in traditional evolutionary computation, there are many choices available to the engineer in designing such systems. The reader is referred to the chapters relating to evolutionary learning for more details.

However, there are some fundamental differences between traditional evolution and coevolution. In coevolution, measuring fitness requires evaluating the interaction between multiple individuals. Interacting individuals may reside in the same population or in different populations; the interactive nature of coevolution evokes notions of cooperation and competition in entirely new ways; the choices regarding how to best conduct evaluation of these interactions for the purposes of selection are particularly important; and there are unique coevolutionary issues surrounding representation. In addition, because of its interactive nature, the dynamics of coevolution can lead to some well-known pathological behaviors, and particularly careful attention to implementation choices to avoid such conditions is generally necessary.

Multiple Versus Single Population Approaches

Coevolution can typically be broadly classified as to whether interacting individuals reside in different populations or in the same population.

In the case of multipopulation coevolution, measuring fitness requires evaluating how individuals in one population interact with individuals in another. For example, individuals in each population may represent potential strategies for particular players of a game, they may represent roles in a larger ecosystem (e.g., predators and prey), or they may represent components that are fitted into a composite assembly with other component then applied to a problem. Though individuals in different populations interact for the purposes of evaluation, they are typically otherwise independent of one another in the coevolutionary search process.

In single population coevolution, an individual in the population is evaluated based on his or her interaction with other individuals in the same population. Such individuals may again represent potential strategies in a game, but evaluation may require them to trade off roles as to which player they represent in that game. Here, individuals interact not only for evaluation, but also implicitly compete with one another as resources used in the coevolutionary search process itself.

There is some controversy in the field as to whether this latter type qualifies as "coevolution." Evolutionary biologists often define coevolution exclusively in terms of multiple populations; however, in biological systems, fitness is always subjective, while the vast majority of computational approaches to evolutionary learning involve objective fitness assessment – and this subjective/objective fitness distinction creates a useful classification.

To be sure, there are fundamental differences between how single population and multipopulation learning systems behave (Ficici, 2004). Still, single population systems that employ subjective fitness assessment behave a lot more like multipopulation coevolutionary systems than like objective fitness based evolution. Moreover, historically, the field has used the term coevolution whenever fitness assessment is based on interactions between individuals, and a large amount of that research has involved systems with only one population.

Competition and Cooperation

The terms *cooperative* and *competitive* have been used to describe aspects of coevolution learning in at least three ways.

First and less commonly, these adjectives can describe qualitatively observed behaviors of potential solutions in coevolutionary systems, the results of some evolutionary process (e.g., "tit-for-tat" strategies, Axelrod, 1984).

Second, problems are sometimes considered to be inherently competitive or cooperative. Indeed, game theory provides some guidance for making such distinctions. However, since in many kinds of problems little may be known about the actual structure of the payoff functions involved, we may not actually be able to classify the problem as definitively competitive or cooperative.

The final and by far most common use of the term is to distinguish algorithms themselves. Cooperative algorithms are those in which interacting individuals succeed or fail together, while competitive algorithms are those in which individuals succeed at the expense of other individuals.

Because of the ambiguity of the terms, some researchers advocate abandoning them altogether, instead focusing distinguishing terminology on the form a potential solution takes. For example, using the term **>**compositional coevolution to describe an algorithm designed to return a solution composed of multiple individuals (e.g., a multiagent team) and using the term **>**test-based coevolution to describe an algorithm designed to return an individual who performs well against an adaptive set of tests (e.g., sorting network). This latter pair of terms is a slightly different, though probably more useful distinction than the cooperative and competitive terms.

Still, it is instructive to survey the algorithms based on how they have been historically classified.

Examples of competitive coevolutionary learning include simultaneously learning sorting networks and challenging data sets in a predator–prey type relationship (Hillis, 1991). Here, individuals in one population representing potential sorting networks are awarded a fitness score based on how well they sort opponent data sets from the other population. Individuals in the second population represent potential data sets whose fitness is based on how well they distinguish opponent sorting networks.

Competitive coevolution has also been applied to learning game-playing strategies (Fogel, 2001; Rosin & Belew, 1996). Additionally, competition has played a vital part in the attempts to coevolve complex agent

C

behaviors (Sims, 1994). Finally, competitive approaches have been applied to a variety of more traditional machine learning problems, for example, learning classifiers in one population and challenging subsets of exemplars in the other (Paredis, 1994).

Potter developed a relatively general framework for cooperative coevolutionary learning, applying it first to static function optimization and later to neural network learning (Potter, 1997). Here, each population contains individuals representing a portion of the network, and evolution of these components occurs almost independently, in tandem with one another, interacting only to be assembled into a complete network in order to obtain fitness. The decomposition of the network can be static and a priori, or dynamic in the sense that components may be added or removed during the learning process.

Moriarty et al. take a different, somewhat more adaptive approach to cooperative coevolution of neural networks (Moriarty & Miikkulainen, 1997). In this case, one population represents potential network *plans*, while a second is used to acquire node information. Plans are evaluated based on how well they solve a problem with their collaborating nodes, and the nodes receive a share of this fitness. Thus, a node is rewarded for participating more with successful plans, and thus receives fitness only indirectly.

Evaluation

Choices surrounding how interacting individuals in coevolutionary systems are evaluated for the purposes of selection are perhaps the most important choices facing an engineer employing these methods. Designing the evaluation method involves a variety of practical choices, as well as a broader eye to the ultimate purpose of the algorithm itself.

Practical concerns in evaluation include determining the number of individuals with whom to interact, how those individuals will be chosen for the interaction, and how the selection will operate on the results of multiple interactions (Wiegand, 2003). For example, one might determine the fitness of an individual by pairing him or her with all other individuals in the other populations (or the same population for single population approaches) and taking the average or maximum value

of such evaluations as the fitness assessment. Alternatively, one may simply use the single best individual as determined by a previous generation of the algorithm, or a combination of those approaches. Random pairings between individuals is also common. This idea can be extended to use tournament evaluation where successful individuals from pairwise interactions are promoted and further paired, assigning fitness based on how far an individual progresses in the tournament. Many of these methods have been evaluated empirically on a variety of types of problems (Angeline & Pollack, 1993; Bull, 1997; Wiegand, 2003).

However, the designing of the evaluation method also speaks to the broader issue of how to best implement the desired ▶ solution concept, (a criterion specifying which locations in the search space are solutions and which are not) (Ficici, 2004). The key to successful application of coevolutionary learning is to first elicit a clear and precise solution concept and then design an algorithm (an evaluation method in particular) that implements such a concept explicitly.

A successful coevolutionary learner capable of achieving reliable progress toward a particular solution concept often makes use of an archive of individuals and an update rule for that archive that insists the distance to a particular solution concept decrease with every change to the archive. For example, if one is interested in finding game strategies that satisfy Nash equilibrium constraints, one might consider comparing new individuals to an archive of potential individual strategies found so far that together represent a potential Nash mixed strategy (Ficici, 2004). Alternatively, if one is interested in maximizing the sum of an individual's outcomes over all tests, one may likewise employ an archive of discovered tests that candidate solutions are able to solve (de Jong, 2004).

It is useful to note that many coevolutionary learning problems are multiobjective in nature. That is, **\rightarrow** underlying objectives may exist in such problems, each creating a different ranking for individuals depending on the set of tests being considered during evaluation (Bucci & Pollack, 2002). The set of all possible underlying objectives (were it known) is sufficient to determine the outcomes on all possible tests. A careful understanding of this can yield approaches that create

ideal and minimal evaluation sets for such problems (de Jong & Pollack, 2004).

By acknowledging the link between multiobjective optimization and coevolutionary learning, a variety of evaluation and selection methods based on notions of multiobjective optimization have been employed. For example, there are selection methods that use Pareto dominance between candidate solutions and their tests as their basis of comparison (Ficici, 2004). Additionally, such methods can be combined with archive-based approaches to ensure monotonicity of progress toward a Pareto dominance solution concept (de Jong & Pollack, 2004).

Representation

Perhaps the core representational question in coevolution is the role that an individual plays. In test-based coevolution, an individual typically represents a potential solution to the problem or a test for a potential solution, whereas in compositional coevolution individuals typically represent a candidate component for a composite or ensemble solution.

Even in test-based approaches, the true solution to the problem may be expressed as a population of individuals, rather than a single individual. The population may represent a mixed strategy while individuals represent potential pure strategies for a game. Engineers using such approaches should be clear of the form of the final solution produced by the algorithm, and that this form is consistent with the prescribed solution concept.

In compositional approaches, the key issues tend to surround about how the problem is decomposed. In some algorithms, this decomposition is performed a priori, having different populations represent explicit components of the problem (Potter, 1997). In other approaches, the decomposition is intended to be somewhat more dynamic (Moriarty & Miikkulainen, 1997; Potter, 1997). Still more recent approaches seek to harness the potential of compositional coevolutionary systems to search open-ended representational spaces by gradually *complexifying* the representational space during the search (Stanley, 2004).

In addition, a variety of coevolutionary systems have successfully dealt with some inherent pathologies by representing populations in spatial topologies, and restricting selection and interaction using geometric constraints defined by those topologies (Pagie, 1999). Typically, these systems involve overlaying multiple grids of individuals, applying selection within some neighborhood in a given grid, and evaluating interactions between individuals in different grids using a similar type of cross-population neighborhood. The benefits of these systems are in part due to their ability to naturally regulate loss of diversity and spread of interaction information by explicit control over the size and shape of these neighborhoods.

Pathologies and Remedies

Perhaps the most commonly cited pathology is the socalled *loss of gradient* problem, in which one population comes to severely dominate the others, thus creating a situation in which individuals cannot be distinguished from one another. The populations become disengaged and evolutionary progress may *stall* or *drift* (Watson & Pollack, 2001). Disengagement most commonly occurs when distinguishing individuals are lost in the evolutionary process (*forgetting*), and the solution to this problem typically involves somehow retaining potentially informative, though possibly inferior quality individuals (e.g., archives).

Intransitivities in the reward system can cause some coevolutionary systems to exhibit cycling dynamics (Watson & Pollack, 2001), where reciprocal changes force the system to orbit some part of a potential search space. The remedy to this problem often involves creating coevolutionary systems that change in response to traits in several other populations. Mechanisms introduced to produce such effects include competitive fitness sharing (Rosin & Belew, 1996).

Another challenging problem occurs when individuals in a coevolutionary systems *overspecialize* on one underlying objective at the expense of other necessary objectives (Watson & Pollack, 2001). In fact, overspecialization can be seen as a form of disengagement on some subset of underlying objectives, and likewise the repair to this problem often involves retaining individuals capable of making distinctions in as many underlying objectives as possible (Bucci & Pollack, 2003).

C

Finally, certain kinds of compositional coevolutionary learning algorithms can be prone to *relative over-generalization*, a pathology in which components that perform reasonably well in a variety of composite solutions are favored over those that are part of an optimal solution (Wiegand, 2003). In this case, it is typically possible to bias the evaluation process toward optimal values by evaluating an individual in a variety of composite assemblies and assigning the best objective value found as the fitness (Panait, 2006).

In addition to pathological behaviors in coevolution, the subjective nature of these learning systems creates difficulty in measuring progress. Since fitness is subjective, it is impossible to determine whether these relative measures indicate progress or stagnation when the measurement values do not change much. Without engaging some kind of external or objective measure, it is difficult to understand what the system is really doing. Obviously, if an objective measure exists then it can be employed directly to measure progress (Watson & Pollack, 2001).

A variety of measurement methodologies have been employed when objective measurement is not possible. One method is to compare current individuals against all ancestral opponents (Cliff & Miller, 1995). Another predator/prey based method holds master tournaments between all the best predators and all the best prey found during the search (Nolfi & Floreano, 1998). A similar approach suggests maintaining the best individuals from each generation in each population in a hall of fame for comparison purposes (Rosin & Belew, 1996). Still other approaches seek to record the points during the coevolutionary search in which a new dominant individual was found (Stanley, 2004). A more recent approach advises looking at the population differential, examining all the information from ancestral generations rather than simply selecting a biased subset (Bader-Natal & Pollack, 2005). Conversely, an alternative idea is to consider how well the dynamics of the best individuals in different populations reflect the fundamental best response curves defined by the problem (Popovici, 2006).

With a clear solution concept, an appropriate evaluation mechanism implementing that concept, and practical progress measures in place, coevolution can be an effective and versatile machine learning tool.

Cross References

► Evolutionary Algorithms

Recommended Reading

- Angeline, P., & Pollack, J. (1993). Competitive environments evolve better solutions for complex tasks. In S. Forest (Ed.), Proceedings of the fifth international conference on genetic algorithms (pp. 264-270). San Mateo, CA: Morgan Kaufmann.
- Axelrod, R. (1984). The evolution of cooperation. New York: Basic Books
- Bader-Natal, A., & Pollack, J. (2005). Towards metrics and visualizations sensitive to Coevolutionary failures. In AAAI technical report FS-05-03 coevolutionary and coadaptive systems. AAAI Fall Symposium, Washington, DC.
- Bucci, A., & Pollack, J. B. (2002). A mathematical framework for the study of coevolution. In R. Poli, et al. (Eds.), Foundations of genetic algorithms VII (pp. 221-235). San Francisco: Morgan Kaufmann.
- Bucci, A., & Pollack, J. B. (2003). Focusing versus intransitivity geometrical aspects of coevolution. In E. Cantú-Paz, et al. (Eds.), Proceedings of the 2003 genetic and evolutionary computation conference (pp. 250-261). Berlin: Springer.
- Bull, L. (1997). Evolutionary computing in multi-agent environments: Partners. In T. Bäck (Ed.), Proceedings of the seventh international conference on genetic algorithms (pp. 370–377). San Mateo, CA: Morgan Kaufmann.
- Cliff, D., & Miller, G. F. (1995). Tracking the red queen: Measurements of adaptive progress in co-evolutionary simulations. In Proceedings of the third European conference on artificial life (pp. 200–218). Berlin: Springer.
- de Jong, E. (2004). The maxsolve algorithm for coevolution. In H. Beyer, et al. (Eds.), Proceedings of the 2005 genetic and evolutionary computation conference (pp. 483-489). New York, NY: ACM Press.
- de Jong, E., & Pollack, J. (2004). Ideal evaluation from coevolution. Evolutionary Computation, 12, 159–192.
- Ficici, S. G. (2004). Solution concepts in coevolutionary algorithms. PhD thesis, Brandeis University, Boston, MA.
- Fogel, D. (2001). Blondie24: Playing at the edge of artificial intelligence. San Francisco: Morgan Kaufmann.
- Hillis, D. (1991). Co-evolving parasites improve simulated evolution as an optimization procedure. *Artificial life II, SFI studies in the sciences of complexity* (Vol. 10, pp. 313–324).
- Moriarty, D., & Miikkulainen, R. (1997). Forming neural networks through efficient and adaptive coevolution. *Evolutionary Computation*, 5, 373–399.
- Nolfi, S., & Floreano, D. (1998). Co-evolving predator and prey robots: Do "arm races" arise in artificial evolution? Artificial Life, 4, 311–335.
- Pagie, L. (1999). Information integration in evolutionary processes.PhD thesis, Universiteit Utrecht, the Netherlands.
- Panait, L. (2006). The analysis and design of concurrent learning algorithms for cooperative multiagent systems. PhD thesis, George Mason University, Fairfax, VA.
- Paredis, J. (1994). Steps towards co-evolutionary classification networks. In R. A. Brooks & P. Maes (Eds.), Artificial life IV,

proceedings of the fourth international workshop on the synthesis and simulation of living systems (pp. 359–365). Cambridge, MA: MIT Press.

Popovici, E. (2006). An analysis of multi-population co-evolution. PhD thesis, George Mason University, Fairfax, VA.

Potter, M. (1997). The design and analysis of a computational model of cooperative co-evolution. PhD thesis, George Mason University, Fairfax, VA.

Rosin, C., & Belew, R. (1996). New methods for competitive coevolution. *Evolutionary Computation*, 5, 1–29.

Sims, K. (1994). Evolving 3D morphology and behavior by competition. In R. A. Brooks & P. Maes (Eds.), Artificial life IV, proceedings of the fourth international workshop on the synthesis and simulation of living systems (pp. 28–39). Cambridge, MA: MIT Press.

Stanley, K. (2004). Efficient evolution of neural networks through complexification. PhD thesis, The University of Texas at Austin, Austin, TX.

Watson, R., & Pollack, J. (2001). Coevolutionary dynamics in a minimal substrate. In L. Spector, et al. (Eds.), Proceedings from the 2001 genetic and evolutionary computation conference (pp. 702–709). San Francisco: Morgan Kaufmann.

Wiegand, R. P. (2003). An analysis of cooperative coevolutionary algorithms. PhD thesis, George Mason University, Fairfax, VA.

Collaborative Filtering

Collaborative Filtering (CF) refers to a class of techniques used in that recommend items to users that other users with similar tastes have liked in the past. CF methods are commonly sub-divided into neighborhood-based and model-based approaches. In neighborhood-based approaches, a subset of users are chosen based on their similarity to the active user, and a weighted combination of their ratings is used to produce predictions for this user. In contrast, model-based approaches assume an underlying structure to users' rating behavior, and induce predictive models based on the past ratings of all users.

Collection



Collective Classification

Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor

University of Maryland, MD, USA

Synonyms

Iterative classification; Link-based classification

Definition

Many real-world Classification problems can be best described as a set of objects interconnected via links to form a network structure. The links in the network denote relationships among the instances such that the class labels of the instances are often correlated. Thus, knowledge of the correct label for one instance improves our knowledge about the correct assignments to the other instances it connects to. The goal of collective classification is to *jointly* determine the correct label assignments of all the objects in the network.

Motivation and Background

Traditionally, a major focus of machine learning is to solve classification problems: given a corpus of documents, classify each according to its topic label; given a collection of e-mails, determine which are spam; given a sentence, determine the part-of-speech tag for each word; given a hand-written document, determine the characters, etc. However, much of the work in machine learning makes an *independent and identically* distributed (IID) assumption, and focuses on predicting the class label of each instance in isolation. In many cases, however, the class labels whose values need to be determined can benefit if we know the correct assignments to related class labels. For example, it is easier to predict the topic of a webpage if we know the topics of the webpages that link to it, the chance of a particular word being a verb increases if we know that the previous word in the sentence is a noun, knowing the rest of the characters in a word can make it easier to identify an unknown character, etc. In the last decade, many researchers have proposed techniques that attempt to classify samples in a joint or collective manner instead of treating each sample in isolation, and reported significant gains in classification accuracy.

Theory/Solution

Collective classification is a combinatorial optimization problem, in which we are given a set of nodes, $V = \{v_1, ..., v_n\}$, and a neighborhood function \mathcal{N} , where $\mathcal{N}_i \subseteq \mathcal{V}\setminus\{v_i\}$, which describes the underlying network structure. Each node in \mathcal{V} is a random variable that can take a value from an appropriate domain, $\mathcal{L} = \{l_1, ..., l_q\}$. \mathcal{V} is further divided into two sets of nodes: \mathcal{X} , the nodes for which we know the correct values (observed variables) and, \mathcal{Y} , the nodes whose values need to be determined. Our task is to label the nodes $y_i \in \mathcal{Y}$ with one of a small number of predefined labels in \mathcal{L} .

Even though it is only in the last decade that collective classification has entered the collective conscience of machine learning researchers, the general idea can be traced further back (Besag, 1986). As a result, a number of approaches have been proposed. The various approaches to collective classification differ in the kinds of information they aim to exploit to arrive at the correct classification, and their mathematical underpinnings. We discuss each in turn.

Relational Classification

Traditional classification concentrates on using the observed attributes of the instance to be classified. Relational classification (Slattery & Craven, 1998) attempts to go a step further by classifying the instance using not only the instance's own attributes but also the instance's neighbors' attributes. For example, in a hypertext classification domain where we want to classify webpages, not only would we use the webpage's own words but we would also look at the webpages linking to this webpage using hyperlinks and their words to arrive at the correct class label. Results obtained using relational classification have been mixed. For example, even though there have been reports of classification accuracy gains using such techniques, in certain cases, these techniques can harm classification accuracy (Chakrabarti, Dom, & Indyk, 1998).

Iterative Collective Classification with Neighborhood Labels

A second approach to collective classification is to use the class labels assigned to the neighbor instead of using the neighbor's observed attributes. For example, going back to our hypertext classification example, instead of using the linking webpage's words we would, in this case, use its assigned labels to classify the current webpage. Chakrabarti et al. (1998) illustrated the use of this approach and reported impressive classification accuracy gains. Neville and Jensen (2000) further developed the approach, and referred to the approach as iterative classification, and studied the conditions under which it improved classification performance (Jensen, Neville, & Gallagher, 2004). Techniques for feature construction from the neighboring labels were developed and studied (Lu & Getoor, 2003), along with methods that make use of only the label information (Macskassy & Provost, 2007), as well as a variety of strategies for when to commit the class labels (McDowell, Gupta, & Aha, 2007).

Algorithm 1 depicts pseudo-code for a simple version of the Iterative Classification Algorithm (ICA). The basic premise behind ICA is extremely simple. Consider a node $Y_i \in \mathcal{Y}$ whose value we need to determine and suppose we know the values of all the other nodes in its neighborhood \mathcal{N}_i (note that \mathcal{N}_i can contain both observed and unobserved variables). Then, ICA assumes that we are given a local classifier f that takes the values of \mathcal{N}_i as arguments and returns a label value for Y_i from the class label set \mathcal{L} . For local classifiers f that do not return a class label but a goodness/likelihood value given a set of attribute values and a label, we

Algorithm 1 Iterative classification algorithm

```
Iterative Classification Algorithm (ICA)
```

```
for each node Y_i \in \mathcal{Y} do {bootstrapping}

{compute label using only observed nodes in \mathcal{N}_i}

compute \vec{a}_i using only \mathcal{X} \cap \mathcal{N}_i

y_i \leftarrow f(\vec{a}_i)

end for

repeat {iterative classification}

generate ordering \mathcal{O} over nodes in \mathcal{Y}

for each node Y_i \in \mathcal{O} do

{compute new estimate of y_i}

compute \vec{a}_i using current assignments to \mathcal{N}_i

y_i \leftarrow f(\vec{a}_i)

end for
```

until all class labels have stabilized or a threshold number of iterations have elapsed

simply choose the label that corresponds to the maximum goodness/likelihood value; in other words, we replace f with $\operatorname{argmax}_{l \in \mathcal{L}} f$. This makes the local classifier f extremely flexible and we can use anything ranging from a decision tree to a \triangleright support vector machine (SVM). Unfortunately, it is rare in practice that we know all values in \mathcal{N}_i , which is why we need to repeat the process iteratively, in each iteration, labeling each Y_i using the current best estimates of \mathcal{N}_i and the local classifier f, and continuing to do so until the assignments to the labels stabilize.

Most local classifiers are defined as functions whose argument consists of a fixed-length vector of attribute values. A common approach to circumvent such a situation is to use an aggregation operator such as count, mode, or prop, which measures the proportion of neighbors with a given label. In Algorithm 1, we use \vec{a}_i to denote the vector encoding the values in \mathcal{N}_i obtained after aggregation. Note that in the first ICA iteration, all labels y_i are undefined and to initialize them we simply apply the local classifier to the observed attributes in the neighborhood of Y_i , this is referred to as "bootstrapping" in Algorithm 1.

Researchers in collective classification (Macskassy & Provost, 2007; McDowell et al., 2007; Neville & Jensen, 2000) have extended the simple algorithm described above, and developed a version of Gibbs sampling that is easy to implement and faster than traditional Gibbs sampling approaches. The basic idea behind this algorithm is to assume, just like in the case of ICA, that we have access to a local classifier f that can sample for the best label estimate for Y_i given all the values for the nodes in \mathcal{N}_i . We keep doing this repeatedly for a fixed number of iterations (a period known as "burnin"). After that, not only do we sample for labels for each $Y_i \in \mathcal{Y}$ but we also maintain count statistics as to how many times we sampled label l for node Y_i . After collecting a predefined number of such samples we output the best label assignment for node Y_i by choosing the label that was assigned the maximum number of times to Y_i while collecting samples.

One of the benefits of both variants of ICA is fairly simple to make use of any local classifier. Some of the classifiers used included the following: naïve Bayes (Chakrabarti et al., 1998; Neville & Jensen, 2000), ▶logistic regression (Lu & Getoor, 2003), ▶decision trees, (Jensen et al., 2004) and weighted-vote relational

neighbor (Macskassy & Provost, 2007). There is some evidence to indicate that discriminately trained local classifiers such as logistic regression tend to produce higher accuracies than others; this is consistent with results in other areas.

Other aspects of ICA that have been the subject of investigation include the ordering strategy to determine in which order to visit the nodes to relabel in each ICA iteration. There is some evidence to suggest that ICA is fairly robust to a number of simple ordering strategies such as random ordering, visiting nodes in ascending order of diversity of its neighborhood class labels, and labeling nodes in descending order of label confidences (Getoor, 2005). However, there is also some evidence that certain modifications to the basic ICA procedure tend to produce improved classification accuracies. For example, both (Neville & Jensen, 2000) and (McDowell et al., 2007) propose a strategy where only a subset of the unobserved variables are utilized as inputs for feature construction. More specifically, in each iteration, they choose the top-k most confident predicted labels and use only those unobserved variables in the following iteration's predictions, thus ignoring the less confident predicted labels. In each subsequent iteration they increase the value of k so that in the last iteration all nodes are used for prediction. McDowell et al. report that such a "cautious" approach leads to improved accuracies.

Collective Classification with Graphical Models

In addition to the approaches described above, which essentially focus on local representations and propagation methods, another approach to collective classification is by first representing the problem with a highlevel global paper graphical model and then using learning and inference techniques for the graphical modeling approach to arrive at the correct classifications. These proposals include the use of both directed paper graphical models (Getoor, Segal, Taskar, & Koller, 2001) and undirected graphical models (Lafferty, McCallum, & Pereira, 2001; Taskar, Abbeel, & Koller, 2002). See statistical relational learning and Getoor and Taskar (2007) for a survey of various graphical models that are suitable for collective classification. In general, these techniques can use both neighborhood labels and observed attributes

of neighbors. On the other hand, due to their generality, these techniques also tend to be less efficient than the iterative collective classification techniques.

One common way of defining such a global model uses a *pairwise Markov random field* (pairwise MRF) (Taskar et al., 2002). Let $G = (\mathcal{V}, E)$ denote a graph of random variables as before where \mathcal{V} consists of two types of random variables, the unobserved variables, \mathcal{Y} , which need to be assigned domain values from label set \mathcal{L} , and observed variables \mathcal{X} whose values we know (see \blacktriangleright Graphical Models). Let \mathcal{V} denote a set of *clique potentials*. \mathcal{V} contains three distinct types of functions:

- For each $Y_i \in \mathcal{Y}$, $\psi_i \in \mathcal{Y}$ is a mapping $\psi_i : \mathcal{L} \to \mathfrak{R}_{\geq 0}$, where $\mathfrak{R}_{\geq 0}$ is the set of nonnegative real numbers.
- For each $(Y_i, X_j) \in E$, $\psi_{ij} \in \Psi$ is a mapping $\psi_{ij} : \mathcal{L} \to \mathfrak{R}_{\geq 0}$.
- For each $(Y_i, Y_j) \in E$, $\psi_{ij} \in \Psi$ is a mapping $\psi_{ij} : \mathcal{L} \times \mathcal{L} \to \Re_{>0}$.

Let **x** denote the values assigned to all the observed variables in \mathcal{V} and let x_i denote the value assigned to X_i . Similarly, let **y** denote any assignment to all the unobserved variables in \mathcal{V} and let y_i denote a value assigned to Y_i . For brevity of notation we will denote by ϕ_i the clique potential obtained by computing $\phi_i(y_i) = \psi_i(y_i) \prod_{(Y_i, X_j) \in E} \psi_{ij}(y_i)$. We are now in a position to define a pairwise MRF.

Definition 1 A pairwise Markov random field (MRF) is given by a pair $\langle G, \Psi \rangle$ where G is a graph and Ψ is a set of clique potentials with ϕ_i and ψ_{ij} as defined above. Given an assignment \mathbf{y} to all the unobserved variables \mathcal{Y} , the pairwise MRF is associated with the probability distribution $P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{Y_i \in \mathcal{Y}} \phi_i(y_i) \prod_{(Y_i, Y_j) \in E} \psi_{ij}(y_i, y_j)$ where \mathbf{x} denotes the observed values of \mathcal{X} and $\mathcal{Z}(\mathbf{x}) = \sum_{\mathbf{y}'} \prod_{Y_i \in \mathcal{Y}} \phi_i(y_i') \prod_{(Y_i, Y_i) \in E} \psi_{ij}(y_i', y_j')$.

Given a pairwise MRF, it is conceptually simple to extract the best assignments to each unobserved variable in the network. For example, we may adopt the criterion that the best label value for Y_i is simply the one corresponding to the highest marginal probability obtained by summing over all other variables from the probability distribution associated with the pairwise MRF. Computationally, however, this is difficult to achieve since computing one marginal probability

requires summing over an exponentially large number of terms, which is why we need approximate inference algorithms. Hence, approximate inference algorithms are typically employed, the two most common being loopy belief propagation (LBP) and mean-field relaxation labeling.

Applications

Due to its general applicability, collective classification has been applied to a number of real-world problems. Foremost in this list is document classification. Chakrabarti et al. (1998) was one of the first to apply collective classification to corpora of patents linked via hyperlinks and reported that considering attributes of neighboring documents actually hurts classification performance. Slattery and Craven (1998) also considered the problem of document classification by constructing features from neighboring documents using an *inductive logic programming* rule learner. Yang, Slattery, & Ghani (2002) conducted an in-depth investigation over multiple datasets commonly used for document classification experiments and identified different patterns. Other applications of collective classification include object labeling in images (Hummel & Zucker, 1983), analysis of spatial statistics (Besag, 1986), iterative decoding (Berrou, Glavieux, & Thitimajshima, 1993), part-of-speech tagging (Lafferty et al., 2001), classification of hypertext documents using hyperlinks (Taskar et al., 2002), link prediction (Getoor, Friedman, Koller, & Taskar, 2002; Taskar, Wong, Abbeel, & Koller, 2003), optical character recognition (Taskar, Guestrin, & Koller, 2003), entity resolution in sensor networks (Chen, Wainwright, Cetin, & Willsky, 2003), predicting disulphide bonds in protein molecules (Taskar, Chatalbashev, Koller, & Guestrin, 2005), segmentation of 3D scan data (Anguelov et al., 2005), and classification of e-mail speech acts (Carvalho & Cohen, 2005). Recently, there have also been attempts to extend collective classification techniques to the semi-supervised learning scenario (Lu & Getoor, 2003b; Macskassy, 2007; Xu, Wilkinson, Southey, & Schuurmans, 2006).

Cross References

- ► Decision Trees
- ► Inductive Logic Programming
- ► Learning From Structured Data

Community Detection 193

- ▶ Relational Learning
- ►Semi-Supervised Learning
- ► Statistical Relational Learning

Recommended Reading

- Anguelov, D., Taskar, B., Chatalbashev, V., Koller, D., Gupta. D., Heitz, G., et al. (2005). Discriminative learning of Markov random fields for segmentation of 3d scan data. In *IEEE computer* society conference on computer vision and pattern recognition. IEEE Computer Society, Washington D.C.
- Berrou, C., Glavieux, A., & Thitimajshima, P. (1993). Near Shannon limit error-correcting coding and decoding: Turbo codes. In *Proceedings of IEEE international communications conference*, Geneva, Switzerland, IEEE.
- Besag, J. (1986). On the statistical analysis of dirty pictures. *Journal of the Royal Statistical Society, B-48*, 259–302.
- Carvalho, V., & Cohen, W. W. (2005). On the collective classification of email speech acts. In Special interest group on information retrieval, Salvador, Brazil, ACM.
- Chakrabarti, S., Dom, B., & Indyk, P. (1998). Enhanced hypertext categorization using hyperlinks. In *International conference on management of data*, Seattle, Washington New York: ACM.
- Chen, L., Wainwright, M., Cetin, M., & Willsky, A. (2003). Multitargetmultisensor data association using the tree-reweighted max-product algorithm. In SPIE Aerosense conference. Orlando, Florida.
- Getoor, L. (2005). Link-based classification. In Advanced methods for knowledge discovery from complex data. New York: Springer.
- Getoor, L., & Taskar, B. (Eds.). (2007). Introduction to statistical relational learning. Cambridge, MA: MIT Press.
- Getoor, L., Segal, E., Taskar, B., & Koller, D. (2001). Probabilistic models of text and link structure fro hypertext classification. In *Proceedings of the IJCAI workshop on text learning: Beyond supervision*, Seattle, WA.
- Getoor, L., Friedman, N., Koller, D., & Taskar, B. (2002). Learning probabilistic models of link structure. *Journal of Machine Learning Research*, 3, 679-707.
- Hummel, R., & Zucker, S. (1983). On the foundations of relaxation labeling processes. IEEE Transactions on Pattern Analysis and Machine Intelligence, 5, 267-287.
- Jensen, D., Neville, J., & Gallagher, B. (2004). Why collective inference improves relational classification. In Proceedings of the 10th ACM SIGKDD international conference on knowledge discovery and data mining, Seattle, WA. ACM.
- Lafferty, J. D., McCallum, A., & Pereira, F. C. N. (2001). conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the international conference on machine learning*, Washington DC. San Francisco, CA: Morgan Kaufmann.
- Lu, Q., & Getoor, L. (2003a). Link based classification. In Proceedings of the international conference on machine learning. AAAI Press, Washington, D.C.
- Lu, Q., & Getoor, L. (2003b). Link-based classification using labeled and unlabeled data. In ICML workshop on the continuum from labeled to unlabeled data in machine learning and data mining. Washington, D.C.
- Macskassy, S., & Provost, F. (2007). Classification in networked data: A toolkit and a univariate case study. *Journal of Machine Learning Research*, 8, 935-983.

- Macskassy, S. A. (2007). Improving learning in networked data by combining explicit and mined links. In *Proceedings of the twenty-second conference on artificial intelligence*. AAAI Press, Vancouver, Canada.
- McDowell, L. K., Gupta, K. M., & Aha, D. W. (2007). Cautious inference in collective classification. In *Proceedings of AAAI*. AAAI Press, Vancouver, Canada.
- Neville, J., & Jensen, D. (2007). Relational dependency networks. Journal of Machine Learning Research, 8, 653-692.
- Neville, J., & Jensen, D. (2000). Iterative classification in relation data. In Workshop on statistical relational learning, AAAI.
- Slattery, S., & Craven, M. (1998). Combining statistical and relational methods for learning in hypertext domains. In *International conferences on inductive logic programming*. Springer-Verlag, London, UK.
- Taskar, B., Abbeel, P., & Koller, D. (2002). Discriminative probabilistic models for relational data. In Proceedings of the annual conference on uncertainty in artificial intelligence. Morgan Kauffman, San Francisco, CA.
- Taskar, B., Guestrin, C., & Koller, D. (2003a). Max-margin markov networks. In *Neural information processing systems*. MIT Press, Cambridge, MA.
- Taskar, B., Wong, M. F., Abbeel, P., & Koller, D. (2003b). Link prediction in relational data. In *Natural information processing systems*. MIT Press, Cambridge, MA.
- Taskar, B., Chatalbashev, V., Koller, D., & Guestrin, C. (2005). Learning structured prediction models: A large margin approach. In Proceedings of the international conference on machine learning. ACM, New York, NY.
- Xu, L., Wilkinson, D., Southey, F., & Schuurmans, D. (2006). Discriminative unsupervised learning of structured predictors. In Proceedings of the international conference on machine learning. ACM, New York, NY.
- Yang, Y., Slattery, S., & Ghani, R. (2002). A study of approaches to hypertext categorization. *Journal of Intelligent Information* Systems. 18(2-3), 219-241.

Commercial Email Filtering

► Text Mining for Spam Filtering

Committee Machines

►Ensemble Learning

Community Detection

▶Group Detection

C

194 Comparable Corpus

Comparable Corpus

A comparable corpus (pl. corpora) is a document collection composed of two or more disjoint subsets, each written in a different language, such that documents in each subset are on a same topic as the documents in the others. The prototypical example of a comparable corpora is a collection of newspaper article written in different languages and reporting about the same events: while they will not be, strictly speaking, the translation of one another, they will share most of the semantic content. Some methods for **\rightarrow cross-language text mining** rely, totally or partially, on the statistical properties of comparable corpora.

Competitive Coevolution

▶ Test-Based Coevolution

Competitive Learning

Competitive learning is an partificial neural network learning process where different neurons or processing elements compete on who is allowed to learn to represent the current input. In its purest form competitive learning is in the so-called winner-take-all networks where only the neuron that best represents the input is allowed to learn. Since all neurons learn to better represent the kinds of inputs they already are good at representing, they become specialized to represent different kinds of inputs. For vector-valued inputs and representations, the input becomes quantized to the unit having the closest representation (model), and the representations are adapted to minimize the representation error using stochastic gradient descent.

Competitive learning networks have been studied as models of how receptive fields and feature detectors, such as orientation-selective visual neurons, develop in neural networks. The same process is at work in online ►K-means clustering, and variants of it in ►Self-Organizing Maps (SOM) and the EM algorithm of mixture models.

Complex Adaptive System

▶ Complexity in Adaptive Systems

Complexity in Adaptive Systems

Jun He Aberystwyth University, Wales, UK

Synonyms

Adaptive system; Complex adaptive system

Definition

An ▶adaptive system, or complex adaptive system, is a special case of complex systems, which is able to adapt its behavior according to changes in its environment or in parts of the system itself. In this way, the system can improve its performance through a continuing interaction with its environment. The concept of ▶complexity in an adaptive system is used to analyze the interactive relationship between the system and its environment, which can be classified into two types: ▶internal complexity for model complexity, and ▶external complexity for data complexity. The internal complexity is defined by the amount of input, information, or energy that the system receives from its environment. The external complexity refers to the complexity of how the system represents these inputs through its internal process.

Motivation and Background

Adaptive systems range from natural systems to artificial systems (Holland, 1992, 1995; Waldrop, 1992). Examples of natural systems include ant colonies, ecosystem, the brain, neural network and immune system, cell and developing embryo; examples of artificial systems include stock market, social system, manufacturing businesses, and human social group-based

Complexity in Adaptive Systems 195

endeavor in a cultural and social system such as political parties or communities. All these systems have a common feature: they can adapt to their environment.

An adaptive system is adaptive in that way it has the capacity to change its internal structure for adapting the environment. It is complex in the sense that it is interactive with its environment. The interaction between an adaptive system and its environment is dynamic and nonlinear. Complexity emerges from the interaction between the system and environment, the elements of the system, where the emergent macroscopic patterns are more complex than the sum of the these low-level (microscopic) elements encompassed in the system. Understanding the evolution and development of adaptive systems still faces many mathematical challenges (Levin, 2003).

The concepts of external and internal complexities are used to analyze the relation between an adaptive system and its environment. The description given below is based on Jürgen Jost's (2004) work, which introduced these two concepts and applied the theoretical framework to the construction of learning models, e.g., to design neural network architectures. In the following, the concepts are mainly applied to analyze the interaction between the system and its environment. The interaction among individual elements of the system is less discussed however, the concepts can be explored in that situation too.

Theory

Adaptive System Environment and Regularities

The environment of an adaptive system is more complex than the system itself and its changes cannot be completely predictable for the system. However, the changes of the environment are not purely random and noisy; there exist regularities in the environment. An adaptive system can recognize these regularities, and depending on these regularities the system will express them through its internal process in order to adapt to the environment.

The input that an adaptive system receives or extracts from its environment usually includes two parts: one is the part with regularities; and another is that appears random to the system. The part of regularities is useful and meaningful. An adaptive system will represent these regularities by internal processes. But

the part of random input is useless, and even at the worst it will be detrimental for an adaptive system. However, it will depend on the adaptive system's internal model of the external environment for how to determine which part of input is meaningful and regular, and which part is random and devoid of meaning and structure.

An adaptive system will translate the external regularities into its internal ones, and only the regularities are useful to the system. The system tries to extract as many regularities as possible, and to represent these regularities as efficiently as possible in order to make optimal use of its capacity.

The notions of external complexity and internal complexity are used to investigate these two complementary aspects conceptually and quantitatively. In terms of these notions, an adaptive system aims to increase its external complexity and reduce its internal complexity.

The two processes operate on their own time scale but are intricately linked and mutually dependent on each other. For example, the internal complexity will be only reduced if the external complexity is fixed. Under fixed inputs received from the external environment, an adaptive system can represent these inputs systems more efficiently and optimize its internal structure. If the external complexity is increased, e.g., if additional new input is required to handle by the system, then it is necessary to increase its internal complexity.

The increase of internal complexity may occur through the creation of redundancy in the existing adaptive system, e.g., to duplicate some internal structures, and then enable the system to handle more external input. Once the input is fixed, the adaptive system then will represent the input as efficiently as possible and reduce the internal input. The decrease of internal complexity can be achieved through discarding some input as meaningless and irrelevant, e.g., leaving some regularities out for the purpose.

Since the inputs relevant to the systems are those which can be reflected in the internal model, the external complexity is not equivalent to the amount of raw data received from the environment. In fact, it is only relevant to the inputs which can be processed in the internal model, or observations in some adaptive systems. Thus the external complexity ultimately is decided by the internal model constructed by the system.

C

196 Complexity in Adaptive Systems

External and Internal Complexities

External complexity means data complexity, which is used to measure the amount of input received from the environment for the system to handle and process. Such a complexity can be measured by entropy in the term of information theory.

Internal complexity is model complexity, which is used to measure the complexity of a model for representing the input or information received by the system.

The aim of the adaptive system is to obtain an efficient model as simple as possible, with the capacity to handle as much input as possible. On one hand, the adaptive system will try to maximize its external complexity and then to adapt to its environment in a maximal way; on the other hand, to minimize its internal complexity and then to construct a model to process the input in a most efficient way.

These two aims sometimes seem conflicting, but such a conflict can be avoided when these two processes operate on different time scales. If given a model, the system will organize the input data and try to increase its ability to deal with the input from its environment, and then increase its external complexity. If given the input, conversely, it tries to simplify its model which represents that input and thus to decrease the internal complexity. The meaning of the input is relevant to the time scale under investigation. On a short time scale, for example, the input may consist of individual signals, but on a long time scale, it will be a sequence of signals which satisfies a probability distribution. A good internal model tries to express regularities in the input sequence, rather than several individual signals. And the decrease of internal complexity will happen on this time scale.

A formal definition of the internal and external complexities concepts is based on the concept of entropy from statistical mechanics and information theory. Given a model θ , the system can model data as with $X(\theta) = (X_1, \ldots, X_k)$, which is assumed to have an internal probability distribution $P(X(\theta))$ so that entropy can be computed. The external complexity is defined by

$$-\sum_{i=1}^{k} P(X_i(\theta)) \log_2 P(X_i(\theta)). \tag{1}$$

An adaptive system tries to maximize the above external complexity.

The probability distribution $P(X(\theta))$ is for quantifying the information value of the data $X(\theta)$. The value

of information can be described in other approaches, e.g., the length of the representation of the data in the internal code of the system (Rissanen, 1989). In this case, the optimal coding is a consequence of the minimization of internal complexity, and then the length of the representation of data $X_i(\theta)$ behaves like $\log_2 P(X(\theta))$ (Rissanen, 1989).

On a short time scale, for a given model θ , the system tries to increase the amount of meaningful input information $X(\theta)$. On a long time scale, when the input is given, e.g., when the system has gathered a set of inputs on a time scale with a stationary probability distribution of input patterns Ξ , then the model should be improved to handle the input as efficiently as possible and reduce the complexity of the model. This complexity, or internal complexity, is defined by

$$-\sum_{i=1}^{k} P(\Xi_i \mid \theta) \log_2 P(\Xi_i \mid \theta) - \log_2 P(\theta), \quad (2)$$

with respect to the model θ .

If Rissanen's (1989) > minimum description length principle is applied to the above formula, then the optimal model will satisfy the variation problem

$$\min_{\theta} \left(-\log_2 P(\Xi \mid \theta) - \log_2 P(\theta) \right). \tag{3}$$

Here in the above minimization problem, there are two objectives to minimize. The first term is to measure how efficiently the model represents or encodes the data; and the second one is to measure how complicated the model is. In computer science, this latter term corresponds to the length of the program required to encode the model.

The concepts of external and internal complexities can be applied into a system divided into subsystems. In this case, some internal part of the original whole system will become external to a subsystem. Thus the internal input of a subsystem consists of original external input and also input from the rest of the system, i.e., other subsystems.

Application: Learning

The discussion of these two concepts, external and internal complexities, can be put into the background of learning. In statistical learning theory (Vapnik, 1998), the criterion for evaluating a learning process is the expected prediction error of future data by the model

Complexity in Adaptive Systems 197

based on training data set with partial and incomplete information. The task is to construct a probability distribution drawn from an a-priori specific class for representing the distribution underlying the input data received. Usually, if a higher error is produced by a model on the training data, then a higher error will be expected on the future data. The error will depend on two factors: one is the accuracy of the model on the training data set, another is the simplicity of the model itself. The description of the data set can be split into two parts, the regular part, which is useful in constructing the model; and the random part, which is a noise to the model.

The learning process fits very well into the theory framework of internal and external complexities. If the model is too complicated, it will bring the risk of overfitting the training data. In this case, some spurious or putative regularity is incorporated into the model, which will not appear in the future data. The model should be constrained within some model class with bounded complexity. This complexity in this context of statistical learning theory is measured by the Vapnik-Chervonenkis dimension (see ►VC Dimension) (Vapnik, 1998). Under the simplest form of statistical learning theory, the system aims at finding a representation with smallest error in a class with given complexity constraints; and then the model should minimize the expected error on future data and also over-fitting error.

The two concepts of over-fitting and leaving out regularities can be distinguished in the following sense. The former is caused by the noise in the data, i.e., the random part of the data, and this leads to putative regularities, which will not appear in the future data. The latter, leaving out regularities, means that the system can forgo some part of regularities in the data, or it is possible to make data compression. Thus, leaving out regularities can be used to simplify the model and reduce the internal complexity. However, a problem is still waiting for answer here, that is, what regularities in the data set are useful for data compression and also meaningful for future prediction; and what parts are random to the model.

The internal complexity is the model complexity. If the internal complexity is chosen too small, then the model does not have enough capacity to represent all the important features of the data set. If the internal complexity is too large, on the other hand, then the model does not represent the data efficiently. The internal complexity is preferably minimized under appropriate constraints on the adequacy of the representation of data. This is consistent with Rissanen's principle of Minimum Description Length (Rissanen, 1989) to represent a given data set in the most efficient way. Thus a good model is both to simplify the model itself and to represent the data efficiently.

The external complexity is the data complexity which should be large to represent the input accurately. This is related to Jaynes' principle of maximizing the ignorance (Jaynes, 1957), where a model for representing data should have the maximal possible entropy under the constraint that all regularities can be reproduced. In this way, putative regularities could be eliminated in the model. However, this principle should be applied with some conditions as argued by Gell-Mann and Lloyd (1996); it cannot eliminate the essential regularities in the data, and an overlying complex model should be avoided.

For some learning system, only a selection of data is gathered and observed by the system. Thus a middle term, observation, is added between model and data. The concept of observation refers to the extraction of value of some specific quantity from a given data or data pool. What a system can observe depends on its internal structure and its general model of the environment. The system does not have direct access to the raw data, but through constructing a model of the environment solely on the basis of the values of its observation.

For such kind of learning system, Jaynes' principle (Jaynes, 1957) is still applicable for increasing the external complexity. For the given observation made on a data set, the maximum entropy representation should be selected. However, this principle is still subject to the modification of Gell-Mann and Lloyd (1996) to a principle where the model should not lose the essential regularities observed in the data.

By contrast, the observations should be selected to reduce the internal complexity. Given a model, if the observation can be made on a given data set, then these observations should be selected so as to minimize the resulting entropy of the model, with the purpose of minimizing the uncertainty left about the data. Thus it leads to reduce the complexity.

In most of the cases, the environment is dynamic, i.e., the data set itself can be varied, then the external

198 Complexity of Inductive Inference

complexity should be maximized again. Thus the observation should be chosen for maximal information gain extracted from the data to increase the external complexity. Jaynes' principle (Jaynes, 1957) can be applied as the same as in previous discussion. But on a longer time scale, when the inputs reach some stationary distribution, the model should be simplified to reduce its internal complexity.

Recommended Reading

Gell-Mann, M., & Lloyd, S. (1996). Information measures, effective complexity, and total information. *Complexity*, 2(1), 44–52.

Holland, J. (1992). Adaptation in natural and artificial systems. Cambridge, MA: MIT Press.

Holland, J. (1995). Hidden order: How adaptation builds complexity. Reading, MA: Addison-Wesley.

Jaynes, E. (1957). Information theory and statistical mechanics. Physical Review, 106(4), 620-630.

Jost, J. (2004). External and internal complexity of complex adaptive systems. Theory in Biosciences, 123(1), 69–88.

Levin, S. (2003). Complex adaptive systems: Exploring the known, the unknown and the unknowable. *Bulletin of the American Mathematical Society*, 40(1), 3-19.

Rissanen, J. (1989). Stochastic complexity in statistical inquiry. Singapore: World Scientific.

Vapnik, V. (1998). Statistical learning theory. New York: John Wiley & Sons.

Waldrop, M. (1992). Complexity: The emerging science at the edge of order and chaos. New York: Simon & Schuster.

Complexity of Inductive Inference

Sanjay Jain, Frank Stephan National University of Singapore, Singapore, Republic of Singapore

Definition

In inductive inference, the complexity of learning can be measured in various ways: by the number of hypotheses issued in the worst case until the correct hypothesis is found; by the number of data items to be consumed or to be memorized in order to learn in the worst case; by the Turing degree of oracles needed to learn the class under a certain criterion; by the intrinsic complexity which is – like the Turing degrees in recursion theory – a way to measure the complexity of classes by using reducibilities between them.

Detail

We refer the reader to the article \triangleright Inductive Inference for basic definitions in inductive inference and the notations used below. Let \mathbb{N} denote the set of natural numbers. Let $\varphi_0, \varphi_1, \ldots$ denote a fixed acceptable programming system (Rogers, 1967). Let $W_i = \text{domain}(\varphi_i)$.

Mind Changes and Anomalies

The first measure of complexity of learning can be considered as the number of mind changes needed before the learner converges to its final hypothesis in the **TxtEx** model of learning. The number of mind changes by a learner M on a text T can be counted as card ($\{m:? \neq M(T[m]) \neq M(T[m+1])\}$). A learner M **TxtEx**_n learns a class \mathcal{L} of languages iff M **TxtEx** learns \mathcal{L} and for all $L \in \mathcal{L}$, for all texts T for L, M makes at most n mind changes on T. **TxtEx**_n is defined as the collection of language classes which can be **TxtEx**_n identified (see Case & Smith (1983) for details).

Consider the class of languages $\mathcal{L}_n = \{L : \operatorname{card}(L) \le n\}$. It can be shown that $\mathcal{L}_{n+1} \in \operatorname{TxtEx}_{n+1} - \operatorname{TxtEx}_n$.

Now consider anomalous learning. A class C is \mathbf{TxtEx}_{h}^{a} -learnable iff there is a learner, which makes at most b mind changes (where b = * denotes that the number of mind changes is finite on each text for a language in the class, but not necessarily bounded by a constant) and whose final hypothesis is allowed to make up to a errors (where a = * denotes finitely many errors). For these learning criteria, we get a twodimensional hierarchy on what can be learnt. Let $C_n =$ $\{f: \varphi_{f(0)} = f\}$. For a total function f, let $L_f = f$ $\{\langle x, f(x) \rangle : x \in \mathbb{N} \}$, where $\langle \cdot, \cdot \rangle$ denotes a computable pairing function: a bijective mapping from $\mathbb{N} \times \mathbb{N}$ to \mathbb{N} . Let $\mathcal{L}_{\mathcal{C}} = \{L_f : f \in \mathcal{C}\}$. Then, one can show that $\mathcal{L}_{\mathcal{C}_{n+1}} \in \mathbf{TxtEx}_0^{n+1} - \mathbf{TxtEx}^n$. Similarly, if we consider the class $S_n = \{f : \operatorname{card}(\{m : f(m) \neq f(m+1)\}) \leq n\}$, then one can show that $\mathcal{L}_{S_{n+1}} \in \mathbf{TxtEx}_{n+1}^0 - \mathbf{TxtEx}_n^*$ (we refer the reader to Case and Smith (1983) for a proof of the above).

Data and Time Complexity

Wiehagen (1986) considered the complexity of number of data needed for learning. Regarding time complexity, one should note the result by Pitt (1989) that any **TxtEx**-learnable class of languages can be **TxtEx**-learnt by a learner that has time complexity (with respect to

Complexity of Inductive Inference 199

the size of the input) bounded by a linear function. This result is achieved by a delaying trick, where the learner just repeats its old hypothesis unless it has enough time to compute its later hypothesis. This seriously effects what one can say about time complexity of learning. One proposal made by Daley and Smith (1986) is to consider the total time used by the learner until its sequence of hypotheses converges, resulting in a possibly more reasonable measure of time in the complexity of learning.

Iterative and Memory-Bounded Learning

Another measure of complexity of learning can be considered when one restricts how much past data a learner can remember. Wiehagen introduced the concept of iterative learning in which the learner cannot remember any past data. Its new hypothesis is based only on its previous conjecture and the new datum it receives. In other words, there exists a recursive function *F* such that M(T[n+1]) = F(M(T[n]), T(n)), for all texts Tand for all n. Here, M(T[0]) is some fixed value, say the symbol '?' which is used by the learner to denote the absence of a reasonable conjecture. It can be shown that being iterative restricts the learning capacity of learners. For example, let $L_e = \{2x : x \in \mathbb{N}\}$ and let $\mathcal{L} =$ $\{L_e\} \cup \{\{S \cup \{2n+1\}\}\} : n \in \mathbb{N}, S \subseteq L_e, \text{ and } \max(S) \le n\};$ then \mathcal{L} can be shown to be TxtEx-learnable but not iteratively learnable.

Memory-bounded learning (see Lange & Zeugmann, 1996) is an extension of memory-limited learning, where the learner is allowed to memorize upto some fixed number of elements seen in the past. Thus, M is an m-memory-bounded learner if there exists a function mem and two computable functions mF and F such that, for all texts T and all n:

```
- mem(T[0]) = \emptyset;

- M(T[n+1]) = F(M(T[n]), mem(T[n]), T(n+1));

- mem(T[n+1]) = mF(M(T[n]), mem(T[n]),

T(n+1));

- mem(T[n+1]) - mem(T[n]) \subseteq \{T(n+1)\};

- card(mem(T[n])) \le m.
```

It can be shown that the criteria of inference based on **TxtEx**-learning by *m*-memory-bounded learners form a proper hierarchy.

Besides memorizing some past elements seen, another way to address this issue is by giving feedback to the learner (see Case, Jain, Lange, & Zeugmann, 1999) on whether some element has appeared in the past data. A feedback learner is an iterative learner, which is additionally allowed to query whether certain elements appeared in earlier data. An *n*-feedback learner is allowed to make *n* such queries at each stage (when it receives the new input datum). Thus, M is an *m*-feedback learner if there exist computable functions Q and a *F* such that, for all texts *T* and all *n*:

- Q(M(T[n]), T(n)) is defined and is a set of m elements;
- If $Q(M(T[n]), T(n)) = (x_1, x_2, ..., x_m)$ then $M(T[n+1]) = F(M(T[n]), T(n), y_1, y_2, ..., y_m)$, where $y_i = 1$ iff $x_i \in \text{ctnt}(T[n])$.

Again, it can be shown that allowing more feedback gives greater learning power, and thus one can get a hierarchy based on the amount of feedback allowed.

Complexity of Final Hypothesis

Another possibility on complexity of learning is to consider the complexity or size of the final grammar output by the learner. Freivalds (1975) considered the case when the final program/grammar output by the learner is minimal: that is, there is no smaller index that accepts/generates the same language. He showed that this severely restricts the learning capacity of learners. Not only that, the learning capacity depends on the acceptable programming system chosen, unlike the case for most other criteria of learning such as TxtEx or TxtBc, which are independent of the acceptable programming system chosen. In particular, there are acceptable programming systems in which only classes containing finitely many infinite languages can be learnt using minimal final grammars (see Freivalds, 1975; Jain and Sharma, 1993). Chen (1982) considered a modification of such a paradigm where one considers convergence to nearly minimal grammars rather than minimal. That is, instead of requiring that the final grammars are minimal, one requires that they are within a recursive function h of minimal. Here h may depend on the class being learnt. Chen showed that this allows one to have the criteria of minimal learnability 200 Complexity of Inductive Inference

to be independent of the acceptable programming system chosen. However, one can show that some simple classes are not minimally learnable. An example of such a class is the class $\mathcal{L}_{\mathcal{C}}$ which is derived from $\mathcal{C} = \{f : \forall^{\infty} \times [f(x) = 0]\}$, the class of all functions which are almost everywhere 0.

Intrinsic Complexity

Another way to consider complexity of learning is to consider relative complexity in a way similar to how one considers Turing reductions in computability theory. Such a notion is called intrinsic complexity of the class. This was first considered by Freivalds et al. (1995) for function learning. Jain and Sharma (1996) considered it for language learning, and the following discussion is from there.

An enumeration operator (see Rogers, 1967), Θ , is an algorithmic mapping from SEQ into SEQ such that the following two conditions are satisfied:

- for all σ , $\tau \in SEQ$, if $\sigma \subseteq \tau$, then $\Theta(\sigma) \subseteq \Theta(\tau)$;
- for all texts T, $\lim_{n\to\infty} |\Theta(T[n])| = \infty$.

By extension, we think of Θ as also mapping texts to texts such that $\Theta(T) = \bigcup_n \Theta(T[n])$. Furthermore, we define $\Theta(L) = \{\operatorname{ctnt}(\Theta(T)) : T \text{ is a text for } L\}$. Intuitively, $\Theta(L)$ denotes the set of languages to whose texts Θ maps texts of L. The reader should note the overloading of this notation because the type of the argument to Θ could be a sequence, a text or a language.

One says that a sequence of grammars $g_0, g_1, ...$ is an acceptable **TxtEx**-sequence for *L* if the sequence of grammars converges to a grammar for *L*.

 $\mathcal{L}_1 \leq_{weak} \mathcal{L}_2$ iff there are two operators Θ and Ψ such that for all $L \in \mathcal{L}_1$, for all texts T for L, $\Theta(T)$ is a text for some $L' \in \mathcal{L}_2$ such that if g_0, g_1, \ldots is an acceptable **TxtEx**-sequence for L' then $\Psi(g_0, g_1, \ldots)$ is an acceptable **TxtEx**-sequence for L.

Note that different texts for the same language L may be mapped by Θ to texts for different languages in \mathcal{L}_2 above. If we require that different texts for L are mapped to texts for the same language L' in \mathcal{L}_2 , then we get a stronger notion of reduction called strong reduction: $\mathcal{L}_1 \leq_{strong} \mathcal{L}_2$ iff $\mathcal{L}_1 \leq_{weak} \mathcal{L}_2$ and for all $L \in \mathcal{L}_1$, $\Theta(L)$ contains only one language, where Θ is as in the definition for \leq_{weak} reduction.

It can be shown that FIN is a complete class for \mathbf{TxtEx} -identification with respect to \leq_{weak} reduction (see Jain & Sharma, 1996). Interestingly it was shown that the class of pattern languages (Angluin, 1980), the class $SD = \{L : W_{\min(L)} = L\}$ and the class $COINIT = \{\{x : x \geq n\} : n \in \mathbb{N}\}$ are all equivalent under \leq_{strong} . Let code be a bijective mapping from non-negative rational numbers to natural numbers. Then, one can show that the class RINIT = $\{\{code(x) : 0 \leq x \leq r, x \text{ is a rational number}\} : 0 \leq r \leq 1, r \text{ is a rational number}\}$ is \leq_{strong} complete for \mathbf{TxtEx} (see Jain, Kinber, & Wiehagen, 2001).

Interestingly every finite directed acyclic graph can be embedded into the \leq_{strong} degree structure (Jain & Sharma, 1997). On the other hand the degree structure is non-dense in the sense that there exist classes \mathcal{L}_1 and \mathcal{L}_2 such that $\mathcal{L}_1 <_{strong} \mathcal{L}_2$, but for any class \mathcal{L} such that $\mathcal{L}_1 \leq_{strong} \mathcal{L}_2$, either $\mathcal{L}_1 \equiv_{strong} \mathcal{L}$ or $\mathcal{L} \equiv_{strong} \mathcal{L}_2$. Similar result holds for \leq_{weak} reducibility (see Jain & Sharma, 1997).

Interesting connections between learning of elementary formal systems (Shinohara, 1994), intrinsic complexity and ordinal mind changes (Freivalds & Smith, 1993) were shown in (Jain & Sharma, 1997).

Learning Using Oracles

Another method to measure complexity of learning is to see how powerful an oracle (given to the learning machine) has to be to make a class learnable. It can be shown that an oracle A permits to explanatorily learn the class of all recursive functions iff A is high (Adleman & Blum, 1991). Furthermore, an oracle is trivial, that is, does not give additional learning power for explanatory learning of function classes iff the oracle has 1-generic Turing degree and is Turing reducible to the halting problem (Slaman & Solovay, 1991). The picture is a bit different in the general case of learning languages. For every oracle A there is an oracle B and a class, which is **TxtEx**-learnable using the oracle B but not using the oracle A (Jain & Sharma, 1993). Note that there are also classes of languages like Gold's class of all finite languages plus the set of natural numbers which are not TxtEx-learnable using any oracle. Furthermore, for oracles above the halting problem, TxtEx-learning and TxtBc-learning using these oracles coincide.

Computational Complexity of Learning 201

Recommended Reading

- Adleman, L., & Blum, M. (1991). Inductive inference and unsolvability. *Journal of Symbolic Logic*, 56, 891–900.
- Angluin, D. (1980). Finding patterns common to a set of strings. Journal of Computer and System Sciences, 21, 46-62.
- Case, J., Jain, S., & Lange, S., & Zeugmann, T. (1999). Incremental concept learning for bounded data mining. *Information and Computation*, 152(1), 74–110.
- Case, J., & Smith, C. H. (1983). Comparison of identification criteria for machine inductive inference. *Theoretical Computer Science*, 25, 193–220
- Daley, R. P., & Smith, C. H. (1986). On the complexity of inductive inference. *Information and Control*, 69, 12–40.
- Chen, K.-J. (1982). Tradeoffs in inductive inference of nearly minimal sized programs. *Information and Control*, 52, 68-86.
- Freivalds, R. (1975). Minimal Gödel numbers and their identification in the limit. Lecture Notes in Computer Science, 32, 219-225
- Freivalds, R., Kinber, E., & Smith, C. H. (1995). On the intrinsic complexity of learning. *Information and Computation*, 123, 64-71
- Freivalds, R., & Smith, C. H. (1993). On the role of procrastination in machine learning. *Information and Computation*, 107(2), 237–271.
- Jain, S., Kinber, E., & Wiehagen, R. (2001). Language learning from texts: Degrees of intrinsic complexity and their characterizations. Journal of Computer and System Sciences, 63, 305-354.
- Jain, S., & Sharma, A. (1993). On the non-existence of maximal inference degrees for language identification. *Information Processing Letters*, 47, 81–88.
- Jain, S., & Sharma, A. (1994). Program size restrictions in computational learning. Theoretical Computer Science, 127, 351-386.
- Jain, S., & Sharma, A. (1996). The intrinsic complexity of language identification. *Journal of Computer and System Sciences*, 52, 393-402.
- Jain, S., & Sharma, A. (1997). The structure of intrinsic complexity of learning. *Journal of Symbolic Logic*, 62, 1187-1201.
- Jain, S., & Sharma, A. (1997). Elementary formal systems, intrinsic complexity and procrastination. *Information and Computation*, 132, 65-84.
- Lange, S., & Zeugmann, T. (1996). Incremental learning from positive data. Journal of Computer and System Sciences, 53, 88-103.
- Pitt, L. (1989). Inductive inference, DFAs, and computational complexity. Analogical and inductive inference, second international workshop, AII 1989, LNAI. (Vol. 397, pp. 18-44) Heidelberg: Springer.
- Rogers, H. (1967). Theory of recursive functions and effective computability. New York: McGraw-Hill (Reprinted, MIT Press 1987).
- Shinohara, T. (1994). Rich classes inferable from positive data: Length-bounded elementary formal systems. *Information and Computation*, 108, 175-186.
- Slaman, T. A. & Solovay, R. (1991). When oracles do not help. Proceedings of the Fourth Annual Workshop on Computational Learning Theory, (pp. 379-383), Morgan Kaufmann.
- Wiehagen, R. (1976). Limes-Erkennung rekursiver Funktionen durch spezielle Strategien. Journal of Information Processing and Cybernetics (EIK), 12, 93-99.

Wiehagen, R. (1986). On the complexity of effective program synthesis. In: K. Jantke (Ed.). Analogical and Inductive Inference. *Proceedings of the International Workshop*, Springer LNCS, (Vol. 265, pp. 209–219).

Compositional Coevolution

Synonyms

Cooperative coevolution

Definition

A coevolutionary system constructed to learn composite solutions in which individuals represent different candidate components and must be evaluated together with other individuals in order to form a complete solution. Though not precisely the same as *cooperative coevolution*, there is a significant overlap.

Cross References

►Coevolutionary Learning

Computational Complexity of Learning

Sanjay Jain, Frank Stephan National University of Singapore, Singapore, Republic of Singapore

Definition

Measures of the complexity of learning have been developed for a number of purposes including ▶Inductive Inference, ▶PAC Learning, and ▶Query-Based Learning. The complexity is usually measured by the largest possible usage of ressources that can occur during the learning of a member of a class. Depending on the context, one measures the complexity of learning either by a single number/ordinal for the whole class or by a function in a parameter n describing the complexity of the target to be learnt. The actual measure can be the number of mind changes, the number of queries submitted to a teacher, the number of wrong conjectures issued, the number of errors made or the number of examples processed until learning succeeds. In addition to this, one can equip the learner with an oracle and determine the complexity of the oracle needed to perform

202 Computational Discovery of Quantitative Laws

the learning process. Alternatively, in complexity theory, instead of asking for an NP-complete oracle to learn a certain class, the result can also be turned into the form "this class is unlearnable unless RP = NP" or something similar. (Here RP is the class of decision problems solvable by a randomized polynomial time algorithm and NP is the class of decision problems solvable by a nondeterministic polynomial time algorithm and both algorithms never give "yes" answer for an instance of the problem with "no" answer.)

Detail

In PAC Learning, one usually asks how many examples are needed to learn the concept, where the number of examples needed mainly depends on the Vapnik Chervonenkis dimension of the class to be learnt, the error permitted, and the confidence required. Furthermore, for certain classes of finite Vapnik Chervonenkis dimension, learnability can still fail when the learner is required to be computable in polynomial time; hence there is, besides the dimension, also a restriction stemming from the computational complexity of problems such as the complexity of finding concepts consistent with all data observed so far.

For \times Query-Based Learning, one common criterion to be looked at is the number of queries made during the learning process. If a class contains 2^n different $\{0,1\}$ -valued functions f and one is required to learn the class using membership-queries, that is, by asking queries of the form whether f(x) = 0 or f(x) = 1, then there is a function f on which the learner needs at least n queries until it knows which of the given functions f is; for some classes consisting of 2^n functions the number of queries needed can be much worse - as much as $2^n - 1$. A well-known result of Angluin is that one can learn the class of all regular languages with polynomially many equivalence and membership queries measured with respect to the number of states of the smallest deterministic finite automaton accepting the language to be learnt. Further research has been done dealing with which query algorithms can be implemented by a polynomial time learner and which need for polynomial time learning, in addition to the teacher informing on the target concept, also some oracle supplying information that cannot be computed in polynomial time. See the entry \to Query-Based Learning for an overview of these results.

For Inductive Inference, most complexity measures are measures applying to the overall class and not just a parameterized version. When learning the class of all sets with up to *n* elements, the learner might first issue the conjecture \emptyset and then revise (up to n times) its hypothesis when a new datum is observed; such a measure is called the mind change complexity of learning. Mind change complexity has been generalized to measure the complexity by recursive ordinals or the notation of these. Furthermore, one can measure the long term memory of past data observed either by a certain number of examples remembered or by the number of bits stored on a tape describing the long-term memory of the learner. Besides these quantitative notions, a further frequently studied question is the following: Which oracles support the learning process in a way that some classes become learnable using the oracle, but are unlearnable without using any oracle? An example of such a type of result is that the class of all recursive functions can be learnt if and only if the learner has access to a high oracle, that is, an oracle that permits to compute a function which dominates (i.e., grows faster than) every recursive function. See the entry Complexity of Inductive Inference for more information.

Computational Discovery of Quantitative Laws

► Equation Discovery

Concept Drift

CLAUDE SAMMUT¹, MICHAEL HARRIES²

¹The University of New South Wales,
Sydney, Australia

²Advanced Products Group, Citrix Labs,
North Ryde, NSW, Australia

Synonyms

Context-sensitive learning; Learning with hidden context

Definition

Concept drift occurs when the values of hidden variables change over time. That is, there is some unknown

Concept Drift 203

C

context for **concept** learning and when that context changes, the learned concept may no longer be valid and must be updated or relearned.

Motivation and Background

Prediction in real-world domains is complicated by potentially unstable phenomena that are not known in advance to the learning system. For example, financial market behavior can change dramatically with changes in contract prices, interest rates, inflation rates, budget announcements, and political and world events. Thus, concept definitions that may have been learned in one context become invalid in a new context. This concept drift can be due to changes in context and is often directly reflected by one or more attributes. When changes in context are not reflected by any known attributes they can be said to be hidden. Hidden changes in context cause problems for any predictive approach that assumes concept stability.

Structure of the Learning System

Machine learning approaches can be broadly categorized as either ▶batch learning or ▶incremental learning. Batch systems learn off-line by examining a large collection of instances *en masse* and form a single concept. Incremental systems evolve and change a concept definition as new observations are processed (Schlimmer & Granger 1986a; Aha et al., 1991; Koltzer & Maloof, 2003).

The most common approach to learning in domains with hidden changes in context has been to use an incremental learning approach in which the importance of older items is progressively decayed. A popular implementation of this, originally presented in Kubet (1989), is to use a window of recent instances from which concept updates are derived. Other examples of this approach include Widmer and Kubat (1996), Kubat and Widmer (1995), Kilander and Jansson (1993), and Salganikoff (1993). Swift adaptation to changes in context can be achieved by dynamically varying the window size in response to changes in accuracy and concept complexity (Widmer & Kubat, 1996).

There are many domains in which the context can be expected not only to change but for earlier contexts to hold again at some time in the future. That is, contexts can repeat in domains such as financial prediction, dynamic control, and underrepresented data mining tasks. In these domains, prediction accuracy can be improved by storing knowledge about past contexts for reuse. FLORA3 (Widmer & Kubat, 1993) addresses domains in which contexts recur by storing and retrieving concepts that appear stable as the learner traverses the series of input data.

In many situations, there is no constraint to learn incrementally. For example, many organizations maintain large data bases of historical data that are suitable for data mining. These data bases may hold instances that belong to a number of contexts but do not have this context explicitly recorded. Many of these data bases may incorporate time as an essential attribute, for example, financial records and stock market price data. Interest in mining datasets of this nature suggests the need for systems that can learn global concepts and are sensitive to changing and hidden contexts. Systems such as FLORA3 also imply that an off-line recognition of stable concepts can be useful for \triangleright on-line prediction.

An alternative to on-line learning for domains with hidden changes in context is to examine the data *en masse* in an attempt to directly identify concepts associated with stable, hidden contexts. Some potential benefits of such an approach are:

- Context specific (known as local) concepts could be used as part of a multiple model on-line predictive system.
- 2. Local concepts could be verified by experts, or used to improve domain understanding.
- 3. A model of the hidden context could be induced using context characteristics such as context duration, order, and stability. The model could also use existing attributes and domain feedback if available.
- Stable contexts identified could be used as target characteristics for selecting additional attributes from the outside world as part of an iterative data mining process.

Splice (Harries, Sammut, & Horn, 1998) is a ▶metalearning system that implements a context sensitive batch learning approach. Splice is designed to identify intervals with stable hidden context, and to induce and refine local concepts associated with hidden contexts.

Identifying Context Change

In many domains with hidden changes in context, time can be used to differentiate hidden contexts. Most 204 Concept Drift

machine learning approaches to these domains do not explicitly represent time as they assume that current context can be captured by focusing on recent examples. The implication is that hidden context will be reflected in contiguous intervals of time. For example, an attempt to build a system to predict changes in the stock market could produce the following becision tree:

```
Year > 2002
Year < 2005
Attribute \ A = true: \ Market \ Rising
Attribute \ A = false: \ Market \ Falling
Year \ge 2005
Attribute \ B = true: \ Market \ Rising
Attribute \ B = false: \ Market \ Falling
```

This tree contains embedded knowledge about two intervals of time: in one of which, 2002–2004, attribute A is predictive; in the other, 2005 onward, attribute B is predictive. As time (in this case, year) is a monotonically increasing attribute, future classification using this decision tree will only use attribute B. If this domain can be expected to have recurring hidden context, information about the prior interval of time could be valuable.

The decision tree in the example above contains information about changes in context. We define context as:

Context is any attribute whose values are largely independent but tend to be stable over contiguous intervals of another attribute known as the environmental attribute.

The ability of decision trees to capture context is associated with the fact that decision tree algorithms use a form of context-sensitive feature selection (CSFS). A number of machine learning algorithms can be regarded as using CSFS including decision tree algorithms (Quinlan, 1993), ▶rule induction algorithms (Clark & Niblett, 1989), and ▶ILP systems (Quinlan, 1990). All of these systems produce concepts containing local information about context.

When contiguous intervals of time reflect a hidden attribute or context, we call time the environmental attribute. The environmental attribute is not restricted to time alone as it could be any ordinal attribute over which instances of a hidden context are liable to be contiguous. There is also no restriction, in principle, to one dimension. Some alternatives to time as environmental attributes are dimensions of space, and space–time combinations.

Given an environmental attribute, we can utilize a CSFS machine learning algorithm to gain information on likely hidden changes in context. The accuracy of the change points found will be dependent upon at least hidden context duration, the number of different contexts, the complexity of each local concept, and noise.

The CSFS identified context change points can be expected to contain errors of the following types:

- 1. Noise or serial correlation errors. These would take the form of additional incorrect change points.
- Errors due to the repetition of tests on time in different parts of the concept. These would take the form of a group of values clustered around the actual point where the context changed.
- 3. Errors of omission, change points that are missed altogether.

The initial set of identified context changes can be refined by contextual **clustering**.

This process combines similar intervals of the dataset, where the similarity of two intervals is based upon the degree to which a partial model is accurate on both intervals.

Recent Advances

With the increasing amount of data being generated by organizations, recent work on concept drift has focused on mining from high volume ▶data streams Hulten, Spencer, & Domingos, 2001; Wang, Fan, Yu, & Han, 2003; Koltzer & Maloof, 2003, Mierswa, Wurst, Klinkenberg, Scholz, & Euler, 2006; Chu & Zaniolo, 2004; Gaber, Zaslavsky, & Krishnaswamy, 2005. Methods such as Hulten *et al* s, combine decision tree learning with incremental methods for efficient updates, thus avoiding relearning large decision trees. Koltzer and Maloof also use incremental methods combined in an ▶ensemble.

Concept Learning 205

C

Cross References

- **▶**Decision Trees
- ►Ensemble Methods
- ►Incremental Learning
- ► Inductive Logic Programming
- ► Lazy Learning

Recommended Reading

- Aha, D. W., Kibler, D., & Albert, M. K. (1991). Instance-based learning algorithms. *Machine Learning*, 6, 37–66.
- Chu, F., & Zaniolo, C. (2004). Fast and light boosting for adaptive mining of data streams. In *Advances in knowledge discovery* and data mining. Lecture notes in computer science (Vol. 3056, pp. 282-292). Springer.
- Clark, P., & Niblett, T. (1989). The CN2 induction algorithm. Machine Learning, 3, 261–283.
- Clearwater, S., Cheng, T.-P., & Hirsh, H. (1989). Incremental batch learning. In *Proceedings of the sixth international workshop on machine learning* (pp. 366–370). Morgan Kaufmann.
- Domingos, P. (1997). Context-sensitive feature selection for lazy learners. *Artificial Intelligence Review*, 11, 227-253. [Aha, D. (Ed.). Special issue on lazy learning.]
- Gaber, M. M., Zaslavsky, A., & Krishnaswamy, S. (2005). Mining data streams: A review. SIGMOD Rec., 34(2), 18–26.
- Harries, M., & Horn, K. (1996). Learning stable concepts in domains with hidden changes in context. In M. Kubat & G. Widmer (Eds.), Learning in context-sensitive domains (workshop notes). 13th international conference on machine learning, Bari, Italy.
- Harries, M. B., Sammut, C., & Horn, K. (1998). Extracting hidden context. *Machine Learning*, 32(2), 101–126.
- Hulten, G., Spencer, L., & Domingos, P. (2001). Mining timechanging data streams. In KDD '01: Proceedings of the seventh ACM SIGKDD international conference on knowledge discovery and data mining (pp. 97-106). New York: ACM.
- Kilander, F., & Jansson, C. G. (1993). COBBIT A control procedure for COBWEB in the presence of concept drift. In P. B. Brazdil (Ed.), European conference on machine learning (pp. 244–261). Berlin: Springer.
- Kolter, J. Z., & Maloof, M. A. (2003). Dynamic weighted majority: A new ensemble method for tracking concept drift. In *Third IEEE international conference on data mining ICDM-2003* (pp. 123-130). IEEE CS Press.
- Kubat, M. (1989). Floating approximation in time-varying knowledge bases. Pattern Recognition Letters, 10, 223-227.
- Kubat, M. (1992). A machine learning based approach to load balancing in computer networks. Cybernetics and Systems Journal.
- Kubat, M. (1996). Second tier for decision trees. In Machine learning: Proceedings of the 13th international conference (pp. 293–301). California: Morgan Kaufmann.
- Kubat, M., & Widmer, G. (1995). Adapting to drift in continuous domains. In Proceedings of the eighth European conference on machine learning (pp. 307-310). Berlin: Springer.
- Mierswa, I., Wurst, M., Klinkenberg, R., Scholz, M., & Euler, T. (2006). Yale: Rapid prototyping for complex data mining tasks. In KDD '06: Proceedings of the 12th ACM SIGKDD international conference on knowledge discovery and data mining (pp. 935-940). New York: ACM.

- Quinlan, J. R. (1990). Learning logical definitions from relations. Machine Learning, 5, 239–266.
- Quinlan, J. R. (1993). C4.5: Programs for machine learning. Morgan Kaufmann: San Mateo.
- Salganicoff, M. (1993). Density adaptive learning and forgetting. In *Machine learning: Proceedings of the tenth international conference* (pp. 276-283). San Mateo: Morgan Kaufmann.
- Schlimmer, J. C., & Granger, R. I., Jr. (1986a). Beyond incremental processing: Tracking concept drift. In *Proceedings AAAI-86* (pp. 502–507). Los Altos: Morgan Kaufmann.
- Schlimmer, J., & Granger, R., Jr. (1986b). Incremental learning from noisy data. Machine Learning, *I*(3), 317–354.
- Turney, P. D. (1993a). Exploiting context when learning to classify. In P. B. Brazdil (Ed.), European conference on machine learning (pp. 402–407). Berlin: Springer.
- Turney, P. D. (1993b). Robust classification with context sensitive features. In Paper presented at the industrial and engineering applications of artificial intelligence and expert systems.
- Turney, P., & Halasz, M. (1993). Contextual normalization applied to aircraft gas turbine engine diagnosis. *Journal of Applied Intelligence*, 3, 109–129.
- Wang, H., Fan, W., Yu, P. S., & Han, J. (2003). Mining concept-drifting data streams using ensemble classifiers. In KDD '03: Proceedings of the ninth ACM SIGKDD international conference on knowledge discovery and data mining (pp. 226-235). New York: ACM.
- Widmer, G. (1996). Recognition and exploitation of contextual clues via incremental meta-learning. In L. Saitta (Ed.), Machine learning: Proceedings of the 13th international workshop (pp. 525–533). San Francisco: Morgan Kaufmann.
- Widmer, G., & Kubat, M. (1993). Effective learning in dynamic environments by explicit concept tracking. In P. B. Brazdil (Ed.), European conference on machine learning (pp. 227–243). Berlin: Springer.
- Widmer, G., & Kubat, M. (1996). Learning in the presence of concept drift and hidden contexts. *Machine Learning*, 23, 69-101.

Concept Learning

CLAUDE SAMMUT

The University of New South Wales, Sydney, NSW, Australia

Synonyms

Categorization; Classification learning

Definition

The term *concept learning* is originated in psychology, where it refers to the human ability to learn categories for object and to recognize new instances of those categories. In machine learning, concept is more formally

206 Concept Learning

defined as "inferring a boolean-valued function from training examples of its inputs and outputs" (Mitchell, 1997).

Background

Bruner, Goodnow, and Austin (1956) published their book *A Study of Thinking*, which became a landmark in psychology and would later have a major impact on machine learning. The experiments reported by Bruner, Goodnow, and Austin were directed toward understanding a human's ability to categorize and how categories are learned.

▶ We begin with what seems a paradox. The world of experience of any normal man is composed of a tremendous array of discriminably different objects, events, people, impressions... But were we to utilize fully our capacity for registering the differences in things and to respond to each event encountered as unique, we would soon be overwhelmed by the complexity of our environment... The resolution of this seeming paradox... is achieved by man's capacity to categorize. To categorize is to render discriminably different things equivalent, to group objects and events and people around us into classes... The process of categorizing involves... an act of invention... If we have learned the class "house" as a concept, new exemplars can be readily recognised. The category becomes a tool for further use. The learning and utilization of categories represents one of the most elementary and general forms of cognition by which man adjusts to his environment.

The first question that they had to deal with was that of representation: what is a concept? They assumed that objects and events could be described by a set of attributes and were concerned with how inferences could be drawn from attributes to class membership. Categories were considered to be of three types: conjunctive, disjunctive, and relational.

...when one learns to categorize a subset of events in a certain way, one is doing more than simply learning to recognise instances encountered. One is also learning a rule that may be applied to new instances. The concept or category is basically, this "rule of grouping" and it is such rules that one constructs in forming and attaining concepts.

The notion of a rule as an abstract representation of a concept influenced research in machine learning. For example, ▶decision tree learning was used as a means of creating a cognitive model of concept learning (Hunt, Martin, & Stone, 1966). This model later inspired Quinlan's development of ID3 (Quinlan, 1983).

The learning experience may be in the form of examples from a trainer or the results of trial and error. In either case, the program must be able to represent its observations of the world, and it must also be able to represent hypotheses about the patterns it may find in those observations. Thus, we will often refer to the bobservation language and the hypothesis language. The observation language describes the inputs and outputs of the program and the hypothesis language describes the internal state of the learning program, which corresponds to its theory of the concepts or patterns that exist in the data.

The input to a learning program consists of descriptions of objects from the universe and, in the case of supervised learning, an output value associated with the example. The universe can be an abstract one, such as the set of all natural numbers, or the universe may be a subset of the real world. No matter which method of representation we choose, descriptions of objects in the real world must ultimately rely on measurements of some properties of those objects. These may be physical properties such as size, weight, and color or they may be defined for objects, for example, the length of time a person has been employed for the purpose of approving a loan. The accuracy and reliability of a learned concept depends on the accuracy and reliability of the

A program is limited in the concepts that it can learn by the representational capabilities of both observation and hypothesis languages. For example, if an attribute/value list is used to represent examples for an induction program, the measurement of certain attributes and not others clearly places bounds on the kinds of patterns that the learner can find. The learner is said to be *biased* by its observation language (see Language Bias). The hypothesis language also places constraints on what may and may not be learned. For

Concept Learning 207

example, in the language of attributes and values, relationships between objects are difficult to represent. Whereas, a more expressive language, such as first-order logic, can easily be used to describe relationships.

Unfortunately, representational power comes at a price. Learning can be viewed as a search through the space of all sentences in a language for a sentence that best describes the data. The richer the language, the larger is the search space. When the search space is small, it is possible to use "brute force" search methods. If the search space is very large, additional knowledge is required to reduce the search.

Rules, Relations, and Background Knowledge

In the early 1960s, there was no discipline called "machine learning." Instead, learning was considered to be part of "pattern recognition," which had not yet split from AI. One of the main problems addressed at that time was how to represent patterns so that they could be recognized easily. Symbolic description languages were developed to be expressive and learnable.

Banerji (1960, 1962) first devised a language, which he called a "description list," which utilized an object's attributes to perform pattern recognition. Pennypacker, a masters student of Banerji at the Case Institute of Technology, implemented the recognition procedure and also used Bruner, Goodnow, and Austin's *Conservative Focussing Strategy* to learn conjunctive concepts (Pennypacker, 1963). Bruner, Goodnow, and Austin describe the strategy as follows:

▶ ... this strategy may be described as finding a positive instance to serve as a focus, then making a sequence of choices each of which alters but one attribute value [of the focus] and testing to see whether the change yields a positive or negative instance. Those attributes of the focus which, when changed, still yield positive instance are *not* part of the concept. Those attributes of the focus that yield negative instances when changed *are* features of the concept.

The strategy is only capable of learning *conjunctive concepts*, that is, the concept description can only consist of a simple conjunction of tests on attribute values. Recognizing the limitations of simple attribute/value representations, Banerji (1964) introduced the use of

predicate logic as a description language. Thus, Banerji was one of the earliest advocates of what would, many years later, become *Inductive Logic Programming*.

In the 1970s, a series of algorithms emerged that developed concept learning further. Winston's ARCH program (Winston, 1970) was influential as one of the first widely known concept learning programs. Michalski (1973, 1983) devised the Aq family of learning algorithms that set some of the early benchmarks for learning programs. Early relational learning programs were developed by Hayes-Roth (1973), Hayes-Roth and McDermott (1977), and Vere (1975, 1977).

Banerji emphasized the importance of a description language that could "grow." That is, its descriptive power should increase as new concepts are learned. These concepts become background knowledge for future learning. A simple example from Banerji (1980) illustrates the use of background knowledge. There is a language for describing instances of a concept and another for describing concepts. Suppose we wish to represent the binary number, 10, by a left-recursive binary tree of digits "0" and "1":

```
[head: [head: 1; tail: nil]; tail: 0]
```

"head" and "tail" are the names of attributes. Their values follow the colon. The concepts of binary digit and binary number are defined as

$$x \in digit \equiv x = 0 \lor x = 1$$

 $x \in num \equiv (tail(x) \in digit \land head(x) = nil)$
 $\lor (tail(x) \in digit \land head(x) \in num)$

Thus, an object belongs to a particular class or concept if it satisfies the logical expression in the body of the description. Note that the concept above is *disjunctive*. Predicates in the expression may test the membership of an object in a previously learned concept and can express *relations* between objects. Cohen and Sammut (1982) devised a learning system based on Banerji's ideas of a growing concept description language and this was further extended by Sammut and Banerji (1986).

Concept Learning and Noise

One of the most severe drawbacks of early concept learning systems was that they assumed that data sets

208 Conditional Random Field

were not noisy. That is, all attribute values and class labels in the training data are assumed to be correct. This is unrealistic in most real applications. Thus, concept learning systems began incorporating statistical measures to minimize the effects of noise and to estimate error rates (Breiman, Friedman, Olshen, & Stone, 1984; Cohen, 1995; Quinlan, 1986, 1993).

Learning to classify objects from training examples has gone on to become one of the central themes of machine learning research. As the robustness of classification systems has increased, they have found many applications, particularly in data mining but in a broad range of other areas.

Cross References

- **▶**Data Mining
- ▶ Decision Tree Learning
- ► Inductive Logic Programming
- ► Learning as Search
- ▶ Relational Learning
- ► Rule Learning

Recommended Reading

- Banerji, R. B. (1960). An information processing program for object recognition. General Systems, 5, 117–127.
- Banerji, R. B. (1962). The description list of concepts. Communications of the Association for Computing Machinery, 5(8), 426-432.
- Banerji, R. B. (1964). A Language for the Description of Concepts. General Systems, 9, 135–141.
- Banerji, R. B. (1980). Artificial intelligence: A theoretical approach. New York: North Holland.
- Breiman, L., Friedman, J. H., Olshen, R. A., & Stone, C. J. (1984). Classification and regression trees. Belmont, CA: Wadsworth.
- Bruner, J. S., Goodnow, J. J., & Austin, G. A. (1956). A study of thinking. New York: Wiley.
- Cohen, B. L., & Sammut, C. A. (1982). Object recognition and concept learning with CONFUCIUS. Pattern Recognition Journal, 15(4), 309-316.
- Cohen, W. W. (1995). In fast effective rule induction. In *Proceedings* of the twelfth international conference on machine learning, Lake Tahoe, California. Menlo Park: Morgan Kaufmann.
- Hayes-Roth, F. (1973). A structural approach to pattern learning and the acquisition of classificatory power. In First international joint conference on pattern recognition (pp. 343-355). Washington, D.C.
- Hayes-Roth, F., & McDermott, J. (1977). Knowledge acquisition from structural descriptions. In Fifth international joint conference on artificial intelligence (pp. 356–362). Cambridge, MA.

- Hunt, E. B., Marin, J., & Stone, P. J. (1966). Experiments in induction. New York: Academic.
- Michalski, R. S. (1973). Discovering classification rules using variable valued logic system VL1. In *Third international joint conference on artificial intelligence* (pp. 162–172). Stanford, CA.
- Michalski, R. S. (1983). A theory and methodology of inductive learning. In R. S. Michalski, J. G. Carbonell, & T. M. Mitchell (Eds.), *Machine learning: An artificial intelligence approach*. Palo Alto: Tioga.
- Mitchell, T. M. (1997). Machine learning. New York: McGraw-Hill.
- Pennypacker, J. C. (1963). An elementary information processor for object recognition. SRC No. 30-I-63-1. Case Institute of Technology.
- Quinlan, J. R. (1983). Learning efficient classification procedures and their application to chess end games. In R. S. Michalski, J. G. Carbonell, & T. M. Mitchell (Eds.), Machine learning: An artificial intelligence approach. Palo Alto: Tioga.
- Quinlan, J. R. (1986). The effect of noise on concept learning. In R. S. Michalski, J. G. Carbonell, & T. M. Mitchell (Eds.), *Machine learning: An artificial intelligence approach* (Vol. 2). Los Altos: Morgan Kaufmann.
- Quinlan, J. R. (1993). C4.5: Programs for machine learning. San Mateo, CA: Morgan Kaufmann.
- Sammut, C. A., & Banerji, R. B. (1986). Learning concepts by asking questions. In R. S. Michalski, J. G. Carbonell, & T. M. Mitchell (Eds.), Machine learning: An artificial intelligence approach (Vol. 2, pp. 167–192). Los Altos, CA: Morgan-Kaufmann.
- Vere, S. (1975). Induction of concepts in the predicate calculus. In Fourth international joint conference on artificial intelligence (pp. 351–356). Tbilisi, Georgia, USSR.
- Vere, S. A. (1977). Induction of relational productions in the presence of background information. In *Fifth international joint conference on artificial intelligence*. Cambridge, MA.
- Winston, P. H. (1970). Learning structural descriptions from examples. Unpublished PhD Thesis, MIT Artificial Intelligence Laboratory.

Conditional Random Field

A Conditional Random Field is a form of ▶Graphical Model for segmenting and ▶classifying sequential data. It is the ▶discriminative learning counterpart to the ▶generative learning Markov Chain model.

Recommended Reading

Lafferty, J., McCallum, A., & Pereira, F. (2001). Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the 18th international conference on machine learning* (pp. 282–289). San Francisco, Morgan Kaufmann.

Confirmation Theory

The branch of philosophy concerned with how (and indeed whether) evidence can confirm a hypothesis, even though typically it does not entail it. A distinction is sometimes drawn between *total confirmation*: how well confirmed a hypothesis is, given all available evidence and *weight-of-evidence*: the amount of extra confirmation added to the total confirmation of a hypothesis by a particular piece of evidence. Confirmation is often measured by the probability of a hypothesis conditional on evidence.

Confusion Matrix

KAI MING TING Monash University, Australia

Definition

A confusion matrix summarizes the classification performance of a \triangleright classifier with respect to some \triangleright test data. It is a two-dimensional matrix, indexed in one dimension by the true class of an object and in the other by the class that the classifier assigns. Table 1 presents an example of confusion matrix for a three-class classification task, with the classes A, B, and C.

The first row of the matrix indicates that 13 objects belong to the class A and that 10 are correctly classified as belonging to A, two misclassified as belonging to B, and one as belonging to C.

A special case of the confusion matrix is often utilized with two classes, one designated the *positive* class and the other the *negative* class. In this context, the four cells of the matrix are designated as ▶ *true positives* (TP), ▶ *false positives* (FP), ▶ *true negatives* (TN), and ▶ *false negatives* (FN), as indicated in Table 2.

A number of measures of classification performance are defined in terms of these four classification outcomes.

Confusion Matrix. Table 1 An example of three-class confusion matrix

		Assigned Class		
		A	В	С
Actual Class	Α	10	2	1
	В	0	6	1
	С	0	3	8

Confusion Matrix. Table 2 The outcomes of classification into positive and negative classes

		Assigned Class		
		Positive	Negative	
Actual Class	Positive	TP	FN	
	Negative	FP	TN	

- ► Positive predictive value = ► Precision = TP/(TP + FP)
 - ► *Negative predictive value* = TN/(TN + FN)

Conjunctive Normal Form

Bernhard Pfahringer University of Waikato, Hamilton, New Zealand

Conjunctive normal form (CNF) is an important normal form for propositional logic. A logic formula is in conjunctive normal form if it is a single conjunction of disjunctions of (possibly negated) literals. No more nesting and no other negations are allowed. Examples are:

$$a$$

$$\neg b$$

$$a \wedge b$$

$$(a \vee \neg b) \wedge (c \vee d)$$

$$\neg a \wedge (b \vee \neg c \vee d) \wedge (a \vee \neg d)$$

210 Connection Strength

Any arbitrary formula in propositional logic can be transformed into conjunctive normal form by application of the laws of distribution, De Morgan's laws, and by removing double negations. It is important to note that this process can lead to exponentially larger formulas which implies that the process in the worst case runs in exponential time. An example for this behavior is the following formula given in bedisjunctive normal form (DNF), which is linear in the number of propositional variables in this form. When transformed into conjunctive normal form (CNF), its size is exponentially larger.

DNF:
$$(a_0 \land a_1) \lor (a_2 \land a_3) \lor ... \lor (a_{2n} \land a_{2n+1})$$

CNF:
$$(a_0 \lor a_2 \lor \ldots \lor a_{2n}) \land (a_1 \lor a_2 \lor \ldots \lor a_{2n})$$

 $\land \ldots \land (a_1 \lor a_3 \lor \ldots \lor a_{2n+1})$

Recommended Reading

Russell, S., & Norvig, P. (2002). Artificial intelligence: A modern approach (p. 215). Prentice Hall

Connection Strength

▶Weight

Connections Between Inductive Inference and Machine Learning

John Case¹, Sanjay Jain²

¹University of Delaware, Newark, USA

²National University of Singapore, Singapore, Republic of Singapore

Definition

Inductive inference is a theoretical framework to model learning in the limit. Here we will discuss some results in inductive inference, which have relevance to machine learning community.

The mathematical/theoretical area called ▶Inductive Inference, is also known as *computability theoretic learning* and *learning in the limit* (Jain, Osherson, Royer, &

Sharma, 1999; Odifreddi, 1999) typically *but*, *as will be seen below, not always* involves a situation depicted in (1) just below.

Data
$$d_0, d_1, d_2, \dots \xrightarrow{\text{In}} M \xrightarrow{\text{Out}} \text{Programs } e_0, e_1, e_2, \dots$$
 (1)

Let $\mathbb{N}=$ the set of nonnegative integers. Strings, including program strings, computer reals, and other data structures, inside computers, are finite bit strings and, hence, can be coded into \mathbb{N} . Therefore, mathematically at least, it is without loss of mathematical generality that we sometimes use the data type \mathbb{N} where standard practice would use a different type.

In (1), d_0, d_1, d_2, \ldots can be, e.g., the successive values of a function $f: \mathbb{N} \to \mathbb{N}$ or the elements of a (formal) language $L \subseteq \mathbb{N}$ in some order; M is a machine; the e_i 's are from some hypothesis space of programs; and, for M's successful learning, later e_i 's exactly or approximately compute the f or L.

Such learning is off-line: in successful cases, one comes away with programs for past and future data. For the related problem of online extrapolation of next values for a function f, suitable e_i 's may be the values of f(i)'s based on having seen strictly prior values of f.

Detail

We will discuss the off-line case until we say otherwise. It is typical in applied machine learning to present to a learner whatever data one has and to obtain one corresponding program hopefully for predicting these data and future data. In inductive inference the case where only one program is output is called one-shot learning. More typically, in inductive inference, one allows for mind-changes, i.e., for a succession of output programs, as one receives successively more input data, with the later programs hopefully eventually being useful for predictions. Typically, one does not get success on one's first conjecture/output program, but rather, one may achieve success eventually, or, as it is said, in the limit after some sequence of trial and error. It is helpful at this juncture to present a problem for which this latter approach makes more sense than the one-shot approach.

We will consider some different criteria of successful learning of f or L by M. For example, **Ex**-style criteria

will require that all but finitely many of the e_i 's are syntactically the same and do a reasonable job of computing the f or L. **Bc**-style criteria are more relaxed, more powerful, but less useful (Bārzdiņš, 1974; Case & Lynes, 1982; Case & Smith, 1983): they do not require almost all e_i 's be the same syntactically.

Here is a well-known regression technique from, e.g., (Hildebrand, 1956), for exactly "curve-fitting" polynomials. It is the method involving calculating *forward differences*. We express it as a learning machine M_0 and illustrate with its being fed an *example* data sequence generated by a cubic polynomial

$$x^3 - 2x^2 + 2x + 3. (2)$$

See (Hildebrand, 1956), for how to recover the polynomials themselves.

 M_0 , fed a finite data sequence of natural numbers, first looks for iterated forward differences to become (apparently) constant, then outputs a rule/program, which uses the (apparent) constant to extrapolate the data sequence for any desired prediction. For example, were M_0 given the data sequence in the top row of Table 1, it would calculate 6 to be the apparent constant after *three* differencings, so M_0 then outputs the following informal rule/program.

▶ To generate the level 0 sequence, at level 0, start with 3; at level 1, start with 1; at level 2, start with 2; add the apparent constant 6 from level 3 to get successive level 2 data items; add successive level 2 items to get successive level 1 data items; finally, add successive level 1 items to get as many successive level 0 data items as needed for prediction.

This program, eventually output by M_0 when its input the whole top row of Table 1, correctly predicts

Connections Between Inductive Inference and Machine Learning. Table 1 Example Sequence and Its Iterated Forward Differences

Sequence:	3		4		7		18		43
1st Diffs:		1		3		11		25	
2nd Diffs:			2		8		14		
3rd Diffs:				6		6			

the elements of the cubic polynomial, on successive values in \mathbb{N} – the whole sequence 3, 4, 7, 18, 43, 88, 159, Along the way, though, just after the first data point, M_0 thinks the apparent constant is 0; just after the second that it is 1; just after the third that it is 2; and only after more of the data points does it converge for this cubic polynomial to the apparent (and, on *this* example, actual) constant 6. In general, M_0 , on a polynomial of degree m, *changes its mind* up to m times until converging to its final program (of course on $f(x) = 2^x$, M_0 never converges, and each level of forward differences is just the sequence f again.).

Hence, M_0 above **Ex**-learns, e.g., the integer polynomials $f: N \to N$, but it does not in general one-shot learn these polynomials – since the data alone do not disclose the degree of a generating polynomial.

In this entry we survey some results from inductive inference but with an eye to topics having something to say regarding or to applied machine learning. In some cases, the theoretical results lend mathematical support to preexisting empirical observations about the efficacy of known machine learning techniques. In other cases, the theoretical results provide some, typically abstract, suggestions for the machine learning practitioner. In some of these cases, some of the suggestions apparently pay off in others, intriguingly, we do not know yet.

Multi-Task or Context Sensitive Learning

In empirical, applied machine learning, multitask or context sensitive learning involves trying to learn Y by first (de Garis, 1990a, b; Fahlman, 1991; Thrun, 1996; Thrun & Sullivan, 1996; Tsung & Cottrell, 1989; Waibel, 1989a, b) or simultaneously (Caruana, 1993, 1996; Dietterich, Hild, & Bakiri, 1995; Matwin & Kubat, 1996; Mitchell, Caruana, Freitag, McDermott, & Zabowski, 1994; Pratt, Mostow, & Kamm, 1991; Sejnowski & Rosenberg, 1986; Bartlmae, Gutjahr, & Nakhaeizadeh, 1997) trying to learn also X – even in cases where there may be no inherent interest in learning X (see also \triangleright Transfer Learning). There is, in many cases, an apparent empirical advantage in doing this for some X, Y. It can happen that *Y* is not apparently or easily learnable by itself, but is learnable if one learns *X* first or simultaneously in some case X itself can be a sequence of tasks X_1, \ldots, X_n . Here the X_i s may need to be learned sequentially or simultaneously to learn Y. For example, to teach a robot to drive

a car, it is useful to train it also to predict the center of the road markings (see, e.g., Baluja & Pomerleau, 1995; Caruana, 1996). For another example: an experimental system to predict the value of German *Daimler* stock performed better when it was modified to track simultaneously the German stock-index DAX (Bartlmae et al., 1997). The value of the Daimler stock here was the primary or target concept and the value of the DAX – a related concept – provided useful auxiliary context.

Angluin, Gasarch, and Smith (1989) shows *mathematically* that, in effect, there are (mathematical) learning scenarios for which it was *provable* that *Y* could not be learned without learning *X first* – and, in other scenarios (Angluin et al., 1989; Kinber, Smith, Velauthapillai, & Wiehagen, 1995), *Y* could not be learned without *simultaneously* learning *X*. These *mathematical* results provide a kind of evidence that the *empirical* observations as to the *apparent* usefulness of multitask or context sensitive learning *may* not be illusionary, luck, or a mere accident of happening to use some data sets but not others.

For illustration, here is a particularly simple theoretical example needing to be learned *simultaneously* and similar to examples in Angluin et al. (1989). Let \mathcal{R} be the set of all computable functions mapping \mathbb{N} to \mathbb{N} . We use numerical names in \mathbb{N} for programs. Let

$$S = \{ (f,g) \in \mathcal{R} \times \mathcal{R} \mid f(0) \text{ is a program for } g \land g(0) \text{ is a program for } f \}.$$
 (3)

We say (p,q) is a program for $(f,g) \in \mathcal{R} \times \mathcal{R}$ iff p is a program for f and q is a program for g.

Consider a machine M which, if, as in (1), M is fed d_0, d_1, \ldots , but where each d_i is (f(i), g(i)), then M outputs each $e_i = (g(0), f(0))$. Clearly, M one-shot learns S. It can be easily shown that the component f's and g's for $(f,g) \in S$ are not separately even \mathbf{Bc} -learnable. It is important to note that, perhaps quite unlike real-world problems, the definition of this example S employs a simple self-referential coding trick: useful programs are coded into values of the functions at argument zero. A number of inductive inference results have been proved by means of (sometimes more complicated) self-referential coding tricks (see, e.g., Case, 1994). Bārzdiņš indirectly (see Zeugmann, 1986) provided a kind of informal robustness idea in his attempt to be rid of such coding tricks in inductive inference.

More formally, Fulk (1990) considered a learnability result involving a witnessing class C of (tuples of) functions to be *robust* iff each computable scrambling of Calso witnesses the learnability result (the allowed computable scramblers are the *general recursive operators* of (Rogers, 1967), but we omit the formal details herein.) Example: A simple shift scrambler converting each *f* to f', where f'(x) = f(x+1), would eliminate the coding tricks just above – since the values of f at argument zero would be lost in this scrambling. Some inductive inference results hold robustly and some not (see, e.g., Fulk, 1990; Jain, 1999; Jain, Smith, & Wiehagen, 2001; Jain et al., 1999; Case, Jain, Ott, Sharma, & Stephan, 2000). Happily, the $S \subseteq \mathcal{R} \times \mathcal{R}$ above (that is, learnable, but its components not) can be replaced by a more complicated class S' that robustly witnesses the same result. This is better theoretical evidence that the empirically noticed efficacy of multitask or context sensitive learning is not just an accident. It is residually important to note that (Jain et al., 2001) shows, though, that the computable scramblers can not get rid of more sophisticated coding tricks they called topological. S' mentioned just above turns out to *employ* this latter kind of coding trick. It is hypothesized in (Case et al., 2000) that nature likely employs some sophisticated coding tricks itself. For a separate informal argument about coding tricks of nature, see (Case, 1999). Ott and Stephan (2002) introduces a finite invariance constraint on top of robustness. This so-called hyperrobustness does destroy all coding tricks, and the result about the theoretical efficacy of multitask or context sensitive learning is *not* hyperrobust. However, hyperrobustness, perhaps, leaves unrealistically sparse structure.

Final note: Machine learning is an engineering endeavor. However, philosophers of science as well as practitioners in classical scientific disciplines should likely be considering the relevance of multitask or context sensitive inductive inference to their endeavors.

Special Cases of Inductive Logic Programming

In this section we discuss some *learning in the limit* results for *elementary formal systems (EFSs)* (Smullyan, 1961). Essentially, EFSs are programs in a string rewriting system. It is well known (Arikawa, Shinohara, & Yamamoto, 1992) that EFSs are essentially (pure) logic

programs over strings. Hence, the results have possible relevance for ▶inductive logic programming (ILP) (Bratko & Muggleton, 1995; Lavrač & Džeroski, 1994; Mitchell, 1997; Muggleton & De Raedt, 1994).

First we will discuss some important special cases based on Angulin's *pattern languages* (Angluin, 1980).

A *pattern language* is (by definition) one generated by all the positive length substitution instances in a *pattern*, such as,

$$abXYcbbZXa$$
 (4)

— where the variables (for substitutions) are depicted in upper case and the constants/terminals in lower case and are from, say, the alphabet {a,b,c}. Just below is an EFS or logic program based on this example pattern.

$$abXYcbbZXa \leftarrow .$$
 (5)

It must be understood, though, that in (5) and in the next example EFS below, only positive length strings are allowed to be substituted for the variables.

Angluin (1980) showed the **Ex**-learnability of the class of pattern languages from positive data. For these results, in the paradigm of (1) above $d_0, d_1, d_2, ...$ is a listing or presentation of some formal language L over a finite nonempty alphabet and the e_i 's are programs that generate languages. In particular, for Angluin's M, for L a pattern language, the e_i 's are patterns, and, for each presentation of L, all but finitely many of the corresponding e_i 's are the same *correct* pattern for L.

Much work has been done on the learnability of pattern languages, e.g., Salomaa (1994a, b); Case, Jain, Kaufmann, Sharma, and Stephan (2001), and on bounded finite unions thereof, e.g., Shinohara (1993); Wright (1989); Kilpeläinen, Mannila, and Ukkonen (1995); Brazma, Ukkonen, and Vilo (1996); Case, Jain, Lange, and Zeugmann (1999).

Regarding bounded finite unions of pattern languages: an n-pattern language is the union of the pattern languages for some n patterns P_1, \ldots, P_n . Each n-pattern language is also Ex-learnable from positive data (see Wright (1989)). An EFS or logic program corresponding to the n-patterns P_1, \ldots, P_n and generating the corresponding n-pattern language is just below.

$$P_1 \leftarrow .$$

$$\vdots$$

$$P_n \leftarrow .$$

Pattern language learning algorithms have been successfully applied toward some problems in molecular biology, see, e.g., Shimozono et al. (1994), Shinohara and Arikawa (1995).

Lange and Wiehagen (1991) presents an interesting iterative (Wiehagen, 1976) algorithm learning the class of pattern languages - from positive data only and with polynomial time constraints. Iterative learners are Ex-learners for which each output depends only on its just prior output (if any) and the input data element currently seen. Their algorithm works in polynomial time (actually quadratic time) in the length of the latest data item and the previous hypothesis. Furthermore, the algorithm has a linear set of good examples, in the sense that if the input data contains these good examples, then the algorithm already converges to the correct hypothesis. The number of good examples needed is at most |P| + 1, where P is a pattern generating the data d_0, d_1, d_2, \dots for the language being learned. This algorithm may be useful in practice due to its fast run time, and being able to converge quickly, if enough good data is available early. Furthermore, due to iterativeness, it does not need to store previous data!

Zeugmann (1998) considers total learning time up to convergence of the algorithm just discussed in the just prior paragraph. Note that, for arbitrary presentations, d_0, d_1, d_2, \ldots , of a pattern language, this time can be unbounded. In the best case it is polynomial in the length of a generating pattern P, where d_0, d_1, d_2, \dots is based on using P to get good examples early - in fact the time taken in the best case is $\Theta(|P|^2 \log_s(s +$ k)), where P is the pattern, s is the alphabet size, and k is the number of variables in P. Much more interesting is the case of average time taken up to convergence. The probability distribution (called uniform by Zeugmann) considered is as follows. A variable X is replaced by a string w with probability $\frac{1}{(2s)^{|w|}}$ (i.e., all strings of length r together have probability 2^{-r} , and the distribution is uniform among strings of length r). Different variables are replaced independently of each other. In this case the average total time up to convergence is $O(2^k k^2 s |P|^2 log_s(ks))$. The main thing is that for average case on probabilistic data (as can be expected in real life, though not necessarily with this kind of uniform distribution), the algorithm converges pretty fast and computations are done efficiently.

A number of papers consider Ex-learning of EFSs (Krishna Rao, 1996; Krishna Rao, 2000, 2004, 2005; Krishna Rao & Sattar, 1998) including with various bounds on the number of mind-changes until syntactic convergence to correct programs (Jain & Sharma, 1997, 2002). The EFSs considered are patterns, n-patterns, those with a constant bound on the length of clauses, and some with constant bounds on search trees. The mind-change bounds are typically more dynamic than those given by constants: they involve counting down from finite representations (called notations) for infinite constructive ordinals. An example of this kind of bound: one can algorithmically, based on some input parameters, decide how many mind-changes will be allowed. In other examples, the decision as to how many mindchanges will be allowed can be algorithmically revised some constant number of times. It is possible that not yet created special cases of some of these algorithms could be made feasible enough for practice.

Learning Drifting Concepts

A drifting concept to be learned is one which is a moving target (see ▶Concept Drift). In some machine learning applications, concept drift must be dealt with (Bartlett, Ben-David, & Kulkarni, 1996; Blum & Chalasani, 1992; Devaney & Ram, 1994; Freund & Mansour, 1997; Helmbold and Long, 1994; Kubat, 1992; Widmer & Kubat, 1996; Wrobel, 1994). An inductive inference contribution is (Case et al., 2001) in which it is shown, for online extrapolation by computable martingale betting strategies, upper bounds on the "speed" of the moving target that permit success at all. Here success is to make unbounded amounts of "money" betting on correctness of ones extrapolations. Here is an illustrative result from (Case et al., 2001). For the pattern languages considered in the previous section, only positive length strings of terminals can be substituted for a variable in an associated pattern. The (difficult to learn) pattern languages with erasing are just the languages obtained by also allowing the substitution of the empty string for variables in a pattern. For our example, we restrict the terminal alphabet to be {0,1}. With each pattern language with erasing L (over this terminal alphabet) we associate its characteristic function χ_L , which is 1 on terminal strings in L and 0 on those not in L. For ε denoting the empty string,

and for the terminal strings in length-lexicographical order, ε , 0, 1, 00, 01, 10, 11, 000, . . . , we would input a χ_L itself to a potential extrapolating machine as the bit string, $\chi_L(\varepsilon)$, $\chi_L(0)$, $\chi_L(1)$, $\chi_L(00)$, $\chi_L(01)$, Let $\mathcal E$ be the class of these characteristic functions. Pick a positive integer constant p. To model *drift with permanence* p, we imagine that a potential extrapolator for $\mathcal E$ receives successive bits from a member of $\mathcal E$ but keeps switching to the next bits of another, etc., *but* it must see at least p bits in a row of each member of $\mathcal E$ it sees before it can see the next bits of another. p is, then, a speed limit on drift. The result is that some suitably clever computable martingale betting strategy is successful at extrapolating $\mathcal E$ with drift permanence (speed limit on drift) of p=7.

Behavioral Cloning

Kummer and Ott (1996); Case, Ott, Sharma, and Stephan (2002) studied learning in the limit of winning control strategies for closed computable games. These games nicely model reactive process-control problems. Included are such example process-control games as regulating temperature of a room to be in a desired interval, forever after no more than some fixed number of moves between the thermostat and processes disturbing the temperature (Roughly, closed computable games are those so that one can tell algorithmically when one has lost. A temperature control game that requires stability forever after some undetermined finite number of moves is *not* a closed computable game. For a more formal treatment, see Cenzer and Remmel (1992); Maler, Pnueli, and Sifakis (1995); Thomas (1995); Kummer and Ott (1996)).

In machine learning, there are cases where one wants to teach a machine some motor skill possessed by human experts and where these human experts do not have access to verbalizable knowledge about how they perform expertly. Piloting an aircraft or expert operation of a swinging shipyard crane provide examples, and machine learning employs, in these cases, ▶ behavioral cloning, which uses direct performance data from the experts (Bain & Sammut, 1999; Bratko, Urbančič, & Sammut, 1998; Šuc, 2003).

Case et al. (2002) studies the effects on learning in the limit *closed computable games* where the learning procedures also had access to the *behavioral performance* (but not the algorithms) of masters/experts at the

games. For example, it is showed that, in some cases, there is better performance cloning n+1 disparate masters over cloning only n. For a while it was not known in machine learning how to clone multiple experts even after Case et al. (2002) was known to some; however, independently of Case et al., 2002, and later, Dorian Šuc (Šuc, 2003) found a way to clone behaviorally more than one human expert simultaneously (for the free-swinging shipyard crane problem) – by having more than one level of feedback control, and he got enhanced performance from cloning the multiple experts!

Learning To Coordinate

Montagna and Osherson (1999) begins the study of learning in the limit to *coordinate* (digital) moves between at least two agents.

The machines of Montagna and Osherson (1999) are, in effect, general extrapolating devices (Montagna & Osherson, 1999; Case et al., 2005). Technically, and without loss of generality of the results, we restrict the moves of each coordinator to bits, i.e., zeros and ones. *Coordination* is achieved between two coordinators iff each, reacting to the bit sequence of the other, eventually (in the limit) matches it bit for bit. Montagna and Osherson (1999) gives an example of two people who show up in a park each day at one of noon (bit 0) or 6pm (bit 1); each *silently* watches the other's past behavior; and each *tries*, based on the past behavior of the other, to show up eventually exactly when the other shows up. If they manage it, they have learned to coordinate.

A *blind* coordinator is one that reacts only to the *presence* of a bit from another process, *not* to *which* bit the other process has played (Montagna and Osherson, 1999).

In Case et al. (2005) is developed and studied the notion of probabilistically correct algorithmic coordinators. Next is a sample of theorems to the effect that just a few random bits can enhance learning to coordinate.

Theorem 1 (Case et al., 2005) Suppose $0 \le p < 1$. There exists a class of deterministic algorithmic coordinators C such that

(1) *No* deterministic algorithmic coordinator can coordinate with all of C; and

- (2) For k chosen so that $1 2^{-k} \ge p$, there exists a blind, *probabilistic* algorithmic coordinator **PM**, such that:
 - (i) For each member of C, **PM** can coordinate with with probability $1 2^{-k} \ge p$; and
 - (ii) **PM** is *k*-memory limited in the sense of (Osherson, Stob, & Weinstein, 1986, P. 66); more specifically, **PM** needs to remember whether it is outputting one of its first *k* bits which are its only random bits (e.g., for $p = \frac{255}{256}$, a mere k = 8 random bits suffice.)

Regarding possible eventual applicability: Maye, Hsieh, Sugihara, and Brembs (2007) cites finding deterministic chaos but *not* randomness in the behavior of *animals*. Hence, animals may not be exploiting random bits in learning anything, including to coordinate. However, one might build artifactual devices to exploit randomness, say, from radioactive decay, including, then, for enhancing learning to coordinate.

Learning Geometric Clustering

Case, Jain, Martin, Sharma, and Stephen (2006) showed that learnability in the limit of ▶ clustering, with or without additional information, depends strongly on geometric constraints on the shape of the clusters. In this approach the hypothesis space of possible clusters is pre-given in each setting. It was hoped to obtain thereby insight into the difficulty of clustering when the clusters are restricted to preassigned geometrically defined classes.

This is interestingly complementary to the *conceptual clustering* approach (see, e.g., Mishra, Ron, & Swaminathan, 2004; Pitt & Reinke, 1988) where one restricts the possible clusters to have good "verbal" descriptions in some language.

Clustering of many of the geometric classes investigated was shown to *require* information *in addition* to a presentation, d_0, d_1, d_2, \ldots , of the set of points to be clustered. For example, for clusters as convex hulls of finitely many points in a rational vector space, clustering can be done – but with the number of clusters as additional information. Let $\mathcal S$ consist of all polygons including their interiors – in the rational two-dimensional plane *without intersections and degenerated angles* (Attention was restricted to spaces of rationals since: 1. computer

reals are rationals, 2. this avoids the uncountability of the set of reals, and 3. this avoids dealing with uncomputable real points.) The class S can be clustered – but with the number of vertices of the polygons of the clusters involved as additional information.

Correspondingly, then, it was shown that the class S' containing S together with all such polygons but with one hole (the nondegenerate differences of two members in S) cannot be clustered with the number of vertices as additional information, yet S' can be clustered with area as additional information – and this even in higher dimensions and with any number of holes (Case et al., 2006).

It remains to be seen if some forms of geometrically constrained clustering can be usefully complementary to, say, conceptually/verbally constrained clustering.

Insights for Limitations of Science

We briefly treat below in some problems regarding parsimonious, refutable, and consistent hypotheses.

It is common wisdom in science that one should fit parsimonious explanations, hypotheses, or programs to data. In machine learning, this has been successfully applied, e.g., (Wallace, 2005; Wallace & Dowe, 1999).

Curiously, though, there are many results in inductive inference in which we see sometimes severe *degradations* of learning power caused by demanding *parsimonious* predictive programs (see, e.g., Freivalds (1975); Kinber (1977); Chen (1982); Case, Jain, and Sharma (1996); Ambainis, Case, Jain, and Suraj (2004)).

It is an interesting problem to resolve the seeming, likely not actual contradiction between the just prior two paragraphs.

Popper's Refutability (Popper, 1962) asserts that hypotheses in science should be subject to refutation. Besides the well-known difficulties of Duhem-Quine (Harding, 1976) of knowing which component hypothesis to throw out when a compound hypothesis badly fails to make correct predictions, inductive inference theorems have provided very different difficulties. Case and Smith (1983) outlines cases of usefully *inc*omplete (hence wrong) hypothesis that cannot be refuted, and Case and Suraj (2007) (see also Case, 2007) provides cases of inductively inferable higher order hypothesis not totally subject to refutation in cases where ordinary hypotheses subject to full refutation cannot be inductively inferred.

While Duhem–Quine may impact machine learning eventually, it remains to be seen about the inductive inference results of the just prior paragraph.

Requiring inductive inference procedures always to output an hypothesis in various senses *consistent* with (e.g., not ignoring) the data on which that hypothesis is based seems like mere common sense. However, from Bārzdiņš (1974a); Blum and Blum (1975); Wiehagen (1976), Case, Jain, Stephan, and Wiehagen (2004) we see that strict adherence to various consistency principles can severely attenuate the learning power of inductive inference machines. Furthermore, interestingly, *even when inductive inference is polytime constrained*, we see similar counterintuitive results to the effect that a kind of consistency can strictly attenuate learning power (Wiehagen & Zeugmann, 1994).

A machine learning analog might be Breiman's bagging (Breiman, 1996) and random forests (Breiman, 2001), where data is purposely ignored. However, in these cases, the purpose of ignoring data is to avoid overfitting to noise.

It remains to be seen, whether, in applied machine learning involving cases of practically noiseless data, one can also obtain some advantage in ignoring some consistency principles. Again the potential lesson from inductive inference is abstract and provides only a hint of something to work out in real machine learning problems.

Cross References

- ▶Behavioural Cloning
- **▶**Clustering
- ►Concept Drift
- ► Inductive Logic Programming
- ►Transfer Learning

Recommended Reading

Ambainis, A., Case, J., Jain, S., & Suraj, M. (2004). Parsimony hierarchies for inductive inference. *Journal of Symbolic Logic*, 69, 287–328.

Angluin, D., Gasarch, W., & Smith, C. (1989). Training sequences. Theoretical Computer Science, 66(3), 255–272.

Angluin, D. (1980). Finding patterns common to a set of strings. Journal of Computer and System Sciences, 21, 46-62.

Arikawa, S., Shinohara, T., & Yamamoto, A. (1992). Learning elementary formal systems. Theoretical Computer Science, 95, 97-113.

Bain, M., & Sammut, C. (1999). A framework for behavioural cloning. In K. Furakawa, S. Muggleton, & D. Michie (Eds.), *Machine intelligence 15*. Oxford: Oxford University Press.

- Baluja, S., & Pomerleau, D. (1995). Using the representation in a neural network's hidden layer for task specific focus of attention. Technical Report CMU-CS-95-143, School of Computer Science, CMU, May 1995. Appears in Proceedings of the 1995 IJCAI.
- Bartlett, P., Ben-David, S., & Kulkarni, S. (1996). Learning changing concepts by exploiting the structure of change. In *Proceedings* of the ninth annual conference on computational learning theory, Desenzano del Garda, Italy. New York: ACM Press.
- Bartlmae, K., Gutjahr, S., & Nakhaeizadeh, G. (1997). Incorporating prior knowledge about financial markets through neural multitask learning. In Proceedings of the fifth international conference on neural networks in the capital markets.
- Bārzdiņš, J. (1974a). Inductive inference of automata, functions and programs. In *Proceedings of the international congress of* mathematicians, Vancouver (pp. 771-776).
- Bārzdiņš, J. (1974b). Two theorems on the limiting synthesis of functions. In *Theory of algorithms and programs* (Vol. 210, pp. 82-88). Latvian State University, Riga.
- Blum, L., & Blum, M. (1975). Toward a mathematical theory of inductive inference. *Information and Control*, 28, 125–155.
- Blum, A., & Chalasani, P. (1992). Learning switching concepts. In Proceedings of the fifth annual conference on computational learning theory, Pittsburgh, Pennsylvania, (pp. 231–242). New York: ACM Press.
- Bratko, I., & Muggleton, S. (1995). Applications of inductive logic programming. *Communications of the ACM*, 38(11), 65–70.
- Bratko, I., Urbančič, T., & Sammut, C. (1998). Behavioural cloning of control skill. In R. S. Michalski, I. Bratko, & M. Kubat (Eds.), Machine learning and data mining: Methods and applications, (pp. 335–351). New York: Wiley.
- Brazma, A., Ukkonen, E., & Vilo, J. (1996). Discovering unbounded unions of regular pattern languages from positive examples. In Proceedings of the seventh international symposium on algorithms and computation (ISAAC'96), Lecture notes in computer science, (Vol. 1178, pp. 95-104), Berlin: Springer-Verlag.
- Breiman, L. (1996). Bagging predictors. Machine Learning, 24(2),
- Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32. Caruana, R. (1993). Multitask connectionist learning. In *Proceedings of the 1993 connectionist models summer school* (pp. 372–379). NJ: Lawrence Erlbaum.
- Caruana, R. (1996). Algorithms and applications for multitask learning. In *Proceedings 13th international conference on machine learning* (pp. 87-95). San Francisco, CA: Morgan Kanfmann
- Case, J. (1994). Infinitary self-reference in learning theory. *Journal of Experimental and Theoretical Artificial Intelligence*, 6, 3–16.
- Case, J. (1999). The power of vacillation in language learning. SIAM Journal on Computing, 28(6), 1941–1969.
- Case, J. (2007). Directions for computability theory beyond pure mathematical. In D. Gabbay, S. Goncharov, & M. Zakharyaschev (Eds.), Mathematical problems from applied logic II. New logics for the XXIst century, International Mathematical Series, (Vol. 5). New York: Springer.
- Case, J., & Lynes, C. (1982). Machine inductive inference and language identification. In M. Nielsen & E. Schmidt, (Eds.), Proceedings of the 9th International Colloquium on Automata, Languages and Programming, Lecture notes in computer science, (Vol. 140, pp. 107–115). Berlin: Springer-Verlag.
- Case, J., & Smith, C. (1983). Comparison of identification criteria for machine inductive inference. *Theoretical Computer Science*, 25, 193–220.

- Case, J., & Suraj, M. (2010). Weakened refutability for machine learning of higher order definitions, 2010. (Working paper for eventual journal submission).
- Case, J., Jain, S., Kaufmann, S., Sharma, A., & Stephan, F. (2001).Predictive learning models for concept drift (Special Issue for ALT'98). Theoretical Computer Science, 268, 323-349.
- Case, J., Jain, S., Lange, S., & Zeugmann, T. (1999). Incremental concept learning for bounded data mining. *Information and Computation*, 152, 74-110.
- Case, J., Jain, S., Montagna, F., Simi, G., & Sorbi, A. (2005). On learning to coordinate: Random bits help, insightful normal forms, and competency isomorphisms (Special issue for selected learning theory papers from COLT'03, FOCS'03, and STOC'03). *Journal of Computer and System Sciences*, 71(3), 308–332.
- Case, J., Jain, S., Martin, E., Sharma, A., & Stephan, F. (2006). Identifying clusters from positive data. SIAM Journal on Computing, 36(1), 28-55.
- Case, J., Jain, S., Ott, M., Sharma, A., & Stephan, F. (2000). Robust learning aided by context (Special Issue for COLT'98). *Journal* of Computer and System Sciences, 60, 234–257.
- Case, J., Jain, S., & Sharma, A. (1996). Machine induction without revolutionary changes in hypothesis size. *Information and Computation*, 128, 73–86.
- Case, J., Jain, S., Stephan, F., & Wiehagen, R. (2004). Robust learning rich and poor. *Journal of Computer and System Sciences*, 69(2), 123–165.
- Case, J., Ott, M., Sharma, A., & Stephan, F. (2002). Learning to win process-control games watching gamemasters. *Information and Computation*, 174(1), 1-19.
- Cenzer, D., & Remmel, J. (1992). Recursively presented games and strategies. *Mathematical Social Sciences*, 24, 117–139.
- Chen, K. (1982). Tradeoffs in the inductive inference of nearly minimal size programs. *Information and Control*, 52, 68-86.
- de Garis, H. (1990a). Genetic programming: Building nanobrains with genetically programmed neural network modules. In *IJCNN: International Joint Conference on Neural Networks*, (Vol. 3, pp. 511–516). Piscataway, NJ: IEEE Service Center.
- deGaris, H. (1990b). Genetic programming: Modular neural evolution for Darwin machines. In M. Caudill (Ed.), IJCNN-90-WASH DC; International joint conference on neural networks (Vol. 1, pp. 194–197). Hillsdale, NJ: Lawrence Erlbaum Associates.
- de Garis, H. (1991). Genetic programming: Building artificial nervous systems with genetically programmed neural network modules. In B. Soušek, & The IRIS group (Eds.), Neural and intelligenct systems integeration: Fifth and sixth generation integerated reasoning information systems (Chap. 8, pp. 207–234). New York: Wiley.
- Devaney, M., & Ram, A. (1994). Dynamically adjusting concepts to accommodate changing contexts. In M. Kubat, G. Widmer (Eds.), Proceedings of the ICML-96 Pre-conference workshop on learning in context-sensitive domains, Bari, Italy (Journal submission).
- Dietterich, T., Hild, H., & Bakiri, G. (1995). A comparison of ID3 and backpropogation for English text-tospeech mapping. *Machine Learning*, 18(1), 51–80.
- Fahlman, S. (1991). The recurrent cascade-correlation architecture. In R. Lippmann, J. Moody, and D. Touretzky (Eds.), Advances in neural information processing systems (Vol. 3, pp. 190–196). San Mateo, CA: Morgan Kaufmann Publishers.
- Freivalds, R. (1975). Minimal Gödel numbers and their identification in the limit. In *Lecture notes in computer science* (Vol. 32, pp. 219–225). Berlin: Springer-Verlag.

- Freund, Y., & Mansour, Y. (1997). Learning under persistent drift. In S. Ben-David, (Ed.), Proceedings of the third European conference on computational learning theory (EuroCOLT'97), Lecture notes in artificial intelligence, (Vol. 1208, pp. 94–108). Berlin: Springer-Verlag.
- Fulk, M. (1990). Robust separations in inductive inference. In Proceedings of the 31st annual symposium on foundations of computer science (pp. 405–410). St. Louis, Missouri. Washington, DC: IEEE Computer Society.
- Harding, S. (Ed.). (1976). Can theories be refuted? Essays on the Duhem-Quine thesis. Dordrecht: Kluwer Academic Publishers.
- Helmbold, D., & Long, P. (1994). Tracking drifting concepts by minimizing disagreements. Machine Learning, 14, 27–46.
- Hildebrand, F. (1956). Introduction to numerical analysis. New York: McGraw-Hill.
- Jain, S. (1999). Robust behaviorally correct learning. Information and Computation, 153(2), 238-248.
- Jain, S., & Sharma, A. (1997). Elementary formal systems, intrinsic complexity, and procrastination. *Information and Computation*, 132, 65-84.
- Jain, S., & Sharma, A. (2002). Mind change complexity of learning logic programs. Theoretical Computer Science, 284(1), 143-160.
- Jain, S., Osherson, D., Royer, J., & Sharma, A. (1999). Systems that learn: An introduction to learning theory (2nd ed.). Cambridge, MA: MIT Press.
- Jain, S., Smith, C., & Wiehagen, R. (2001). Robust learning is rich. Journal of Computer and System Sciences, 62(1), 178-212.
- Kilpeläinen, P., Mannila, H., & Ukkonen, E. (1995). MDL learning of unions of simple pattern languages from positive examples. In P. Vitányi (Ed.), Computational learning theory, second European conference, EuroCOLT'95, Lecture notes in artificial intelligence, (Vol. 904, pp. 252–260). Berlin: Springer-Verlag.
- Kinber, E. (1977). On a theory of inductive inference. In *Lecture notes* in computer science (Vol. 56, pp. 435-440). Berlin: Springer-Verlag.
- Kinber, E., Smith, C., Velauthapillai, M., & Wiehagen, R. (1995). On learning multiple concepts in parallel. *Journal of Computer and System Sciences*, 50, 41-52.
- Krishna Rao, M. (1996). A class of prolog programs inferable from positive data. In A. Arikawa & A. Sharma (Eds.), Seventh international conference on algorithmic learning theory (ALT' 96), Lecture notes in artificial intelligence (Vol. 1160, pp. 272-284). Berlin: Springer-Verlag.
- Krishna Rao, M. (2000). Some classes of prolog programs inferable from positive data (Special Issue for ALT'96). Theoretical Computer Science A, 241, 211-234.
- Krishna Rao, M. (2004). Inductive inference of term rewriting systems from positive data. In S. Ben-David, J. Case, & A. Maruoka (Eds.), Algorithmic learning theory: Fifteenth international conference (ALT' 2004), Lecture notes in artificial intelligence (Vol. 3244, pp. 69–82). Berlin: Springer-Verlag.
- Krishna Rao, M. (2005). A class of prolog programs with nonlinear outputs inferable from positive data. In S. Jain, H. U. Simon, & E. Tomita (Eds.), Algorithmic learning theory: Sixteenth international conference (ALT' 2005), Lecture notes in artificial intelligence, (Vol. 3734, pp. 312–326). Berlin: Springer-Verlag.
- Krishna Rao, M., & Sattar, A. (1998). Learning from entailment of logic programs with local variables. In M. Richter, C. Smith, R. Wiehagen, & T. Zeugmann (Eds.), Ninth international conference on algorithmic learning theory (ALT' 98), Lecture notes in

- artificial intelligence (Vol. 1501, pp. 143–157). Berlin: Springer-Verlag.
- Kubat, M. (1992). A machine learning based approach to load balancing in computer networks. Cybernetics and Systems, 23, 389-400.
- Kummer, M., & Ott, M. (1996). Learning branches and learning to win closed recursive games. In *Proceedings of the ninth annual conference on computational learning theory*, Desenzano del Garda, Italy. New York: ACM Press.
- Lange, S., & Wiehagen, R. (1991). Polynomial time inference of arbitrary pattern languages. New Generation Computing, 8, 361–370.
- Lavrač, N., & Džeroski, S. (1994). *Inductive logic programming: Techniques and applications*. New York: Ellis Horwood.
- Maler, O., Pnueli, A., & Sifakis, J. (1995). On the synthesis of discrete controllers for timed systems. In *Proceedings of the annual sym*posium on the theoretical aspects of computer science, LNCS (Vol. 900, pp. 229–242). Berlin: Springer-Verlag.
- Matwin, S., & Kubat, M. (1996). The role of context in concept learning. In M. Kubat & G. Widmer (Eds.), *Proceedings of the ICML-96 pre-conference workshop on learning in context-sensitive domains*, Bari, Italy, (pp. 1-5).
- Maye, A., Hsieh, C., Sugihara, G., & Brembs, B. (2007). Order in spontaneous behavior. *PLoS One*, May, 2007. See: http://brembs.net/spontaneous/
- Mishra, N., Ron, D., & Swaminathan, R. (2004). A new conceptual clustering framework. *Machine Learning*, 56(1-3), 115-151.
- Mitchell, T. (1997). Machine learning. New York: McGraw Hill.
- Mitchell, T., Caruana, R., Freitag, D., McDermott, J., & Zabowski, D. (1994). Experience with a learning, personal assistant. Communications of the ACM, 37, 80–91.
- Montagna, F., & Osherson, D. (1999). Learning to coordinate: A recursion theoretic perspective. *Synthese*, 118, 363–382.
- Muggleton, S., & De Raedt, L. (1994). Inductive logic programming: Theory and methods. *Journal of Logic Programming*, 19/20, 669–679
- Odifreddi, P. (1999). Classical recursion theory (Vol. II). Amsterdam:
- Osherson, D., Stob, M., & Weinstein, S. (1986). Systems that learn: An introduction to learning theory for cognitive and computer scientists. Cambridge, MA: MIT Press.
- Ott, M., & Stephan, F. (2002). Avoiding coding tricks by hyperrobust learning. *Theoretical Computer Science*, 284(1), 161–180.
- Pitt, L., & Reinke, R. (1988). Criteria for polynomial-time (conceptual) clustering. *Machine Learning*, 2, 371–396.
- Popper, K. (1992). Conjectures and refutations: The growth of scientific knowledge. New York: Basic Books.
- Pratt, L., Mostow, J., & Kamm, C. (1991). Direct transfer of learned information among neural networks. In *Proceedings of the 9th national conference on artificial intelligence (AAAI-91)*, Anaheim, California. Menlo Park, CA: AAAI press.
- Rogers, H. (1987). Theory of recursive functions and effective computability. New York: McGraw Hill (Reprinted, MIT Press, 1987).
- Salomaa, A. (1994a). Patterns (The formal language theory column). EATCS Bulletin, 54, 46-62.
- Salomaa, A. (1994b). Return to patterns (The formal language theory column). *EATCS Bulletin*, 55, 144–157.
- Sejnowski, T., & Rosenberg, C. (1986). NETtalk: A parallel network that learns to read aloud. *Technical Report JHU-EECS-86-01*, Johns Hopkins University.

219

- Shimozono, S., Shinohara, A., Shinohara, T., Miyano, S., Kuhara, S., & Arikawa, S. (1994). Knowledge acquisition from amino acid sequences by machine learning system BONSAI. *Transactions* of *Information Processing Society of Japan*, 35, 2009–2018.
- Shinohara, T. (1983). Inferring unions of two pattern languages. Bulletin of Informatics and Cybernetics, 20, 83–88.
- Shinohara, T., & Arikawa, A. (1995). Pattern inference. In K. P. Jantke & S. Lange (Eds.), Algorithmic learning for knowledge-based systems, Lecture notes in artificial intelligence (Vol. 961, pp. 259-291). Berlin: Springer-Verlag.
- Smullyan, R. (1961). Theory of formal systems. In Annals of Mathematics Studies (Vol. 47). Princeton, NJ: Princeton University Press.
- Šuc, D. (2003). Machine reconstruction of human control strategies. Frontiers in artificial intelligence and applications (Vol. 99). Amsterdam: IOS Press.
- Thomas, W. (1995). On the synthesis of strategies in infinite games. In *Proceedings of the annual symposium on the theoretical aspects of computer science*, *LNCS* (Vol. 900, pp. 1–13). Berlin: Springer-Verlag.
- Thrun, S. (1996). Is learning the n-th thing any easier than learning the first? In *Advances in neural information processing systems*, 8. San Mateo, CA: Morgan Kaufmann.
- Thrun, S., & Sullivan, J. (1996). Discovering structure in multiple learning tasks: The TC algorithm. In *Proceedings of the thirteenth international conference on machine learning (ICML-96)* (pp. 489–497). San Francisco, CA: Morgan Kaufmann.
- Tsung, F., & Cottrell, G. (1989). A sequential adder using recurrent networks. In *IJCNN-89-WASHINGTON DC: International joint conference on neural networks* June 18–22 (Vol. 2, pp. 133–139). Piscataway, NJ: IEEE Service Center.
- Waibel, A. (1989a). Connectionist glue: Modular design of neural speech systems. In D. Touretzky, G. Hinton, & T. Sejnowski (Eds.), Proceedings of the 1988 connectionist models summer school (pp. 417–425). San Mateo, CA: Morgan Kaufmann.
- Waibel, A. (1989b). Consonant recognition by modular construction of large phonemic time-delay neural networks. In D. S. Touretzky (Ed.), Advances in neural information processing systems I (pp. 215–223). San Mateo, CA: Morgan Kaufmann.
- Wallace, C. (2005). Statistical and inductive inference by minimum message length. (Information Science and Statistics). New York: Springer (Posthumously published).
- Wallace, C., & Dowe, D. (1999). Minimum message length and kolmogorov complexity (Special Issue on Kolmogorov Complexity). Computer Journal, 42(4), 123–155. http://comjnl. oxfordjournals.org/cgi/reprint/42/4/270.
- Widmer, G., & Kubat, M. (1996). Learning in the presence of concept drift and hidden contexts. *Machine Learning*, 23, 69–101.
- Wiehagen, R. (1976). Limes-Erkennung rekursiver Funktionen durch spezielle Strategien. Electronische Informationverarbeitung und Kybernetik, 12, 93-99.
- Wiehagen, R., & Zeugmann, T. (1994). Ignoring data may be the only way to learn efficiently. Journal of Experimental and Theoretical Artificial Intelligence, 6, 131–144.
- Wright, K. (1989). Identification of unions of languages drawn from an identifiable class. In R. Rivest, D. Haussler, & M. Warmuth (Eds.), Proceedings of the second annual workshop on computational learning theory, Santa Cruz, California, (pp. 328–333). San Mateo, CA: Morgan Kaufmann Publishers.

- Wrobel, S. (1994). Concept formation and knowledge revision. Dordrecht: Kluwer Academic Publishers.
- Zeugmann, T. (1986). On Bārzdiņš' conjecture. In K. P. Jantke (Ed.), Analogical and inductive inference, Proceedings of the international workshop, Lecture notes in computer science, (Vol. 265, pp. 220–227). Berlin: Springer-Verlag.
- Zeugmann, T. (1998). Lange and Wiehagen's pattern language learning algorithm: An average case analysis with respect to its total learning time. Annals of Mathematics and Artificial Intelligence, 23, 117–145.

Connectivity

▶Topology of a Neural Network

Consensus Clustering

Synonyms

Clustering aggregation; Clustering ensembles

Definition

In Consensus Clustering we are given a set of n objects V, and a set of m clusterings $\{C_1, C_2, \ldots, C_m\}$ of the objects in V. The aim is to find a single clustering C that *disagrees* least with the input clusterings, that is, C minimizes

$$D(C) = \sum_{C_i} d(C, C_i),$$

for some metric d on clusterings of V. Meilă (2003) proposed the principled *variation of information* metric on clusterings, but it has been difficult to analyze theoretically. The Mirkin metric is the most widely used, in which d(C, C') is the number of pairs of objects (u, v) that are clustered together in C and apart in C', or vice versa; it can be calculated in time O(mn).

We can interpret each of the clusterings C_i in Consensus Clustering as evidence that pairs ought be put together or separated. That is, w_{uv}^+ is the number of C_i in which $C_i[u] = C_i[v]$ and w_{uv}^- is the number of C_i in which $C_i[u] \neq C_i[v]$. It is clear that $w_{uv}^+ + w_{uv}^- = m$ and

220 Constrained Clustering

that Consensus clustering is an instance of Corre-LATION CLUSTERING in which the w_{uv}^- weights obey the triangle inequality.

Constrained Clustering

Kiri L. Wagstaff Pasadena, CA, USA

Definition

Constrained clustering is a semisupervised approach to clustering data while incorporating domain knowledge in the form of constraints. The constraints are usually expressed as pairwise statements indicating that two items must, or cannot, be placed into the same cluster. Constrained clustering algorithms may enforce every constraint in the solution, or they may use the constraints as guidance rather than hard requirements.

Motivation and Background

▶Unsupervised learning operates without any domainspecific guidance or preexisting knowledge. Supervised learning requires that all training examples be associated with labels. Yet it is often the case that existing knowledge for a problem domain fits neither of these extremes. Semisupervised learning methods fill this gap by making use of both labeled and unlabeled data. Constrained clustering, a form of semisupervised learning, was developed to extend clustering algorithms to incorporate existing domain knowledge, when available. This knowledge may arise from labeled data or from more general rules about the concept to be learned.

One of the original motivating applications was noun phrase coreference resolution, in which noun phrases in a text must be clustered together to represent distinct entities (e.g., "Mr. Obama" and "the President" and "he", separate from "Sarah Palin" and "she" and "the Alaska governor"). This problem domain contains several natural rules for when noun phrases should (such as appositive phrases) or should not (such as a mismatch on gender) be clustered together. These rules can be translated into a collection of pairwise constraints on the data to be clustered.

Constrained clustering algorithms have now been applied to a rich variety of domain areas, including hyperspectral image analysis, road lane divisions from

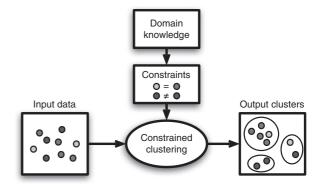
GPS data, gene expression microarray analysis, video object identification, document clustering, and web search result grouping.

Structure of the Learning System

Constrained clustering arises out of existing work with unsupervised clustering algorithms. In this description, we focus on clustering algorithms that seek a partition of the data into disjoint clusters, using a distance or similarity measure to place similar items into the same cluster. Usually, the desired number of clusters, k, is specified as an input to the algorithm. The most common clustering algorithms are k-means (MacQueen, 1967) and expectation maximization or EM (Dempster, Laird, & Rubin, 1977) (Fig. 1).

A constrained clustering algorithm takes the same inputs as a regular (unsupervised) clustering algorithm and also accepts a set of pairwise constraints. Each constraint is a must-link or cannot-link constraint. The must-link constraints form an equivalence relation, which permits the inference of additional transitively implied must-links as well as additional entailed cannot-link constraints between items from distinct must-link cliques. Specifying a significant number of pairwise constraints might be tedious for large data sets, so often they may be generated from a manually labeled subset of the data or from domain-specific rules.

The algorithm may interpret the constraints as hard constraints that must be satisfied in the output or as soft preferences that can be violated, if necessary. The former approach was used in the first constrained clustering algorithms, COP-COBWEB (Wagstaff & Cardie,



Constrained Clustering. Figure 1. The constrained clustering algorithm takes in nine items and two pairwise constraints (one must-link and one cannot-link). The output clusters respect the specified constraints

Basu, S., Davidson, I., & Wagstaff, K. (Eds.). (2008). Constrained Clustering: Advances in Algorithms, Theory, and Applications.

Boca Raton, FL: CRC Press.

Basu, S., Bilenko, M., & Mooney, R. J. (2004). A probabilistic framework for semi-supervised clustering. In *Proceedings of the Tenth*

ery and Data Mining (pp. 59-68). Seattle, WA.

ACM SIGKDD International Conference on Knowledge Discov-

Bilenko, M., Basu, S., & Mooney, R. J. (2004). Integrating constraints and metric learning in semi-supervised clustering. In *Proceedings of the Twenty-first International Conference on Machine Learning* (pp. 11–18). Banff, AB, Canada.

Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society*, 39(1), 1–38.

Kamvar, S., Klein, D., & Manning, C. D. (2003). Spectral learning. In Proceedings of the International Joint Conference on Artificial Intelligence (pp. 561-566). Acapulco, Mexico.

Klein, D., Kamvar, S. D., & Manning, C. D. (2002). From instance-level constraints to space-level constraints: Making the most of prior knowledge in data clustering. In *Proceedings of the Nineteenth International Conference on Machine Learning* (pp. 307–313). Sydney, Australia.

Lu, Z. & Leen, T. (2005). Semi-supervised learning with penalized probabilistic clustering. In Advances in Neural Information Processing Systems (Vol. 17, pp. 849–856). Cambridge, MA: MIT Press.

MacQueen, J. B. (1967). Some methods for classification and analysis of multivariate observations. In *Proceedings of the* Fifth Symposium on Math, Statistics, and Probability (Vol. 1, pp. 281–297). California: University of California Press.

Shental, N., Bar-Hillel, A., Hertz, T., & Weinshall, D. (2004). Computing Gaussian mixture models with EM using equivalence constraints. In Advances in Neural Information Processing Systems (Vol. 16, pp. 465-472). Cambridge, MA: MIT Press.

Wagstaff, K. & Cardie, C. (2000). Clustering with instance-level constraints. In Proceedings of the Seventeenth International Conference on Machine Learning (pp. 1103-1110). San Francisco: Morgan Kaufmann.

Wagstaff, K., Cardie, C., Rogers, S., & Schroedl, S. (2001). Constrained k-means clustering with background knowledge. In Proceedings of the Eighteenth International Conference on Machine Learning (pp. 577-584). San Francisco: Morgan Kaufmann.

Xing, E. P., Ng, A. Y., Jordan, M. I., & Russell, S. (2003). Distance metric learning, with application to clustering with side-information. In Advances in Neural Information Processing Systems (Vol. 15, pp. 505–512). Cambridge, MA: MIT Press.

2000) and COP-kmeans (Wagstaff, Cardie, Rogers, & Schroedl, 2001). COP-kmeans accommodates the constraints by restricting item assignments to exclude any constraint violations. If a solution that satisfies the constraints is not found, COP-kmeans terminates without a solution. Later, algorithms such as PCK-means and MPCK-means (Bilenko, Basu, & Mooney, 2004) permitted the violation of constraints when necessary by introducing a violation penalty. This is useful when the constraints may contain noise or internal inconsistencies, which are especially relevant in real-world domains. Constrained versions of other clustering algorithms such as EM (Shental, Bar-Hillel, Hertz, & Weinshall, 2004) and spectral clustering (Kamvar, Klein, & Manning, 2003) also exist. Penalized probabilistic clustering (PPC) is a modified version of EM that interprets the constraints as (soft) probabilistic priors on the relationships between items (Lu & Leen, 2005).

In addition to constraining the assignment of individual items, constraints can be used to learn a better distance metric for the problem at hand (Bar-Hillel, Hertz, Shental, & Weinshall, 2005; Klein, Kamvar, & Manning, 2002; Xing, Ng, Jordan, & Russell, 2003). Must-link constraints hint that the effective distance between those items should be low, while cannot-link constraints suggest that their pairwise distance should be high. Modifying the metric accordingly permits the subsequent application of a regular clustering algorithm, which need not explicitly work with the constraints at all. The MPCK-means algorithm fuses these approaches together, providing both constraint satisfaction and metric learning simultaneously (Basu, Bilenko, & Mooney, 2004; Bilenko et al., 2004).

More information about subsequent advances in constrained clustering algorithms, theory, and novel applications can be found in a compilation edited by Basu, Davidson, and Wagstaff (2008).

Programs and Data

The MPCK-means algorithm is available in a modified version of the Weka machine learning toolkit (Java) at http://www.cs.utexas.edu/users/ml/risc/code/.

Recommended Reading

Bar-Hillel, A., Hertz, T., Shental, N., & Weinshall, D. (2005). Learning a Mahalanobis metric from equivalence constraints. *Journal of Machine Learning Research*, 6, 937–965.

Constraint-Based Mining

SIEGFRIED NIJSSEN Katholieke Universiteit Leuven, Leuven, Belgium

Definition

Constraint-based mining is the research area studying the development of data mining algorithms that search C

through a pattern or model space restricted by constraints. The term is usually used to refer to algorithms that search for patterns only. The most well-known instance of constraint-based mining is the mining of patterns. Constraints are needed in pattern mining algorithms to increase the efficiency of the search and to reduce the number of patterns that are presented to the user, thus making knowledge discovery more effective and useful.

Motivation and Background

Constraint-based pattern mining is a generalization of frequent itemset mining. For an introduction to frequent itemset mining, see Frequent Patterns. A constraint-based mining problem is specified by providing the following elements:

- A database D, usually consisting of independent transactions (or instances)
- A hypothesis space \mathcal{L} of patterns
- A constraint q(θ, D) expressing criteria that a pattern θ in the hypothesis space should fulfill on the database

The general constraint-based mining problem is to find the set

Th(
$$\mathcal{D}$$
, \mathcal{L} , q) = { $\theta \in \mathcal{L} | q(\theta, \mathcal{D})$ = true}.

Alternative problem settings are obtained by making different choices for \mathcal{D} , \mathcal{L} and q. For instance,

- If the database and the hypothesis space consist of graphs or trees instead of itemsets, a graph mining or a tree mining problem is obtained. For more information about these topics, see ►Graph Mining and
 Tree Mining
- Additional syntactic constraints can be imposed

An overview of important types of constraints is given below.

One can generalize the constraint-based mining problem beyond pattern mining. Also models, such as

Decision Trees, could be seen as languages of interest. In the broadest sense, topics such as ▶Constrained Clustering, ▶Cost-Sensitive Learning, and even learning ▶Support Vector Machines (SVMs) may be seen as constraint-based mining problems. However, it is currently not common to categorize these topics as constraint-based mining; in practice, the term refers to constraint-based pattern mining.

From the perspective of constraint-based mining, the knowledge discovery process can be seen as a process in which a user repeatedly specifies constraints for data mining algorithms; the data mining system is a solver that finds patterns or models that satisfy the constraints.

This approach to data mining is very similar to querying relational databases. Whereas relational databases are usually queried using operations such as projections, selections, and joins, in the constraint-based mining framework data is queried to find patterns or models that satisfy constraints that cannot be expressed in these primitives. A database which supports constraint-based mining queries, stores patterns and models, and allows later reuse of patterns and models, is sometimes also called an *inductive database* (Imielinski & Mannila, 1996).

Structure of the Learning System

Constraints

Frequent pattern mining algorithms can be generalized along several dimensions.

One way to generalize pattern mining algorithms is to allow them to deal with arbitrary **coverage relations**, which determine when a pattern matches a transaction in the data. In the example of mining itemsets, the subset relation determines the coverage relation. The coverage relation is at the basis of constraints such as minimum support; an alternative coverage relation would be the superset relation.

From the coverage relation follows a generality relationship. A pattern θ_1 is defined to be more specific than a pattern θ_2 (denoted by $\theta_1 > \theta_2$) if any transaction that is covered by θ_1 is also covered by θ_2 (see \triangleright Generalization). In frequent itemset mining, itemset I_1 is more general than itemset I_2 if and only $I_1 \subseteq I_2$.

Generalization and coverage relationships can be used to identify the following types of constraints.

Monotonic and Anti-Monotonic Constraints An essential property which is exploited in ▶frequent pattern mining, is that all subsets of a frequent pattern are also frequent. This is a property that can be generalized:

- A constraint is called *monotonic* if any generalization of a pattern that satisfies the constraint, also satisfies the constraint
- A constraint is called *anti-monotonic* if any specialization of a pattern that satisfies the constraint, also satisfies the constraint

In some publications, the definitions of monotonic and anti-monotonic are used reversely.

The following are examples of monotonic constraints:

- Minimum support
- Syntactic constraints, for instance: a constraint that requires that patterns specializing a given pattern *x* are excluded a constraint requiring patterns to be small given a definition of pattern size
- Disjunctions or conjunctions of monotonic constraints
- Negations of anti-monotonic constraints

The following are examples of anti-monotonic constraints:

- Maximum support
- Syntactic constraints, for instance, a constraint that requires that patterns generalizing a given pattern x are excluded
- Disjunctions or conjunctions of anti-monotonic constraints
- Negations of monotonic constraints

Succinct Constraints Constraints that can be pushed in the mining process by adapting the pattern space or data, are called succinct constraints. An example of a succinct constraint is the monotonic constraint that an itemset should contain the item *A*. This constraint could be dealt with by deleting all transactions that do not contain *A*. For any frequent itemset found in the new dataset, it is now known that the item *A* can be added to it.

Convertible Constraints Some constraints that are not monotonic, can still be *convertible* monotonic (Pei &

Han, 2002). A constraint is convertible monotonic if for every pattern θ one least general generalization θ' can be identified such that if θ satisfies the constraint, then θ' also satisfies the constraint. An example of a convertible constraint is a maximum average cost constraint. Assume that every item in an itemset has a cost as defined by a function c(i). The constraint $c(I) = \sum_{i \in I} c(i)/|I| \le maxcost$ is not monotonic. However, for every itemset I with $c(I) \le maxcost$, if an item i is removed with $c(i) = \max_{i \in I} c(i)$, an itemset with $c(I - \{i\}) \le c(I) \le maxcost$ is obtained.

Maximum average cost has the desirable property that no access to the data is needed to identify the generalization that should satisfy the constraints. If it is not possible to identify the necessary least general generalization before accessing the data, the convertible constraint is also sometimes called weak (anti-)monotone (Zhu, Yan, Han, & Yu, 2007).

Boundable Constraints Constraints on non-monotonic measures for which a monotonic bound exist, are called boundable. An example of such a constraint is a minimum accuracy constraint in a database with binary class labels. Assume that every itemset is interpreted as a rule if I then 1 else 2 (thus, class label 1 is predicted if a transaction contains itemset I, or class label 2 otherwise; see ightharpoonup Supervised Descriptive Rule Discovery). A minimum accuracy constraint can be formalized by the formula $(\operatorname{fr}(I, D_1) + |D_2| - \operatorname{fr}(I, D_2))/|D| \ge \minacc$, where D_k is the database containing only the examples labeled with class label k. It can be derived from this that

$$\operatorname{fr}(I, D_1) \ge |D| \min \operatorname{acc} - |D_2| + \operatorname{fr}(I, D_2) \ge |D| \min \operatorname{acc} - |D_2|.$$

In other words, if a high accuracy is desirable, a minimum number of examples of class 1 is required to be covered, and a minimum frequency constraint can thus be derived. Therefore, minimum support can be used as a bound for minimum accuracy.

The principle of deriving bounds for non-monotonic measures can be applied widely (Bayardo, Agrawal, & Gunopulos, 1999; Morishita & Sese, 2000).

Borders If constraints are not restrictive enough, the number of patterns can be huge. Ignoring statistics about patterns such as their exact frequency, the set of patterns can be represented more compactly only by

listing the patterns in the *border*(*s*) (Mannila & Toivonen, 1997), similar to the idea of version spaces. An example of a border is the set of maximal frequent itemsets (see Frequent Patterns). Borders can be computed for other types of both monotonic and antimonotonic constraints as well. There are several complications compared to the simple frequent pattern mining setting:

- If there is an anti-monotonic constraint, such as maximum support, not only is it needed to compute a border for the most specific elements in the set (S-Set), but also a border for the least general elements in the set (G-Set)
- If the formula is a disjunction of conjunctions, the result of a query becomes a union of version spaces, which is called a multi-dimensional version space (see Fig. 1) (De Raedt, Jaeger, Lee, & Mannila, 2002); the G-Set of one version space may be more general than the G-Set of another version space

Both the S-Set and the G-Set can be represented by listing elements just within the version space (the positive border), or elements just outside the version space (the negative border). For instance, the positive border of the G-Set consists of those patterns which are part of the version space, and for which no generalizations exist which are part of the version space.

Similarly, there may exist several representations of multi-dimensional version spaces; optimizing the representation of multi-dimensional version spaces is analogous to optimizing queries in relational databases (De Raedt et al., 2002).

Borders form a *condensed representations*, that is, they compactly represent the solution space; see Frequent Patterns.

Algorithms For many of the constraints specified in the previous section specialized algorithms have been developed in combination with specific hypothesis spaces. It is beyond the scope of this chapter to discuss all these algorithms; only the most common ideas are provided here.

The main idea is that Apriori can easily be updated to deal with general monotonic constraints in arbitrary hypothesis spaces. The concept of a specialization refinement operator is essential to operate on

other hypothesis spaces than itemsets. A specialization operator $\rho(\theta)$ computes a set of specializations in the hypothesis space for a given input pattern. In pattern mining, this operator should have the following properties:

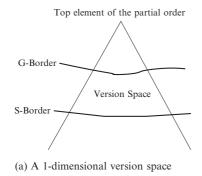
- Completeness: every pattern in the hypothesis space should be reachable by repeated application of the refinement operator starting from the most general pattern in the hypothesis space
- Nonredundancy: every pattern in the hypothesis space should be reachable in only one way starting from the most general pattern in the hypothesis space

In itemset mining, optimal refinement is usually obtained by first ordering the items (for instance, alphabetically, or by frequency), and then adding items that are higher in the chosen order to a set than the items already in the set. For instance, for the itemset $\{A, C\}$, the specialization operator returns $\rho(\{A, C\}) = \{\{A, C, D\}, \{A, C, E\}\}$, assuming that the domain of items $\{A, B, C, D, E\}$ is considered. Other refinement operators are needed while dealing with other hypothesis spaces, such as in \triangleright graph mining.

The search in Apriori proceeds breadth-first. Each level, the specialization operator is applied on patterns satisfying the monotonic constraints to generate candidates for the next level. For every new candidate it is checked whether its generalizations satisfy the monotonic constraints. To create a set of generalizations, a generalization refinement operator can be used. In frequent itemset mining, usually single items are removed from the itemset to generate generalizations.

More changes are required to deal with *anti-monotonic* constraints. A simple way of dealing with both monotonic and anti-monotonic constraints is to first compute all patterns that satisfy the monotonic constraints, and then to prune the patterns that fail to satisfy the anti-monotonic constraints. More challenging is to "push" anti-monotonic constraints in the mining process. An observation which is often exploited is that generalizations of patterns that do not satisfy the anti-monotonic constraint need not be considered. Well-known strategies are:

Constructive Induction 225



Top element of the partial order

G-Border (1)

Version
Space (1)

G-Border (2)

Version Space (2)

Version Space (2)

S-Border (2)

Version Space (2)

(b) A 2-dimensional version space

Constraint-Based Mining. Figure 1. Version spaces

- In a breadth-first setting: traverse the lattice in reverse order for monotonic constraints, after the patterns have been determined satisfying the antimonotonic constraints (De Raedt et al., 2002)
- In a depth-first setting: during the search for patterns, try to guess the largest pattern that can still be reached, and prune a branch in the search if the pattern does not satisfy the monotonic constraint on this pattern (Bucila, Gehrke, Kifer, & White, 2003; Kifer, Gehrke, Bucila, & White, 2003)

It is beyond the scope of this chapter to discuss how to deal with other types of constraints; however, it should be pointed out that not all combinations of constraints and hypothesis spaces have been studied; it is not obvious whether all constraints can be pushed usefully in a pattern search for any hypothesis space, for instance, when boundable constraints in more complex hypothesis spaces (such as graphs) are involved. Research in this area is ongoing.

Cross References

- ► Constrained Clustering
- ▶Frequent Pattern Mining
- ► Graph Mining
- ►Tree Mining

Recommended Reading

Bayardo, R. J., Jr., Agrawal, R., & Gunopulos, D. (1999). Constraint-based rule mining in large, dense databases. In *Proceedings of the 15th international conference on data engineering (ICDE)* (pp. 188–197). Sydney, Australia.

Bucila, C., Gehrke, J., Kifer, D., & White, W. M. (2003). DualMiner: A dual-pruning algorithm for itemsets with constraints. *Data Mining and Knowledge Discovery*, 7(3), 241–272. De Raedt, L., Jaeger, M., Lee, S. D., & Mannila, H. (2002). A theory of inductive query answering (extended abstract). In *Proceedings of the second IEEE international conference on data mining (ICDM)* (pp. 123–130). Los Alamitos, CA: IEEE Press.

Imielinski, T., & Mannila, H. (1996). A database perspective on knowledge discovery. Communications of the ACM, 39, 58-64.

Kifer, D., Gehrke, J., Bucila, C., & White, W. M. (2003). How to quickly find a witness. In *Proceedings of the twenty-second ACM SIGACT-SIGMOD-SIGART symposium on principles of database systems* (pp. 272–283). San Diego, CA: ACM Press.

Mannila, H., & Toivonen, H. (1997). Levelwise search and borders of theories in knowledge discovery. *Data Mining and Knowledge Discovery*, 1(3), 241–258.

Morishita, S., & Sese, J. (2000). Traversing itemset lattices with statistical metric pruning. In *Proceedings of the nineteenth ACM SIGACT-SIGMOD-SIGART symposium on database systems (PODS)* (pp. 226–236). San Diego, CA: ACM Press.

Pei, J., & Han, J. (2002). Constrained frequent pattern mining: A pattern-growth view. SIGKDD Explorations, 4(1), 31–39.

Zhu, F., Yan, X., Han, J., & Yu, P. S. (2007). gPrune: A constraint pushing framework for graph pattern mining. In *Proceedings* of the sixth Pacific-Asia conference on knowledge discovery and data mining (PAKDD). Lecture notes in computer science (Vol. 4426, pp. 388-400). Berlin: Springer.

Constructive Induction

Constructive induction is any form of ▶induction that generates new descriptors not present in the input data (Dietterich & Michalski, 1983).

Recommended Reading

Dietterich, T. G., & Michalski, R. S. (1983). A comparative review of selected methods for learning from examples. In Michalski, R. S., Carbonell, J. G., & Mitchell, T. M. (Eds.). *Machine learning:* An artificial intelligence approach, pp. 41–81. Tioga.

C

226 Content Match

Content Match

► Text Mining for Advertising

Content-Based Filtering

Synonyms

Content-based recommending

Definition

Content-based filtering is prevalent in Information Retrieval, where the text and multimedia content of documents is used to select documents relevant to a user's query. In the context this refers to content-based recommenders, that provide recommendations by comparing representations of content describing an item to representations of content that interests a user.

Content-Based Recommending

► Content-Based Filtering

Context-Sensitive Learning

▶Concept Drift

Contextual Advertising

► Text Mining for Advertising

Continual Learning

Synonyms

Life-Long Learning

Definition

A learning system that can continue adding new data without the need to ever stop or freeze the updating. Usually continual learning requires incremental and bonline learning as a component, but not every incremental learning system has the ability to achieve continual learning, i.e., the learning may deterioate after some time.

Cross References

▶Cumulative Learning

Continuous Attribute

A **continuous attribute** can assume all values on the number line within the value range. See ►Attribute and ►Measurement Scales.

Contrast Set Mining

Definition

Contrast set mining is an area of supervised descriptive rule induction. The contrast set mining problem is defined as finding contrast sets, which are conjunctions of attributes and values that differ meaningfully in their distributions across groups (Bay & Pazzani, 2001). In this context, groups are the properties of interest.

Recommended Reading

Bay, S.D., & Pazzani, M. J. (2001). Detecting group differences: Mining contrast sets. Data Mining and Knowledge Discovery, 5(3), 213-246.

Cooperative Coevolution

► Compositional Coevolution

Co-Reference Resolution

► Entity Resolution

ANTHONY WIRTH
The University of Melbourne, Victoria, Australia

Synonyms

Clustering with advice; Clustering with constraints; Clustering with qualitative information; Clustering with side information

Definition

In its rawest form, *correlation clustering* is graph optimization problem. Consider a \triangleright clustering C to be a mapping from the elements to be clustered, V, to the set $\{1, \ldots, |V|\}$, so that u and v are in the same cluster if and only if C[u] = C[v]. Given a collection of items in which each pair (u, v) has two weights w_{uv}^+ and w_{uv}^- , we must *find* a clustering C that minimizes

$$\sum_{C[u]=C[v]} w_{uv}^{-} + \sum_{C[u]\neq C[v]} w_{uv}^{+}, \tag{1}$$

or, equivalently, maximizes

$$\sum_{C[u]=C[v]} w_{uv}^{+} + \sum_{C[u]\neq C[v]} w_{uv}^{-}. \tag{2}$$

Note that although w_{uv}^+ and w_{uv}^- may be thought of as positive and negative evidence towards coassociation, the actual weights are nonnegative.

Motivation and Background

The notion of *clustering with advice*, that is nonmetric-driven relations between items, had been studied in other communities (Ferligoj & Batagelj, 1982) prior to its appearance in theoretical computer science. Traditional clustering problems, such as *k*-median and *k*-center, assume that there is some type of distance measure (metric) on the data items, and often specify the number of clusters that should be formed. In the clustering with advice framework, however, the number of clusters to be built need not be specified in advance: it can be an outcome of the objective function. Furthermore, instead of, or in addition to, a distance function, we are given advice as to which pairs of

items are similar. The two weights w_{uv}^+ and w_{uv}^- correspond to external advice about whether the pair should be clustered together or separately. Bansal, Blum, and Chawla (2002) introduced the problem to the theoretical computer science and machine-learning communities. They were motivated by database consistency problems, in which the same entity appeared in different forms in various databases. Given a collection of such records from multiple databases, the aim is to cluster together the records that appear to correspond to the same entity. From this viewpoint, the log odds ratio from some classifier.

$$\log\left(\frac{\Pr(\text{same})}{\Pr(\text{different})}\right)$$
,

corresponds to a label w_{uv} for the pair. In many applications only one of the + and – weights for the pair is nonzero, that is

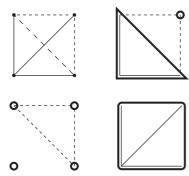
$$(w_{uv}^+, w_{uv}^-) = \begin{cases} (w_{uv}, 0) & \text{for } w_{uv} \ge 0 \\ (0, -w_{uv}) & \text{for } w_{uv} \le 0. \end{cases}$$

In addition, if every pair has weight ± 1 , then the instance is called *complete*, otherwise it is referred to as *general*. Demaine, Emanuel, Fiat, and Immorlica (2006) suggest the following motivation. Suppose we have a set of guests at a party. Each guest has preferences for whom they would like to sit with, and for whom they would like to avoid. We must group the guests into tables in a way that enhances the amicability of the party.

The notion of producing good clusterings when given inconsistent advice first appeared in the work of Ben-Dor, Shamir, and Yakhini (1999). A canonical example of inconsistent advice is this: items u and v are similar, items v and y are similar, but u and y are dissimilar. It is impossible to find a clustering that satisfies all the advice. Figure 1 shows a very simple example of inconsistent advice. In addition, although Correlation clustering is an NP-hard problem, recent algorithms for clustering with advice *guarantee* that their solutions are only a specified factor worse than the optimal: that is, they are *approximation algorithms*.

Theory

In setting out the correlation clustering framework, Bansal et al. (2002) noted that the following algorithm



Correlation Clustering. Figure 1. Top left is a toy *clustering with advice* example showing three similar pairs (*solid edges*) and three dissimilar pairs (*dashed edges*). Bottom left is a clustering solution for this example with four singleton clusters, while bottom right has one cluster. Top right is a partitioning into two clusters that appears to best respect the advice

produces a 2-approximation for the maximization problem:

If the total of the positive weights exceeds the total of the negative weights then, place all the items in a single cluster; otherwise, make each item a singleton cluster.

They then showed that complete instances are NP-hard to optimize, and how to minimize the penalty (1) with a constant factor approximation. The constant for this combinatorial algorithm was rather large. The algorithm relied heavily on the completeness of the instance; it iteratively *cleans* clusters until every cluster is δ -clean. That is, for each item at most a fraction δ (0 < δ < 1) of the other items in its cluster have a negative relation with it, and at most δ outside its cluster a positive relation. Bansal et al. also demonstrated that the minimization problem on general instances is APX-hard: there is some constant, larger than 1, below which approximation is NP-hard. Finally, they provided a polynomial time approximation scheme (PTAS) for maximizing (2) in complete instances.

The constant factor for minimizing (1) on complete instances was improved to 4 by Charikar, Guruswami, and Wirth (2003). They employed a region-growing

type procedure to round the solution of a linear programming relaxation of the problem:

minimize
$$\sum_{ij} w_{ij}^{+} \cdot x_{ij} + w_{ij}^{-} \cdot (1 - x_{ij})$$
subject to
$$x_{ik} \leq x_{ij} + x_{jk} \quad \text{for all } i, j, k$$

$$x_{ij} \in [0, 1] \quad \text{for all } i, j$$

In this setting, $x_{ij} = 1$ implies i and j's separation, while $x_{ij} = 0$ implies coclustering, with values in between representing partial evidence. In practice solving this linear program is very slow and has huge memory demands (Bertolacci & Wirth, 2007). Charikar et al. also showed that this version of problem is APX-hard.

For the maximization problem (2), they showed that instances with general weights were APX-hard and provided a rounding of the following semidefinite program (SDP) that yields a 0.7664 factor approximation algorithm.

maximize
$$\sum_{+(ij)} w_{ij} (v_i \cdot v_j) + \sum_{-(ij)} w_{ij} (1 - v_i \cdot v_j)$$
 subject to
$$v_i \cdot v_i = 1 \quad \text{for all } i$$

$$v_i \cdot v_j \ge 0 \quad \text{for all } i, j$$

In this case we interpret $v_i \cdot v_j = 1$ as evidence that i and j are in the same cluster, but $v_i \cdot v_j = 0$ as evidence toward separation.

Emanuel and Fiat (2003) extended the work of Bansal et al. by drawing a link between Correlation Clustering and the Minimum Multicut problem. This reduction to Multicut provided an $O(\log n)$ approximation algorithm for minimizing general instances of Correlation Clustering. Interestingly, Emanuel and Fiat also showed that there was reduction in the opposite direction: an optimal solution to Correlation Clustering induced an optimal solution to Minimum Multicut.

Demaine and Immorlica (2003) also drew the link from Correlation Clustering to Minimum multicut and its $O(\log n)$ approximation algorithm. In addition, they described an $O(r^3)$ -approximation algorithm for graphs that exclude the complete bipartite graph $K_{r,r}$ as a minor.

C

Swamy (2004), using the same SDP (4) as Charikar et al., but different rounding techniques, showed how to maximize (2) within factor 0.7666 in general instances.

The factor 4 approximation for minimization (1) of complete instances was lowered to 2.5 by Ailon, Charikar, and Newman (2005). Using the *distances* obtained by solving the linear program (3), they repeat the following steps:

▶ form a cluster around random item i by including each (unclustered) j with probability $1 - x_{ij}$; set the cluster aside.

Since solving the linear program is highly resource hungry, Ailon et al. provided a combinatorial alternative: add j to i's cluster if $w_{ij}^+ > w_{ij}^-$. Not only is this algorithm very fast, it is actually a factor 3 approximation.

Recently, Tan (2007) has shown that the $79/80 + \epsilon$ inapproximability for maximizing (2) on general weighted graphs extends to general unweighted graphs.

A further variant in the Correlation Clustering family of problems is the maximization of (2)–(1), known as *maximizing correlation*. Charikar and Wirth (2004) proved an $\Omega(1/\log n)$ approximation for the general problem of maximizing

$$\sum_{i=1}^{n} \sum_{j=1}^{n} a_{ij} x_i x_j, \quad \text{s.t. } x_i \in \{-1, 1\} \text{ for all } i,$$
 (5)

for a matrix A with null diagonal entries, by rounding the canonical SDP relaxation. This effectively maximized correlation with the requirement that two clusters be formed; it was not hard to extend this to general instances. The gap between the vector SDP solution and the integral solution to maximizing the quadratic program (5) was in fact shown to be $\Theta(1/\log n)$ in general (Alon, Makarychev, Makarychev, & Naor, 2006). However, in other instances such as those with a bounded number of nonzero weights for each item, a constant factor approximation was possible. Arora, Berger, Hazan, Kindler, and Safra (2005) went further and showed that it is *quasi*-NP-hard to approximate the maximization to a factor better than $\Omega(1/\log^{\gamma} n)$ for some $\gamma > 0$.

Shamir, Sharan, and Tsur (2004) showed that \blacktriangleright Cluster Editing and p-Cluster Editing, in which p clusters must be formed, are NP-complete (for $p \ge 2$). Gramm, Guo, Hüffner, and Niedermeier (2004) took

an innovative approach to solving the Clustering Editing problem exactly. They had previously produced an $O(2.27^k + n^3)$ time hand-made search tree algorithm, where k is the number of edges that need to be modified. This "awkward and error-prone work" was then replaced with a computer program that itself designed a search tree algorithm, involving automated case analysis, that ran in $O(1.92^k + n^3)$ time.

Kulis, Basu, Dhillon, and Mooney (2005) unify various forms of clustering, correlation clustering, spectral clustering, and clustering with constraints in their kernel-based approach to k-means. In this, they have a general objective function that includes penalties for violating pairwise constraints and for having points spread far apart from their cluster centers, where the *spread* is measured in some high-dimensional space.

Applications

The work of Demaine and Immorlica (2003) on Correlation Clustering was closely linked with that of Bejerano et al. on Location Area Planning. This problem is concerned with the allocation of cells in a cellular network to clusters known as *location areas*. There are costs associated with traffic between the location areas (cuts between clusters) and with the size of clusters themselves (related to paging phones within individual cells). These costs drive the clustering solution in opposite directions, on top of which there are constraints on cells that must (or cannot) be in the same cluster. The authors show that the same $O(\log n)$ region-growing algorithm for minimizing Correlation Clustering and Multicut applies to Location Area Planning.

Correlation clustering has been directly applied to the coreference problem in natural language processing and other instances in which there are multiple references to the same object (Daume, 2006; McCallum & Wellner, 2005). Assuming some sort of undirected graphical model, such as a Conditional Random Field, algorithms for correlation clustering are used to partition a graph whose edge weights corresponding to logpotentials between node pairs. The machine learning community has applied some of the algorithms for Correlation clustering to problems such as email clustering and image segmentation. With similar applications in mind, Finley and Joachims (2005) explore the idea of adapting the pairwise input information to fit example

clusterings given by a user. Their objective function is the same as Correlation Clustering (2), but their main tool is the Support Vector Machine.

There has been considerable interest in the ▶ consensus clustering problem, which is an excellent application of Correlation clustering techniques. Gionis, Mannila, and Tsaparas (2005) note several sources of motivation for the Consensus Clustering; these include identifying the correct number of clusters and improving clustering robustness. They adapt Charikar et al.'s region-growing algorithm to create a three-approximation that performs reasonably well in practice, though not as well as local search techniques. Gionis et al. also suggest using sampling as a tool for handling large data sets. Bertolacci and Wirth (2007) extended this study by implementing Ailon et al.'s algorithms with sampling, and therefore a variety of ways of developing a full clustering from the clustering of the sample. They noted that LP-based methods performed best, but placed a significant strain on resources.

Applications of Clustering with Advice

The $\triangleright k$ -means clustering algorithm is perhaps the most-used clustering technique: Wagstaff et al. incorporated constraints into a highly cited k-means variant called COP-KMEANS. They applied this algorithm to the task of identifying lanes of traffic based on input GPS data.

In the constrained-clustering framework, the constraints are usually assumed to be consistent (noncontradictory) and hard. In addition to the usual must- and cannot-link constraints, Davidson and Ravi (2005) added constraints enforcing various requirements on the distances between points in particular clusters. They analyzed the computational feasibility of the problem of establishing the (in) feasibility of a set of constraints, for various constraint types. Their constrained *k*-means algorithms were used to help a robot discover objects in a scene.

Recommended Reading

Ailon, N., Charikar, M., & Newman, A. (2005). Aggregating inconsistent information: Ranking and clustering. In Proceedings of the Thirty-Seventh ACM Symposium on the Theory of Computing (pp. 684-693). New York: ACM Press.

- Alon, N., Makarychev, K., Makarychev, Y., & Naor, A. (2006). Quadratic forms on graphs. *Inventiones Mathematicae*, 163(3), 499-522.
- Arora, S., Berger, E., Hazan, E., Kindler, G., & Safra, S. (2005).
 On non-approximability for quadratic programs. In Proceedings of Forty-Sixth Symposium on Foundations of Computer Science. (pp. 206–215). Washington DC: IEEE Computer Society.
- Bansal, N., Blum, A., & Chawla, S. (2002). Correlation clustering. In Correlation clustering (pp. 238-247). Washington, DC: IEEE Computer Society.
- Ben-Dor, A., Shamir, R., & Yakhini, Z. (1999). Clustering gene expression patterns. *Journal of Computational Biology*, 6, 281–297.
- Bertolacci, M., & Wirth, A. (2007). Are approximation algorithms for consensus clustering worthwhile? In *Proceedings of Seventh SIAM International Conference on Data Mining*. (pp. 437–442). Philadelphia: SIAM.
- Charikar, M., Guruswami, V., & Wirth, A. (2003). Clustering with qualitative information. In *Proceedings of forty fourth FOCS* (pp. 524-533).
- Charikar, M., & Wirth, A. (2004). Maximizing quadratic programs: Extending Grothendieck's inequality. In *Proceedings of forty fifth FOCS* (pp. 54–60).
- Daume, H. (2006). Practical structured learning techniques for natural language processing. PhD thesis, University of Southern California.
- Davidson, I., & Ravi, S. (2005). Clustering with constraints: Feasibility issues and the k-means algorithm. In *Proceedings of Fifth* SIAM International Conference on Data Mining.
- Demaine, E., Emanuel, D., Fiat, A., & Immorlica, N. (2006). Correlation clustering in general weighted graphs. *Theoretical Computer Science*, 361(2), 172–187.
- Demaine, E., & Immorlica, N. (2003). Correlation clustering with partial information. In *Proceedings of Sixth Workshop on Approximation Algorithms for Combinatorial Optimization Problems.* (pp. 1-13).
- Emanuel, D., & Fiat, A. (2003). Correlation clustering minimizing disagreements on arbitrary weighted graphs. In *Proceedings of Eleventh European Symposium on Algorithms* (pp. 208–220).
- Ferligoj, A., & Batagelj, V. (1982). Clustering with relational constraint. *Psychometrika*, 47(4), 413–426.
- Finley, T., & Joachims, T. (2005). Supervised clustering with support vector machines. In *Proceedings of Twenty-Second International Conference on Machine Learning*.
- Gionis, A., Mannila, H., & Tsaparas, P. (2005). Clustering aggregation. In *Proceedings of Twenty-First International Conference on Data Engineering*. To appear.
- Gramm, J., Guo, J., Hüffner, F., & Niedermeier, R. (2004).
 Automated generation of search tree algorithms for hard graph modification problems. Algorithmica, 39(4), 321–347.
- Kulis, B., Basu, S., Dhillon, I., & Mooney, R. (2005). Semi-supervised graph clustering: A kernel approach. In Proceedings of Twenty-Second International Conference on Machine Learning (pp. 457-464).
- McCallum, A., & Wellner, B. (2005). Conditional models of identity uncertainty with application to noun coreference. In L. Saul,

Cost-Sensitive Learning 231

C

Y. Weiss, & L. Bottou, (Eds.), Advances in neural information processing systems 17 (pp. 905-912). Cambridge, MA: MIT Press.

Meilă, M. (2003). Comparing clusterings by the variation of information. In *Proceedings of Sixteenth Conference on Learning Theory* (pp. 173-187).

Shamir, R., Sharan, R., & Tsur, D. (2004). Cluster graph modification problems. *Discrete Applied Mathematics*, 144, 173–182.

Swamy, C. (2004). Correlation Clustering: Maximizing agreements via semidefinite programming. In Proceedings of Fifteenth ACM-SIAM Symposium on Discrete Algorithms (pp. 519–520).

Tan, J. (2007). A Note on the inapproximability of correlation clustering. Technical Report 0704.2092, eprint arXiv, 2007.

Correlation-Based Learning

► Biological Learning: Synaptic Plasticity, Hebb Rule and Spike Timing Dependent Plasticity

Cost

In \blacktriangleright Markov decision processes, negative rewards are often expressed as costs. A reward of -x is expressed as a cost of x. In \blacktriangleright supervised learning, cost is used as a synonym for \blacktriangleright loss.

Cross References

►Loss

Cost Function

►Loss Function

Cost-Sensitive Classification

► Cost-Sensitive Learning

Cost-Sensitive Learning

CHARLES X. LING, VICTOR S. SHENG
The University of Western Ontario, Canada

Synonyms

Cost-sensitive classification; Learning with different classification costs

Definition

Cost-Sensitive Learning is a type of learning that takes the misclassification costs (and possibly other types of cost) into consideration. The goal of this type of learning is to minimize the total cost. The key difference between cost-sensitive learning and cost-insensitive learning is that cost-sensitive learning treats different misclassifications differently. That is, the cost for labeling a positive example as negative can be different from the cost for labeling a negative example as positive. Cost-insensitive learning does not take misclassification costs into consideration.

Motivation and Background

Classification is an important task in inductive learning and machine learning. A classifier, trained from a set of training examples with class labels, can then be used to predict the class labels of new examples. The class label is usually discrete and finite. Many effective classification algorithms have been developed, such as Inaïve Bayes, decision trees, neural networks, and support vector machines. However, most classification algorithms seek to minimize the error rate: the percentage of the incorrect prediction of class labels. They ignore the difference between types of misclassification errors. In particular, they implicitly assume that all misclassification errors have equal cost.

In many real-world applications, this assumption is not true. The differences between different misclassification errors can be quite large. For example, in medical diagnosis of a certain cancer (where having cancer is regarded as the positive class, and non-cancer (healthy) as negative), misdiagnosing a cancer patient as healthy (the patient is actually positive but is classified as negative; thus it is also called "false negative") is much more serious (thus expensive) than a false-positive error. The patient could lose his/her life because of a delay in correct diagnosis and treatment. Similarly, if carrying a bomb is positive, then it is much more expensive to miss a terrorist who carries a bomb onto a flight than searching an innocent person.

Cost-sensitive learning takes costs, such as the misclassification cost, into consideration. Turney (2000) provides a comprehensive survey of a large variety of different types of costs in data mining and machine 232 Cost-Sensitive Learning

learning, including misclassification costs, data acquisition cost (instance costs and attribute costs), ractive learning costs, computation cost, human-computer interaction cost, and so on. The misclassification cost is singled out as the most important cost, and it has received the most attention in recent years.

Theory

The theory of cost-sensitive learning (Elkan, 2001; Zadrozny and Elkan, 2001) describes how the misclassification cost plays its essential role in various cost-sensitive learning algorithms.

Without loss of generality, binary classification is assumed (i.e., positive and negative class) in this paper. In cost-sensitive learning, the costs of false positive (actual negative but predicted as positive; denoted as FP), false negative (FN), true positive (TP), and true negative (TN) can be given in a cost matrix, as shown in Table 1. In the table, the notation C(i, j) is also used to represent the misclassification cost of classifying an instance from its actual class *j* into the predicted class *i* (1 is used for positive, and 0 for negative). These misclassification cost values can be given by domain experts, or learned via other approaches. In cost-sensitive learning, it is usually assumed that such a cost matrix is given and known. For multiple classes, the cost matrix can be easily extended by adding more rows and more columns.

Note that C(i,i) (*TP* and *TN*) is usually regarded as the "benefit" (i.e., negated cost) when an instance is predicted correctly. In addition, cost-sensitive learning is often used to deal with datasets with very imbalanced class distributions (see \triangleright Class Imbalance Problem) (Japkowicz & Stephen, 2002). Usually (and without loss of generality), the minority or rare class is regarded as the positive class, and it is often more expensive to misclassify an actual positive example into negative,

Cost-Sensitive Learning. Table 1 An Example of Cost Matrix for Binary Classification

	Actual negative	Actual positive
Predict negative	C(0,0), or TP	C(0,1), or FN
Predict positive	C(1,0), or FP	C(1,1), or TP

than an actual negative example into positive. That is, the value of FN = C(0,1) is usually larger than that of FP = C(1,0). This is true for the cancer example mentioned earlier (cancer patients are usually rare in the population, but predicting an actual cancer patient as negative is usually very costly) and the bomb example (terrorists are rare).

Given the cost matrix, an example should be classified into the class that has the minimum expected cost. This is the minimum expected cost principle. The expected cost R(i|x) of classifying an instance x into class i (by a classifier) can be expressed as:

$$R(i|x) = \sum_{j} P(j|x) C(j,i), \qquad (1)$$

where P(j|x) is the probability estimation of classifying an instance into class j. That is, the classifier will classify an instance x into positive class if and only if:

$$P(0|x) C(1,0) + P(1|x) C(1,1) \le P(0|x) C(0,0) + P(1|x) C(0,1)$$

This is equivalent to:

$$P(0|x)(C(1,0) - C(0,0)) \le P(1|x)$$

 $(C(0,1) - C(1,1))$

Thus, the decision (of classifying an example into positive) will not be changed if a constant is added into a column of the original cost matrix. Thus, the original cost matrix can always be converted to a simpler one by subtracting C(0,0)to the first column, and C(1,1) to the second column. After such conversion, the simpler cost matrix is shown in Table 2. Thus, any given cost-matrix can be converted to one with C(0,0) = C(1,1) = 0. (Here it is assumed that the misclassification cost is the same for

Cost-Sensitive Learning. Table 2 A Simpler Cost Matrix with an Equivalent Optimal Classification

	True negative	True positive
Predict negative	0	C(0,1) – C(1,1)
Predict positive	C(1,0) - C(0,0)	0

all examples. This property is a special case of the one discussed in Elkan (2001).) In the rest of the paper, it will be assumed that C(0,0) = C(1,1) = 0. Under this assumption, the classifier will classify an instance x into positive class if and only if:

$$P(0|x)C(1,0) \leq P(1|x)C(0,1)$$

As P(0|x) = 1 - P(1|x), a threshold p^* can be obtained for the classifier to classify an instance x into positive if $P(1|x) \ge p^*$, where

$$p^* = \frac{C(1,0)}{C(1,0) + C(0,1)}. (2)$$

Thus, if a cost-insensitive classifier can produce a posterior probability estimation p(1|x) for each test example x, one can make the classifier cost-sensitive by simply choosing the classification threshold according to (2), and classify any example to be positive whenever $P(1|x) \ge p^*$. This is what several cost-sensitive metalearning algorithms, such as *Relabeling*, are based on (see later for details). However, some cost-insensitive classifiers, such as C4.5, may not be able to produce accurate probability estimation; they return a class label without a probability estimate. *Empirical Thresholding* (Sheng & Ling, 2006) does not require accurate estimation of probabilities – an accurate ranking is sufficient. It simply uses \triangleright cross-validation to search for the best probability value p^* to use as a threshold.

Traditional cost-insensitive classifiers are designed to predict the class in terms of a default, fixed threshold of 0.5. Elkan (2001) shows that one can "rebalance" the original training examples by sampling, such that the classifiers with the 0.5 threshold is equivalent to the classifiers with the p^* threshold as in (2), in order to achieve cost-sensitivity. The rebalance is done as follows. If all positive examples (as they are assumed as the rare class) are kept, then the number of negative examples should be multiplied by C(1,0)/C(0,1) = FP/FN. Note that as usually FP < FN, the multiple is less than 1. This is, thus, often called "under-sampling the majority class." This is also equivalent to "proportional sampling," where positive and negative examples are sampled by the ratio of:

$$p(1) FN : p(0) FP \tag{3}$$

where p(1) and p(0) are the prior probability of the positive and negative examples in the original training set. That is, the prior probabilities and the costs are interchangeable: doubling p(1) has the same effect as doubling FN, or halving FP (Drummond & Holte, 2000). Most sampling meta-learning methods, such as costing (Zadrozny, Langford, & Abe, 2003), are based on (3) above (see later for details).

Almost all meta-learning approaches are either based on (2) or (3) for the thresholding- and sampling-based meta-learning methods, respectively, to be discussed in the next section.

Structure of Learning System

Broadly speaking, cost-sensitive learning can be categorized into two categories. The first one is to design classifiers that are cost-sensitive in themselves. They are called the direct method. Examples of direct cost-sensitive learning are ICET (Turney, 1995) and cost-sensitive decision tree (Drummond & Holte, 2000; Ling, Yang, Wang, & Zhang, 2004). The other category is to design a "wrapper" that converts any existing cost-insensitive (or cost-blind) classifiers into cost-sensitive ones. The wrapper method is also called cost-sensitive metalearning method, and it can be further categorized into thresholding and sampling. Here is a hierarchy of the cost-sensitive learning and some typical methods. This paper will focus on cost-sensitive meta-learning that considers the misclassification cost only.

Cost-Sensitive learning

- Direct methods
 - ICET (Turney, 1995)
 - Cost-sensitive decision trees (Drummond & Holte, 2000; Ling et al., 2004)
- Meta-learning
 - Thresholding
 - MetaCost (Domingos, 1999)
 - CostSensitiveClassifier (CSC in short) (Witten & Frank, 2005)
 - Cost-sensitive naïve Bayes (Chai, Deng, Yang, & Ling, 2004)
 - Empirical Thresholding (ET in short) (Sheng & Ling, 2006)
 - Sampling
 - Costing (Zadrozny et al., 2003)
 - Weighting (Ting, 1998)

234 Cost-Sensitive Learning

Direct Cost-Sensitive Learning

The main idea of building a direct cost-sensitive learning algorithm is to directly introduce and utilize misclassification costs into the learning algorithms. There are several works on direct cost-sensitive learning algorithms, such as ICET (Turney, 1995) and cost-sensitive decision trees (Ling et al., 2004).

ICET (Turney, 1995) incorporates misclassification costs in the fitness function of genetic algorithms. On the other hand, cost-sensitive decision tree (Ling et al., 2004), called CSTree here, uses the misclassification costs directly in its tree building process. That is, instead of minimizing entropy in attribute selection as in C4.5, CSTree selects the best attribute by the expected total cost reduction. That is, an attribute is selected as a root of the (sub) tree if it minimizes the total misclassification cost.

Note that as both ICET and CSTree directly take costs into model building, they can also take easily attribute costs (and perhaps other costs) directly into consideration, while meta cost-sensitive learning algorithms generally cannot.

Drummond and Holte (2000) investigate the costsensitivity of the four commonly used attribute selection criteria of decision tree learning: accuracy, Gini, entropy, and DKM. They claim that the sensitivity of cost is highest with the accuracy, followed by Gini, entropy, and DKM.

Cost-Sensitive Meta-Learning

Cost-sensitive meta-learning converts existing costinsensitive classifiers into cost-sensitive ones without modifying them. Thus, it can be regarded as a middleware component that preprocesses the training data, or post-processes the output, from the cost-insensitive learning algorithms.

Cost-sensitive meta-learning can be further classified into two main categories: *thresholding* and *sampling*, based on (2) and (3) respectively, as discussed in the theory section.

Thresholding uses (2) as a threshold to classify examples into positive or negative if the cost-insensitive classifiers can produce probability estimations. *MetaCost* (Domingos, 1999) is a *thresholding* method. It first uses bagging on decision trees to obtain reliable probability estimations of training examples, relabels the classes of training examples according to (2), and then uses the

relabeled training instances to build a cost-insensitive classifier. *CSC* (Witten & Frank, 2005) also uses (2) to predict the class of test instances. More specifically, *CSC* uses a cost-insensitive algorithm to obtain the probability estimations $P(j \mid x)$ of each test instance. (CSC is a meta-learning method and can be applied to any classifiers.) Then it uses (2) to predict the class label of the test examples. Cost-sensitive naïve Bayes (Chai et al., 2004) uses (2) to classify test examples based on the posterior probability produced by the naïve Bayes.

As seen, all *thresholding*-based meta-learning methods rely on accurate probability estimations of p(1|x) for the test example x. To achieve this, Zadrozny and Elkan (2001) propose several methods to improve the calibration of probability estimates. ET (Empirical Thresholding) (Sheng and Ling, 2006) is a thresholding-based meta-learning method. It does not require accurate estimation of probabilities – an accurate ranking is sufficient. ET simply uses cross-validation to search the best probability from the training instances as the threshold, and uses the searched threshold to predict the class label of test instances.

On the other hand, *sampling* first modifies the class distribution of the training data according to (3), and then applies cost-insensitive classifiers on the sampled data directly. There is no need for the classifiers to produce probability estimations, as long as they can classify positive or negative examples accurately. Zadrozny et al. (2003) show that proportional sampling with replacement produces duplicated cases in the training, which in turn produces overfitting in model building. Instead, Zadrozny et al. (2003) proposes to use "rejection sampling" to avoid duplication. More specifically, each instance in the original training set is drawn once, and accepted into the sample with the accepting probability C(j,i)/Z, where C(j,i) is the misclassification cost of class i, and Z is an arbitrary constant such that $Z \ge max C(j,i)$. When $Z = max_{ij}C(j,i)$, this is equivalent to keeping all examples of the rare class, and sampling the majority class without replacement according to C(1,0)/C(0,1) – in accordance with (3). Bagging is applied after rejection sampling to improve the results further. The resulting method is called

Weighting (Ting, 1998) can also be viewed as a sampling method. It assigns a normalized weight to each instance according to the misclassification costs specified in (3). That is, examples of the rare class (which carries a higher misclassification cost) are assigned, proportionally, high weights. Examples with high weights can be viewed as example duplication – thus oversampling. Weighting then induces cost-sensitivity by integrating the instances' weights directly into C4.5, as C4.5 can take example weights directly in the entropy calculation. It works whenever the original cost-insensitive classifiers can accept example weights directly. (Thus, it can be said that Weighting is a semi meta-learning method.) In addition, Weighting does not rely on bagging as Costing does, as it "utilizes" all examples in the training set.

Recommended Reading

Chai, X., Deng, L., Yang, Q., & Ling, C. X. (2004). Test-cost sensitive naïve Bayesian classification. In Proceedings of the fourth IEEE international conference on data mining. Brighton: IEEE Computer Society Press.

Domingos, P. (1999). MetaCost: A general method for making classifiers cost-sensitive. In *Proceedings of the fifth international conference on knowledge discovery and data mining, San Diego* (pp. 155–164). New York: ACM.

Drummond, C., & Holte, R. (2000). Exploiting the cost (in)sensitivity of decision tree splitting criteria. In Proceedings of the 17th international conference on machine learning (pp. 239–246).

Elkan, C. (2001). The foundations of cost-sensitive learning. In Proceedings of the 17th international joint conference of artificial intelligence (pp. 973–978). Seattle: Morgan Kaufmann.

Japkowicz, N., & Stephen, S. (2002). The class imbalance problem: A systematic study. *Intelligent Data Analysis*, 6(5), 429-450.

Ling, C. X., Yang, Q., Wang, J., & Zhang, S. (2004). Decision trees with minimal costs. InProceedings of 2004 international conference on machine learning (ICML'2004).

Sheng, V. S., & Ling, C. X. (2006). Thresholding for making classifiers cost-sensitive. In *Proceedings of the 21st national conference on artificial intelligence* (pp. 476–481), 16–20 July 2006, Boston, Massachusetts.

Ting, K. M. (1998). Inducing cost-sensitive trees via instance weighting. In *Proceedings of the second European symposium on principles of data mining and knowledge discovery* (pp. 23–26). Heidelberg: Springer.

Turney, P. D. (1995). Cost-sensitive classification: Empirical evaluation of a hybrid genetic decision tree induction algorithm. *Journal of Artificial Intelligence Research*, 2, 369–409.

Turney, P. D. (2000). Types of cost in inductive concept learning. In Proceedings of the workshop on cost-sensitive learning at the 17th international conference on machine learning, Stanford University, California.

Witten, I. H., & Frank, E. (2005). Data mining - Practical machine learning tools and techniques with Java implementations. San Francisco: Morgan Kaufmann.

Zadrozny, B., & Elkan, C. (2001). Learning and making decisions when costs and probabilities are both unknown. In *Proceedings*

of the seventh international conference on knowledge discovery and data mining (pp. 204–213).

Zadrozny, B., Langford, J., & Abe, N. (2003). Cost-sensitive learning by cost-proportionate instance weighting. In *Proceedings of the third International conference on data mining.*

Cost-to-Go Function Approximation

► Value Function Approximation

Covariance Matrix

XINHUA ZHANG Australian National University, Canberra, Australia

Definition

It is convenient to define a covariance matrix by using multi-variate random variables (mrv): $\mathbf{X} = (X_1, \dots, X_d)^{\mathsf{T}}$. For univariate random variables X_i and X_j , their covariance is defined as:

$$Cov(X_i, X_j) = \mathbb{E}\left[(X_i - \mu_i)(X_j - \mu_j)\right],$$

where μ_i is the mean of X_i : $\mu_i = \mathbb{E}[X_i]$. As a special case, when i = j, then we get the variance of X_i , $Var(X_i) = Cov(X_i, X_i)$. Now in the setting of mrv, assuming that each component random variable X_i has finite variance under its marginal distribution, the covariance matrix $Cov(\mathbf{X}, \mathbf{X})$ can be defined as a d-by-d matrix whose (i, j)-th entry is the covariance:

$$(Cov(\mathbf{X}, \mathbf{X}))_{ij} = Cov(X_i, X_j) = \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)].$$

And its inverse is also called *precision matrix*.

Motivation and Background

The covariance between two univariate random variables measures how much they change together, and as a special case, the covariance of a random variable with itself is exactly its variance. It is important to note that covariance is an unnormalized measure of the correlation between the random variables.

As a generalization to multi-variate random variables $\mathbf{X} = (X_1, \dots, X_d)^{\mathsf{T}}$, the covariance matrix is a

236 Covariance Matrix

d-by-*d* matrix whose (i,j)-th component is the covariance between X_i and X_i .

In many applications, it is important to characterize the relations between a set of factors, hence the covariance matrix plays an important role in practice, especially in machine learning.

Theory

It is easy to rewrite the element-wise definition into the matrix form:

$$Cov(\mathbf{X}, \mathbf{X}) = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^{\mathsf{T}}],$$
 (1)

which naturally generalizes the variance of univariate random variables: $Var(X) = \mathbb{E}[(X - \mathbb{E}[X])^2]$.

Moreover, it is also straightforward to extend the covariance of a single $mrv \ \mathbf{X}$ to two $mrv \$'s $\mathbf{X} \ (d \ dimensional)$ and $\mathbf{y} \ (s \ dimensional)$, under the name cross-covariance. It quantifies how much the component random variables in \mathbf{X} and \mathbf{y} change together. The cross-covariance matrix is defined as a $d \times s$ matrix $Cov(\mathbf{X}, \mathbf{y})$ whose (i, j)-th entry is

$$(Cov(\mathbf{X}, \mathbf{y}))_{ij} = Cov(X_i, Y_j)$$

= $\mathbb{E} [(X_i - \mathbb{E}[X_i])(Y_j - \mathbb{E}[Y_j])].$

Cov(X, y) can also be written in the matrix form as

$$Cov(X, y) = \mathbb{E} [(X - \mathbb{E}[X])(y - \mathbb{E}[y])^{T}],$$

where the expectation is with respect to the joint distribution of (X, y). Obviously, Cov(X, y) becomes Cov(X, X) when y = X.

Properties

Covariance Cov(X, X) has the following properties:

- 1. Positive semi-definiteness. It follows from (1) that $Cov(\mathbf{X}, \mathbf{X})$ is positive semi-definite. $Cov(\mathbf{X}, \mathbf{X}) = \mathbf{0}$ if, and only if, \mathbf{X} is a constant almost surely, i.e., there exists a constant \mathbf{x} such that $Pr(\mathbf{X} \neq \mathbf{x}) = 0$. $Cov(\mathbf{X}, \mathbf{X})$ is not positive definite if, and only if, there exists a constant $\boldsymbol{\alpha}$ such that $\langle \boldsymbol{\alpha}, \mathbf{X} \rangle$ is constant almost surely.
- 2. Relating cumulant to moments: $Cov(X, X) = \mathbb{E}[XX^{T}] \mathbb{E}[X]\mathbb{E}[X]^{T}$.
- 3. Linear transform: If $\mathbf{y} = \mathbf{A}\mathbf{X} + \mathbf{b}$ where $\mathbf{A} \in \mathbb{R}^{s \times d}$ and $\mathbf{b} \in \mathbb{R}^{s}$, then $Cov(\mathbf{y}, \mathbf{y}) = \mathbf{A}Cov(\mathbf{X}, \mathbf{X})\mathbf{A}^{\mathsf{T}}$.

Cross-covariance Cov(X,y) has the following properties.

- 1. Symmetry: Cov(X, y) = Cov(y, X).
- 2. Linearity: $Cov(\mathbf{X}_1 + \mathbf{X}_2, \mathbf{y}) = Cov(\mathbf{X}_1, \mathbf{y}) + Cov(\mathbf{X}_2, \mathbf{y})$.
- 3. Relating to covariance: If X and y have the same dimension, then Cov(X + y, X + y) = Cov(X, X) + Cov(y, y) + 2Cov(y, X).
- 4. Linear transform: Cov(AX, By) = ACov(X, y)B.

It is highly important to note that Cov(X, y) = 0 is a necessary but not sufficient condition for X and y to be independent.

Correlation Coefficient

Entries in the covariance matrix are sometimes presented in a normalized form by dividing each entry by its corresponding standard deviations. This quantity is called the correlation coefficient, represented as ρ_{X_i,X_j} , and defined as

$$\rho_{X_i,X_j} = \frac{\text{Cov}(X_i,X_j)}{\text{Cov}(X_i,X_i)^{1/2}\text{Cov}(X_j,X_j)^{1/2}}.$$

The corresponding matrix is called the correlation matrix, and for Γ_X set to $Cov(\mathbf{X}, \mathbf{X})$ with all non-diagonal entries zeroed, and Γ_Y likewise, then the correlation matrix is given by

$$\operatorname{Corr}(\mathbf{X}, \mathbf{y}) = \Gamma_X^{-1/2} \operatorname{Cov}(\mathbf{X}, \mathbf{y}) \Gamma_Y^{-1/2}.$$

The correlation coefficient takes on values between [-1,1].

Parameter Estimation

Given observations $\mathbf{x}_1, \dots, \mathbf{x}_n$ of a *mrv* \mathbf{X} , an unbiased estimator of $Cov(\mathbf{X}, \mathbf{X})$ is:

$$S = \frac{1}{n-1} \sum_{i=1}^{n} (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^{\mathsf{T}},$$

where $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_{i}$. The denominator n-1 reflects the fact that the mean is unknown and the sample mean is used in place. Note the maximum likelihood estimator in this case replaces the denominator n-1 by n.

Covariance Matrix 237

Conjugate Priors

A covariance matrix is used to define the Gaussian distribution. In this case, the inverse Wishart distribution is the conjugate prior for the covariance matrix. Since the Gamma distribution is a 1-D version of the Wishart distribution, in the 1-D case the Gamma is the conjugate prior for precision matrix.

Applications

Several key uses of the covariance matrix are reviewed here.

Correlation and Kernel Methods

In many machine learning problems, we often need to quantify the correlation of two mrv s which may be from two different spaces. For example, we may want to study how much the image stream of a movie is correlated with the comments it receives. For simplicity, we consider a r-dimensional mrv X and a s-dimensional mrv y. To study their correlation, suppose we have npairs of observations $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ drawn *iid* from certain underlying joint distribution of (\mathbf{X}, \mathbf{y}) . Let $\bar{\mathbf{x}} =$ $\frac{1}{n}\sum_{i=1}^{n}\mathbf{x}_{i}$ and $\bar{\mathbf{y}} = \frac{1}{n}\sum_{i=1}^{n}\mathbf{y}_{i}$, and stack $\{\mathbf{x}_{i}\}$ and $\{\mathbf{y}_{i}\}$ into $\tilde{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^{\mathsf{T}}$ and $\tilde{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)^{\mathsf{T}}$ respectively. Then the cross-covariance matrix Cov(X, y) can be estimated by $\frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{y}_i - \bar{\mathbf{y}})^{\mathsf{T}}$. To quantify the crosscorrelation by a real number, we need to apply some norm of the cross-covariance matrix, and the simplest one is the Frobenius norm, $||A||_F^2 = \sum_{ij} A_{ij}^2$. Therefore, we obtain a measure of cross-correlation,

$$\left\| \frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{y}_i - \bar{\mathbf{y}})^{\top} \right\|_{\Gamma}^2 = \frac{1}{n} H \tilde{\mathbf{x}} \tilde{\mathbf{x}}^{\top} H \tilde{\mathbf{Y}} \tilde{\mathbf{Y}}^{\top}, \quad (2)$$

where $H_{ij} = \delta_{ij} - \frac{1}{n}$, and $\delta_{ij} = 1$ if i = j and 0 otherwise.

It is important to notice that (1) in this measure, inner product is performed only in the space of \mathbf{X} and \mathbf{y} separately, i.e., no transformation between \mathbf{X} and \mathbf{y} is required, (2) the data points affect the measure only via inner products $\mathbf{x}_i^{\mathsf{T}}\mathbf{x}_j$ as the (i,j)-th entry of $\tilde{x}\tilde{x}^{\mathsf{T}}$ (and similarly for \mathbf{y}_i). Hence we can endow new inner products on \mathbf{X} and \mathbf{y} , which eventually allows us to apply kernels, e.g., Gretton, Herbrich, Smola, Bousquet, & Schölkopf (2005). In a nutshell, kernel methods (Schölkopf & Smola, 2002) redefine the inner product $\mathbf{x}_i^{\mathsf{T}}\mathbf{x}_j$ by mapping \mathbf{x}_i to a richer feature space via $\phi(\mathbf{x}_i)$ and then compute the inner product

there: $k(\mathbf{x}_i, \mathbf{x}_j) := \phi(\mathbf{x}_i)^{\top} \phi(\mathbf{x}_j)$. Since the measure in (2) only needs inner products, one can even directly define k(,) without explicitly specifying ϕ . This allows us to

- Implicitly use a rich feature space whose dimension can be infinitely high.
- Apply this measure of cross correlation to non-Euclidean spaces as long as a kernel $k(\mathbf{x}_i, \mathbf{x}_j)$ can be defined on it.

Correlation and Least Squares Approximation

The measure of (2) can be equivalently motivated by least square \triangleright linear regression. That is, we look for a linear transform $T : \mathbb{R}^d \to \mathbb{R}^s$ which minimizes

$$\frac{1}{n}\sum_{i=1}^n \left\| (\mathbf{y}_i - \bar{\mathbf{y}}) - T(\mathbf{x}_i - \bar{\mathbf{x}}) \right\|^2.$$

And one can show that its minimum objective value is exactly equal to (2) up to a constant, as long as all $\mathbf{y}_i - \bar{\mathbf{y}}$ and $\mathbf{x}_i - \bar{\mathbf{x}}$ have unit length. In practice, this can be achieved by normalization. Or, the measure in (2) itself can be normalized by replacing the covariance matrix with the correlation matrix.

Principal Component Analysis

The covariance matrix plays a key role in principal component analysis (PCA). Assume that we are given n *iid* observations $\mathbf{x}_1, \dots, \mathbf{x}_n$ of a mrv \mathbf{X} , and let $\bar{\mathbf{x}} = \frac{1}{n} \sum_i \mathbf{x}_i$. PCA tries to find a set of orthogonal directions $\mathbf{w}_1, \mathbf{w}_2, \dots$, such that the projection of \mathbf{X} to the direction $\mathbf{w}_1, \mathbf{w}_1^T \mathbf{X}$, has the highest variance among all possible directions in the d-dimensional space. After subtracting from \mathbf{X} the projection to $\mathbf{w}_1, \mathbf{w}_2$ is chosen as the highest variance projection direction for the remainder. This procedure goes on for the required number of components.

To find $\mathbf{w}_1 := \operatorname{argmax}_{\mathbf{w}} \operatorname{Var}(\mathbf{w}^{\mathsf{T}} \mathbf{X})$, we need an empirical estimate of $\operatorname{Var}(\mathbf{w}^{\mathsf{T}} \mathbf{X})$. Estimating $\mathbb{E}[(\mathbf{w}^{\mathsf{T}} \mathbf{X})^2]$ by $\mathbf{w}^{\mathsf{T}} \left(\frac{1}{n} \sum_i \mathbf{x}_i \mathbf{x}_i^{\mathsf{T}}\right) \mathbf{w}$, and $\mathbb{E}[\mathbf{w}^{\mathsf{T}} \mathbf{X}]$ by $\frac{1}{n} \sum_i \mathbf{w}^{\mathsf{T}} \mathbf{x}_i$, we get

$$\mathbf{w}_1 = \operatorname{argmax}_{\mathbf{w}} : \|\mathbf{w} = 1\| \mathbf{w}^{\mathsf{T}} S \mathbf{w},$$

where $S = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^{\mathsf{T}},$

i.e., *S* is $\frac{n}{n-1}$ times the unbias empirical estimate of the covariance of **X**, based on samples $\mathbf{x}_1, \dots, \mathbf{x}_n$. \mathbf{w}_1 turns

238 Covering Algorithm

out to be exactly the eigenvector of *S* corresponding to the greatest eigenvalue.

Note that PCA is independent of the distribution of **X**. More details on PCA can be found at Jolliffe (2002).

Gaussian Processes

Gaussian processes are another important framework in machine learning that rely on the covariance matrix. It is a distribution over functions $f(\cdot)$ from certain space \mathcal{X} to \mathbb{R} , such that for any $n \in \mathbb{N}$ and any n points $\{\mathbf{x}_i \in \mathcal{X}\}_{i=1}^n$, the set of values of f evaluated at $\{\mathbf{x}_i\}_i$, $\{f(x_1),\ldots,f(x_n)\}$, will have an *n*-dimensional Gaussian distribution. Different choices of the covariance matrix of the multi-variate Gaussian lead to different stochastic processes such as Wiener process, Brownian motion, Ornstein-Uhlenbeck process, etc. In these cases, it makes more sense to define a covariance function $C: \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$, such that given any set $\{\mathbf{x}_i \in \mathcal{X}\}_{i=1}^n$ for any $n \in \mathbb{N}$, the *n*-by-*n* matrix $(C(\mathbf{x}_i, \mathbf{x}_j))_{ii}$ is positive semi-definite and can be used as the covariance matrix. This further allows straightforward kernelization of a Gaussian process by using the kernel function as the covariance function.

Although the space of functions is infinite dimensional, the marginalization property of multi-variate Gaussian distributions guarantees that the user of the model only needs to consider the observed \mathbf{x}_i , and ignore all the other possible $\mathbf{x} \in \mathcal{X}$. This important property says that for a $mrv \ \mathbf{X} = (\mathbf{X}_1^{\mathsf{T}}, \mathbf{X}_2^{\mathsf{T}})^{\mathsf{T}} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, the marginal distribution of \mathbf{X}_1 is $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$, where $\boldsymbol{\Sigma}_{11}$ is the submatrix of $\boldsymbol{\Sigma}$ corresponding to \mathbf{X}_1 (and similarly for $\boldsymbol{\mu}_1$). So taking into account the random variable \mathbf{X}_2 will not change the marginal distribution of \mathbf{X}_1 .

For a complete treatment of covariance matrix from a statistical perspective, see Casella and Berger (2002), and Mardia, Kent, and Bibby (1979) provides details for the multi-variate case. PCA is comprehensively discussed in Jolliffe (2002), and kernel methods are introduced in Schölkopf and Smola (2002). Williams & Rasmussen (2006) gives the state of the art on how Gaussian processes can be utilized for machine learning.

Cross References

- ► Gaussian Distribution
- ►Gaussian Processes
- ►Kernel Methods

Recommended Reading

Casella, G., & Berger, R. (2002). Statistical inference (2nd ed.). Pacific Grove, CA: Duxbury.

Gretton, A., Herbrich, R., Smola, A., Bousquet, O., & Schölkopf, B. (2005). Kernel methods for measuring independence. *Journal of Machine Learning Research*, 6, 2075–2129.

Jolliffe, I. T. (2002) Principal component analysis (2nd ed.). Springer series in statistics. New York: Springer.

Mardia, K. V., Kent, J. T., & Bibby, J. M. (1979). *Multivariate analysis*. London: Academic Press.

Schölkopf, B., & Smola, A. (2002). Learning with kernels. Cambridge, MA: MIT Press.

Williams, C. K. I., & Rasmussen, C. E. (2006). Gaussian processes for regression. Cambridge, MA: MIT Press.

Covering Algorithm

►Rule Learning

Credit Assignment

CLAUDE SAMMUT
The University of New South Wales

Synonyms

Structural credit assignment; Temporal credit assignment

Definition

When a learning system employs a complex decision process, it must assign credit or blame for the outcomes to each of its decisions. Where it is not possible to directly attribute an individual outcome to each decision, it is necessary to apportion credit and blame between each of the combinations of decisions that contributed to the outcome. We distinguish two cases in the credit assignment problem. *Temporal credit assignment* refers to the assignment of credit for outcomes to actions. *Structural credit assignment* refers to the assignment of credit for actions to internal decisions. The first subproblem involves determining when the actions that deserve credit were taken and the second involves assigning credit to the internal structure of actions (Sutton, 1984).

Credit Assignment 239

Motivation

Consider the problem of learning to balance a pole that is hinged on a cart (Michie & Chambers, 1968, Anderson & Miller, 1991). The cart is constrained to run along a track of finite length and a fixed force can be applied to push the cart left or right. A controller for the pole and cart system must make a decision whether to push left or right at frequent, regular time intervals, for example, 20 times a second. Suppose that this controller is capable of learning from trial-and-error. If the pole falls over, then it must determine which actions it took helped or hurt its performance. Determining that action is the problem of temporal credit assignment. Although the actions are directly responsible for the outcome of a trial, the internal process for choosing the action indirectly affects the outcome. Assigning credit or blame to those internal processes that lead to the choice of action is the structural credit assignment problem. In the case of pole balancing, the learning system will typically keep statistics such as how long, on average, the pole remained balanced after taking a particular action in a particular state, or after a failure, it may count back and determine the average amount of time to failure after taking a particular action in a particular state. Using these statistics, the learner attempts to determine the best action for a given state.

The above example is typical of many problems in reinforcement learning (Sutton & Barto, 1998), where an agent interacts with its environment and through that interaction, learns to improve its performance in a task. Although Samuel (1959) was the first to use a form of reinforcement learning in his checkers playing program, Minksy (1961) first articulated the credit assignment, as follows:

▶ Using devices that also learn which events are associated with reinforcement, i.e., reward, we can build more autonomous "secondary reinforcement" systems. In applying such methods to complex problems, one encounters a serious difficulty – in distributing credit for success of a complex strategy among the many decisions that were involved.

The BOXES algorithm of Michie and Chambers (1968) learned to control a pole balancer and performed credit assignment but the problem of credit assignment later became central to reinforcement learning, particularly following the work of Sutton (1984). Although credit

assignment has become most strongly identified with reinforcement learning, it may appear in any learning system that attempts to assess and revise its decisionmaking process.

Structural Credit Assignment

The setting for our learning system is that we have an agent that interacts with an environment. The environment may be a virtual one, as in game playing, or it may be physical, as in a robot performing some task. The agent receives input, possibly through sensing devices, that allows it to characterize the state of the world. Somehow, the agent must map these inputs to appropriate responses. These responses may change the state of the world. In reinforcement learning, we assume that the agent will receive some reward signal after an action or sequence of actions. Its job is to maximize these rewards over time.

Structural credit assignment is associated with generalization over the input space of the agent. For example, a game player may have to respond to a very large number of potential board positions or a robot may have to respond to a stream of camera images. It is infeasible to learn a complete mapping from every possible input to every possible output. Therefore, a learning agent will typically use some means of grouping input signals. In the case of the BOXES pole balancer, Michie and Chambers discretized the state space. The state is characterized by the cart's position and velocity and the pole's angle and angular velocity. These parameters create a four-dimensional space, which was broken into three regions (left, center, right) for the pole angle, five for the angular velocity, and three for the cart position and velocity. These choices were arbitrary and other combinations also worked.

Having divided the input space into non-overlapping regions, Michie and Chambers associated a push-left and push-right action with each region, or box. The learning algorithm maintains a score for each action and chooses the next action based on that score. BOXES was an early, and simple example, of creating an internal representation for mapping inputs to outputs. The problem with this method is that the structure of the decision-making system is fixed at the start and the learner is incapable of changing the representation. This may be needed if, for example, the subdivisions

240 Credit Assignment

that were chosen do not correspond to a real decision boundary. A learning system that could adapt its representation has an advantage, in this case.

The BOXES representation can be thought of as a lookup table that implements a function that maps an input to an output. The fixed lookup table can be replaced by a ▶function approximator that, given examples from the desired function, generalizes from them to construct an approximation of that function. Different function approximation techniques can be used. For example, Moore's (1990) function approximator was a ▶nearest-neighbor algorithm, implemented using **\rightarrow**kd-tree to improve efficiency. Other function approximation methods may also be used, e.g., Albus' CMAC algorithm (1975), >locally weighted regression (Atkeson, Schaal, & Moore, 1997), ▶perceptrons (Rosenblatt, 1962), ▶multi-layer networks (Hinton, Rumelhart, & Williams, 1985), ▶radial basis functions, etc. Structural credit assignment is also addressed in the creation of hierarchical representations. See ▶hierarchical reinforcement learning. Other approaches to structural credit assignment include ► Value function approximation (Bertsekas & Tsitsiklis, 1996) and automatic basis generation (Mahadevan, 2009). See the entry on ▶Gaussian Processes for examples of recent Bayesian and kernel method based approaches to solving the credit assignment problem.

Temporal Credit Assignment

In the pole balancing example described above, the learning system receives a signal when the pole has fallen over. How does it know which actions leading up to the failure contributed to the fall? The system will receive a high-level punishment in the event of a failure or a reward in tasks where there is a goal to be achieved. In either case, it makes sense to assign the greatest credit or blame to the most recent actions and assign progressively less to the preceding actions. Each time a learning trial is repeated, the value of an action is updated so that if it leads to another action of higher value, its weight is increased. Thus, the reward or punishment propagates back through the sequence of decisions taken by the system. The credit assignment problem was addressed by Michie and Chambers, in the BOXES, algorithm but many other solutions

have subsequently been proposed. See the entries on • Q-learning (Watkins, 1989; Watkins & Dayan, 1992) and • temporal difference learning (Barto, Sutton, & Anderson, 1983; Sutton, 1984).

Although temporal credit assignment is usually associated with reinforcement learning, it also appears in other forms of learning. In ▶learning by imitation or ▶behavioral cloning, an agent observes the actions of another agent and tries to learn from traces of behaviors. In this case, the learner must judge which actions of the other agent should receive credit or blame. Plan learning also encounters the same problem (Benson & Nilsson, 1995; Wang, Simon, & Lehman, 1996), as does ▶explanation-based learning (Mitchell, Keller, & Kedar-Cabelli, 1986; Dejong & Mooney, 1986; Laird, Newell, & Rosenbloom, 1987).

To illustrate the connection with explanation-based learning, we use one of the earliest examples of this kind of learning, Mitchell and Utgoff's, LEX program (Mitchell, Utgoff, & Banerji, 1983). The program was intended to learn heuristics for performing symbolic integration. Given a mathematical expression that included an integral sign, the program tried to transform the expression into one they did not. The standard symbolic integration operators were known to the program but not when it is best to apply them. The task of the learning system was to learn the heuristics for when to apply the operators. This was done by experimentation. If no heuristics were available, the program attempted a brute force search. If the search was successful, all the operators applied, leading to the success were assumed to be positive examples for a heuristic, whereas operators applied during a failed attempt became negative examples. Thus, LEX performed a simple form of credit assignment, which is typical of any system that learns how to improve sequences of decisions.

▶Genetic algorithms can also be used to evolve rules that perform sequences of actions (Holland, 1986). When situation-action rules are applied in a sequence, we have a credit assignment problem that is similar to when we use a reinforcement learning. That is, how do we know which rules were responsible for success or failure and to what extent? Grefenstette (1988) describes a bucket brigade algorithm in which rules are given strengths that are adjusted to reflect credit or blame.

Credit Assignment 241

This is similar to temporal difference learning except that in the bucket brigade the strengths apply to rules rather than states. See Classifier Systems and for a more comprehensive survey of bucket brigade methods, see Goldberg (1989).

Transfer Learning

After a person has learned to perform some task, learning a new, but related, task is usually easier because knowledge of the first learning episode is transferred to the new task. Transfer Learning is particularly useful for acquiring new concepts or behaviors when given only a small amount for training data. It can be viewed as a form of credit assignment because successes or failures in previous learning episodes bias future learning. Reid (2004, 2007) identifies three forms of ▶inductive bias involved in transfer learning for rules: language bias, which determines what kinds of rules can be constructed by the learner; the search bias, which determines the order in which rules will be searched; and the evaluation bias, which determines how the quality of the rules will be assessed. Note that learning language bias is a form of structural credit assignment. Similarly, where rules are applied sequentially, evaluation bias becomes temporal credit assignment. Taylor and Stone (2009) give a comprehensive survey of transfer in ▶reinforcement learning, in which they describe a variety of techniques for transferring the structure of an RL task from one case to another. They also survey methods for transferring evaluation bias.

Transfer learning can be applied in many different settings. Caruana (1997) developed a system for transferring inductive bias in ▶neural networks performing multitask learning and more recent research has been directed toward transfer learning in ▶Bayesian Networks (Niculescu-mizil & Caruana, 2007).

See ▶ Transfer Learning and Silver et al. (2005) and Banerjee, Liu, and Youngblood (2006) for recent work on transfer learning.

Cross References

- ► Bayesian Network
- ► Classifier Systems
- ►Genetic Algorithms

- ► Hierarchical Reinforcement Learning
- ►Inductive Bias
- ▶kd-Trees
- ► Locally Weighted Regression
- ▶Nearest-Neighbor
- **▶**Perceptrons
- ► Radial Basis Function
- ▶ Reinforcement Learning
- ▶ Temporal Difference Learning
- ► Transfer Learning

Recommended Reading

- Albus, J. S. (1975). A new approach to manipulator control: The cerebellar model articulation controller (CMAC). *Journal of Dynamic Systems, Measurement and Control, Transactions ASME, 97*(3), 220–227.
- Anderson, C. W., & Miller, W. T. (1991). A set of challenging control problems. In W. Miller, R. S. Sutton, & P. J. Werbos (Eds.), Neural Networks for Control. Cambridge: MIT Press.
- Atkeson, C., Schaal, S., & Moore, A. (1997). Locally weighted learning. AI Review, 11, 11–73.
- Banerjee, B., Liu, Y., & Youngblood, G. M. (Eds.), (2006). Proceedings of the ICML workshop on "Structural knowledge transfer for machine learning." Pittsburgh, PA.
- Barto, A., Sutton, R., & Anderson, C. (1983). Neuron-like adaptive elements that can solve difficult learning control problems. IEEE Transactions on Systems, Man, and Cybernetics, SMC-13, 834-846.
- Benson, S., & Nilsson, N. J. (1995). Reacting, planning and learning in an autonomous agent. In K. Furukawa, D. Michie, & S. Muggleton (Eds.), Machine Intelligence 14. Oxford: Oxford University Press.
- Bertsekas, D. P., & Tsitsiklis, J. (1996). Neuro-dynamic programming. Nashua, NH: Athena Scientific.
- Caruana, R. (1997). Multitask learning. *Machine Learning*, 28, 41–75.
- Dejong, G., & Mooney, R. (1986). Explanation-based learning: An alternative view. Machine Learning, 1, 145-176.
- Goldberg, D. E. (1989). Genetic algorithms in search, optimization and machine learning. Boston: Addison-Wesley Longman Publishing.
- Grefenstette, J. J. (1988). Credit assignment in rule discovery systems based on genetic algorithms. Machine Learning, 3(2-3), 225-245
- Hinton, G., Rumelhart, D., & Williams, R. (1985). Learning internal representation by back-propagating errors. In D. Rumelhart, J. McClelland, & T. P. R. Group (Eds.), Parallel distributed computing: Explorations in the microstructure of cognition (Vol. 1., pp. 31–362). Cambridge: MIT Press.

242 Cross-Language Document Categorization

Holland, J. (1986). Escaping brittleness: The possibilities of general-purpose learning algorithms applied to parallel rule-based systems. In R. S. Michalski, J. G. Carbonell, & T. M. Mitchell (Eds.), Machine learning: An artificial intelligence approach (Vol. 2). Los Altos: Morgan Kaufmann.

- Laird, J. E., Newell, A., & Rosenbloom, P. S. (1987). SOAR: An architecture for general intelligence. Artificial Intelligence, 33(1), 1-64
- Mahadevan, S. (2009). Learning representation and control in Markov decision processes: New frontiers. Foundations and Trends in Machine Learning, 1(4), 403-565.
- Michie, D., & Chambers, R. (1968). Boxes: An experiment in adaptive control. In E. Dale & D. Michie (Eds.), Machine Intelligence 2. Edinburgh: Oliver and Boyd.
- Minsky, M. (1961). Steps towards artificial intelligence. Proceedings of the IRE, 49, 8–30.
- Mitchell, T. M., Keller, R. M., & Kedar-Cabelli, S. T. (1986). Explanation based generalisation: A unifying view. Machine Learning, 1, 47–80.
- Mitchell, T. M., Utgoff, P. E., & Banerji, R. B. (1983). Learning by experimentation: Acquiring and refining problem-solving heuristics. In R. Michalski, J. Carbonell, & T. Mitchell (Eds.), Machine kearning: An artificial intelligence approach. Palo Alto: Tioga.
- Moore, A. W. (1990). Efficient memory-based learning for robot control. Ph.D. Thesis, UCAM-CL-TR-209, Computer Laboratory, University of Cambridge, Cambridge.
- Niculescu-mizil, A., & Caruana, R. (2007). Inductive transfer for Bayesian network structure learning. In *Proceedings of the 11th International Conference on AI and Statistics (AISTATS 2007)*. San Juan, Puerto Rico.
- Reid, M. D. (2004). Improving rule evaluation using multitask learning. In Proceedings of the 14th International Conference on Inductive Logic Programming (pp. 252-269). Porto, Portugal.
- Reid, M. D. (2007). DEFT guessing: Using inductive transfer to improve rule evaluation from limited data. Ph.D. thesis, School of Computer Science and Engineering, The University of New South Wales, Sydney, Australia.
- Rosenblatt, F. (1962). Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanics. Washington, DC: Spartan Books.
- Samuel, A. (1959). Some studies in machine learning using the game of checkers. IBM Journal on Research and Development, 3(3), 210-229.
- Silver, D., Bakir, G., Bennett, K., Caruana, R., Pontil, M., Russell, S., et al. (2005). NIPS workshop on "Inductive transfer: 10 years later". Whistler, Canada.
- Sutton, R. (1984). Temporal credit assignment in reinforcement learning. Ph.D. thesis, Department of Computer and Information Science, University of Massachusetts, Amherst, MA.
- Sutton, R., & Barto, A. (1998). Reinforcement learning: An introduction. Cambridge: MIT Press.
- Taylor, M. E., & Stone, P. (2009). Transfer learning for reinforcement learning domains: A survey. *Journal of Machine Learning Research*, 10, 1633–1685.
- Wang, X., Simon, H. A., Lehman, J. F., & Fisher, D. H. (1996). Learning planning operators by observation and practice. In Proceedings of the Second International Conference on AI Planning Systems, AIPS-94 (pp. 335-340). Chicago, IL.

Watkins, C. (1989). Learning with delayed rewards. Ph.D. thesis,
Psychology Department, University of Cambridge, Cambridge.
Watkins, C., & Dayan, P. (1992). Q-learning. Machine Learning,
8(3-4), 279-292.

Cross-Language Document Categorization

Document Categorization is the task consisting in assigning a document to zero, one or more categories in a predefined taxonomy. *Cross-language document categorization* describes the specific case in which one is interested in automatically categorize a document in a same taxonomy regardless of the fact that the document is written in one of several languages. For more details on the methods used to perform this task see cross-lingual text mining.

Cross-Language Information Retrieval

Cross-language information retrieval (CLIR) is the task consisting in recovering the subset of a document collection D relevant to a query q, in the special case in which D contains documents written in more than one language. Generally, it is additionally assumed that the subset of relevant documents must be returned as an ordered list, in decreasing order of relevance. For more details on methods and applications see ightharpoonup cross-lingual text mining.

Cross-Language Question Answering

Question answering is the task consisting in finding in a document collection the answer to a question. CLCat is the specific case in which the question and the documents can be in different languages. For more details on the methods used to perform this task see **crosslingual text mining**.

Cross-Lingual Text Mining

NICOLA CANCEDDA, JEAN-MICHEL RENDERS Xerox Research Centre Europe, Meylan, France

Definition

Cross-lingual text mining is a general category denoting tasks and methods for accessing the information in sets of documents written in several languages, or whenever the language used to express an information need is different from the language of the documents. A distinguishing feature of cross-lingual text mining is the necessity to overcome some language translation barrier.

Motivation and Background

Advances in mass storage and network connectivity make enormous amounts of information easily accessible to an increasingly large fraction of the world population. Such information is mostly encoded in the form of running text which, in most cases, is written in a language different from the native language of the user. This state of affairs creates many situations in which the main barrier to the fulfillment of an information need is not technological but linguistic. For example, in some cases the user has some knowledge of the language in which the text containing a relevant piece of information is written, but does not have a sufficient control of this language to express his/her information needs. In other cases, documents in many different languages must be categorized in a same categorization schema, but manually categorized examples are available for only one language.

While the automatic translation of text from a natural language into another (machine translation) is one of the oldest problems on which computers have been used, a palette of other tasks has become relevant only more recently, due to the technological advances mentioned above. Most of them were originally motivated by needs of government Intelligence communities, but received a strong impulse from the diffusion of the World-Wide Web and of the Internet in general.

Tasks and Methods

A number of specific tasks fall under the term of Crosslingual text mining (CLTM), including:

- Cross-language information retrieval
- Cross-language document categorization
- Cross-language document clustering
- Cross-language question answering

These tasks can in principle be performed using methods which do not involve any ▶Text Mining, but as a matter of fact all of them have been successfully approached relying on the statistical analysis of multilingual document collections, especially *parallel corpora*. While CLTM tasks differ in many respect, they are all characterized by the fact that they require to reliably measure the similarity of two text spans written in different languages. There are essentially two families of approaches for doing this:

- 1. In *translation-based* approaches one of the two text spans is first translated into the language of the other. Similarity is then computed based on any measure used in mono-lingual cases. As a variant, both text spans can be translated in a third *pivot* language.
- 2. In *latent semantics* approaches, an abstract vector space is defined based on the statistical properties of a *parallel corpus* (or, more rarely, of a *comparable corpus*). Both text spans are then represented as vectors in such *latent semantic* space, where any similarity measure for vector spaces can be used.

The rest of this entry is organized as follows: first Translation-related approaches will be introduced, followed by Latent-semantic approaches. Finally, each of the specific CLTM tasks will be discussed in turn.

Translation-Based Approaches

The simplest approach consists in using a manually-written machine-readable bilingual dictionary: words from the first span are looked up and replaced with words in the second language (see e.g., Zhang & Vines, 2005). Since typically dictionaries contain entries for "citation forms" only (e.g., the singular for nouns, the infinitive for verbs etc.), words in both spans are preliminarily *lemmatized*, i.e., replaced with the corresponding

C

citation form. In all cases when the lexica and morphological analyzers required to perform lemmatization are not available, a frequently adopted crude alternative consists in *stemming* (i.e., truncating by taking away a suffix) both the words in the span to be translated and in the corresponding side in the lexicon. Some languages (e.g., Germanic languages) are characterized by a very productive compounding: simpler words are connected together to form complex words. Compound words are rarely in dictionaries as such: in order to find them it is first necessary to break compounds into their elements. This can be done based on additional linguistic resources or by means of heuristics, but in all cases it is a challenging operation in itself. If the method used afterward to compare the two spans in the target language can take weights into account, translations are "normalized" in such a way that the cumulative weight of all translations of a word is the same regardless of the number of alternative translations. Most often, the weight is simply distributed uniformly among all alternative translations. Sometimes, only the first translation for each word is kept, or the first two or three.

A second approach consists in extracting a bilingual lexicon from a *parallel corpus* instead of using a manually-written one. Methods for extracting probabilistic lexica look at the frequencies with which a word s in one language was translated with a word t to estimate the translation probability p(t|s). In order to determine which word is the translation of which other word in the available examples, these examples are preliminarily aligned, first at the sentence level (to know what sentence is the translation of what other sentence) and then at the word level. Several methods for aligning sentences at the word level have been proposed, and this problem is a lively research topic in itself (see Brown, Della Pietra, Della Pietra, & Mercer, 1993 for a seminal paper).

Once a probabilistic bilingual dictionary is available, it can be used much in the same way as human-written dictionaries, with the notable difference that the estimated conditional probabilities provide a natural way to distribute weight across translations. When the example documents used for extracting the bilingual dictionaries are of the same style and domain as the text spans to be translated, this can result in a significant increase in accuracy for the final task, whatever this is.

It is often the case that a parallel corpus sufficiently similar in topic and style to the spans to be translated is unavailable, or it is too small to be used for reliably estimating translation probabilities. In such cases, it can be possible to replace or complement the parallel corpus with a "comparable" corpus. A comparable corpus is a pair of collections of documents, one in each of the languages of interest, which are known to be similar in content, although not the translation of one another. A typical case might be two sets of articles from corresponding sections of different newspapers collected during a same period of time. If some additional bilingual seed dictionary (human-written or extracted from a parallel corpus) is also available, then the comparable corpus can be leveraged as well: a word t is likely to be the translation of a word s if it turns out that the words often appearing near s are translations of the words often appearing near t. Using this observation it is thus possible to estimate the probability that t is a valid translation of s even though they are not contained in the original dictionary. Most approaches proceed by associating with s a context vector. This vector, with one component for each word in the source language, can simply be formed by summing together the count histograms of the words occurring within a fixed window centered in all occurrences of s in the corpus, but is often constructed using statistically more robust association measures, such as mutual information. After a possible normalization step, the context vector CV(s) is translated using the seed dictionary into the target language. A context vector is also extracted from the corpus for all target words t. Eventually, a translation score between s and t is computed as $\langle Tr(CV(s)), CV(t) \rangle$:

$$S(s,t) = \langle CV(s), Tr(CV(t)) \rangle$$

=
$$\sum_{(s',t') \in \mathcal{D}} a(s,s')a(t,t'),$$

where *a* is the association score used to construct the context vector. While effective in many cases, this approach can provide inaccurate similarity values when polysemous words and synonyms appear in the corpus. To deal with this problem, Gaussier, Renders, Matveeva, Goutte, and Déjean (2004) propose the following extension:

$$S(s,t) = \sum_{(s',t')\in\mathcal{D}} \left(\sum_{s'} a(s',s'')a(s,s'')\right)$$
$$\left(\sum_{t''} a(t',t'')a(t,t'')\right),$$

which is more robust in cases when the entries in the seed bilingual dictionary do not cover all senses

actually present in the two sides of the comparable corpus.

Although these methods for building bilingual dictionaries can be (and often are) used in isolation, it can be more effective to combine them.

Using a bilingual dictionary directly is not the only way for translating a span from one language into another. A second alternative consists in using a machine translation (MT) system. While the MT system, in turn, relies on a bilingual dictionary of some sort, it is in general in the position of leveraging contextual clues to select the correct words and put them in the right order in the translation. This can be more or less useful depending on the specific task. MT systems fall, broadly speaking, into two classes: rule-based and statistical. Systems in the first class rely on sets of hand-written rules describing how words and syntactic structures should be translated. Statistical machine translation (SMT) systems learn this mapping by performing a statistical analysis of a parallel corpus. Some authors (e.g., Savoy & Berger, 2005) also experimented with combining translation from multiple machine translation systems.

Latent Semantic Approaches

In CLTM, *Latent Semantic* approaches rely on some interlingua (language-independent) representation. Most of the time, this interlingua representation is obtained by linear or non-linear statistical analysis techniques and more specifically ▶dimensionality reduction methods with ad-hoc optimization criterion and constraints. But, others adopt a more manual approach by exploiting multilingual thesauri or even multilingual ontologies in order to map textual objects towards a list − possibly weighted − of interlingua concepts.

For any textual object (typically a document or a section of document), the *interlingua* concept representation is derived from a sequence of operations that encompass:

1. Linguistic preprocessing (as explained in previous sections, this step amounts to extract the relevant, normalized "terms" of the textual objects, by tokenisation, word segmentation/decompounding, lemmatisation/stemming, part-of-speech tagging, stopword removal, corpus-based term filtering, Noun-phrase extractions, etc.).

- 2. Semantic enrichment and/or monolingual dimensionality reduction.
- 3. Interlingua semantic projection.

A typical semantic enrichment method is the generalized vector space model, that adds related terms or neighbour terms - to each term of the textual object, neighbour terms being defined by some cooccurrence measures (for instance, mutual information). Semantic enrichment can alternatively be achieved by using (monolingual) thesaurus, exploiting relationships such as synonymy, hyperonymy and hyponymy. Monolingual dimensionality reduction consists typically in performing some latent semantic analysis (LSA), some form of principal component analysis on the textual object/term matrix. Dimensionality reduction techniques such as LSA or their discrete/probabilistic variants such as probabilistic semantic analysis (PLSA) and latent dirichlet allocation (LDA) offer to some extent a semantic robustness to deal with the effects of polysemy/synonymy, adopting a languagedependent concept representation in a space of dimension much smaller than the size of the vocabulary in a language.

Of course, steps (1) and (2) are highly languagedependent. Textual objects written in different languages will not follow the same linguistic processing or semantic enrichment/ dimensionality reduction. The last step (3), however, aims at projecting textual objects in the same language-independent concept space, for any source language. This is done by first extracting these common concepts, typically from a parallel corpus that offers a natural multiple-view representation of the same objects. Starting from these multiple-view observations, common factors are extracted through the use of canonical correlation analysis (CCA), crosslanguage latent semantic analysis, their kernelized variants (eg. Kernel-CCA) or their discrete, probabilistic extensions (cross-language latent dirichlet allocation, multinomial CCA, ...). All these methods try to discover latent factors that simultaneously explain as much as possible the "intra-language" variance and the "inter-language" correlation. They differ in the choice of the underlying distributions and how they precisely define and combine these two criteria. The following subsections will describe them in more details.

As already emphasized, CLTM mainly relies on defining appropriate similarities between textual objects

expressed in different languages. Numerous categorization, clustering and retrieval algorithms focus on defining efficient and powerful measures of similarity between objects, as strengthened recently by the development of kernel methods for textual information access. We will see that the (linear) statistical algorithms used for performing steps (2) and (3) can most of the time be embedded into one valid (Mercer) kernel, so that we can very easily obtain non-linear variants of these algorithms, just by adopting some standard non-linear kernels.

Cross-Language Semantic Analysis

This amounts to concatenate the vectorial representation of each view of the objects of the parallel collection (typically, objects are aligned sentences), and then to perform standard singular value decomposition of the global object/term matrix. Equivalently, defining the kernel similarity matrix between all pairs of multiview objects as the sum of the mono-lingual textual similarity matrices, this amounts to perform the eigenvalue decomposition of the corresponding kernel Gram matrix, if a dual formulation is adopted. The number of eigenvalues/eigenvectors that are retained to define the latent factors and the corresponding projections is typically from several hundreds of components to several thousands, still much fewer than the original sizes of the vocabulary. Note that this process does not really control the formation of interlingua concepts: nothing prevents the method from extracting factors that are linear combination of terms in one language only.

Cross-Language Latent Dirichlet Allocation

The extraction of *interlingua components* is realised by using LDA to model the set of parallel objects, by imposing the same proportion of components (topics) for all views of the same object. This is represented in Fig. 1.

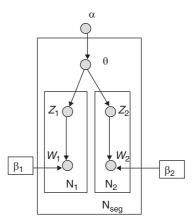
LDA is performing some form of clustering, with a predefined number of components (*K*) and with the constraint that the two views of the same object belongs to the clusters with the same membership values. This results in 2.*K* component profiles that are then used for "folding in" (projecting) new documents by launching some form of EM to derive their posterior probabilities to belong to each of the language-independent component. The similarity between two documents written in

different languages is obtained by comparing their posterior distribution over these latent classes. Note that this approach could easily integrate supervised topic information and provides a nice framework for semisupervised *interlingua* concept extraction.

Cross-Language Canonical Correlation Analysis

The Primal Formulation CCA is a standard statistical method to perform multi-block multivariate analysis, the goal being to find linear combinations of variables for each block (i.e., each language) that are maximally correlated. In other words, CCA is able to enforce the commonality of latent concept formations by extracting maximally correlated projections. Starting from a set of paired views of the same objects (typically, aligned sentences of a parallel corpus) in languages L1 and L2, the algebraic formulation of this optimization problem leads to a generalized eigenvalue problem of size $(n_1 + n_2)$, where n_1 and n_2 are the sizes of the vocabularies in L1 and L2 respectively. For obvious scalability reasons, the dual – or kernel – formulation (of size N, the number of paired objects in the training set) is often preferred.

Kernel Canonical Correlation Analysis Basically, Kernel Canonical Correlation Analysis amounts to do CCA on some implicit, but more complex feature space and to express the projection coefficients as linear combination of the training paired objects. This results in the dual formulation, which is a generalized eigenvalue/vector



Cross-Lingual Text Mining. Figure 1. Latent allocation of a parallel corpus

dirichlet

Cross-Lingual Text Mining 247

problem of size 2N, that involves only the monolingual kernel gram matrices K_1 and K_2 (matrices of monolingual textual similarities between all pairs of objects in the training set in language L1 and L2 respectively). Note that it is easy to show that the eigenvalues go by pairs: we always have two symmetrical eigenvalues $+\lambda$ and $-\lambda$. This kernel formulation has the advantage to include any text specific prior properties in the kernel (e.g., use of N-gram kernels, word-sequence kernels, and any semantically-smoothed kernel). After extraction of the first *k* generalized eigenvalues/eigenvectors, the similarity between any pair of test objects in languages L1 and L2 can be computed by using projection matrices composed of extracted eigenvector as well as the (monolingual) kernels of the test objects with the training objects.

Regularization and Partial Least Squares Solution When the number of training examples (N) is less than n_1 and n_2 (the dimensions of the monolingual feature spaces), the eigenvalue spectrum of the KCCA problem has generally two null eigenvalues (due to data centering), (N-1) eigenvalues in +1 and (N-1) eigenvalues in -1, so that, as such, the KCCA problem only results in trivial solutions and is useless. When using kernel methods, the case $(N < n_1, n_2)$ is frequent, so that some regularization scheme is needed. One way of realizing this regularization is to resort to finding the directions of maximum covariance (instead of correlation): this can be considered as a partial least squares (PLS) problem, whose formulation is very similar to the CCA problem. Adopting a mixed criterion CCA/PLS (trying to maximize a combination of covariance and correlation between projections) turns out to both avoid overfitting (or spurious solutions) and to enhance numerical stability.

Approximate Solutions Both CCA and KCCA suffer from a lack of scalability, due to the fact the complexity of generalized eigenvalue/vector decomposition is $O(N^3)$ for KCCA or $O(\min(n_1, n_2)^3)$ for CCA. As it can be shown that performing a complete KCCA (or KPLS) analysis amounts to do first complete PCA's, and then a linear CCA (or PLS) on the resulting new projections, it is obvious that we could reduce the complexity by working on a reduced-rank approximation (incomplete

KPCA) of the kernel matrices. However, the implicit projections derived from incomplete KPCA may be not optimal with respect to cross-correlation or covariance criteria. Another idea to decrease the complexity is to perform some incomplete Cholesky decomposition of the (monolingual) kernel matrices K_1 and K_2 (that is equivalent to partial Gram-Schmit orthogonalisation in the feature space): $K_1 = G_1.G_1^t$ and $K_2 = G_2.G_2^t$, with G_i of rank $k \ll N$. Considering G_i as the new representation of the training data, KCCA now reduces to solving a generalized eigenvalue problem of size 2.k.

Specific Applications

The previous sections illustrated a number of different ways of solving the core problem of cross-language text mining: quantifying the similarity between two spans of text in different languages. In this section we turn to describing some actual applications relying on these methods.

Cross-Language Information Retrieval (CLIR)

Given a collection of documents in several languages and a single query, the CLIR problem consists in producing a single ranking of all documents according to their relevance to the query. CLIR is in particular useful whenever a user has some knowledge of the languages in which documents are written, but not enough to express his/her information needs in those languages by means of a precise query. Sometimes CLIR engines are coupled with translation tools to help the user access the content of relevant documents written in languages unknown to him/her. In this case document collections in an even larger number of languages can be effectively queried.

It is probably fair to say that the vast majority of the CLIR systems use a translation-based approach. In most cases it is the query which is translated in all languages before being sent to monolingual search engines. While this limits the amount of translation work that needs be done, it requires doing it on-line at query time. Moreover, when queries are short it can be difficult to translate them correctly, since there is little context to help identifying the correct sense in which words are used. For these reasons several groups also proposed translating all documents at indexing time instead. Regardless of whether queries or documents

248 Cross-Lingual Text Mining

are translated, whenever similarity scores between (possibly translated) queries and (possibly translated) documents are not directly comparable, all methods then face the problem of merging multiple monolingual rankings in a single multilingual ranking.

Research in CLIR and cross-language question answering (see below) has been significantly stimulated by at least three government-sponsored evaluation campaigns:

- The NII Test Collection for IR Systems (NTCIR) (http://research.nii.ac.jp/ntcir/), running yearly since 1999, focusing on Asian languages (Japanese, Chinese, Korean) and English.
- The Cross-Language Evaluation Forum (CLEF) (http://www.clef-campaign.org), running yearly since 2000, focusing on European languages.
- A cross-language track at the Text Retrieval Conference (TREC) (http://trec.nist.gov/), which was run until 2002, focused on querying documents in Arabic using queries in English.

The respective websites are ideal starting points for any further exploration on the subject.

Cross-Language Question Answering (CLQA)

Question answering is the task of automatically finding the answer to a specific question in a document collection. While in practice this vague description can be instantiated in many different ways, the sense in which the term is mostly understood is strongly influenced by the task specification formulated by the National Institute of Science and Technology (NIST) of the United States for its TREC evaluation conferences (see above). In this sense, the task consists in identifying a *text snippet*, i.e., a substring, of a predefined maximal length (e.g., 50 characters, or 200 characters) within a document in the collection containing the answer. Different classes of questions are considered:

- Questions around facts and events.
- Questions requiring the definition of people, things and organizations.
- Questions requiring as answer lists of people, objects or data.

Most proposals for solving the QA problem proceed by first identifying promising documents (or document segments) by using information retrieval techniques treating the question as a query, and then performing some finer-grained analysis to converge to a sufficiently short snippet. Questions are classified in a hierarchy of possible "question types." Also, documents are preliminarily indexed to identify elements (e.g., person names) that are potential answers to questions of relevant types (e.g., "Who" questions).

Cross-language question answering (CLQA) is the extension of this task to the case where the collection contains documents in a language different than the language of the question. In this task a CLIR step replaces the monolingual IR step to shortlist promising documents. The classification of the question is generally done in the source language.

Both CLEF and NTCIR (see above) organize crosslanguage question answering comparative evaluations on an annual basis.

Cross-Language Categorization (CLCat) and Clustering (CLCLu)

Cross-language categorization tackles the problem of categorizing documents in different languages in a same categorization scheme.

The vast majority of document categorization systems rely on machine learning techniques to automatically acquire the necessary knowledge (often referred to as a model) from a possibly large collection of manually categorized documents. Most often the model is based on frequency counts of words, and is thus intrinsically language-dependent. The most direct way to perform categorization in different languages would consist in manually categorizing a sufficient amount of documents in all languages of interest and then train a set of independent categorizer. In some cases, however, it is impractical to manually categorize a sufficient number of documents to ensure accurate categorization in all languages, while it can be easier to identify bilingual dictionaries or parallel (or comparable) corpora for the language pairs and in the application domain of interest. In such cases it is then preferable to obtain manually categorized documents only for a single language A and use them to train a monolingual categorizer. Any of the translation-based approaches described above can then be used to translate a document originally in language B – or most often its representation as a bag of

words- into language A. Once the document is translated, it can be categorized using the monolingual A system.

As an alternative, latent-semantics approaches can be used as well. An existing parallel corpus can be used to identify an abstract vector space common to *A* and *B*. The manually categorized documents in *A* can then be represented in this space, and a model can be learned which operates directly on this latent-semantic representation. Whenever a document in *B* needs be categorized, it is first projected in the common semantic space and then categorized using the same model.

All these considerations carry unchanged to the cross-language clustering task, which consists in identifying subsets of documents in a multilingual document collection which are mutually similar to one another according to some criterion. Again, this task can be effectively solved by either translating all documents into a single language or by learning a common semantic space and performing the clustering task there.

While CLCat and Clustering are relevant tasks in many real-world situations, it is probably fair to say that less effort has been devoted to them by the research community than to CLIR and CLQA.

Recommended Reading

Brown, P. E., Della Pietra, V. J., Della Pietra, S. A., & Mercer, R. L. (1993). The mathematics of statistical machine translation: Parameter estimation. *Computational Linguistics*, 12(2), 263–311.

Gaussier, E., Renders, J.-M., Matveeva, I., Goutte, C., & Déjean, H. (2004). A geometric view on bilingual lexicon extraction from comparable corpora. In Proceedings of the 42nd annual meeting of the association for computational linguistics, Barcelona, Spain. Morristown, NJ: Association for Computational Linguistics.

Savoy, J., & Berger, P. Y. (2005). Report on CLEF-2005 evaluation campaign: Monolingual, bilingual and GIRT information retrieval. In *Proceedings of the cross-language evaluation forum (CLEF)* (pp. 131–140). Heidelberg: Springer.

Zhang, Y., & Vines, P. (2005). Using the web for translation disambiguation. In Proceedings of the NTCIR-5 workshop meeting, Tokyo, Japan.

Cross-Validation

Definition

Cross-validation is a process for creating a distribution of pairs of ▶training and ▶test sets out of a single

▶ data set. In cross validation the data are partitioned into k subsets, $S_1...S_k$, each called a *fold*. The folds are usually of approximately the same size. The learning algorithm is then applied k times, for i = 1 to k, each time using the union of all subsets other than S_i as the ▶ training set and using S_i as the ▶ test set.

Cross References

- ► Algorithm Evaluation
- ▶Leave-One-Out Cross-Validation

Cumulative Learning

Pietro Michelucci¹, Daniel Oblinger²
¹Strategic Analysis, Inc., Arlington, VA, USA
²DARPA/IPTO, Arlington, VA, USA

Synonyms

Continual learning; Lifelong learning; Sequential inductive transfer

Definition

Cumulative learning (CL) exploits knowledge acquired on prior tasks to improve learning performance on subsequent related tasks. Consider, for example, a CL system that is learning to play chess. Here, one might expect the system to learn from prior games concepts (e.g., favorable board positions, standard openings, end games, etc.) that can be used for future learning. This is in contrast to base learning (Vilalta & Drissi, 2002) in which a fixed learning algorithm is applied to a single task and performance tends to improve only with more exemplars. So, in CL there tends to be explicit reuse of learned knowledge to constrain new learning, whereas base learning depends entirely upon new external inputs.

Relevant techniques for CL operate over multiple tasks, often at higher levels of abstraction, such as new problem space representations, task-based selection of learning algorithms, dynamic adjustment of learning parameters, and iterative analysis and modification of the learning algorithms themselves. Though actual usage of this term is varied and evolving, CL typically connotes sequential inductive transfer. It should be noted that the word "inductive" in this connotation

C

qualifies the transfer of knowledge to new tasks, not the underlying learning algorithms.

Related Terminology

The terms "meta-learning" and "learning to learn" are sometimes used interchangeably with CL. However each of these concepts has a specific relationship to CL.

▶ Meta-learning (Brazdil et al., 2009; Vilalta & Drissi, 2002) involves the application of learning algorithms to meta-data, which are abstracted representations of input data or learning system knowledge. In the case that abstractions of system knowledge are themselves learning algorithms, meta-learning involves assessing the suitability of these algorithms for previous tasks and, on that basis, selecting algorithms for new tasks (see entry on "meta-learning"). In general, the sharing of abstracted knowledge across tasks in a CL system implies the use of meta-learning techniques. However, the converse is not true. Meta-learning can and does occur in learning systems that do not accumulate and transfer knowledge across tasks.

Learning to learn is a synonym for inductive transfer. Thus, learning to learn is more general than CL. Though it specifies the application of knowledge learned in one domain to another, it does not stipulate whether that knowledge is accumulated and applied sequentially or shared in a parallel learning context.

Motivation and Background

Traditional supervised learning approaches require large datasets and extensive training in order to generalize to new inputs in a single task. Furthermore, traditional (non-CL) reinforcement learning approaches require tightly constrained environments to ensure a

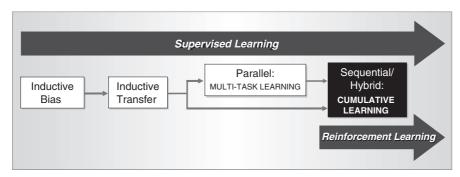
tractable state space. In contrast, humans are able to generalize across tasks in dynamic environments from brief exposure to small datasets. The human advantage seems to derive from the ability to draw upon prior task and context knowledge to constrain hypothesis development for new tasks. Recognition of this disparity between human learning and traditional machine learning had led to the pursuit of methods that seek to emulate the accumulation and exploitation of task-based knowledge that is observed in humans. A coarse evolution of this work is depicted in Fig. 1.

History

Advancements in CL have resulted from two classes of innovation: the development of techniques for
inductive transfer and the integration of those techniques into autonomous learning systems.

Alan Turing (1950) was the first to propose a cumulative learning system. His 1950 paper is best remembered for the imitation game, later known as the Turing test. However, the final sections of the paper address the question of how a machine could be made sufficiently complex to be able to pass the test. He posited that programming it would be too difficult a task. Therefore, it should be instructed as one might teach a child, starting with simple concepts and working up to more complex ones.

Banerji (1964) introduced the use of predicate logic as a description language for machine learning. Thus, Banerji was one of the earliest advocates of what would later become ILP. His concept description language allowed the use of background knowledge and therefore was an extensible language. The first implementation of a cumulative learning system based on Banerji's ideas was Cohen's CONFUCIUS (Cohen, 1978;



Cumulative Learning. Figure 1. Evolution of cumulative learning

Cohen & Sammut, 1982). In this work, an instructor teaches the system concepts that are stored in a long-term memory. When examples of a new concept are seen, their descriptions are matched against stored concepts, which allow the system to re-describe the examples in terms of the background knowledge. Thus, as more concepts are accumulated, the system is capable of describing complex objects more compactly than if it had not had the background knowledge. Compact representations generally allow complex concepts to be learned more efficiently. In many cases, learning would be intractable without the prior knowledge. See the entries on Inductive Logic Programming, which describe the use of background knowledge further.

Independent of the research in symbolic learning, much of the binductive transfer research that underlies CL took root in >artificial neural network research, a traditional approach to supervised learning. For example, Abu-Mostafa (1990) introduced the notion of reducing the hypothesis space of a neural network by introducing "hints" either as hard-wired additions to the network or via examples designed to teach a particular invariance. The task of a neural network can be thought of as the determination of a function that maps exemplars into a classification space. So, in this context, hints constitute an articulation of some aspect of the target mapping function. For example, if a neural network is tasked with mapping numbers into primes and composites, one "hint" would be that all even numbers (besides 2) are composite. Leveraging such a priori knowledge about the mapping function may facilitate convergence on a solution. An inherent limitation to neural networks, however, is their immutable architecture. which does not lend itself to the continual accumulation of knowledge. Consequently, Ring (1991) introduced a neural network that constructs new nodes on demand in a reinforcement learning context in order to support ongoing hierarchical knowledge acquisition and transfer. In this model, nodes called "bions" correspond simultaneously to the enactment and perception of a single behavior. If two bions are activated in sequence repeatedly, a new bion is created to join the coincident pair and represent their collective functionality.

Contemporaneously, Pratt, Mostow, and Kamm (1991) investigated the hypothesis that knowledge

acquired by one neural network could be used to assist another neural network learn a related task. In the speech recognition domain, they trained three separate networks, each corresponding to speech segments of a different length, such that each network was optimized to learn certain types of phonemes. They then demonstrated that a direct transfer of information encoded as network weights from these three specialized networks to a single, combined speech recognition network resulted in a tenfold reduction in training epochs for the combined network compared with the number of training epochs required when no knowledge was transferred. This was one of the first empirical results in neural network-based transfer learning. Caruana (1993) extended this work to demonstrate the performance benefits associated with the simultaneous transfer of ▶inductive bias in a "Multitask Learning" (MTL) methodology. In this work, Caruana hypothesized that training the same neural network simultaneously on related tasks would naturally induce additional constraints on learning for each individual task. The intuition was that converging on a mapping in support of multiple tasks with shared representations might best reveal aspects of the input that are invariant across tasks, thus obviating within-task regularities, which might be less relevant to classification. Those empirical results are supported by Baxter (1995) who proved that the number of examples required by a representation learner for learning a single task is an inverse linear function of the number of simultaneous tasks being learned.

Though the innovative underpinnings of inductive transfer that critically underlie CL evolved in a supervised learning context, it was the integration of those methods with classical reinforcement learning that has led to current models of CL. Early integration of this type comes from Thrun and Mitchell (1995), who applied an extension of explanation-based learning (EBL), called explanation-based neural networks (EBNN) (Mitchell & Thrun, 1993), to an agent-based "lifelong learning framework." This framework provides for the acquisition of different control policies for different environments and reward functions. Since the robot actuators, sensors, and the environment (largely) remain invariant, this framework supports the use of knowledge acquired from one control problem to be applied to another. By using EBNN to allow learning C

from previous control problems to constrain learning on new control problems, learning is accelerated over the lifetime of the robot.

More recently, Silver and Mercer (2002) introduced a hybrid model that involves a combination of parallel and sequential inductive transfer in an autonomous agent framework. The so-called task rehearsal method (TRM) uses MTL to combine new training inputs with relevant exemplars that are generated from prior task knowledge. Thus, inductive bias is achieved by training the neural networks on new tasks while simultaneously rehearsing learned task knowledge.

Structure of the Learning System

CL is characterized by systems that use prior knowledge to bias future learning. The canonical interpretation is that knowledge transfer occurs at the task level. Although this description encompasses a broad research space, it is not boundless. In particular, CL systems must be able to (1) retain knowledge and (2) use that knowledge to restrict the hypothesis space for new learning. Nonetheless, learning systems can vary widely across numerous orthogonal dimensions and still meet these criteria.

Toward a CL Specification

Recognizing the empirical utility of a more specific delineation of CL systems, Silver and Poirier (2005) introduced a set of functional requirements, classification criteria, and performance specifications that characterize more precisely the scope of machines capable of lifelong learning. Any system that meets these requirements is considered a machine lifelong learning (ML3) system. A general CL architecture that conforms to the ML3 standard is depicted in Fig. 2.

Two basic memory constructs are typical of CL systems. Long term memory (LTM) is required for storing domain knowledge (DK) that can be used to bias new learning. Short term memory (STM) provides a working memory for building representations and testing hypotheses associated with new task learning. Most of the ML3 requirements specify the interplay of these constructs.

LTM and STM are depicted in Fig. 2, along with a comparison process, an assessment process, and the learning environment. In this model, the comparison

process evaluates the training input in the context of LTM to determine the most relevant domain knowledge that can be used to constrain short term learning. The comparison process also determines the weight assigned to domain knowledge that is used to bias short term learning. Once the rate of performance improvement on the primary task falls below a threshold the assessment process compares the state of STM to the environment to determine which domain knowledge to extract and store in LTM.

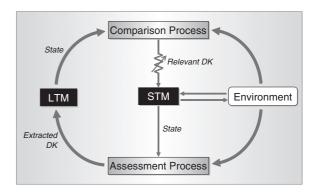
Classification of CL Systems

The simplicity of the architecture shown in Fig. 2 belies the richness of the feature space for CL systems. The following classification dimensions are derived largely from the ML3 specification. This list includes both qualitative and quantitative dimensions. They are presented in three overlapping categories: architectural features, characteristics of the knowledge base, and learning capabilities.

Architecture

The following architectural dimensions for a CL system range from paradigm choices to low-level interface considerations.

Learning paradigm – The learning paradigm(s) may include supervised learning (e.g., neural network, SVM, ILP, etc.), unsupervised learning (e.g., clustering), reinforcement learning (e.g., automated agent), or some combination thereof. Figure 2 depicts a general architecture with processes that are common across these



Cumulative Learning. Figure 2. Typical CL system

laborated to Data type – The input and output data types can be fixed or variable. A type-flexible system can produce sks sequen- both categorical and scalar predictions.

Scalability – CL systems may or may not scale on a variety of dimensions including inputs, outputs, training examples, and tasks.

learning paradigms, and which could be elaborated to reflect the details of each.

Task order – CL systems may learn tasks sequentially (Thrun & Mitchell, 1995), in parallel (e.g., MTL (Caruana, 1993)), or via a hybrid methodology (e.g., TRM (Silver & Mercer, 2002)). One hybrid approach is to engage in practice (i.e., revisiting prior learned tasks). Transferring knowledge between learned tasks through practice may serve to improve generalization accuracy. Task order would be reflected in the sequence of events within and among process arrows in the Fig. 2 architecture. For example, a system may alternate between processing new exemplars and "practicing" with old, stored exemplars.

Transfer method - Knowledge transfer can also be representational or functional. Functional transfer provides implicit pressure from related training exemplars. For example, the environmental input in Fig. 2 may take the form of training exemplars drawn randomly from data representing two related tasks, such that learning to classify exemplars from one task implicitly constrains learning on the other task. Representational knowledge transfer involves the direct or indirect (Pratt et al., 1991) assignment of a hypothesis representation. A direct inductive transfer entails the assignment of an original hypothesis representation, such as a vector of trained neural network activation weights. This might take the form of a direct injection to LTM in Fig. 2. Indirect transfer implies that some level of abstraction analysis has been applied to the hypothesis representation prior to assignment.

Learning stages – A learning system may implement learning in a single stage or in a series of stages. An example of a two-stage system is one that waits to initiate the long-term storage of domain knowledge until after primary task learning in short-term memory is complete. Like task order, learning stages would be reflected in the sequence of events within and among process arrows in the Fig. 2 architecture. But in this case, ordering pertains to the manner in which learning is staged across encoding processes.

Interface cardinality – The interface cardinality can be fixed or variable. Fixing the number of inputs and outputs has the advantage of providing a consistent interface without posing restrictions on the growth of the internal representation.

Knowledge

This category pertains to the long-term storage of learned knowledge. Thus, the following CL dimensions characterize knowledge representation, storage, and retrieval.

Knowledge representation – Stored knowledge can manifest as functional or representational. Functional knowledge retention involves the storage of specific exemplars or parameter values, which tends to be more accurate, whereas representational knowledge retention involves the storage of hypotheses derived from training on exemplars, which has the advantage of storage economy.

Retention efficacy – The efficacy of long term retention varies across CL systems. Effective retention implies that only domain knowledge with an acceptable level of accuracy is retained so that errors aren't propagated to future hypotheses. A related consideration is whether or not the consolidation of new domain knowledge degrades the accuracy of current or prior hypotheses.

Retention efficiency – The retention efficiency of long term memory can vary according to both economy of representation and computationally efficiency.

Indexing method – The input to the comparison process used to select appropriate knowledge for biasing new learning may simply be exemplars (as provided by LTM in Fig. 2) or may take a representational form (e.g., a vector of neural network weights).

Indexing efficiency – CL systems vary in terms of the speed and accuracy with which they can identify related prior knowledge that is suitable for inductive transfer during short term learning. The input to this selection process is the indexing method.

Meta-knowledge – CL systems differentially exhibit the ability to abstract, store, and utilize meta-knowledge, such as characteristics of the input space, learning system parameter values, etc. C

Cumulative Learning. Table 1 CL System Dimensions

Category	Dimension	Values (ML3 guidance is indicated by √)
Architecture	Learning paradigm	Supervised learning
		Reinforcement learning
		Unsupervised learning
		√ Hybrid
	Task order	Sequential
		Parallel
		✓ Revisit (practice)
		Hybrid
	Transfer method	Functional
		Representational – direct
		Representational – indirect
	Learning stages	✓ Single (computational retention efficiency)
		Multiple
	Interface cardinality	√ Fixed
		Variable
	Data type	Fixed
		Variable
	Scalability	✓ Inputs
		✓ Outputs
		✓ Exemplars
		√ Tasks
Knowledge	Representation	Functional
		Representational – disjoint
		✓ Representational – continuous
	Retention efficacy	✓ Improves prior task performance
		✓ Improves new task performance
	Retention efficiency	✓ Space (memory usage)
		✓ Time (computational processing)
	Indexing method	✓ Deliberative – functional
		✓ Deliberative – representational
		Reflexive

C

Cumulative Learning. Table 1 (Continued)

Category	Dimension	Values (ML3 guidance is indicated by √)
	Indexing efficiency	\checkmark Time $< O(n^c)$, $c > 1$ ($n = $ tasks)
	Meta-knowledge	✓ Probability distribution of input space
		Learning curve
		Error rate
Learning	Agency	Single learning method
		Task-based selection of learning method
	Utility	Single learning method
		Task-based selection of learning method
	Task awareness	Task boundary identification (begin/end)
	Bias modulation	✓ Estimated sample complexity
		✓ Number of task exemplars
		✓ Generalization accuracy of retained knowledge
		√ Relatedness of retained knowledge
	Learning efficacy	\checkmark Generalization \mid bias \ge generalization \mid no bias
	Learning efficiency	✓ Time bias ≤ time no bias

Learning

While all of the dimensions listed herein impact learning, the following dimensions correspond to specific learning capabilities or learning performance metrics.

Agency – The agency of a learning system is the degree of sophistication exhibited by its top-level controller. For example a learning system may be on the low end of the agency continuum if it always applies one predetermined learning method to one task or on the high end if it selects among many learning methods as a function of the learning task. One might imagine, for example, two process diagrams such as the one depicted in Fig. 2, that share the same LTM, but are otherwise distinct and differentially activated by a governing controller as a function of qualitative aspects of the input.

Utility – Domain knowledge acquisition can be deliberative in the sense that the learning system decides which hypotheses to incorporate based upon their estimated utility, or reflexive, in which case all

hypotheses are stored irrespective of utility considerations.

Task awareness – Task awareness characterizes the system's ability to identify the beginning and end of a new task.

Bias modulation – A CL system may have the ability to determine the extent to which short-term learning would benefit from inductive transfer and, on that basis, assign a relevant weight. The depth of this analysis can vary and might consider factors such as the estimated sample complexity, number of exemplars, the generalization accuracy of retained knowledge, and relatedness of retained knowledge.

Learning efficacy – A measure of learning efficacy is derived by comparing generalization performance in the presence and absence of an inductive bias. Learning is considered effective when the application of an inductive bias results in greater generalization performance on the primary task than when the bias is absent.

Learning efficiency – Similarly, learning efficiency is assessed by comparing the computational time needed to generate a hypothesis in the presence and absence of an inductive bias. Lower computational time in the presence of bias signifies greater learning efficiency.

The Research Space

Table 1 summarizes the classification dimensions, providing an overview of the research space, an evaluative framework for assessing and contrasting CL approaches, and a generative framework for identifying new areas of exploration. In addition, checked items in the Values column indicate ML3 guidance. Specifically, an ideal ML3 system would correspond functionally to the called-out items and performance criteria. However, Silver and Poirier (2005) allude to the fact that it would be nigh impossible to generate a strictly compliant ML3 system since some of the recommended criteria do not coexist easily. For example, effective and efficient learning are mutually incompatible because they require different forms of knowledge transfer. Nonetheless, a CL system that falls within scope of the majority of the ML3 criteria would be well-positioned to exhibit lifelong learning behavior.

Future Directions

Emergent work (Oblinger, 2006; Swarup, Lakkaraju, Ray, & Gasser, 2006) in instructable computing has given rise to a new CL paradigm that is largely ML3 compliant and involves high degrees of task awareness and agency sophistication. Swarup et al. (2006) describe an approach in which domain knowledge is represented in the form of structured graphs. Short term (primary task) learning occurs via a genetic algorithm, after which domain knowledge is extracted by mining frequent subgraphs. The accumulated domain knowledge forms an ontology to which the learning system grounds symbols as a result of structured interactions with instructional agents. Subsequent interactions occur using the symbol system as a shared lexicon for communication between the instructor and the learning system. Knowledge acquired from these interactions bootstrap future learning.

The Bootstrapped Learning framework proposed by Oblinger (2006) provides for hierarchical, domainindependent learning that, like the effort described above, is also premised on a model of building concepts from structured lessons. In this case, however, there is no a priori knowledge acquisition. Instead, some "common" knowledge about the world is provided explicitly to the learning system, and then lessons are taught by a human teacher using the same natural instruction methods that would be used to teach another human. Rather than requiring a specific learning algorithm, this framework provides a context for evaluating and comparing learning algorithms. It includes a knowledge representation language that supports syntactic, logical, procedural, and functional knowledge, an interaction language for communication among the learning system, instructor, and environment, and an integration architecture that evaluates, processes, and responds to interaction language communiqués in the context of existing knowledge and through the selective utilization of available learning algorithms.

The learning performance advantages anticipated by these proposals for instructional computing seem to stem from the economy of representation afforded by hierarchical knowledge combined with the tremendous learning bias imposed by explicit instruction.

Recommended Reading

Abu-Mostafa, Y. (1990). Learning from hints in neural networks (invited). *Journal of Complexity*, 6(2), 192-198.

Banerji, R. B. (1964). A Language for the Description of Concepts. General Systems, 9, 135–141.

Baxter, J. (1995). Learning internal representations. In (COLT): Proceeding of the workshop on computational learning theory, Santa Cruz, California. Morgan Kaufmann.

Brazdil P., Giraud-Carrier, C., Soares, C., & Vilalta, R. (2009).

Metalearning - Applications to Data Mining, Springer.

Caruana, R. (1993). Multitask learning: A knowledge-based source of inductive bias. In *Proceedings of the tenth international conference on machine learning*, University of Massachusetts, Amherst (pp. 41–48).

Caruana, R. (1996). Algorithms and applications for multitask learning. In Machine learning: Proceedings of the 13th international conference on machine learning (ICML 1996), Bari, Italy (pp. 87-95). Morgan Kauffmann.

Cohen, B. L. (1978). A Theory of Structural Concept Formation and Pattern Recognition. Ph.D. Thesis, Department of Computer Science, The University of New South Wales.

Cohen, B. L., & Sammut, C. A. (1982). Object Recognition and Concept Learning with CONFUCIUS. *Pattern Recognition Journal*, 15(4), 309-316.

Mitchell, T. (1980). The need for biases in learning generalizations. Rutgers TR CBM-TR-117.

Mitchell, T. M., & Thrun, S. B. (1993). Explanation-based neural network learning for robot control. In Hanson, Cowan, &

Curse of Dimensionality 257

C

- Giles (Eds.), Advances in neural information processing systems 5 (pp. 287-294). San Francisco, CA: Morgan-Kaufmann.
- Nilsson, N. J. (1996). Introduction to machine learning: An early draft of a proposed textbook (p. 12). Online at http://ai.stanford.edu/ \$\sim\$nilsson/MLBOOK.pdf. Accessed on July 22, 2010.
- Oblinger, D. (2006). Bootstrapped learning proposer information pamphlet for broad agency announcement 07-04. Online at http://fsl.fbo.gov/EPSData/ODA/Synopses/4965/BAA07-04/BLPIPfinal.pdf.
- Pratt, L. Y., Mostow, J., & Kamm, C. A. (1991). Direct transfer of learned information among neural networks. In *Proceedings of* the ninth national conference on artificial intelligence (AAAI-91), Anaheim, CA (pp. 584-589).
- Ring, M. (1991). Incremental development of complex behaviors through automatic construction of sensory-motor hierarchies. In Proceedings of the eighth international workshop (ML91), San Mateo, California.
- Silver, D., & Mercer, R. (2002). The task rehearsal method of lifelong learning: Overcoming impoverished data. In R. Cohen & B. Spencer (Eds.), Advances in artificial intelligence, 15th conference of the Canadian society for computational studies of intelligence (AI 2002), Calgary, Canada, May 27-29, 2002. Lecture notes in computer science (Vol. 2338, pp. 90-101). London: Springer.
- Silver, D., & Poirier, R. (2005). Requirements for machine lifelong learning. JSOCS Technical Report TR-2005-009, Acadia University.
- Swarup, S., Lakkaraju, K., Ray, S. R., & Gasser, L. (2006). Symbol grounding through cumulative learning. In P. Vogt et al. (Eds.), Symbol grounding and beyond: Proceedings of the third international workshop on the emergence and evolution of linguistic communication, Rome, Italy (pp. 180–191). Berlin: Springer.
- Swarup, S., Mahmud, M. M. H., Lakkaraju, K., & Ray, S. R. (2005). Cumulative learning: Towards designing cognitive architectures for artificial agents that have a lifetime. Tech. Rep. UIUCDCS-R-2005-2514.
- Thrun, S. (1998). Lifelong learning algorithms. In S. Thrun & L. Y.
 Pratt (Eds.), Learning to learn. Norwell, MA: Kluwer Academic.
 Thrun, S., & Mitchell, T. (1995). Lifelong robot learning. Robotics and Autonomous Systems, 15, 25-46.
- Turing, A. M. (1950). Computing Machinery and Intelligence. *Mind Mind*, 59(236), 433–460.
- Vilalta, R., & Drissi, Y. (2002). A perspective view and survey of meta-learning. Artificial Intelligence Review, 18, 77-95.

Curse of Dimensionality

EAMONN KEOGH, ABDULLAH MUEEN University California-Riverside, Riverside, CA, USA

Definition

The curse of dimensionality is a term introduced by Bellman to describe the problem caused by the exponential increase in volume associated with adding extra dimensions to Euclidean space (Bellman, 1957).

For example, 100 evenly-spaced sample points suffice to sample a unit interval with no more than 0.01 distance between points; an equivalent sampling of a 10-dimensional unit hypercube with a grid with a spacing of 0.01 between adjacent points would require 10^{20} sample points: thus, in some sense, the 10D hypercube can be said to be a factor of 10^{18} "larger" than the unit interval.

Informally, the phrase *curse of dimensionality* is often used to simply refer to the fact that one's intuitions about how data structures, similarity measures, and algorithms behave in low dimensions do typically generalize well to higher dimensions.

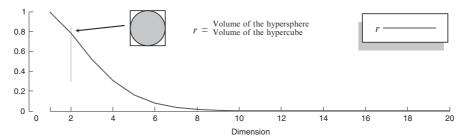
Background

Another way to envisage the vastness of high-dimensional Euclidean space is to compare the size of the unit sphere with the unit cube as the dimension of the space increases: as the dimension increases. As we can see in Fig. 1, the unit sphere becomes an insignificant volume relative to that of the unit cube. In other words, almost all of the high-dimensional space is *far away* from the center.

In research papers, the phrase *curse of dimensionality* is often used as shorthand for one of its many implications for machine learning algorithms. Examples of these implications include:

- Nearest neighbor searches can be made significantly faster for low-dimensional data by indexing the data with an R-tree, a KD-tree, or a similar spatial access method. However, for high-dimensional data all such methods degrade to the performance of a simple linear scan across the data.
- For machine learning problems, a small increase in dimensionality generally requires a large increase in the numerosity of the data, in order to keep the same level of performance for regression, clustering, etc.
- In high-dimensional spaces, the normally intuitive concept of proximity or similarity may not be qualitatively meaningful. This is because the ratio of an object's nearest neighbor over its farthest neighbor approaches one for high-dimensional spaces (Aggarwal, Hinneburg, & Keim, 2001). In other

258 Curse of Dimensionality



Curse of Dimensionality. Figure 1. The ratio of the volume of the hypersphere enclosed by the unit hypercube. The most intuitive example, the unit square and unit circle, are shown as an inset. Note that the volume of the hypersphere quickly becomes irrelevant for higher dimensionality

words, all objects are approximately equidistant from each other.

There are many ways to attempt to mitigate the *curse* of dimensionality, including ▶feature selection and ▶dimensionality reduction. However, there is no single solution to the many difficulties caused by the effect.

Recommended Reading

The major database (SIGMOD, VLDB, PODS), data mining (SIGKDD, ICDM, SDM), and machine learning (ICML, NIPS)

conferences typically feature several papers which explicitly address the curse of dimensionality each year.

Aggarwal, C. C., Hinneburg, A., & Keim, D. A. (2001). On the surprising behavior of distance metrics in high dimensional spaces. In *ICDT* (pp. 420–434). London, England.

Bellman, R. E. (1957). *Dynamic programming*. Princeton, NJ: Princeton University Press.

Ding, H., Trajcevski, G., Scheuermann, P., Wang, X., & Keogh, E. (2008). Querying and mining of time series data: Experimental comparison of representations and distance measures. In Proceedings of the VLDB endowment (Vol. 1, pp. 1542–1552). Auckland, NewZealand.