#### Course 395: Machine Learning – Lectures

- Lecture 1-2: Concept Learning (M. Pantic)
- Lecture 3-4: Decision Trees & CBC Intro (M. Pantic & S. Petridis)
- Lecture 5-6: Evaluating Hypotheses (S. Petridis)
- Lecture 7-8: Artificial Neural Networks I (S. Petridis)
- Lecture 9-10: Artificial Neural Networks II (S. Petridis)
- Lecture 11-12: Instance Based Learning (*M. Pantic*)
- Lecture 13-14: Genetic Algorithms (M. Pantic)



#### Instance Based Learning – Lecture Overview

- Lazy learning
- K-Nearest Neighbour learning
- Locally weighted regression
- Case-based reasoning (CBR)
- Advantages and disadvantages of lazy learning



#### Eager vs. Lazy Learning

- Eager learning methods construct general, explicit description of the target function based on the provided training examples.
- Lazy learning methods simply store the data and generalizing beyond these data is postponed until an explicit request is made.



- 1. Search the memory for similar instances
- 2. Retrieve the related solutions
- 3. Adapt the solutions to the current instance
- 4. Assign the value of the target function estimated for the current instance



#### Eager vs. Lazy Learning

- Eager learning methods construct general, explicit description of the target function based on the provided training examples.
- Lazy learning methods simply store the data and generalizing beyond these data is postponed until an explicit request is made.
- Lazy learning methods can construct a different approximation to the target function for each encountered query instance.
- Eager learning methods use the same approximation to the target function, which must be learned based on training examples and before input queries are observed.
- Lazy learning is very suitable for complex and incomplete problem domains, where a complex target function can be represented by a collection of less complex local approximations.



The main idea behind *k*-NN learning is so-called *majority voting*.



Imperial College London

Maja Pantic

• Given the target function  $V: X \rightarrow C$  and a set of *n* already observed instances  $(x_i, c_j)$ , where  $x_i \in X$ , i = [1..n],  $c_j \in C$ , j = [1..m],  $V(x_i) = c_j$ , *k*-NN algorithm will decide the class of the new query instance  $x_q$  based on its *k* nearest neighbours (previously observed instances)  $x_r$ , r = [1..k], in the following way:

$$V(x_q) \leftarrow c_l \in C \iff (\forall j \neq l) \sum_r E(c_l, V(x_r)) > \sum_r E(c_j, V(x_r))$$
 where  
 $E(a, b) = 1$  if  $a = b$  and  $E(a, b) = 0$  if  $a \neq b$ 

• The nearest neighbours of a query instance  $x_q$  are usually defined in terms of standard Euclidean distance:

$$d_{e}(x_{i}, x_{q}) = \sqrt{\{\sum_{g} (a_{g}(x_{i}) - a_{g}(x_{q}))^{2}\}}$$

where the instances  $x_i, x_q \in X$  are described with a set of g = [1..p] arguments  $a_g$ 

• Distance between two instances  $x_i, x_q \in X$ , described with a set of g = [1..p] arguments  $a_g$ , can be calculated as:



**Maja Pantic** 

ndor

For k = 1, the decision surface is a set of polygons (*Voronoi diagram*), completely defined by previously observed instances (training examples).





Maja Pantic

- The nearest neighbours (previously observed instances)  $x_r$ , r = [1..k], of a query instance  $x_q$  are defined based on a distance  $d(x_r, x_q)$  such as the Euclidian distance.
- A refinement of the *k*-NN algorithm: assign a weight  $w_r$  to each neighbour  $x_r$  of the query instance  $x_q$  based on the distance  $d(x_r, x_q)$  such that  $(d(x_r, x_q) \downarrow \leftrightarrow w_r \uparrow)$
- **Distance-weighted** *k*-NN algorithm: Given the target function  $V: X \to C$  and a set of *n* already observed instances  $(x_i, c_j)$ , where  $x_i \in X$ , i = [1..n],  $c_j \in C$ , j = [1..m],  $V(x_i) = c_j$ , distance weighted *k*-NN algorithm will decide the class of the query instance  $x_q$  based on its *k* nearest neighbours  $x_r$ , r = [1..k], in the following way:



## k-Nearest Neighbour Learning: Remarks

- By the distance-weighted *k*-NN algorithm, the value of *k* is of minor importance as distant examples will have very small weight and will not greatly affect the value of  $V(x_q)$ .
- If k = n, where *n* is the total number of previously observed instances, we call the algorithm a global method. Otherwise, if k < n, the algorithm is called a local method.
- *Advantage* Distance-weighted *k*-NN algorithm is robust to noisy training data: it calculates  $V(x_q)$  based on a weighted  $V(x_r)$  values of *all k* nearest neighbours  $x_r$ , effectively smoothing out the impact of isolated noisy training data.
- **Disadvantage** All *k*-NN algorithms calculate the distance between instances based on *all* attributes → if there are many irrelevant attributes, instances that belong together may still be distant from one another.
- *Remedy* weight each attribute differently when calculating the distance between two instances



## Locally Weighted Regression

- Locally weighted regression is a most general form of *k*-NN learning.
  It constructs an explicit approximation to target function *V* that fits the training examples in the local neighbourhood of the query instance *x*<sub>q</sub>.
- Local V is approximated based only on the data (neighbours) near xq.
  Weighted contribution of a datum is weighted by its distance from xq.
  Regression refers to the problem of approximating a real-valued target function.
- Locally weighted regression:

target function:  $V: X \to C$ , target function approximation near  $x_q$ :  $V'(x_q) = w_0 + w_1 a_1(x_q) + ... + w_n a_n(x_q)$ , where  $x_q \in X$  is described with a set of g = [1..n] attributes  $a_j$ training examples: set of *k* nearest neighbours  $x_r$ , r = [1..k], of the query instance  $x_q$ , learning problem: learn the most optimal weights *w* given the set of training examples

learning algorithm (distance-weighted gradient descent training rule):

$$\Delta w_j = \eta \cdot \sum_r \{ k(d(x_r, x_q)) \cdot (V(x_r) - V'(x_r)) \cdot a_j(x_r) \}$$

 $rac{1}{5}$  function of distance that determines weights of  $x_r$ 



#### Instance Based Learning – Lecture Overview

- Lazy learning
- K-Nearest Neighbour learning
- Locally weighted regression
- Case-based reasoning (CBR)
- Advantages and disadvantages of lazy learning
- (Example: CBR-based system for facial expression interpretation)



# Case Based Reasoning (CBR) – Schank's Theory

- The work of Roger Schank, inspired by findings in cognitive sciences on human reasoning and memory organization, is held to be the origin of CBR.
- Human knowledge about the world is organized in memory packets holding similar concepts and/or episodes that one experienced.
- If a memory packet contains a situation when a problem was successfully solved and the person experiences a similar situation, the previous experience is recollected and the same steps are followed to reach a solution.
- Rather than following a general set of ruls, reapplying previously successful solution schemes in a new but similar context solves the newly encountered problems.



Lazy learning is much closer to human reasoning model than this is the case with eager learning

Machine Learning (course 395)

# Case Based Reasoning (CBR) – Schank's Theory

Schank's memory-based reasoning model:

based on similarity→ of cases

- The memory of experiences is derived from enumaration of the observed cases, which are stored further in memory organization packets.
- If problems occur to which no specific case can match exactly, reason from more general similarities to come up with solutions.
  *1*-NN, *otherwise k*-NN Note: the retrieval is almost never full breadth (exhaustive).

distance

measure – The basis of memory-based model is automatic (online) learning.
 Memory of experiences is augmented by each novel experience (case).
 I.e., the process of learning never ceases.

opposite of offline learning (typical for eager learning methods), where the process of learning ceases when the training is completed

Imperial Colle

Maja Pantic

## Case Based Reasoning (CBR)

- Schank's memory-based reasoning model is the underlying reasoning model of CBR.
- CBR is *reasoning by remembering*: previously solved cases are used to suggest solutions for novel but similar problems.



- 1. Search the memory for similar instances
- 2. Retrieve the related solutions (1 / k NN)
- **3.** Adapt the solutions to the current instance
- 4. Store the new case in the memory of experiences



#### Case Based Reasoning (CBR) – Working Cycle

CBR working cycle:

- 1. RETRIEVE the most similar case(s).
- 2. REUSE the case(s) to suggest the solution for the current case.
- 3. REVISE the suggested solution.
- 4. RETAIN the case by storing it in the memory of experiences.

case base





# Case Based Reasoning (CBR) – System Design

- How the cases will be represented?
- How the case base should be organized?
- How the indexing (assigning indexes to cases to facilitate their retrieval) should be defined?
- Which retrieval algorithm is to be used?
- Which (case base) adaptation algorithm is to be used?



# Case Based Reasoning (CBR) – Cases

- Cases contain knowledge about previous experiences (solved problems).
- A case is typically composed of the problem description and the problem solution.
- The classic guideline 'the more information it stores, the more useful the case is', should be applied cautiously.
- Problem description should contain enough data for an accurate and efficient case retrieval. Useful info: retrieval statistics.
- Problem solution can be either atomic (e.g., an action) or compound (e.g., a sequence of actions).
- Cases can be either monolithic (e.g., observation → action) or compound (e.g., a set of observations → a sequence of actions; Note: parts can be processed separately).
- Cases can be represented in various ways: feature vectors, semantic nets, objects, frames, rules...



ondon

Cases should be such that an accurate and efficient retrieval is facilitated.

**Maja Pantic** 



# Case Based Reasoning (CBR) – Organisation

- *Flat* Case Base Organisation
  - The simplest case base organisation without any specific structure. Case retrieval is based on case-by-case search.
- *Clustered* Case Base Organisation
  - Cases are stored in clusters of similar cases (as originally proposed by Schank). Case retrieval includes finding the appropriate cluster(s) and searching through it for similar cases.

Case addition / deletion algorithm is more complex than by flat organisation.

- Hierarchical Case Base Organisation
  - Cases that share *features* are grouped together.
  - A semantic network containing interlinked features and categories is used.
  - Cases are associated with categories.
  - Case retrieval is feature based. It is fast and accurate.
  - Reorganisation of the case base may be very complex and difficult.



# Case Based Reasoning (CBR) – Indexing

- Case indexing: assigning indexes to cases to facilitate efficient and accurate retrieval of cases from the case base.
- Indexes are defined in terms of features / attributes of cases.
- Indexes should be:
  - *predictive* of the case relevance
  - recognisable it should be clear why they are used
  - *abstract* enough to allow for widening of the case base
  - *discriminative* enough to facilitate efficient and accurate case retrieval

# trade-off between the generality and specificity of the hypotheses (set of features) to be used for indexing



most informative

features

# Case Based Reasoning (CBR) – Retrieval

- Retrieval algorithm should retrieve case(s) most similar to the currently presented problem / situation. *preferred as it results in faster*
- 1-NN (k-NN) search retrieval and more accurate solutions A case-by-case search. Search is accurate but highly time consuming.
- *1-NN (k-NN) search* through *preselected* cases

Uses the indexing structure of the case base to preselect the cases. Then, applies *I*-NN or *k*-NN search. Faster than simple case-by-case search. It can happen that the best match is not in the preselected cases.

• *1-NN (k-NN) search* through (*preselected* and) *ranked* cases

Uses the retrieval statistics to rank the cases. Then applies the *1*-NN or *k*-NN search (through preselected cases). Search is faster than in the above mentioned cases but not necessarily more accurate.



Good retrieval algorithm: the best compromise between accuracy and efficiency.

Maja Pantic

# Case Based Reasoning (CBR) – Adaptation

- Adaptation algorithm adapts the solutions associated with the retrieved cases to the currently presented problem / situation.
- *Structural* Adaptation

Applies a set of adaptation rules directly to the retrieved solutions. Adaptation rules can include, e.g., modifying certain attributes through interpolating between relevant attributes of the retrieved cases.

• *Derivational* Adaptation

Uses algorithms / rules that have been used to generate the original solution. Can be used only for problem domains that are completely transparent. ↔ Not used very often.

• *Manual (User-driven)* Adaptation

If no exact match is found, asks the user for a feedback. Adapts the solutions accordingly. Faulty adaptations cannot be encountered. Used very often.



## Lazy Learning – Advantages

- *Incremental (online) learning*: The problem-solving ability is increased with each newly presented case.
- *Suitability for complex and incomplete problem domains*: A complex target function can be described as a collection of less complex local approximations and unknown classes can be learned.
- *Suitability for simultaneous application to multiple problems*: Examples are simply stored and can be used for multiple problem-solving purposes.
- *Ease of maintenance*: A lazy learner adapts automatically to changes in the problem domain.



### Lazy Learning – Disadvantages

- *Handling very large problem domains*: This implies high memory / storage requirements and time-consuming search for similar examples.
- Handling highly dynamic problem domains: In CBR, this involves continuous reorganisation of the case base, which may introduce errors in the case base.
  Overall, the set previously encountered examples may become outdated if a sudden large shift in the problem domain occurs.
- *Handling overly noisy data*: Such data may result in storing same problems numerous times because of the differences in cases due to noise. In turn, this implies high memory / storage requirements and time-consuming search for similar examples.
- Achieving fully automatic operation: Only for complete problem domains a fully automatic operation of a lazy learner can be expected. Otherwise, user feedback is needed for situations for which the learner has no solution.



- Tom Mitchell's book –chapter 8
- Relevant exercises from chapter 8: 8.1, 8.2, 8.3
- Case-Based Reasoning Syllabus



#### Course 395: Machine Learning – Lectures

- Lecture 1-2: Concept Learning (M. Pantic)
- Lecture 3-4: Decision Trees & CBC Intro (M. Pantic & S. Petridis)
- Lecture 5-6: Evaluating Hypotheses (S. Petridis)
- Lecture 7-8: Artificial Neural Networks I (S. Petridis)
- Lecture 9-10: Artificial Neural Networks II (S. Petridis)
- Lecture 11-12: Instance Based Learning (M. Pantic)
- Lecture 13-14: Genetic Algorithms (*M. Pantic*)

