# **Chapter ML:I**

- I. Introduction
  - Examples of Learning Tasks
  - □ Specification of Learning Problems

ML:I-8 Introduction © STEIN/LETTMANN 2005-2018

## Car Shopping Guide













What criteria form the basis of a decision?

## Risk Analysis for Credit Approval

Customer 1	
house owner	yes
income (p.a.)	51 000 EUR
repayment (p.m.)	1 000 EUR
credit period	7 years
SCHUFA entry	no
age	37
married	yes

Customer n
house owner no
income (p.a.) 55 000 EUR
repayment (p.m.) 1 200 EUR
credit period 8 years
SCHUFA entry no
age ?
married yes
...

## Risk Analysis for Credit Approval

Customer 1	
house owner	yes
income (p.a.)	51 000 EUR
repayment (p.m.)	1 000 EUR
credit period	7 years
SCHUFA entry	no
age	37
married	yes

Customer n	
house owner	no
income (p.a.)	55 000 EUR
repayment (p.m.)	1 200 EUR
credit period	8 years
SCHUFA entry	no
age	?
married	yes

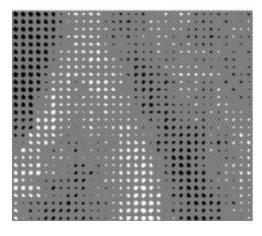
### Learned rules:

ML:I-11 Introduction

Image Analysis [Mitchell 1997]







[1992]

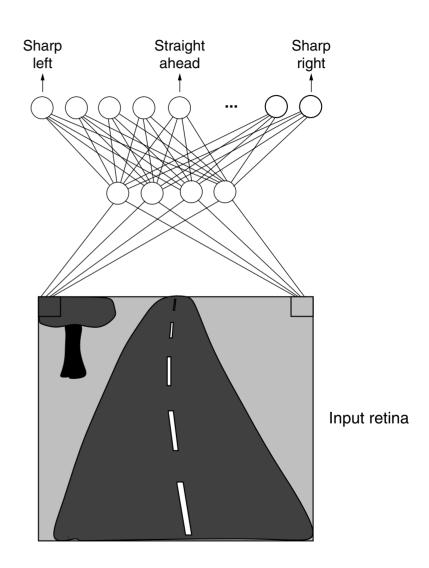
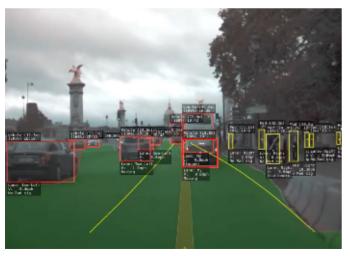
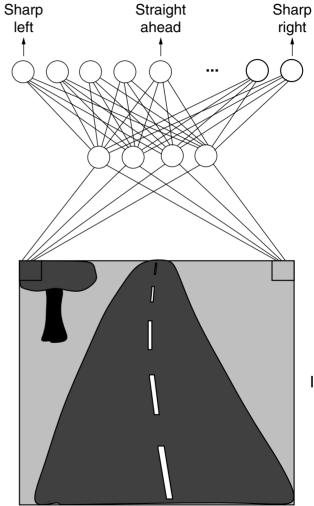


Image Analysis [Mitchell 1997]







Input retina

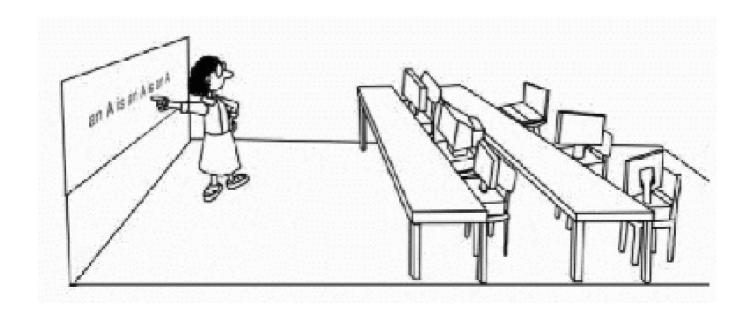
[2018]

### **Definition 1 (Machine Learning** [Mitchell 1997])

A computer program is said to learn

- from experience
- with respect to some class of tasks and
- a performance measure,

if its performance at the tasks improves with the experience.



#### Remarks:

- Example: chess
  - task = playing chess
  - performance measure = number of games won during a world championship
  - experience = possibility to play against itself
- □ Example: optical character recognition
  - task = isolation and classification of handwritten words in bitmaps
  - performance measure = percentage of correctly classified words
  - experience = collection of correctly classified, handwritten words
- □ A data set (a corpus) with labeled examples forms a kind of "compiled experience".
- Consider the different data sets that are developed and exploited for different learning tasks in the webis group. [webis.de/data.html]

ML:I-15 Introduction ©STEIN/LETTMANN 2005-2018

## Learning Paradigms

### 1. Supervised learning

Learn a function from a set of input-output-pairs. An important branch of supervised learning is automated classification. Example: optical character recognition

## 2. Unsupervised learning

Identify structures in data. Important subareas of unsupervised learning include automated categorization (e.g. via cluster analysis), parameter optimization (e.g. via expectation maximization), and feature extraction (e.g. via factor analysis).

### 3. Reinforcement learning

Learn, adapt, or optimize a behavior strategy in order to maximize the own benefit by interpreting feedback that is provided by the environment. Example: development of behavior strategies for agents in a hostile environment.

### Learning Paradigms

### 1. Supervised learning

Learn a function from a set of input-output-pairs. An important branch of supervised learning is automated classification. Example: optical character recognition

### 2. Unsupervised learning

Identify structures in data. Important subareas of unsupervised learning include automated categorization (e.g. via cluster analysis), parameter optimization (e.g. via expectation maximization), and feature extraction (e.g. via factor analysis).

### 3. Reinforcement learning

Learn, adapt, or optimize a behavior strategy in order to maximize the own benefit by interpreting feedback that is provided by the environment. Example: development of behavior strategies for agents in a hostile environment.

Example Chess: Kinds of Experience [Mitchell 1997]

### 1. Feedback

- direct: for each board configuration the best move is given.
- indirect: only the final result is given after a series of moves.

Example Chess: Kinds of Experience [Mitchell 1997]

### 1. Feedback

- direct: for each board configuration the best move is given.
- indirect: only the final result is given after a series of moves.

### 2. Sequence and distribution of examples

- A teacher presents important example problems along with a solution.
- The learner chooses from the examples; e.g., picks a board for which the best move is unknown.

The selection of examples to learn from should follow the (expected) distribution of future problems.

Example Chess: Kinds of Experience [Mitchell 1997]

### 1. Feedback

- direct: for each board configuration the best move is given.
- indirect: only the final result is given after a series of moves.

### 2. Sequence and distribution of examples

- A teacher presents important example problems along with a solution.
- The learner chooses from the examples; e.g., picks a board for which the best move is unknown.

The selection of examples to learn from should follow the (expected) distribution of future problems.

### 3. Relevance under a performance measure

- How far can we get with experience?
- Can we master situations in the wild?
   (playing against itself will be not enough to become world class)

Example Chess: Ideal Target Function  $\gamma$  [Mitchell 1997]

(a)  $\gamma$  : Boards  $\rightarrow$  Moves

(b)  $\gamma: \textit{Boards} \to \mathbf{R}$ 

Example Chess: Ideal Target Function  $\gamma$  [Mitchell 1997]

- (a)  $\gamma$  : Boards  $\rightarrow$  Moves
- (b)  $\gamma$  : *Boards*  $\rightarrow$  R

A recursive definition of  $\gamma$ , following a kind of *means-ends analysis*:

Let be  $o \in Boards$ .

- 1.  $\gamma(o) = 100$ , if o represents a final board state that is won.
- 2.  $\gamma(o) = -100$ , if o represents a final board state that is lost.
- 3.  $\gamma(o) = 0$ , if o represents a final board state that is drawn.
- 4.  $\gamma(o) = \gamma(o^*)$  otherwise.

 $o^*$  denotes the best final state that can be reached if both sides play optimally. Related: minimax strategy,  $\alpha$ - $\beta$  pruning. [Course on Search Algorithms, Stein 1998-2018]

Example Chess: From the Real World  $\gamma$  to a Model World y

$$\gamma(o) \leadsto y(\alpha(o)) \equiv y(\mathbf{x})$$

Example Chess: From the Real World  $\gamma$  to a Model World y

$$\gamma(o) \leadsto y(\alpha(o)) \equiv y(\mathbf{x})$$
  
 $y(\mathbf{x}) = w_0 + w_1 \cdot x_1 + w_2 \cdot x_2 + w_3 \cdot x_3 + w_4 \cdot x_4 + w_5 \cdot x_5 + w_6 \cdot x_6$ 

### where

- $x_1$  = number of black pawns on board o
- $x_2$  = number of white pawns on board o
- $x_3$  = number of black pieces on board o
- $x_4$  = number of white pieces on board o
- $x_5$  = number of black pieces threatened on board o
- $x_6$  = number of white pieces threatened on board o

Example Chess: From the Real World  $\gamma$  to a Model World y

$$\gamma(o) \leadsto y(\alpha(o)) \equiv y(\mathbf{x})$$
  
 $y(\mathbf{x}) = w_0 + w_1 \cdot x_1 + w_2 \cdot x_2 + w_3 \cdot x_3 + w_4 \cdot x_4 + w_5 \cdot x_5 + w_6 \cdot x_6$ 

### where

- $x_1$  = number of black pawns on board o
- $x_2$  = number of white pawns on board o
- $x_3$  = number of black pieces on board o
- $x_4$  = number of white pieces on board o
- $x_5$  = number of black pieces threatened on board o
- $x_6$  = number of white pieces threatened on board o

### Other approaches to formulate y:

- case base
- set of rules
- neural network
- polynomial function of board features

#### Remarks:

- The *ideal target function*  $\gamma$  interprets the real world, say, a real-world object o, to "compute"  $\gamma(o)$ . This "computation" can be operationalized by a human or by some other (even arcane) mechanism of the real world.
- To simulate the interesting aspects of the real world by means of a computer, we define a model world. This model world is restricted to particular—typically easily measurable—features  $\mathbf{x}$  that are derived from o, with  $\mathbf{x} = \alpha(o)$ . In the model world,  $y(\mathbf{x})$  is the abstracted and formalized counterpart of  $\gamma(o)$ .
- $\Box$  The key difference between an ideal target function  $\gamma$  and a model function y lies in the complexity and the representation of their respective domains. Examples:
  - A chess grand master assesses a board o in its entirety, both intuitively and analytically; a chess program is restricted to particular features  $\mathbf{x}$ ,  $\mathbf{x} = \alpha(o)$ .
  - A human mushroom picker assesses a mushroom o with all her skills (intuitively, analytically, by tickled senses); a classification program is restricted to a few surface features  $\mathbf{x}$ ,  $\mathbf{x} = \alpha(o)$ .

#### Remarks (continued):

- □ For automated chess playing a real-valued assessment function is needed; such kind of problems form regression problems. If only a small number of values are to be considered (e.g. school grades), we are given a classification problem. A regression problem can be transformed into a classification problem via domain discretization.
- Regression problems and classification problems often differ with regard to assessing the achieved accuracy or goodness of fit. For regression problems the sum of the squared residuals may be a sensible criterion; for classification problems the number of misclassified examples may be more relevant.
- $\supset$  For classification problems, the ideal target function  $\gamma$  is also called ideal *classifier*; similarly, the model function y is called classifier.
- Decision problems are classification problems with two classes.
- ☐ The halting problem for Turing machines is an undecidable classification problem.

ML:I-27 Introduction © STEIN/LETTMANN 2005-2018

# Specification of Learning Problems [model world]

How to Build a Classifier y

### Characterization of the real world:

- □ *O* is a set of objects.
- $\Box$  *C* is a set of classes.
- $\neg \gamma: O \to C$  is the ideal classifier for O.

(example: mails)

(example: spam versus ham)

( $\gamma$  is a human expert)

# Specification of Learning Problems [model world]

How to Build a Classifier y

### Characterization of the real world:

- □ O is a set of objects.
  - $\Box$  C is a set of classes. (example: spam versus ham)
- $\neg \gamma: O \to C$  is the ideal classifier for O.  $(\gamma \text{ is a human expert})$

### Classification problem:

 $\Box$  Given some  $o \in O$ , determine its class  $\gamma(o) \in C$ .

(example: is a mail spam?)

(example: mails)

### Acquisition of classification knowledge:

- 1. Collect real-world examples of the form  $(o, \gamma(o)), o \in O$ .
- 2. Abstract the objects towards feature vectors  $\mathbf{x} \in X$ , where  $\mathbf{x} = \alpha(o)$ .
- 3. Construct examples as  $(\mathbf{x}, c(\mathbf{x}))$ , where  $\mathbf{x} = \alpha(o)$  and  $c(\mathbf{x}) \equiv \gamma(o)$ .

## Specification of Learning Problems [real world]

How to Build a Classifier *y* (continued)

### Characterization of the model world:

- $\Box$  X is a set of feature vectors, called feature space. (example: word frequencies)
- $\Box$  C is a set of classes. (as before: spam versus ham)
- $c: X \to C$  is the ideal classifier for X. (c is unknown)
- $D = \{(\mathbf{x}_1, c(\mathbf{x}_1)), \dots, (\mathbf{x}_n, c(\mathbf{x}_n))\} \subseteq X \times C$  is a set of examples.

## Specification of Learning Problems [real world]

How to Build a Classifier *y* (continued)

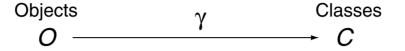
### Characterization of the model world:

- $\Box$  X is a set of feature vectors, called feature space. (example: word frequencies)
- $\Box$  C is a set of classes. (as before: spam versus ham)
- $\neg c: X \to C$  is the ideal classifier for X. (c is unknown)
- $D = \{(\mathbf{x}_1, c(\mathbf{x}_1)), \dots, (\mathbf{x}_n, c(\mathbf{x}_n))\} \subseteq X \times C \text{ is a set of examples.}$

### Machine learning problem:

- $\Box$  Approximate the ideal classifier c (implicitly given via D) by a function y:
- ightharpoonup Formulate a model function  $y: X \to C$ ,  $\mathbf{x} \mapsto y(\mathbf{x})$  (y needs to be fitted)
- → Apply statistics, search, theory, and algorithms from the field of machine learning to maximize the goodness of fit between c and y.

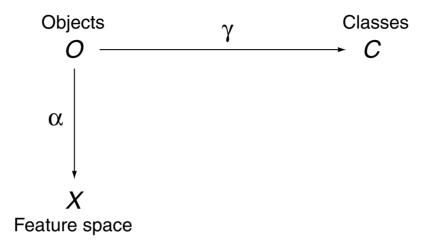
How to Build a Classifier y (continued)



### Semantics:

 $\gamma$  Ideal classifier (a human) for real-world objects.

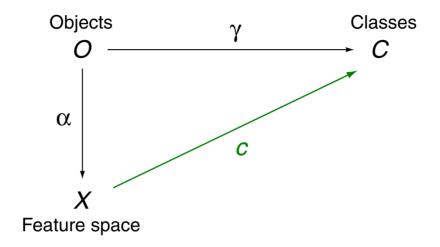
How to Build a Classifier *y* (continued)



### Semantics:

- $\gamma$  Ideal classifier (a human) for real-world objects.
- $\alpha$  Model formation function.

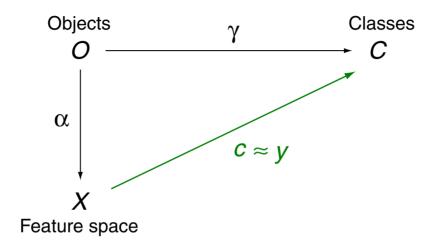
How to Build a Classifier *y* (continued)



### Semantics:

- $\gamma$  Ideal classifier (a human) for real-world objects.
- $\alpha$  Model formation function.
- c Unknown ideal classifier for vectors from the feature space.

How to Build a Classifier *y* (continued)



### Semantics:

- $\gamma$  Ideal classifier (a human) for real-world objects.
- $\alpha$  Model formation function.
- c Unknown ideal classifier for vectors from the feature space.
- y Classifier to be learned.
- $c \approx y$  c is approximated by y (based on a set of examples).

#### Remarks:

- The feature space X comprises vectors  $\mathbf{x}_1, \mathbf{x}_2, \ldots$ , which can be considered as abstractions of real-world objects  $o_1, o_2, \ldots$ , and which have been computed according to our view of the real world.
- The model formation function  $\alpha$  determines the level of abstraction between o and  $\mathbf{x}$ ,  $\mathbf{x} = \alpha(o)$ . I.e.,  $\alpha$  determines the representation fidelity, exactness, quality, or simplification.
- Though  $\alpha$  models an object  $o \in O$  only imperfectly as  $\mathbf{x} = \alpha(o)$ , c must be considered as *ideal* classifier, since  $c(\mathbf{x})$  is defined as  $\gamma(o)$ , the true real-world class. I.e., c and  $\gamma$  have different domains each, but they return the same images.
- $\Box$  The function c is only implicitly given, in the form of the example set D. The explicit form of c is unknown, it is approximated by y.
- $c(\mathbf{x})$  is often termed "ground truth" (for  $\mathbf{x}$  and the underlying classification problem). Observe that this term is justified by the fact that  $c(\mathbf{x}) \equiv \gamma(o)$ .

LMS Algorithm for Fitting y [IGD Algorithm]

Algorithm: LMS Least Mean Squares.

Input: D Training examples of the form  $(\mathbf{x}, c(\mathbf{x}))$  with target function value  $c(\mathbf{x})$  for  $\mathbf{x}$ .

 $\eta$  Learning rate, a small positive constant.

Internal: y(D) Set of y(x)-values computed from the elements x in D given some w.

Output: w Weight vector.

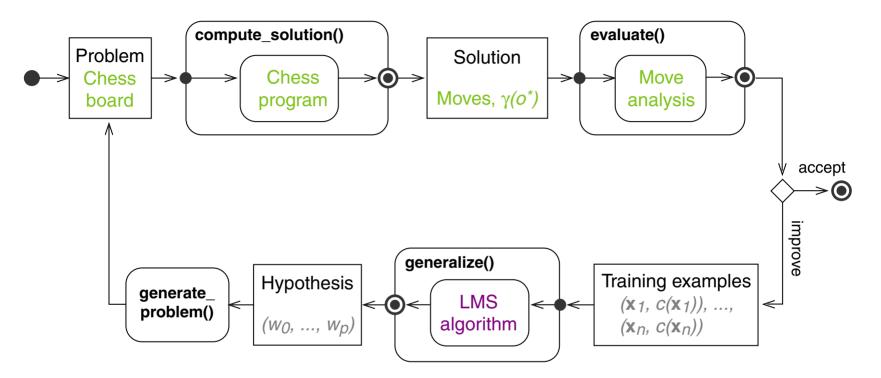
```
LMS(D, \eta)
```

- 1. initialize\_random\_weights $((w_0, w_1, \dots, w_p))$
- 2. REPEAT
- 3.  $(\mathbf{x}, c(\mathbf{x})) = random\_select(D)$
- 4.  $y(\mathbf{x}) = w_0 + w_1 \cdot x_1 + \ldots + w_p \cdot x_p = \mathbf{w}^T \mathbf{x}$  //  $\forall_{\mathbf{x} \in D} : \mathbf{x}|_{x_0} \equiv 1$
- 5.  $error = c(\mathbf{x}) y(\mathbf{x})$
- 6.  $\Delta \mathbf{w} = \eta \cdot \textit{error} \cdot \mathbf{x}$
- 7.  $\mathbf{w} = \mathbf{w} + \Delta \mathbf{w}$
- 8.  $\mathbf{UNTIL}(\mathbf{convergence}(D, y(D)))$
- 9.  $return((w_0, w_1, \dots, w_p))$

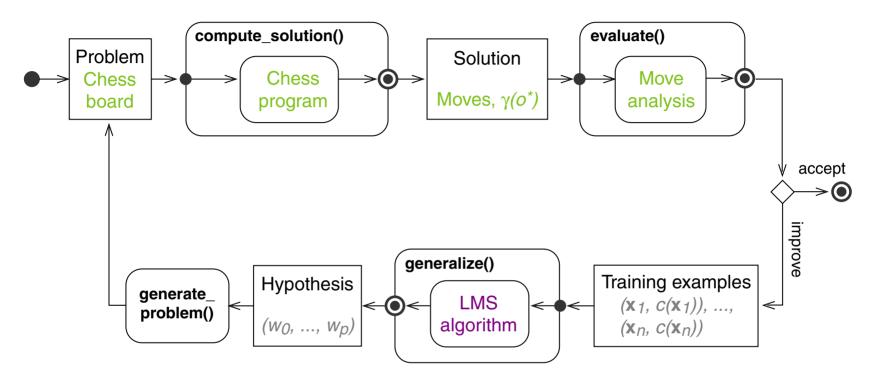
#### Remarks:

- □ The LMS weight adaptation corresponds to the incremental gradient descend [IGD] algorithm, and approximates the global direction of steepest error descent as used by the batch gradient descent [BGD] algorithm, for which more rigorous statements on convergence are possible.
- The *convergence* function may compute the global error quantified as the sum of the squared residuals,  $\sum_{(\mathbf{x},c(\mathbf{x}))\in D} (c(\mathbf{x})-y(\mathbf{x}))^2$ , or employ an upper bound on the number of iterations.

Design of Learning Systems [p.12, Mitchell 1997]



Design of Learning Systems [p.12, Mitchell 1997]



## Important design aspects:

- 1. kind of experience
- 2. fidelity of the model formation function  $\alpha: O \to X$
- 3. class or structure of the model function y
- 4. learning method for fitting y

ML:I-40 Introduction

### **Related Questions**

### Model functions *y*:

- What are important classes of model functions?
- □ What are methods to fit (= learn) model functions?
- What are measures to assess the goodness of fit?
- How does the example number affect the learning process?
- How does noise affect the learning process?

Related Questions (continued)

### Generic learnability:

- What are the theoretical limits of learnability?
- How can we use nature as a model for learning?

### Knowledge acquisition:

- How can we integrate background knowledge into the learning process?
- □ How can we integrate human expertise into the learning process?