

# Chapter ML:I

## I. Introduction

- ❑ Examples of Learning Tasks
- ❑ Specification of Learning Tasks
- ❑ Elements of Machine Learning
- ❑ Notation Overview
- ❑ Classification Approaches Overview

# Examples of Learning Tasks

## (1) Car Shopping Guide



?

What criteria form the basis of a decision?

# Examples of Learning Tasks

## (2) 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
...	

# Examples of Learning Tasks

## (2) Risk Analysis for Credit Approval (continued)

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:

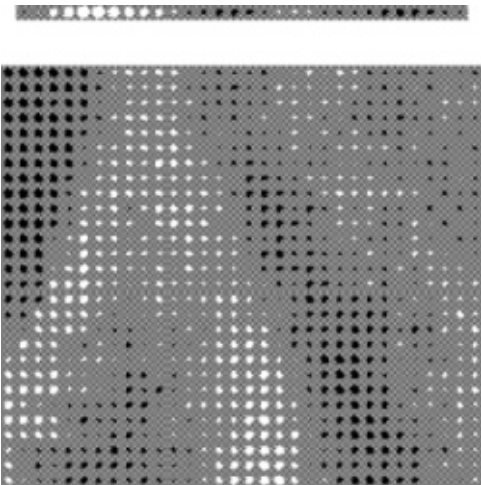
**IF** ( income>40 000 **AND** credit\_period<3 ) **OR** house\_owner=yes  
**THEN** credit\_approval=yes

**IF** SCHUFA\_entry=yes **OR** ( income<20 000 **AND** repayment>800 )  
**THEN** credit\_approval=no

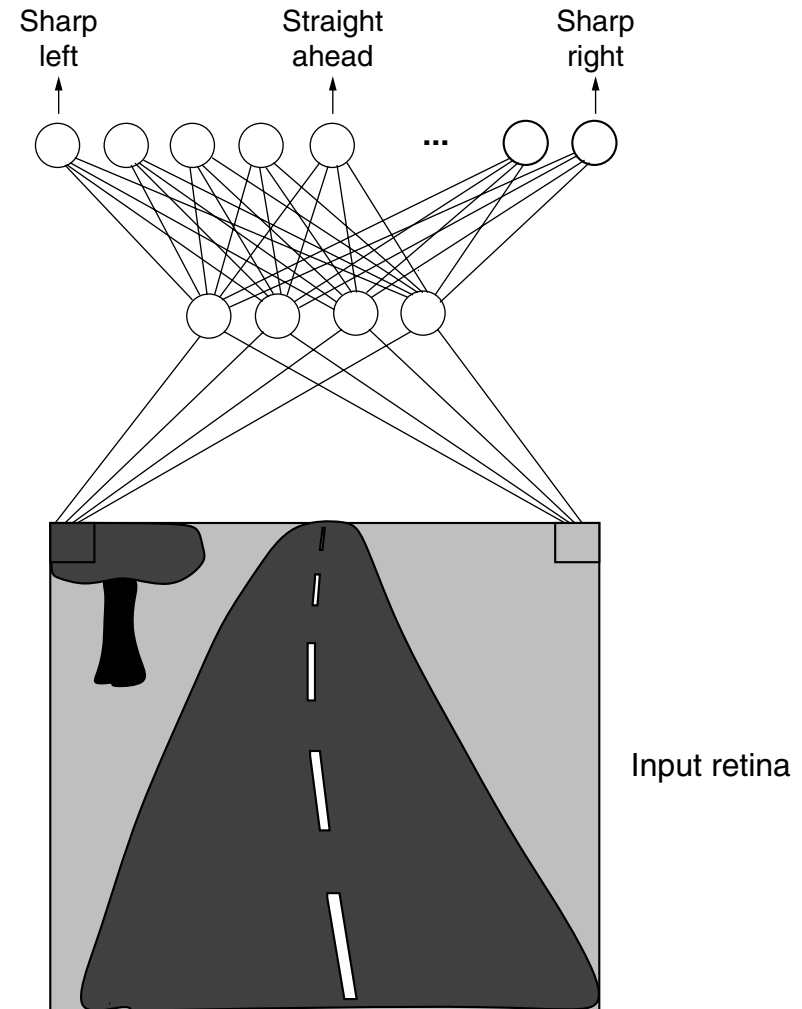
...

# Examples of Learning Tasks

## (3) Image Analysis [Mitchell 1997]

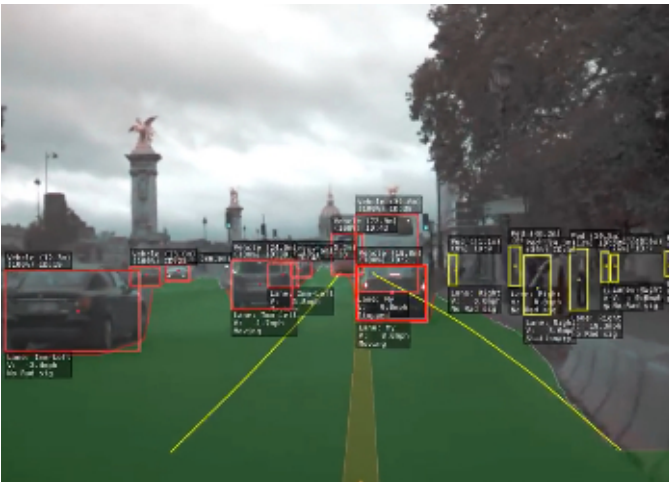


[1992]

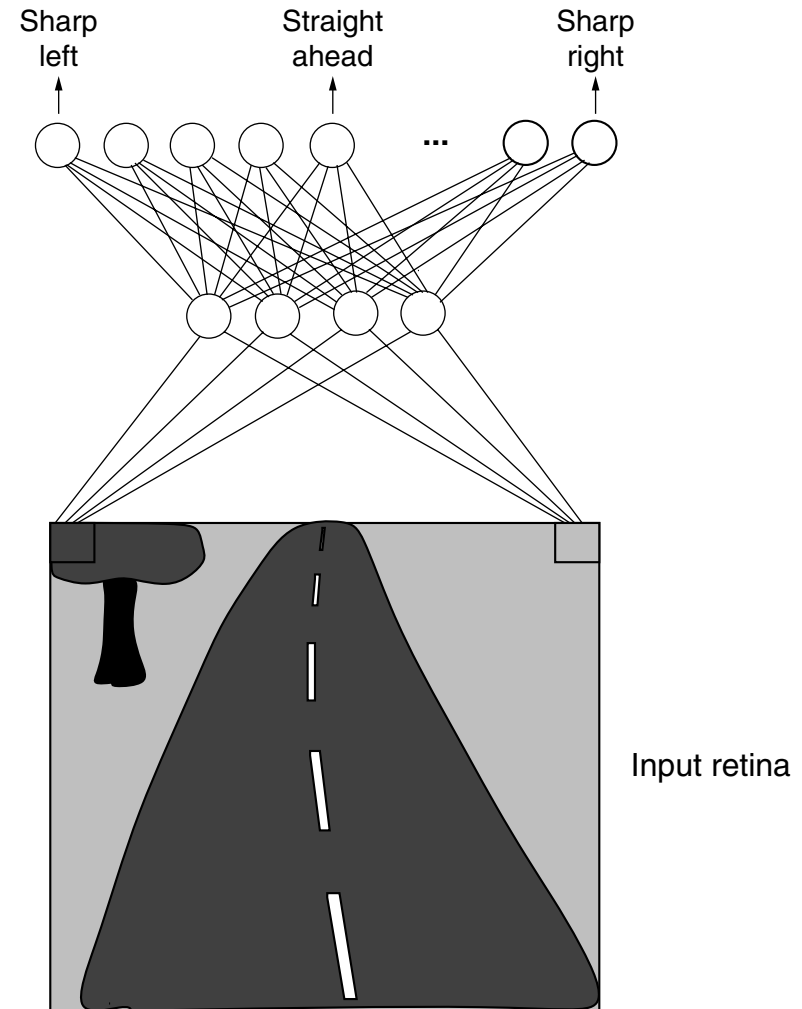


# Examples of Learning Tasks

## (3) Image Analysis (continued)



[2018]



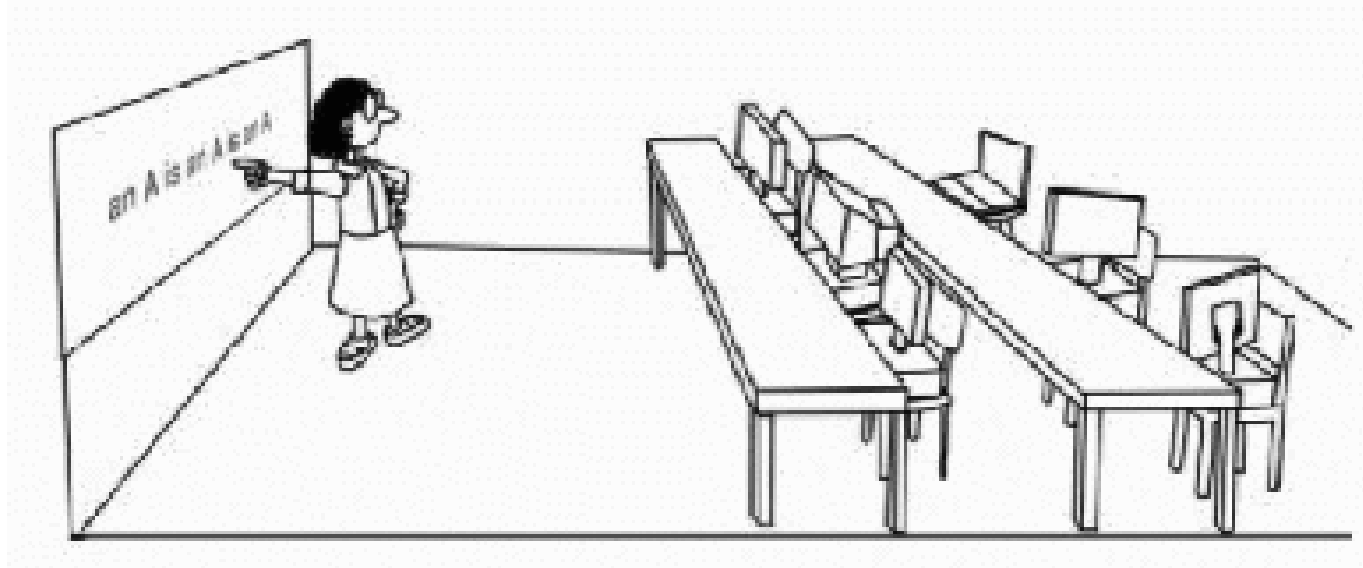
# Specification of Learning Tasks

## 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 do self play
- ❑ 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](http://webis.de/data.html)]



# Specification of Learning Tasks

## 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).

Example: intrusion detection in a network data stream

### 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 in a hostile environment.

# Specification of Learning Tasks

## 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).

Example: intrusion detection in a network data stream

### 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 in a hostile environment.

# Specification of Learning Tasks

## (4) 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.

# Specification of Learning Tasks

## (4) Example Chess: Kinds of Experience (continued)

### 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 problem instances.

# Specification of Learning Tasks

## (4) Example Chess: Kinds of Experience (continued)

### 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 problem instances.

### 3. Relevance under a performance measure

- How far can we get with experience?
- Can we master situations in the wild?  
(self play may not be enough to become world class)

# Specification of Learning Tasks

## (4) Example Chess: Ideal Target Function $\gamma()$ [Mitchell 1997]

a)  $\gamma : \text{Boards} \rightarrow \text{Moves}$

b)  $\gamma : \text{Boards} \rightarrow \mathbb{R}$

# Specification of Learning Tasks

## (4) Example Chess: Ideal Target Function $\gamma()$ (continued)

a)  $\gamma : \text{Boards} \rightarrow \text{Moves}$

b)  $\gamma : \text{Boards} \rightarrow \mathbb{R}$

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

Let be  $o \in \text{Boards}$ .

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. [\[Course on Search Algorithms\]](#)

# Specification of Learning Tasks

## (4) Example Chess: Real World $\rightarrow$ Model World

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

“model function”  
“model formation function”



# Specification of Learning Tasks

## (4) Example Chess: Real World $\rightarrow$ Model World (continued)

$$\gamma(o) \rightarrow y(\alpha(o)) \equiv y(\mathbf{x}) \quad \begin{array}{l} \text{“model function”} \\ \text{“model formation function”} \end{array}$$

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

where

$x_{1, 2}$  = number of black / white pawns on board  $o$

$x_{3, 4}$  = number of black / white pieces on board  $o$

$x_{5, 6}$  = number of black / white pieces threatened on board  $o$

$D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ , a set of board descriptions  $\mathbf{x}_i$  with scores  $y_i, y_i \in [-100; 100]$ .

# Specification of Learning Tasks

## (4) Example Chess: Real World $\rightarrow$ Model World (continued)

$$\gamma(o) \rightarrow y(\alpha(o)) \equiv y(\mathbf{x}) \quad \text{"model function"}$$

"model formation function"

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

where

$x_{1,2}$  = number of black / white pawns on board  $o$

$x_{3,4}$  = number of black / white pieces on board  $o$

$x_{5,6}$  = number of black / white pieces threatened on board  $o$

$D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ , a set of board descriptions  $\mathbf{x}_i$  with scores  $y_i, y_i \in [-100; 100]$ .

Other approaches to specify a model function:

- ❑ case base
- ❑ set of rules
- ❑ neural network
- ❑ higher order polynomial of board features

## Remarks:

- ❑  $\gamma : O \rightarrow \mathbf{R}, o \mapsto \gamma(o)$ , is called *ideal target function*. The ideal target function interprets the real world, say, a real-world object  $o$ , to “compute”  $\gamma(o)$ . This “computation” can be done 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 consider a model world. This model world is restricted to particular and typically easily measurable features  $\mathbf{x}$ , which are derived from  $o$ , with  $\mathbf{x} = \alpha(o)$ . In the model world,  $y(\mathbf{x})$  is the abstracted and formalized counterpart of  $\gamma(o)$ .
- ❑  $y : \mathbf{X} \rightarrow \mathbf{R}, \mathbf{x} \mapsto y(\mathbf{x})$ , is called *model function* or *model*. The choice and computation of a suitable model function is a central aspect in the field of machine learning.
- ❑  $\alpha : O \rightarrow \mathbf{X}, o \mapsto \alpha(o)$ , is called *model formation function*. The development of a suitable model formation function is often not treated as part of machine learning but “outsourced” to respective domain experts. However, tackling advanced learning tasks such as autonomous driving, automated debating, or playing chess requires a tight cooperation between the developers of  $\alpha()$  and  $y()$ .
- ❑ The 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 is needed; such kind of tasks form regression problems. If only a small set of values, here denoted as  $C$ , are to be considered (e.g. school grades), we are given a classification problem. A regression problem can be transformed into a classification problem by domain discretization.
- ❑ Regression problems and classification problems differ in the way how an achieved accuracy or goodness of fit is assessed. For regression problems the sum of the squared residuals, RSS, may be a sensible criterion; for classification problems the number of misclassified examples may be more relevant. Keywords: *regression loss* versus *classification loss*
- ❑ For classification problems, the ideal target function,  $\gamma : \mathcal{O} \rightarrow C$ , is also called *ideal classifier*; analogously, the model function,  $y : \mathbf{X} \rightarrow C$ , is called classifier.
- ❑ Decision problems are classification problems with two classes.
- ❑ The halting problem for Turing machines is an undecidable classification problem.

# Specification of Learning Tasks [model world]

Real World  $\rightarrow$  Model World

Setting of the **real world**:

- $O$  is a set of objects. (example: emails)
- $C$  is a set of classes. (example: spam versus ham)
- $\gamma : O \rightarrow C$  is the ideal classifier. ( $\gamma$  is a human expert)

Classification task:

- Given some  $o \in O$ , determine its class  $\gamma(o) \in C$ . (example: is an email spam?)

# Specification of Learning Tasks [model world]

Real World  $\rightarrow$  Model World (continued)

Setting of the **real world**:

- $O$  is a set of objects. (example: emails)
- $C$  is a set of classes. (example: spam versus ham)
- $\gamma : O \rightarrow C$  is the ideal classifier. ( $\gamma$  is a human expert)

Classification task:

- Given some  $o \in O$ , determine its class  $\gamma(o) \in C$ . (example: is an email spam?)

Acquisition of classification knowledge  $D$ :

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

# Specification of Learning Tasks [real world]

Real World  $\rightarrow$  Model World (continued)

Setting of the **model world**:

- $X$  is a multiset of feature vectors. (example: word frequency vectors)
- $C$  is a set of classes. (as before: spam versus ham)
- $\underline{D} = \{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_n, c_n)\} \subseteq X \times C$  is a multiset of examples.

Learning task:

- Approximate  $D$  by a model function  $y()$ .

# Specification of Learning Tasks [real world]

Real World  $\rightarrow$  Model World (continued)

Setting of the **model world**:

- $X$  is a multiset of feature vectors. (example: word frequency vectors)
- $C$  is a set of classes. (as before: spam versus ham)
- $D = \{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_n, c_n)\} \subseteq X \times C$  is a multiset of examples.

Learning task:

- Approximate  $D$  by a model function  $y()$ .

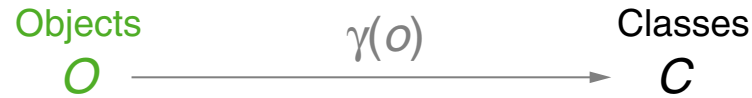
Machine learning:

1. Formulate a model function  $y : \mathbf{X} \rightarrow C$ ,  $\mathbf{x} \mapsto y(\mathbf{x})$  ( $y$  needs to be fitted)
2. Apply statistics, theory, and algorithms from the field of machine learning to maximize the goodness of fit between  $(\mathbf{x}, c)$  and  $(\mathbf{x}, y(\mathbf{x}))$ ,  $(\mathbf{x}, c) \in D$ .



# Specification of Learning Tasks

Real World  $\rightarrow$  Model World (continued)



Semantics:

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

# Specification of Learning Tasks

Real World  $\rightarrow$  Model World (continued)

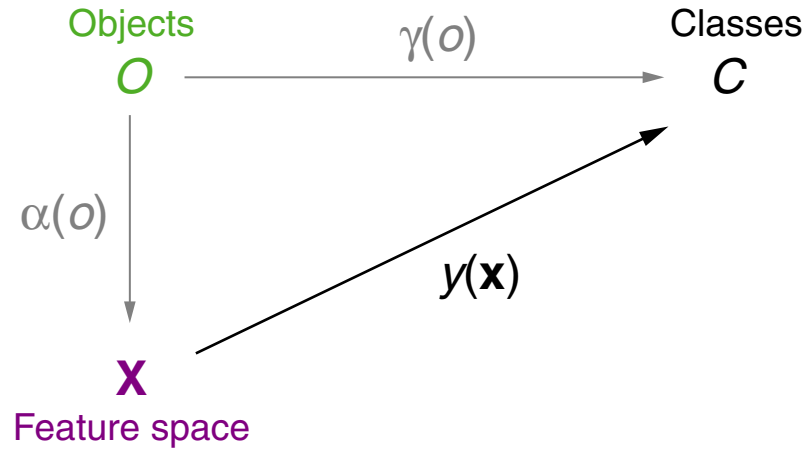


Semantics:

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

# Specification of Learning Tasks

Real World  $\rightarrow$  Model World (continued)



Semantics:

- $\gamma(o)$  Ideal classifier (a human) for real-world objects.
- $\alpha(o)$  Model formation function.
- $y(\mathbf{x})$  Classifier (model function) to be learned.

## Remarks:

- ❑ Feature vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots$ , are abstractions of real-world objects  $o_1, o_2, \dots$ . While the real-world objects can be considered as unique, the feature vectors are typically not since similar objects can be mapped onto the same feature vector. I.e, the set  $X$  actually is a multiset, comprised of feature vectors resulting from the analysis of a sequence of unique real-world objects.
- ❑ To be distinguished from the multiset  $X$  of feature vectors is the space  $\mathbf{X}$ , which is the Cartesian product of the domains of the different features (dimensions) of a feature vector  $\mathbf{x}$ . In the [chess example](#),  $\mathbf{X}$  is the  $\mathbb{N}^6$ , and, in the upcoming regression settings,  $\mathbf{X}$  is the  $\mathbb{R}^p$ .  
With respect to the support of the multiset  $X$  holds  $\text{Supp}(X) \subseteq \mathbf{X}$ , although typically  $\text{Supp}(X) \subset \mathbf{X}$  with  $|\text{Supp}(X)| \ll |\mathbf{X}|$ .
- ❑  $\mathbf{X}$  is called input space, feature space, feature domain, domain, or similar. The [structure of the feature space](#) (consider a pre-Hilbert space, a  $\sigma$ -algebra, or even no structure) is an important determinant of the learning problem.
- ❑ Heads-up: We will distinguish carefully between the two terms input space and feature space if an explicit feature transformation step is done such as [basis expansion](#) (regression), [multilayer perceptron processing](#) (neural networks), or non-linear kernelization (support vector machines).

## Remarks (continued) :

- ❑ 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. A suitable “representation”  $\mathbf{x}$  for an object  $o$  can also be algorithmically computed, with embedding approaches for example. Keyword: *representation learning*
- ❑ The value  $c$  in an example  $(\mathbf{x}, c)$ ,  $(\mathbf{x}, c) \in D$ , is termed “target value”, “ground truth”, or “observation” (for  $\mathbf{x}$  and the underlying classification problem). Observe that these terms are justified by the fact that  $c := \gamma(o)$ .
- ❑ Note that in the [chess example](#)  $\gamma()$  defines the scores  $y_i \in \mathbb{R}$  only for three types of boards (won, lost, drawn);  $\gamma()$  is unknown for all boards  $o$  that fall under Point (4) of the recursive definition. I.e., for most chess boards  $o$  we *cannot provide the ground truth*  $\gamma(o)$ , say, we can neither give a statement whether  $o$  leads to a final board state that is won or lost if both sides play optimally nor provide the next optimum move.

# Specification of Learning Tasks

## The LMS Algorithm for Fitting $y(\mathbf{x})$ [algorithms: LMS BGD <sub>$\sigma$</sub> PT BGD IGD]

Algorithm: LMS Least Mean Squares.

Input:  $D$  Training examples  $(\mathbf{x}, c)$  with  $\mathbf{x} \in \mathbb{R}^p$  and target class  $c$ .

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

Output:  $\mathbf{w}$  Weight vector from  $\mathbb{R}^{p+1}$ . (= hypothesis)

$\text{LMS}(D, \eta)$

1. *initialize\_random\_weights* $((w_0, w_1, \dots, w_p))$ ,  $t = 0$
2. **REPEAT**
3.    $t = t + 1$
4.    $(\mathbf{x}, c) = \text{random\_select}(D)$
5.    $y(\mathbf{x}) = w_0 + w_1 \cdot x_1 + \dots + w_p \cdot x_p \stackrel{(*)}{=} \mathbf{w}^T \mathbf{x}$
6.    $\delta = c - y(\mathbf{x})$
7.    $\Delta \mathbf{w} \stackrel{(*)}{=} \eta \cdot \delta \cdot \mathbf{x}$    //  $-\delta \cdot \mathbf{x}$  is the derivative of loss  $l_2(c, y(\mathbf{x}))$  wrt.  $\mathbf{w}$ .
8.    $\mathbf{w} = \mathbf{w} + \Delta \mathbf{w}$
9. **UNTIL** $(\text{convergence}(D, y(\cdot), t))$
10. *return* $((w_0, w_1, \dots, w_p))$

# Specification of Learning Tasks

## The LMS Algorithm for Fitting $y(\mathbf{x})$ [algorithms: LMS BGD <sub>$\sigma$</sub> PT BGD IGD]

Algorithm: LMS Least Mean Squares.

Input:  $D$  Training examples  $(\mathbf{x}, c)$  with  $\mathbf{x} \in \mathbb{R}^p$  and target class  $c$ .

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

Output:  $\mathbf{w}$  Weight vector from  $\mathbb{R}^{p+1}$ . (= hypothesis)

$\text{LMS}(D, \eta)$

1. *initialize\_random\_weights* $((w_0, w_1, \dots, w_p))$ ,  $t = 0$
2. **REPEAT**
3.  $t = t + 1$
4.  $(\mathbf{x}, c) = \text{random\_select}(D)$
5.  Model function evaluation.
6.  Calculation of residual.
7.  Calculation of derivative.
8.  Parameter vector update = one gradient step down.
9. **UNTIL** $(\text{convergence}(D, y(\cdot), t))$
10. *return* $((w_0, w_1, \dots, w_p))$

## Remarks:

- (★) We consider the feature vector  $\mathbf{x}$  in its extended form when used as operand in a scalar product with the weight vector,  $\mathbf{w}^T \mathbf{x}$ , and consequently, when noted as argument of the model function,  $y(\mathbf{x})$ . I.e.,  $\mathbf{x} = (1, x_1, \dots, x_p)^T \in \mathbf{R}^{p+1}$ , and  $x_0 = 1$ .
- Line 7: The derivative of the pointwise squared loss  $l_2(c, y(\mathbf{x}))$  is used to update the weight vector. See section Gradient Descent of part Neural Networks for a derivation of the update term.
- Line 9: The function *convergence*( $\cdot$ ) can analyze the global loss quantified as sum of squared residuals (RSS),  $\sum_{(\mathbf{x}, c) \in D} (c - y(\mathbf{x}))^2$ , or the norm of the loss gradient,  $\|\nabla \text{RSS}\|$ , and compare it to a small positive bound  $\varepsilon$ . Consider in this regard the vectors of observed and computed classes,  $D|_c$  and  $y(D|_x)$  respectively. In addition, the function may check via  $t$  an upper bound on the number of iterations.



## Remarks (continued):

- ❑ The LMS algorithm approximates the solution of a least squares problem. It is a specialization of the SGD algorithm (stochastic gradient descent)—more specifically, an SGD that uses a squared loss (cost) function.

The SGD algorithm, in turn, is a stochastic approximation of gradient descent optimization; it replaces the true gradient (calculated from the entire data set) by an estimate thereof (calculated from a randomly selected subset of the data).

- ❑ The LMS algorithm is hence very similar to the IGD algorithm (incremental gradient descent) equipped with a squared loss (cost) function. They differ in their selection of examples:

[LMS] Line 4:  $(\mathbf{x}, c) = \text{random\_select}(D)$

[IGD] Line 4: **FOREACH**  $(\mathbf{x}, c) \in D$  **DO**

The IGD algorithm approximates the global direction of steepest loss descent as used by the BGD algorithm (batch gradient descent), for which more rigorous statements on convergence are possible.

- ❑ The LMS algorithm was proposed in 1959 by Bernard Widrow and Ted Hoff.