# A Short  Introduction to Machine Learning

*Introduction to Machine Learning*
## Lect.s 2 and 3

# Alessio Micheli

**micheli@di.unipi.it**

Dipartimento di Informatica
Università di Pisa - Italy

**Computational Intelligence &
Machine Learning Group**

# About ML

- **Machine Learning (ML)**
  - *Master Programme in Computer Science*
  - *Master Programme in Data Science and Business Informatics*
  - *Master Programme in Digital Humanities*

- **Code: 654AA Credits (ECTS): 9 Semester: 1**

- **Lecturer: Alessio Micheli:** micheli@di.unipi.it

# Practical information

In class:

- Please, max **silence** during the lecture (to avoid noise) *recording in progress!* (of course, you can make questions)

Connect to ML:

- **Material: Moodle https://elearning.di.unipi.it/**
- **Streaming & recordings of lectures: Teams platform**
  - See lecture 1 and Moodle: «FAQ and general information»
  - The enrolling students mechanism for attendance "in presence" (and **to connect to Teams**) is through the App "Didactic Agenda" ("Agenda Didattica") for 654AA 24/25
  - Please, remember to fill the poll (see INTRO-curricula22)

# Introduction to ML: plan of the next lectures

- Introduction aims:
  - Critical contextualization of the ML in comp. science [lect 1 and 2]

  - **Overview  and <u>Terminologies</u>** [lect 2, 3, 4]
    - the relevant concepts will be developed later in the course

  - First basic models and learning algorithms  [lect  5, 6, 7]

  - Then, we will start with Neural Networks!

*See the "Course structure" slide!*

# Learning

*The problem of **learning** is arguably at the very core of the problem of **intelligence**, both biological and artificial*

Poggio, Shelton, *AI Magazine 1999*

i.e. Learning as a major challenge and a strategic way  to provide intelligence into the systems

# Machine Learning (I)

We restrict to the *computational* framework:

- Principles, methods and algorithms for learning and prediction:
  - Learning by the system of the experience (known data) to approach a <u>defined computational task</u>
  - Build a model  (or hypothesis) to be used for predictions
    - ❖ (see examples on email-spam or face recognition)

*Most common specific framework :*

- Infer a model / *function* from a set of examples which allows the **generalization** (to provide accurate response on new data)

# Machine Learning (II): When?

*Opportunity* (if useful) and *awareness* (needs and limits)

- Utility of predictive models: (in the following cases)
    - **no (or poor) theory** (or knowledge to explain the phenomenon)
    - **uncertain, noisy or incomplete data** (which hinder formalization of solutions)
- Requests:
    - source of training experience (representative data)
    - tolerance on the precision of results

# Machine Learning (III): When?

- Models to solve real-world problems that are difficult to be treated with traditional techniques (*complementary* to analytical models based on previous knowledge, algorithms and imperative programming, classical AI, …)

- Examples of appropriate applications versus standard programming:

  - Knowledge is too difficult (to be formalized by 'hand-made' algorithm)
    - e.g. face recognition: humans can do it but cannot describe how they do it
    - e.g. voice automatic telephone answering service

  - Not enough human knowledge
    - e.g., predicting binding strength of molecules to proteins

  - Personalized behavior
    - scoring email messages or web pages according to user preferences
    - individualized (intelligent) human-computer interfaces

- Due to this flexibility ML applicative area is very large:  see lecture 1

# General challenges

- Build autonomous Intelligent/learning systems:
  - Robotics, HRI, search engines, …

- Build powerful tools for emerging challenges in intelligent data analysis
  - Tools for  the "data scientist"

- Open new areas of  applications in CS: innovative interdisciplinary open problems (more in general, "machine learning scientist")
  - Fantasy is your limitation !
  - ML in the era of  "changing of paradigm in science, in which scientific advances are becoming more and more *data-driven*"
  - *Growing data sources opens up a huge application area for ML and related areas (Web, Social Net., IoT, BioMed, …)*

# An useful framework:
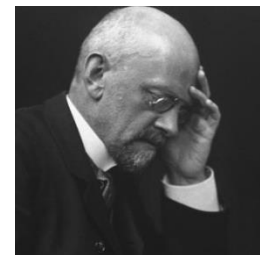# Learning as an approximation of an unknown function from examples

**Specific vision but widespread in ML**
### For us:

- **Different tasks seen in uniform framework**
- **Enables a rigorous formulation**
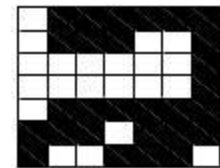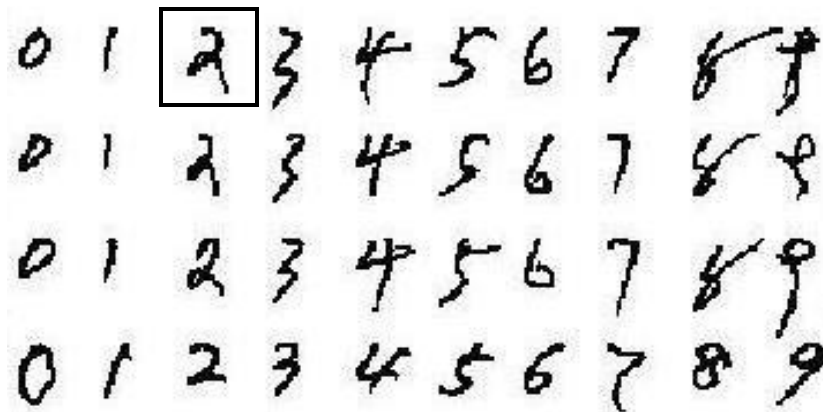
→ **Intro guided by intuitive examples**

Hilbert spaces

**Please, note that the following example was already introduced in Lect 1**

# An Example

- A pilot example: *recognition of handwritten digits*

- Input: collection of images of handwritten digits (arrays/matrix of values)
- Problem: build model that receives in input an image of handwritten digit and "predict" the digits

8 x 8

0  1  2  3  4  5  6  7  8  9

**Image
8 x 8**

$f$

**Output class**

# Handwritten Digits Recognition

Image
8 x 8

$f$

Output class

0 1 2 3 4 5 6 7 8 9
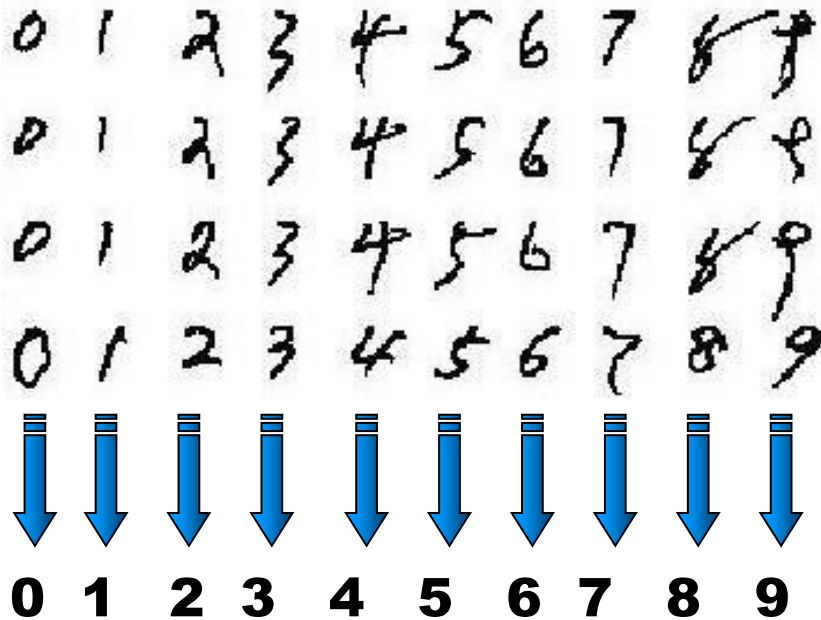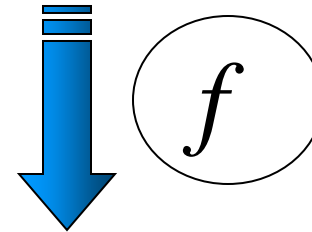
**Classification problem**

- Difficult to formalize exactly the solution of the problem:
  Possible presence of noise and ambiguous data;
- Relatively easy to collect a set of labeled examples

  => Example of successful application of the ML!

# **Machine Learning**

A new extended definition (looking to the pilot example)

- The ML studies and proposes methods to build (infer) dependencies / **functions** / hypotheses from examples of observed data
    - that *fits* the know examples
    - able to *generalize*, with reasonable accuracy for new data
        - According to verifiable results
        - Under statistical and computational conditions and criteria
    - Considering the expressiveness and algorithmic complexity of the models and learning algorithms

# Examples of $x$ - $f(x)$

*Inferring general functions from know data:*

- Handwriting Recognition
  - *x*: Data from pen motion.
  - *f(x)*: Letter of the alphabet.
- Disease diagnosis (from database of past medical records)
  - *x*: Properties of patient (symptoms, lab tests)
  - *f(x)*: Disease (or maybe, recommended therapy)
  - TR Training Set: *<x,f(x)>*: database of past medical records
- Face recognition
  - *x*: Bitmap picture of person's face
  - *f(x)*: Name of the person.
- Spam Detection
  - *x*: Email message
  - *f(x)*: Spam or not spam.

# **Complex data**

- Protein folding

    - *x*: sequence of amino acids

    - *f(x)*: sequence of atoms' 3D coordinates

    - TR <*x,f(x)*>: known proteins

    - Type of *x*: string (variable length)

    - Type of *f(x)*: sequence of 3D vectors

- Drug design

    - *x*: a molecule

    - *f(x)*: binding strength to HIV protease

    - TR <*x,f(x)*>: molecules already tested

    - Type of **x**: a graph or a relational description of atoms/chemical bonds

    - Type of *f(x)*: a real number

# Overview of a ML (predictive) System

Build or improve the agent/model/hypothesis by learning from data (world observations)

**DATA** → **MODEL** → **Prediction**

world observations

**TASK**

**LEARNING ALG.**

**VALIDATION**

Drive the model building by tuning the system parameters to the problem at hand

Also as a guide to the **key design choices** **(**ML system **"ingredients")**

# DATA

- The data *represent* the available facts (*experience*).
  - Representation problem: to capture the structure of the analyzed objects

**Types**: Flat, Structured, …

- **Flat** (*attribute-value language*):

fixed-size vectors of properties (*features*), single table of tuple (measurements of the objects)

| Fruits | Weight | Cost $ | Color | Bio |
|--------|--------|--------|-------|-----|
| Fruit 1 (lemon) | 2.1 | 0.5 | y | 1 |
| Fruit 2 (apple) | 3.5 | 0.6 | r | ? |

→ Attributes (categorical/discrete or continuous)

→ missing data

Data can be subject to
*preprocessing*: e.g. Variable scaling, encoding*, feature selection…

# Examples and terminologies

**Medical records**     $i$

| Patients | Age | Smoke | Sex | Lab Test |
|----------|-----|-------|-----|----------|
| Pat 1 | 101 | 0.8 | M | 1 |
| Pat 2 | 30 | 0.0 | F | ? |

$p$ (Pat 2 row) → $x_p$

→ Attributes (discrete/continuous)

- Each row ($x$, vector in bold): example, pattern, instance, sample,….

- Dimension of data set: number of examples $l$

- Dimension (of the input $x$): number of features $n$

- If we will index the features/inputs/variables by $i$ or $j$ : variable $x_i$ is (typically) the *i-th* feature/property/attribute/element/component of $x$.

(but may be to simplify we need to use subscript index for other meanings)

- $x_p$ (or $x_i$) is (typically) the *p-th* (or *i-th)* pattern/example/raw (vector)

- $x_{p,i}$ (for example) can be the attribute $i$ of the pattern $p$

A. Micheli

# DATA Encoding

## Flat case:

- **Numerical encoding for categories: e.g.**
  - 0/1 (or −1/+1) for 2 classes
  - More classes:
    - 1,2,3… Warning: grade of similarity (1 vs 2 or 3): useful for "order categorical" variables (e.g small, medium, large)
    - *1-of-k* (or *1-hot*) encoding: useful for symbols

| A | **1** | 0 | 0 |
|---|---|---|---|
| B | 0 | **1** | 0 |
| C | 0 | 0 | **1** |

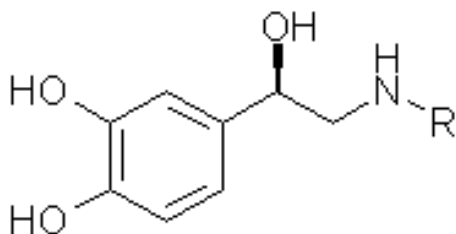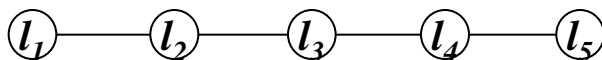**It will be useful for the project !**

Useful both for input or output variables
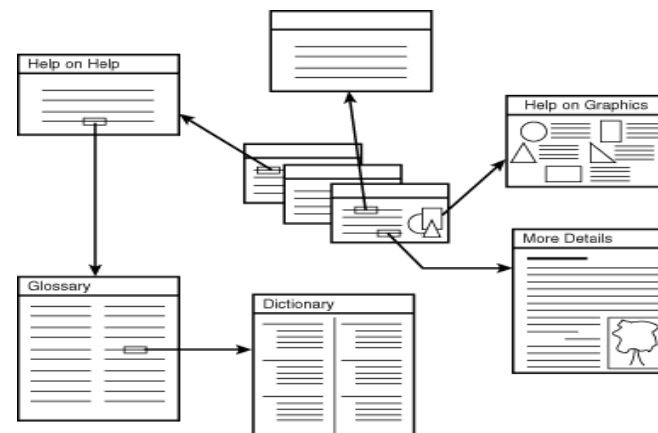
# DATA : Structures

- **Structured:** Sequences (lists), trees, graphs, Multi-relational data (table) (in DB)

Examples: images, microarray, temporal data, strings of a language, DNA e proteins, hierarchical relationships, molecules, hyperlink connectivity in web pages, ...

Which natural representation?

Graph/network data

- **Noise**: addition of external factors to the stream of (target) information (*signal*); due to randomness in the measurements, not due to the underlying law: e.g. Gaussian noise

SIGNAL + NOISE    SIGNAL

NOISE

- **Outliers**: are unusual data values that are not consistent with most observations (e.g. due to abnormal measurements errors)
  - outlier detection – preprocessing: removal
  - Robust modeling methods

- **Feature selection**: selection of a small number of informative features: it can provide an optimal input representation for a learning problem

# TASKS

- The task defines the purpose of the application:
  - Knowledge that we want to achieve? (e.g. pattern in DM or model in ML)
  - Which is the helpful nature of the result?
  - What information are available?

Mainly in the ML course

- **Predictive** (Classification, Regression): function approximation

$$x$$

Input space

$$f$$

Categories o real values $(R)$

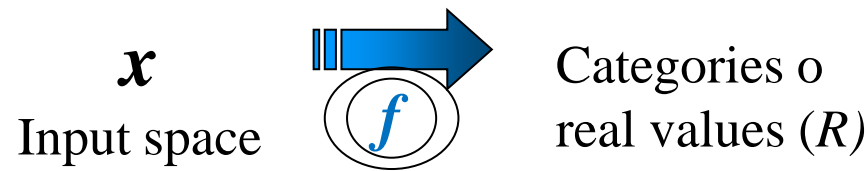E.g. recall the *"pilot" example on handwritten digits:* Build a *function* from examples

- **Descriptive** (Cluster Analysis, Association Rules): find subsets or groups of unclassified data

# Tasks: Supervised Learning

- <u>Given</u>: Training examples as $<input,output>=<\boldsymbol{x},d>$ (**labeled** examples)

  *Def* for an unknown function $f$ (known only at the given points of example)

  - Target value: desiderate value $d$ or $t$ or $y$ ... is given by the teacher according to $f(\boldsymbol{x})$ to label the data

- <u>Find</u>: A *good* approximation to $f$ (a <u>hypothesis</u> $h$ that can used for prediction on unseen data $\boldsymbol{x'}$, i.e. that is able to generalize)

$$\boldsymbol{x}$$
$$f$$

Input space     Categories o real values (*R*)

- Target $d$ *(or t or y)*: a categorical or numerical *label*

  - ***Classification:*** discrete value outputs:

    $$f(\boldsymbol{x}) \in \{1,2,...,K\} \quad classes \quad (discrete\text{-}valued\ function)$$

  - ***Regression:*** real continuous output values (approximate a real-valued target function, in $R$ or $R^K$)

> Unified vision thanks to the formalism of a **function approximation** task

A. Micheli

28

## Unsupervised Learning: No teacher!

- TR  (Training Set)= set of  unlabeled data $<x>$

- E.g. to find *natural groupings* in a set of data
  - Clustering
  - Dimensionality reduction/ Visualization/Preprocessing
  - Modeling the data density



◎ Centroids

■ Clustering:

Partition of data into clusters (subsets of "similar" data)

# Tasks: Classification

(Supervised) Classification: Patterns (features vectors) are seen as members of a class and the goal is to assign the patterns observed classes (label)

- *Classification*: $f(\boldsymbol{x})$ return the correct class for $\boldsymbol{x}$
- Number of classes:

    - **=2** : $f(\boldsymbol{x})$ is a Boolean function: binary classification, **concept learning**  (T/F or 0/1 or −1/+1 or negative/positive),

    - **> 2**:  multi-class problem ($C_1, C_2, C_3 \ldots C_K$)

# Example

From DATA to TASK (e.g. classification)

| Patients | Age | Smoke | Sex | Lab Test |
|----------|-----|-------|-----|----------|
| Pat 1 | 101 | 0.8 | M | 1 |
| Pat 2 | 30 | 0.0 | F | ? |

| Target: diagnose |
|------------------|
| + |
| - |

$f$

$x$ : Input space

Terminology in statistics:
• Inputs are the "independent variables"
• Outputs are the "dependent variables" or "responses"

A. Micheli

31

The classification may be viewed as the allocation of the input space in decision regions (e.g. **0/1**)

<u>Example</u>: graphical illustration of a linear separator on a instance space $x^T=(x_1,x_2)$ in $IR^2$, $f(x)=0/1$ (or $-1/+1$)

Point belonging to class 1

Separating (hyper)plane : $x$ s.t.

*PREVIEW*

$$w^T x + w_0 = w_1 x_1 + w_2 x_2 + w_0 = 0$$

$$h(x) = \begin{cases} 1 & \_if \quad w^T x + w_0 \geq 0 \\ 0 & _____ otherwise \end{cases}$$

or

$$h(x) = sign(w^T x + w_0)$$

Linear threshold unit (LTU)

Indicator functions

A. Micheli

How many? (H): set of dichotomies induced by hyperplanes

The 0/1 classification function in 3D
(on a 2D input space)

Region where the
output of the
classifier is 1

# Tasks: Regression: example

- Process of estimating of a real-value function  on the basis of finite set of noisy samples (supervised task)
  - known pairs $(x, f(x)+random\ noise)$

  Task (exercise): find $f$ for the data in the following table:

| $x$ | $target$ |
|-----|----------|
| 1 | 2.1 |
| 2 | 3.9 |
| 3 | 6.1 |
| 4 | 8.4 |
| 5 | 9.8 |
| ... | ... |

Via Neural Network ?
or by …
Guessing   $f(x)=2x$

Small errors at the points!

- *Regression:* $x$ = variables (e.g. real values), $f(x)$ *real values*: curve fitting  ($x$ is 1-dim in the example but it becomes  *k-dim* in general)

- Process of estimating of a real-value function  on the basis of finite set of noisy samples

  - known pairs $(x, f(x)+random\ noise)$



Point where we know the value of  $f(x)$
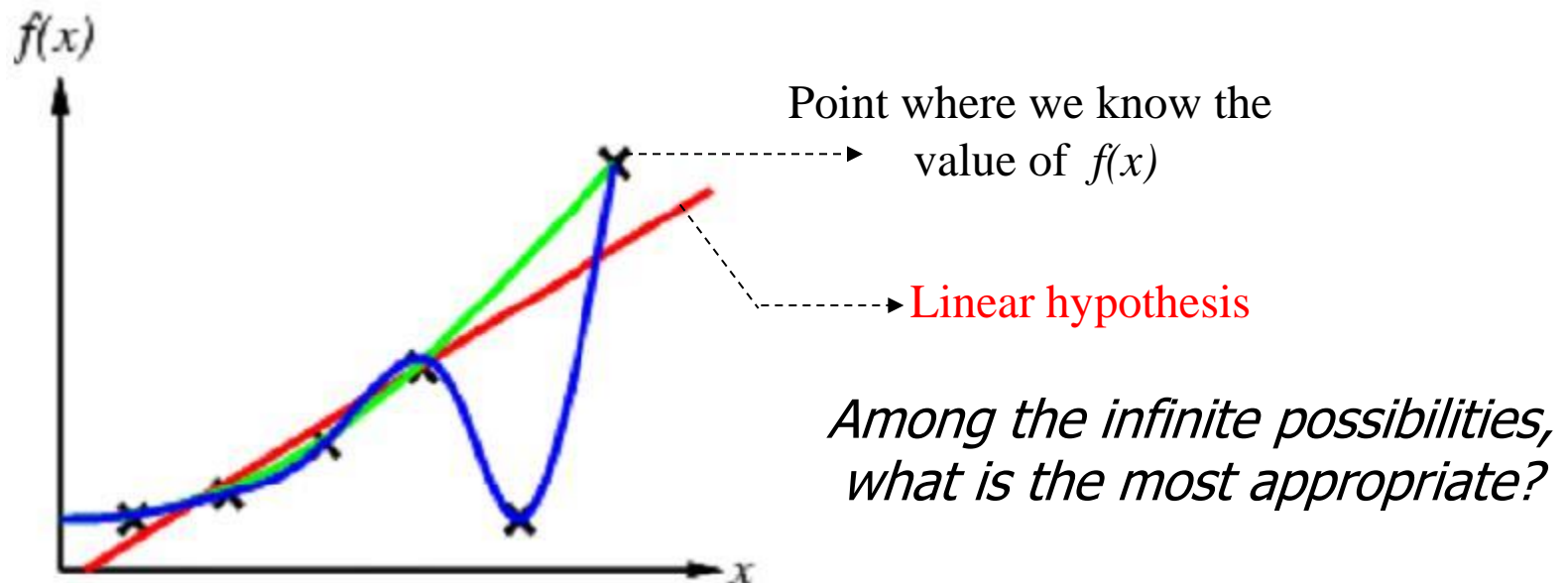
Linear hypothesis

*Among the infinite possibilities, what is the most appropriate?*

An example (linear hypothesis): $h_w(x)=w_1 x+w_0=0.2\ x\ -0.4$

- **Semi-supervised learning**
  - combines both labeled and unlabeled examples to generate an appropriate function or classifier.

- **Reinforcement Learning** (learning with right/wrong critic).
  - Adaptation in *autonomous systems*
  - "the algorithm learns a *policy* of how to act given an observation of the world. Every action has some impact in the environment, and the environment provides feedback that guides the learning algorithm".
  - Not step by step examples
  - Toward decision-making aims
  - Useful in modern AI



Reinforcement Learning Setup

# Models
# and survey of useful concepts

- MODEL:
  - Aim: to capture/describes the relationships among the data (on the basis of the task) by a "language" (numerical, symbolic, …)
  - The "language" is related to the *representation* used to get knowledge
  - The model defines the class of functions that the learning machine can implement (*hypotheses space*)
    - E.g. set of functions $h(\boldsymbol{x},\boldsymbol{w})$, where $\boldsymbol{w}$ is the (abstract) parameter

- **Training example** (superv.): An example of the form ($\boldsymbol{x}$, $f(\boldsymbol{x})+noise$)
  $\boldsymbol{x}$ is usually an input vector of features, *(d or t or) $y=f(\boldsymbol{x})+noise$* is called the target value
- **Target function**: The true function $f$
- **Hypothesis**: A proposed function $h$ believed to be similar to $f$. An expression in a given *language* that describes the relationships among data
- **Hypotheses space** H: The space of all hypotheses (specific models) that can, in principle be output by the learning algorithm

A. Micheli

# Models:
# few trivial examples....

Just to have a preview of different *representation* of hypothesis
(because you already know the *language* of equations, logic, probability):

- **Linear models** (representation of H defines a **continuously** parameterized space of potential hypothesis);

  each assignment of $w$ is a different hypothesis, e.g:

  – $h(\mathbf{x}) = sign(\mathbf{w}^T \mathbf{x} + w_0)$

    *binary classifier*

    $h_{\mathbf{w}}(x) = w_1 x + w_0$     *E.g.* $h_{\mathbf{w}}(x) = 2x + 150$

    *simple linear regression*

- **Symbolic Rules**: (hypothesis space is based on **discrete** representations); different rules are possible , e.g:

  – *if $(x_1=0)$ and $(x_2=1)$ then $h(\mathbf{x})=1$*

    *binary classifier*

  – *else $h(\mathbf{x})=0$*

- **Probabilistic models**: estimate $p(x,y)$

- K Nearest neighbor regression: Predict mean y value of nearest neighbors (memory-based)

A. Micheli

42

# Neural Networks (just a look)

An example: we will see a **neural networks**, beyond the *neurobiological inspiration*, as a computational model for the treatment of data, capable of approximating complex (*non-linear*) relationships between inputs and outputs

$x$

Input space

Categories or $IR$ *(real)* values

$f$

Age

Smoke

Alcool

Again,
a class of functions !!!

# Paradigms and methods (Languages for H)

- Symbolic and Rule-based (or *discrete* H)
  - <u>Conjuction of literals</u>*, Decision trees (propositional rules)
  - Inductive grammars, Evolutionary algorithms, …
  - Inductive Logic Programming (first order logic rules)
- Sub-symbolic (or *continuous* H)
  - Linear discriminant analysis, <u>Multiple Linear Regression</u>*, <u>LTU</u>
  - <u>Neural networks</u>
  - Kernel methods (<u>SVMs</u>, gaussian kernels, spectral kernels, etc)
- Probabilistic/Generative
  - <u>Traditional parametric models</u> (density estimation, discriminant analysis, polynomial regression,…)
  - <u>Graphical models</u>: <u>Bayesian networks</u>, <u>Naïve Bayes</u>, PLSA, Markov models, Hidden Markov models, …
- Instance-based
  - <u>Nearest neighbor</u>*

Note:  Underlined – >ML
1. Some models can be expressed by different languages
2. * Next lectures

# How many models?

- Theory (*No Free Lunch Theorem*) : *there is no universal "best" learning method* (*without any knowledge, for any problems,...*):

   *if an algorithm achieves superior results on some problems, it must pay with inferiority on other problems. In this sense there is no free lunch.*

   E.g. Devroye (1982), Wolpert and Macready (1997), and others

- → The course provides a
   - set of models and the
   - critical instruments to compare them

- However, not all the models are equivalent:
   - Important differences are for the **flexibility** of the approaches, toward models that can in principle approximate arbitrary functions (e.g. no just linear approximation seen in the examples)
   - Important differences are for the **control of the complexity** (we will see later)
   - Use of flexible models and principia for the control of the complexity are the core of ML

A. Micheli

- LEARNING ALGORITHM      Basing on data, task and model

- (Heuristic) *search through the hypothesis space* **H** *of the* **best hypothesis**
    - i.e. the best approximation to the (unknown) target function
    - Typically searching for the *h* with the *minimum "error"*

    - E.g. free parameters of the model are fitted to the task at hand:
    - Examples: best *w* in linear models, best rules for symbolic models, ….
    - Remember the regression example, we proposed $h(x)=2x$, for $h_w(x)=w_1x+w_0$ assuming $w_1=2$ and $w_0=0$ as the best parameter value: *how?*

- **H** may not coincide with the set of all possible functions and the search can not be exhaustive: we need to make assumptions → (we will see the role of) *Inductive bias*

Initial solution

Set of solutions chosen a priori

Hypotheses space **H**

Each point represents a different hypothesis (function)

Optimal solution

(minimum "error")

Set of solutions compatible with training set

Typically local search approaches

# Learning (terminologies)

According to the different paradigms/contexts "learning" can be differently termed or have different acceptations:

- Inference (statistics)
- Inference: Abduction/Induction (logic)
- Adapting (biology, systems)
- Optimizing (mathematics)
- Training  (e.g. Neural Networks)
- Function approximations  (mathematics)

Can be more specifically found in other sub-fields:

- Regression analysis (statistics), curve fitting (math, CS), …
- Or using other terminologies e.g. "*Fitting* a multivariate function"

# Recap and next topics

After the introduction of the first four ingredients (Data, Task, Model and Learning Alg.), we need to focus on three mentioned *relevant concepts* not yet discussed so far:

1. The *inductive bias* (examples in discrete hypothesis spaces)
2. The *loss,* used to measure the quality of our approximation
3. The concept of *generalization* and *validation* (next lecture)

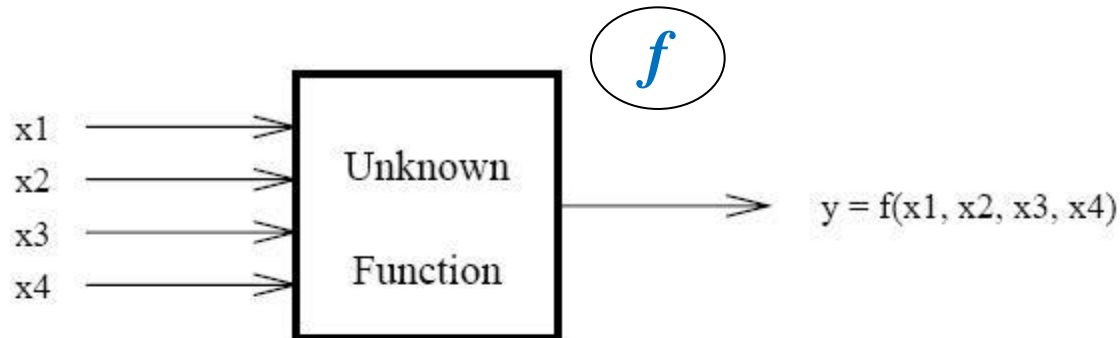A. Micheli

# 1. The Role of the Inductive Bias

In order to set up a model and a learning algorithm we can make *assumptions* (about the nature of the target function) concerning either

- Constraints in the model (in the hypothesis space H, due to the set of hypotheses that we can express or consider) (**Language Bias**)

- Constraints or preferences in learning algorithm/search strategy (**Search Bias**)

- Or **Both**.

- We will see that such assumptions are strictly need to obtain an useful model for the ML aims, i.e. a model with generalization capabilities

- We start to discuss it within examples in *discrete hypotheses spaces (rules)*, **learning a concept** (a Boolean function) [Mitchell chapt. 2]

  - E.g. $x$ is a "cat" if $h_{cat}(x) = 1$, otherwise is 0 for $x$ in "animals"

# An example: Learning Boolean functions

$f$

x1
x2
x3
x4

Unknown Function

$y = f(x1, x2, x3, x4)$

Find the function s.t.

| Example | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|---------|-------|-------|-------|-------|-----|
| 1 | 0 | 0 | 1 | 0 | 0 |
| 2 | 0 | 1 | 0 | 0 | 0 |
| 3 | 0 | 0 | 1 | 1 | 1 |
| 4 | 1 | 0 | 0 | 1 | 1 |
| 5 | 0 | 1 | 1 | 0 | 0 |
| 6 | 1 | 1 | 0 | 0 | 0 |
| 7 | 0 | 1 | 0 | 1 | 0 |

This is an **ill posed** (inverse) problem:
We may violate either existence, **uniqueness**, stability of the solution or solutions

Table 1

A. Micheli

59

- There are $2^{16} = 2^{2^4} = 65536$ possible Boolean functions over four input features. We can not figure out which one is correct until we have seen every possible input-output pair.

- After 7 examples, we still have $2^9$ possibilities.

- In the general case, in this discrete hypothesis space H: $|H| = 2^{\#\text{-input-instances}} = 2^{2^n}$

    for binary inputs/outputs, *n= input dimension*

    **Lookup table** *model* $\longrightarrow$

- I.e. a rote learner: Store/memorize examples, classify **x** if and only if it matches a previously observed example (else "no answer").

    – No inductive bias $\rightarrow$ *no generalization!*

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | ? |
| 0 | 0 | 0 | 1 | ? |
| 0 | 0 | 1 | 0 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 1 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | ? |
| 1 | 0 | 0 | 0 | ? |
| 1 | 0 | 0 | 1 | 1 |
| 1 | 0 | 1 | 0 | ? |
| 1 | 0 | 1 | 1 | ? |
| 1 | 1 | 0 | 0 | 0 |
| 1 | 1 | 0 | 1 | ? |
| 1 | 1 | 1 | 0 | ? |
| 1 | 1 | 1 | 1 | ? |

# Another discrete H space: Conjutive rules

- As second example of discrete H, we can image to learn a discrete function with discrete inputs assuming **conjunctive rules** (propositions with AND among literals, a language bias)

- i.e. using a *language bias* to work with a restricted hypothesis space

- E.g. $h_1 = l_2$, $h_2 = (l_1 \text{ and } l_2)$, $h_3 = true$, $h_4 = not(l_1) \text{ and } l_2, \ldots$

  - Rules such as if $l_2(=true)$ then $h(x)=true$, else $h(x)=false$
    or equivalently $if (x_2=1)$ then $h(x)=1$, else $h(x)=0$

- With $n$ binary inputs we had $|H| = 2^{\text{\#-input-instances}} = 2^{2^n}$

- With only conjunctive rules:

  #semantically distinct hypotheses (conjunctions):

  $3^n$ (for each of the $n$ positions we can have $l_i$, $not(l_i)$, $don't\ care$) + 1
       (+1 because all h with ($l_i\ AND\ not(l_i)$) are equivalent to "*false*")

  (e.g. from 65536 to just $3^4+1=82$ in the example with $n=4$)

# Find the Version Space

- Given the def.: a hypothesis h is **consistent** with the TR, if *h(**x**)=d(**x**)* for each training example *<**x**,d(**x**)>* in TR.

- It  is possible to perform a *complete search* (finding the set of *all* h consistent with the TR set) *in an efficient way* in this reduced space (of conjunctive rules) by cleverer algorithms  (Mitchell chap. 2)
  - Instead of searching enumerating all the possible combination of literals, i.e. every h in H

- We are only interested to say that these algorithms find the VS:
- Call the **version space**, $VS_{H,TR}$ ,  with respect to hypothesis space H, and training set TR, the *subset of hypotheses* from H *consistent* with all training examples

- Hence, this conjunctive assumption for H leads to an efficient solution in finding a VS.
  However, using only conjunctive rules may be **too restrictive**: if the target concept is not in H, it *cannot be* represented in H.

  - e.g. *if ($x_1$=1) **or** ($x_2$=1) then h($\boldsymbol{x}$)=1, else h($\boldsymbol{x}$)=0*

- **Idea**: Choose H that expresses every teachable concept (among propositions), that means H is the set of all possible subsets of X (*instance* or *input* space):  the power set P(X)

- E.g. *n*=10 binary inputs |X|= $2^{10}$=1024, |P(X)|=$2^{1024}$ ~ $10^{308}$ distinct concepts (much more than the  num. of atoms in the universe)

- H = disjunctions, conjunctions, negations

- H surely contains the target concept.

- *What for generalitazion ?*

# Unbiased Learner II (formal)

Recall that the **version space**, $VS_{H,TR}$, with respect to hypothesis space H, and training set TR, is the subset of hypotheses from H consistent with all training examples

The only examples that are unambiguolsy classified by an **unbiased learner** represented with the VS are the training examples themselves I.e. the *lookup table* !

---

**Property:** An <u>unbiased learner is *unable to generalize*</u> (on new instances):
<u>*Proof*</u>: Each unobserved instance will be classified 1 (positive) by precisely half the hypothesis in VS and 0 (or negative) by the other half (*rejection:* no answer is made by the VS for new input instances).
Indeed:
$\forall$h consistent with $x_i$ (test), $\exists$ $h'$ identical to $h$ except $h'(x_i) <> h(x_i)$,
$h \in VS \rightarrow h' \in VS$ (because they are identical on TR)

---

# Futility of Bias-Free Learning

- A learner that makes no prior assumptions regarding the identity of the target function/concept has no rational basis for classifying any unseen instances.

- (Restriction, preference) bias *not* only assumed for *efficiency*, *it is needed for the generalization capability*
  - However, it does no tell us (quantify) which one is the best solution for generalization yet
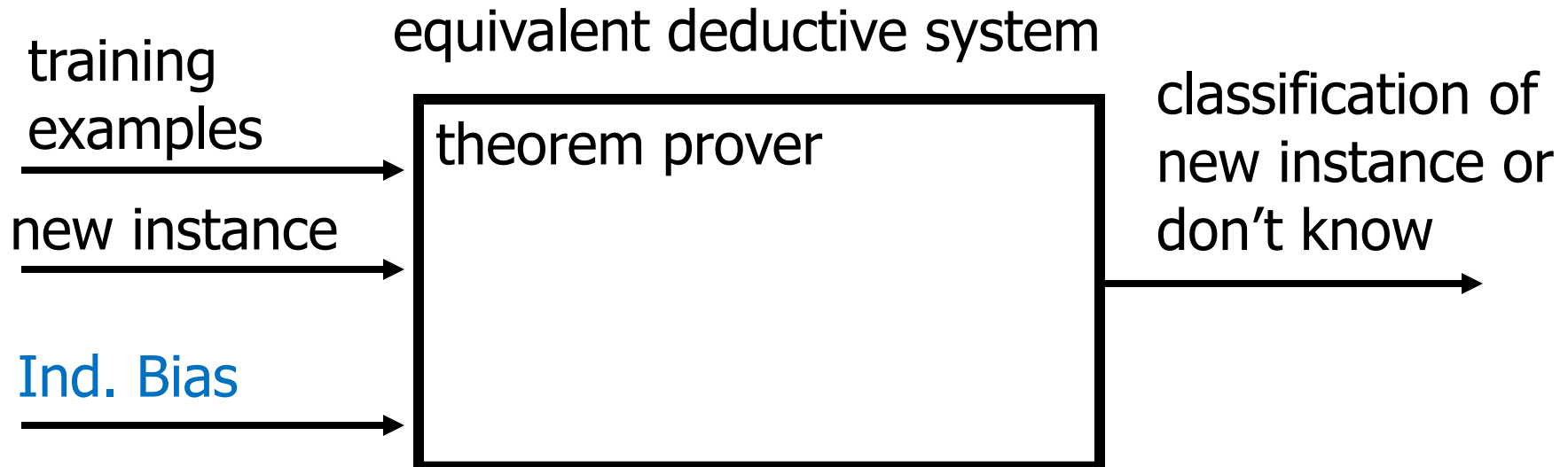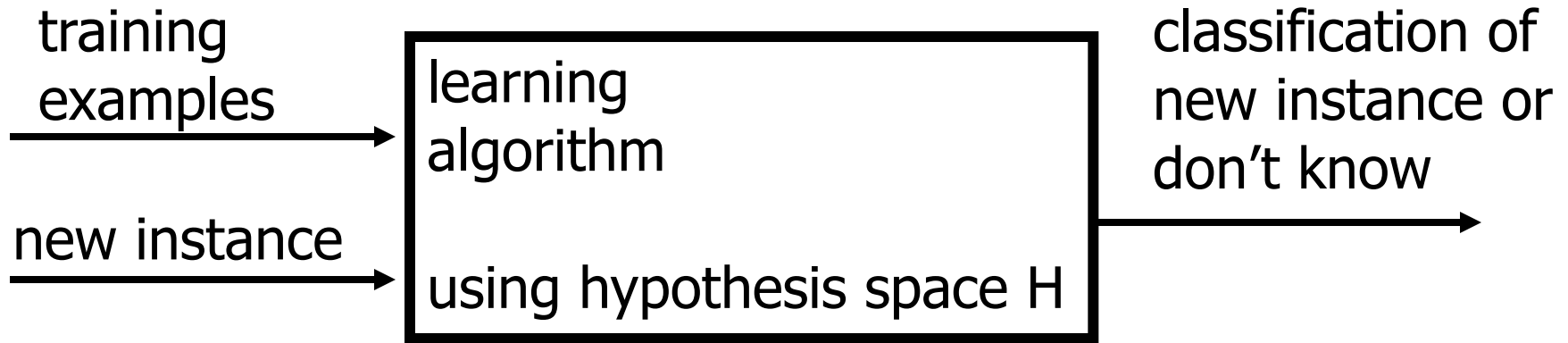
- **Trivial Example** (TR= Training Set, TS= Test Set): **:**

  X  *d(x)*    H={*x*, not(*x*), **0**, **1**}

TR   0    0        VS={*x*,**0**}

TS   1    ?   →   <span style="color:orange">Can be 1 or 0</span> ... Unless you use all X as TR set.

  In other words, in order to learn  the target concept,  one would have to present every single instance in X as a training example (lookup table)

# Inductive Systems and Equivalent Deductive Systems

training examples →

new instance →

**learning algorithm**

**using hypothesis space H**

→ classification of new instance or don't know

equivalent deductive system

training examples →

new instance →

Ind. Bias →

**theorem prover**

→ classification of new instance or don't know

# Language or search bias?

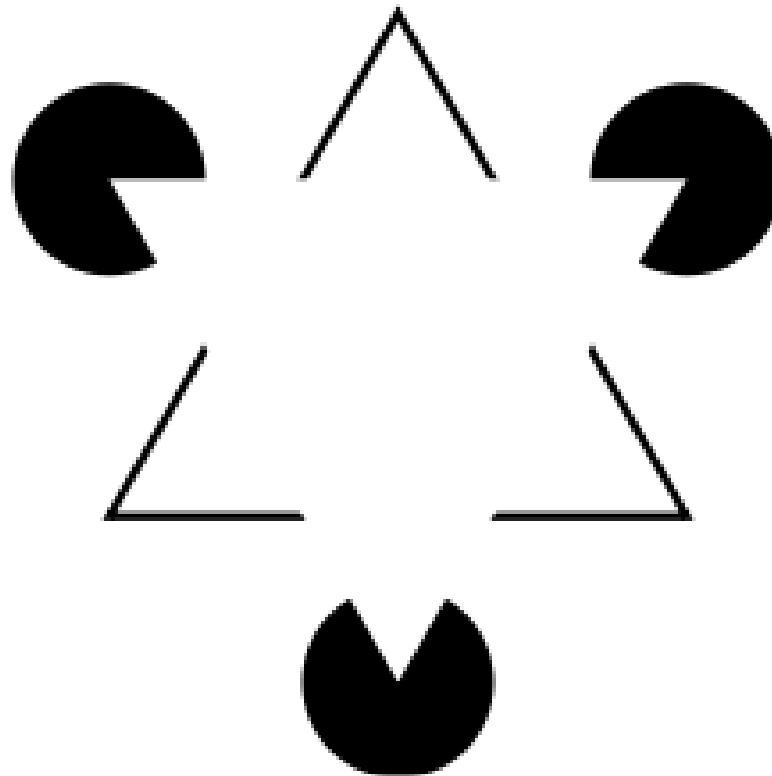Why the *search bias* can be preferred over the *language bias*?

- In ML typically use **flexible** *approaches* (expressive hypothesis spaces, universal capability of the models, e.g. Neural Networks, DT)

- avoiding the *language bias*, hence <u>without excluding a priori the unknown target function</u>,

- retaining an inductive bias but focusing on the *search bias* (which is ruled by the learning algorithm).

  - In practice using an *incomplete* search strategy.

**Conclusions:**

- Learning without bias cannot extract any regularities from data (lookup-table: no generalization capabilities)

- Every state-of-the-art ML approach shows an *inductive bias*

- Issue: characterize the bias for different models/learning approaches

Dip. Informatica
University of Pisa

We said ... A "**good**" approximation to $f$ from examples.

How to <u>measure</u> the quality of the approximation?

- Recall that we produce $h(\boldsymbol{x})$ value (output of the model for input $\boldsymbol{x}$)

- We want to measure the "distance" between $h(\boldsymbol{x})$ and $d$
  (objective function for minimization of errors in training, check of errors in test)

We use a ("inner") *loss function/measure*: $L(h_w(\boldsymbol{x}), d)$ (for a pattern $\boldsymbol{x}$)

e.g. high value → poor approximation

The *Error* (or *Risk or Loss*) is an expected value of this $L$

e.g. a "sum" or mean of the inner loss $L$ over the set of samples

$$Loss(h_{\boldsymbol{w}}) = E(\boldsymbol{w}) = \frac{1}{l} \sum_{p=1}^{l} L(h_{\boldsymbol{w}}(\boldsymbol{x}_p), d_p)$$

Note:
index $p$ is used for the
samples $p=1..l$

We will change $L$ for different tasks

Note: at moment Error, Risk and Loss are considered equivalent, we will specify
differences later through the course

A. Micheli

69

# **Tasks: Common Tasks review**

I will show a short survey of common learning _tasks_  by specifying the (changing of the) nature

- of the output and hypothesis space
- of the _loss function (in particular of  $L)$_,

_i.e._ Examples of loss functions: use it for future reference

# Regression

- Regression: *predicting a numerical value*

- **Output**: $d_p=f(\boldsymbol{x}_p) + e$   *(real value function + random error)*
- **H**:  a set of real-valued functions

- **Loss function** $L$ : measures the approximation accuracy/error
- A *common* loss function for regression: the squared error

$$L(h_{\boldsymbol{w}}(\boldsymbol{x}_p), d_p) = (d_p - h_{\boldsymbol{w}}(\boldsymbol{x}_p))^2$$

- The mean over the data set provide the *Mean Square Error (MSE)*
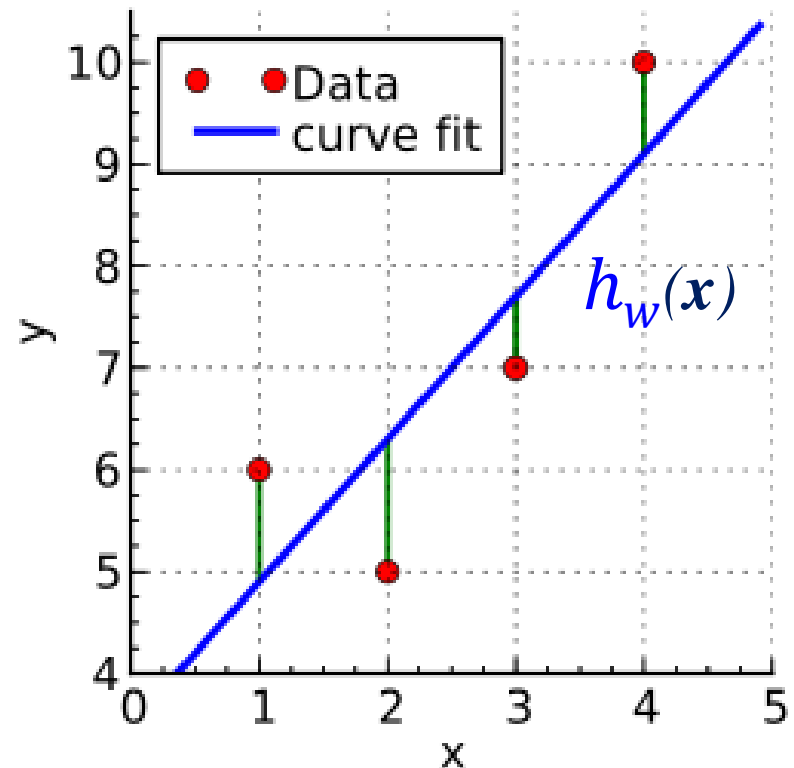
A. Micheli

In the example we have
$h(x)=w_1x+w_0$ as the blue line

and in green the errors at the data
points $(x_i\,y_i)$ (in red), where the
target $d_i$ for $x_i$ is denoted $y_i$ in the
example



$h_w(x)$

The Mean Square Error (MSE)

is the mean of the square of the
green errors:

$$E(w)= \frac{1}{l}\sum_{p=1}^{l}\left(y_p - h_w(x_p)\right)^2$$

$w$ are the free parameters of
the linear model

Note: this plot is taken elsewhere, I used
different colors before: here the line is in blue.
Also, the $y$ are therein the desidered (target $d$)
values

A. Micheli

72

# Classification

- Classification of *data into discrete classes*

- **Output**: e.g. {0,1}
- **H**: a set of indicator functions

- **Loss function** $L$ : measures the classification error

$$L(h_{\boldsymbol{w}}(\boldsymbol{x}_p), d_p) = \begin{cases} 0 & if \quad h_{\boldsymbol{w}}(\boldsymbol{x}_p) = d_p \\ 1 & otherwise \end{cases}$$

*0/1 Loss*

Def

- The mean over the data set provide the *number/percentage of misclassified patterns*
- *E.g. 20 out of 100 are misclassified → 20% errors, i.e. 80% of **accuracy***
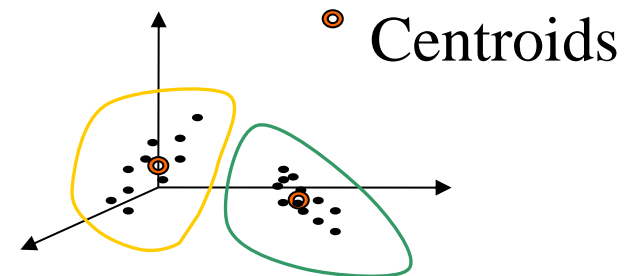
# Clustering and Vector Quantization*<sub>preview</sub>

- **Goal:** optimal partitioning of unknown distribution in **x**-space into regions (clusters) approximated by a cluster center or *prototype*.


° Centroids

- **H**: a set of vector quantizers $x \rightarrow c(x)$

  *continuos space → discrete space*

- Loss function $L$ : measures the vector quantizer optimality

- A *common* **loss function** *would be* the *squared error distortion*:

$$L(h(\mathbf{x}_p)) = (\mathbf{x}_p - h(\mathbf{x}_p)) \bullet (\mathbf{x}_p - h(\mathbf{x}_p))$$
$$\bullet = inner\_product$$

→ We'll see later

Proximity of the pattern to the centroid of its cluster

# Density estimation* *preview*

- Density estimation (generative, "parametric methods")
  from an assumed class of density

- **Output**: a density e.g. normal distribution with mean $m$ and variance $sigma^2$ : $p(x \mid m, sigma^2 )$

- **H**: a set of densities (e.g. $m$ and $sigma^2$ are the two unknown *parameters*)

- A *common* **loss function** $L$ for density estimation:

$$L(h(\mathbf{x}_p)) = -\ln( h(\mathbf{x}_p))$$

→ We'll see later

- Related to "maximizing the (log) likelihood function". [not hear]
- E.g. $P(x_1, x_2, x_3, \ldots \mid m, sigma^2 )$

# 3. Machine Learning & generalization

This is a fundamental concept of the course

- *Learning*: search for a *good function* in a function space from known data *(typically minimizing an Error/Loss)*

- Good w.r.t. generalization error: it measures how accurately the model predicts over novel samples of data (*Error/Loss measured over new data*)

Generalization: crucial point of ML!!!

Easy to **use** ML tools *versus* **correct/good use** of ML

# Generalization

- **Learning** phase **(training, fitting)**: build the model from know data  – *training data* (and bias)

- **Predictive** or **Test** phase (deployment/ *Inference* use of the ML built model): apply the model  to new examples:
  - we take the new inputs *x'* and we compute the response by the model
  - we compare with its target *d'* that the model has never seen
  - i.e. we make evaluation of the <span style="color:orange">generalization capability</span> of  our predictive hypothesis

Note: *performance*  in ML = generalization accuracy/ *predictive accuracy*
> estimated by the error computed on the (hold out) <span style="color:red">Test Set</span>

- **Theory**: E.g. Statistical Learning Theory [Vapnik] :

  - *under what (mathematical) conditions is a model  able to generalize?* → see next lecture (just basic notions)

# **Validation**

- Evaluation of performances for ML systems =
    Generalization/Predictive accuracy evaluation, i.e.:

- Validation !
- Validation !!

# • Validation !!!

- In the following (next lecture) we will discuss some **validation techniques**
    – to evaluate (model assessment) and
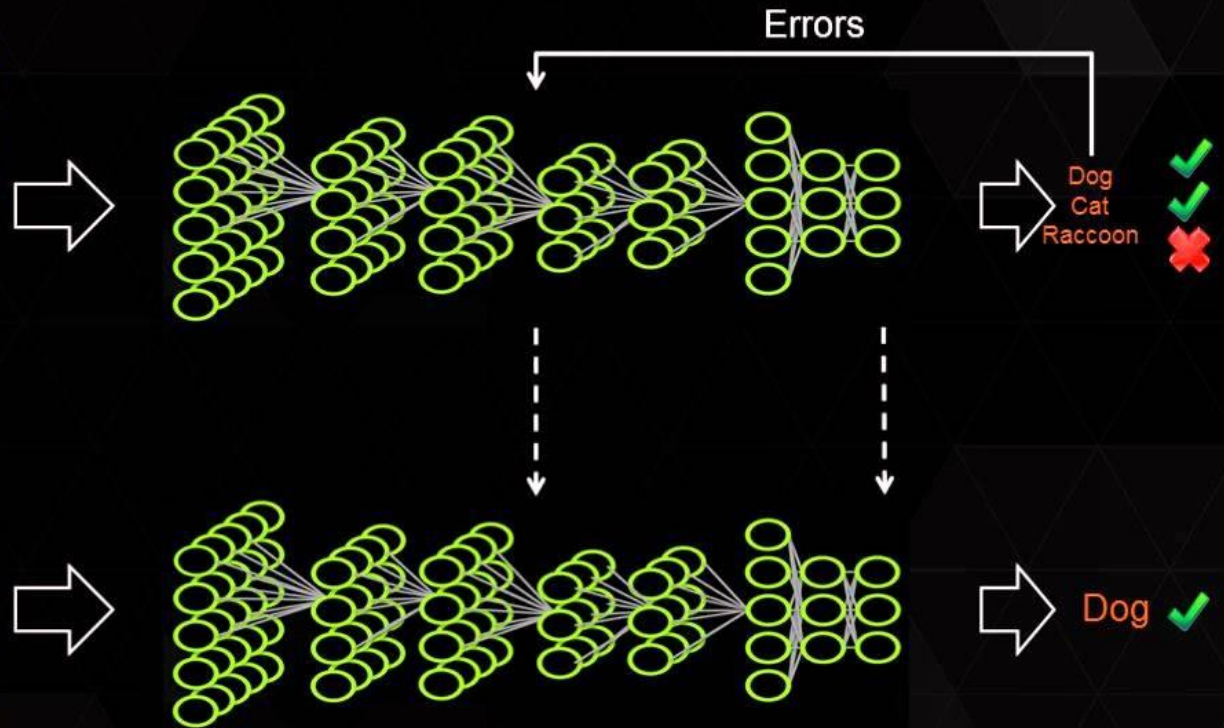    – to manage the generalization capability (model  selection).

# Exemplification of the Deployment/ *Inference* use

Even the inference part can be costly if you have millions of requests
(e.g. at google)

A Google server rack containing multiple **Tensor Processing Units**, a special-
purpose chip designed specifically for machine learning
The original TPU was designed specifically to work best with Google's TensorFlow.

**Just for inference (mapping) !!!!**

# Summary of the Intro to ML

- Part I (now)
  - Motivations, contextualization in CS
  - Course info

- Part II (in Lect.s 2 and 3)
  - Utility of ML
  - Learning as function approximation (pilot example)
  - Design components of a ML system, including
    - Learning tasks
    - Hypothesis space (and first overview)
    - Inductive bias (examples in discrete hypothesis spaces)
    - Loss and learning tasks
    - Generalization (first part)

- Part III (in Lect. 4)
  - Generalization and Validation

***Aim:** overview and terminology
before starting to study models and learning algorithms*

# For information

**Alessio Micheli**
**micheli@di.unipi.it**

**http://ciml.di.unipi.it**

Dipartimento di Informatica
Università di Pisa - Italy

**Computational Intelligence &**
**Machine Learning Group**