# LEARNING PROBLEM

Bùi Tiến Lên

01/09/2019

KHOA CÔNG NGHỆ THÔNG TIN
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN

# Contents

Learning
Components
A Simple
Learning Model
Type of
Learnings
Feasibility Of
Learning
Probability to the
rescue

## Notation

🧠

| symbol | meaning |
|---|---|
| $a, b, c, N \ldots$ | scalar number |
| $\boldsymbol{w}, \boldsymbol{v}, \boldsymbol{x}, \boldsymbol{y} \ldots$ | column vector |
| $\boldsymbol{X}, \boldsymbol{Y} \ldots$ | matrix |
| $\mathbb{R}$ | set of real numbers |
| $\mathbb{Z}$ | set of integer numbers |
| $\mathbb{N}$ | set of natural numbers |
| $\mathbb{R}^D$ | set of vectors |
| $\mathcal{X}, \mathcal{Y}, \ldots$ | set |
| $\mathcal{A}$ | algorithm |

| operator | meaning |
|---|---|
| $\boldsymbol{w}^\top$ | transpose |
| $\boldsymbol{XY}$ | matrix multiplication |
| $\boldsymbol{X}^{-1}$ | inverse |

3

# Learning Components

# Credit Approval

- Suppose that a bank receives thousands of credit card applications every day, and it wants to automate the process of evaluating them.
- Applicant information

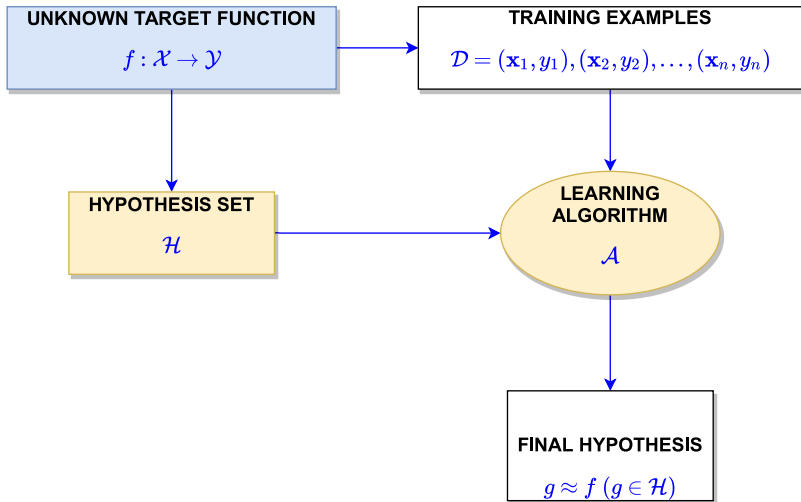| | |
|---|---|
| age | 23 years |
| gender | male |
| annual salary | $30000 |
| years in residence | 1 year |
| years in job | 1 year |
| current debt | $15000 |
| ... | ... |

- Approve credit?

## Problem Statement

Formalization

- Input: $\boldsymbol{x}$ (*customer application*)
- Output: $y$ (*good/bad customer?* or $\{1, -1\}$)
- Data $(\boldsymbol{x}_1, y_1), (\boldsymbol{x}_2, y_2), ...(\boldsymbol{x}_N, y_N)$ (*historical records*)
- Target function: $f : \mathcal{X} \to \mathcal{Y}$ (*ideal credit approval formula*)
- Best approximate function $g : \mathcal{X} \to \mathcal{Y}$ (*formula to be used*)

Learning
Components

A Simple
Learning Model

Type of
Learnings

Feasibility Of
Learning

Probability to the
rescue

# Components of Learning



UNKNOWN TARGET FUNCTION
$f : \mathcal{X} \to \mathcal{Y}$

TRAINING EXAMPLES
$\mathcal{D} = (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \ldots, (\mathbf{x}_n, y_n)$

HYPOTHESIS SET
$\mathcal{H}$

LEARNING
ALGORITHM
$\mathcal{A}$

FINAL HYPOTHESIS
$g \approx f \ (g \in \mathcal{H})$

**Learning
Components**

A Simple
Learning Model

Type of
Learnings

Feasibility Of
Learning

Probability to the
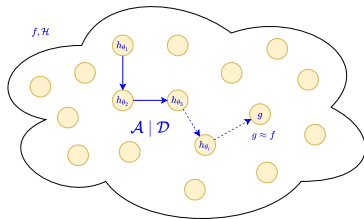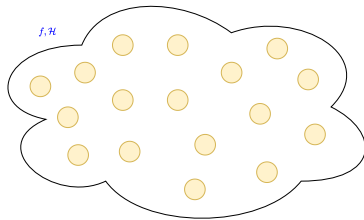rescue

## Solution components

The 2 solution components are referred as
the **learning model**

- The **hypothesis set** $\mathcal{H}$ built up
  from the problem

$$\mathcal{H} = \{h_{\theta_1}, h_{\theta_2}, ...\}$$

- The **learning algorithm** $\mathcal{A}$ is a
  **search algorithm** which finds
  $g \in \mathcal{H}$ such that

$$g \overset{best}{\approx} f$$

Learning
Components

**A Simple
Learning Model**

Type of
Learnings

Feasibility Of
Learning

Probability to the
rescue

# A Simple Hypothesis Set

We starts with the simple model (**the perceptron model**)

- For input $x = (x_1, ..., x_d)$ (*attributes of a customer*)

$$\text{Approve credit if} \sum_{i=1}^{d} w_i x_i \geq \textit{threshold}$$

$$\text{Deny credit if} \sum_{i=1} w_i x_i < \textit{threshold}$$

- This linear formula $h \in \mathcal{H}$ can be written as

$$h(x) = h_{\textbf{w}, threshold}(x) = sign\left(\left(\sum_{i=1}^{d} w_i x_i\right) - \textit{threshold}\right)$$

Learning
Components

**A Simple
Learning Model**

Type of
Learnings

Feasibility Of
Learning

Probability to the
rescue

## A Simple Hypothesis Set (cont.)

- Set $w_0 = -threshold$

$$h(x) = h_{\mathbf{w}}(x) = sign\left(\left(\sum_{i=1}^{d} w_i x_i\right) + w_0\right)$$

- Introduce an artfificial coordinate $x_0 = 1$

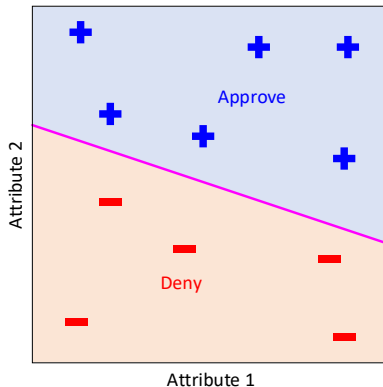$$h(x) = h_{\mathbf{w}}(x) = sign\left(\sum_{i=0}^{d} w_i x_i\right)$$

- In vector form, the perceptron implements

$$h(x) = h_{\mathbf{w}}(x) = sign\left(\mathbf{w}^\mathsf{T} \mathbf{x}\right)$$

Learning
Components

**A Simple
Learning Model**

Type of
Learnings

Feasibility Of
Learning
Probability to the
rescue

## 2D Model

- **Decision boundaries**: line
- **Decision regions**: approve and deny regions

Learning
Components

**A Simple
Learning Model**

Type of
Learnings

Feasibility Of
Learning

Probability to the
rescue

## A Simple Learning Algorithm

We uses the simple learning algorithm (**perceptron learning algorithm** - **PLA**) to implements

$$h(x) = h_{\mathbf{w}}(x) = sign\left(\mathbf{w}^\mathsf{T}\mathbf{x}\right)$$

- Given the training set

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), ...(\mathbf{x}_N, y_N)\}$$

- pick a *misclassified* point $(\mathbf{x}_i, y_i)$

$$sign(\mathbf{w}^\mathsf{T}\mathbf{x}_i) \neq y_i$$

- and update the weight vector

$$\mathbf{w} \leftarrow \mathbf{w} + y_n\mathbf{x}_n$$

Learning
Components

**A Simple
Learning Model**

Type of
Learnings

Feasibility Of
Learning
Probability to the
rescue

## Iterations of PLA

- At iteration $t = 1, 2, 3, ...$ pick a misclassified point from

$$\mathcal{D} = \{(\boldsymbol{x}_1, y_1), (\boldsymbol{x}_2, y_2), ...(\boldsymbol{x}_N, y_N)\}$$

  and run a PLA iteration on it
- That's it

Learning
Components

**A Simple
Learning Model**

Type of
Learnings

Feasibility Of
Learning
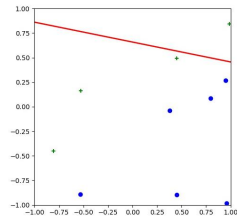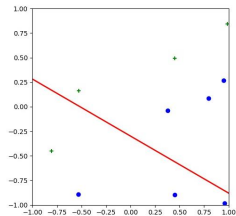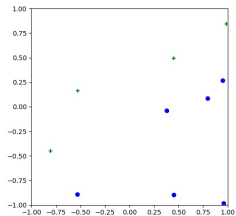Probability to the
rescue

## Is It Learning Algorithm?

Learning
Components

A Simple
Learning Model

Type of
Learnings

Feasibility Of
Learning
Probability to the
rescue

# A Learning Puzzle



y = -1

y = +1

y = ?

Learning
Components

A Simple
Learning Model

**Type of
Learnings**

Feasibility Of
Learning

Probability to the
rescue

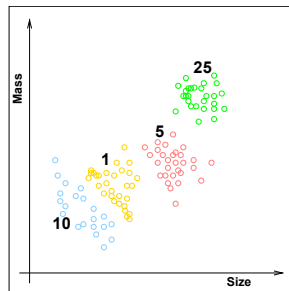# Basic Premise of Learning

"**using a set of observations to uncover an underlying process**"

broad premise $\implies$ many variations

- Supervise learning
- Unsupervised learning
- Reinforcement learning

Learning
Components

A Simple
Learning Model

Type of
Learnings

Feasibility Of
Learning

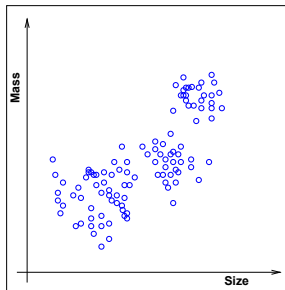Probability to the
rescue

## Supervised Learning

- We get data $\mathcal{D}$: (**input**, **correct ouput**)
  - When the **output** is one of *a finite set of values*, the learning problem is called **classification**
  - When the **output** is a *number*, the learning problem is called **regression**
- Example from vending machine - **coin classification**



19

Learning
Components

A Simple
Learning Model

Type of
Learnings

Feasibility Of
Learning
Probability to the
rescue

## Unsupervised Learning

- Instead of (**input**, **correct input**), we get (**input**, ?)

Learning
Components

A Simple
Learning Model

Type of
Learnings

Feasibility Of
Learning

Probability to the
rescue

# Reinforcement Learning

- Instead of (**input**, **correct input**),
  we get (**input**, *some* **ouput**, **grade** *for this output*)

Learning
Components

A Simple
Learning Model

Type of
Learnings

Feasibility Of
Learning

Probability to the
rescue

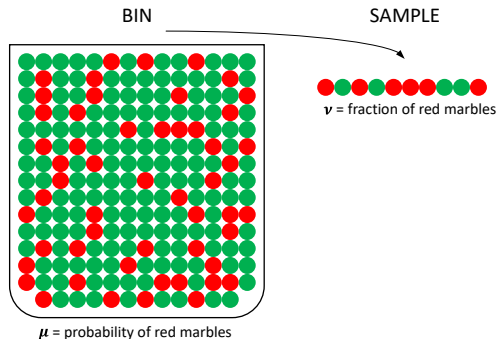# A Related Experiment - Bin Problem

- Consider a **BIN** with red and green marbles

    $P[\text{picking a red marble}] = \mu$

    $P[\text{picking a green marble}] = 1 - \mu$

- The value of $\mu$ is unknown to us
- We pick $N$ marbles independently
- The fraction of red marbles in **SAMPLE** $= \nu$



BIN

SAMPLE

$\nu$ = fraction of red marbles

$\mu$ = probability of red marbles

Learning
Components

A Simple
Learning Model

Type of
Learnings

Feasibility Of
Learning

**Probability to the
rescue**

# Does $\nu$ say anything about $\mu$?

- **No!** (certain answer)
  - Sample can be mostly red while bin is mostly red
- **Yes!** (uncertain answer)
  - Sample frequency $\nu$ is likely close to bin frequency $\mu$

Learning
Components

A Simple
Learning Model

Type of
Learnings

Feasibility Of
Learning

**Probability to the
rescue**

# What does $\nu$ say about $\mu$?

- In a big sample (large $N$), $\nu$ is probably close $\mu$ (within $\epsilon$)
- Formally,
$$P[|\nu - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N} \text{ for any } \epsilon > 0$$

  This is called **Hoeffding's Inequality**
- **Bound** does not depend on $\mu$; tradeoff: $N, \epsilon$ and the bound
- We have
$$\nu \approx \mu \implies \mu \approx \nu$$

- In other words, the statement "$\mu = \nu$" is **probably approximately correct** (P.A.C)

Learning
Components

A Simple
Learning Model

Type of
Learnings

Feasibility Of
Learning

**Probability to the
rescue**

## Connection to Learning

- **Bin problem**: The unknown is a number $\mu$
- **Learning problem**: The unknown is a function $f : \mathcal{X} \to \mathcal{Y}$
- Each marble ● is a point $\boldsymbol{x} \in \mathcal{X}$

| Bin problem | Learning problem |
|:---:|:---:|
| ● | hypothesis got it right $h(x) = f(x)$ |
| ● | hypothesis got it wrong $h(x) \neq f(x)$ |

Learning
Components

A Simple
Learning Model

Type of
Learnings

Feasibility Of
Learning

Probability to the
rescue

## Connection to Learning (cont.)

- The error rate within the sample, which corresponds to $\nu$ in the bin model, will be called the *in-sample error*, (domain $\mathcal{D}$)

$$E_{in}(h) = \text{fraction of } \mathcal{D} \text{ where } f \text{ and } h \text{ disagre}$$

$$= \frac{1}{N} \sum_{n=1}^{N} [\![ h(\mathbf{x}_n) \neq f(\mathbf{x}_n) ]\!]$$

where $[\![ statement ]\!] = 1$ if the statement is true, and $= 0$ if the statement is false

- In the same way, we define the *out-of-sample error* , (domain $\mathcal{X}$)

$$E_{out}(h) = P(h(\mathbf{x}) \neq f(\mathbf{x})), \mathbf{x} \in \mathcal{X}$$

which corresponds to $\mu$ in the bin model.

Learning
Components
A Simple
Learning Model
Type of
Learnings
Feasibility Of
Learning
**Probability to the
rescue**

## Connection to Learning (cont.)
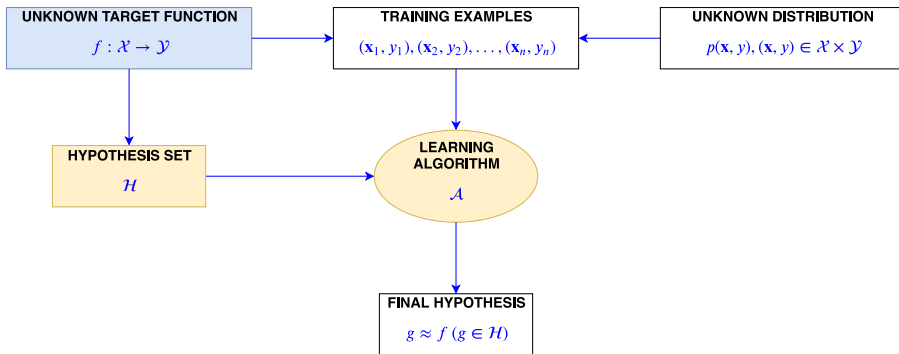
- The Hoeffding inequality becomes:

$$P[|E_{in}(h) - E_{out}(h)| > \epsilon] \le 2e^{-2\epsilon^2 N} \text{ for any } \epsilon > 0$$

**Inductive Learning Hypothesis**

Generalization is possible.

- If a machine performs well on most **training data** AND it is not too complex, it will probably do well on **similar test data**.

Learning
Components

A Simple
Learning Model

Type of
Learnings

Feasibility Of
Learning

Probability to the
rescue

# Back to Learning Diagram

# References

📄 Goodfellow, I., Bengio, Y., and Courville, A. (2016).
*Deep learning.*
MIT press.

📄 Lê, B. and Tô, V. (2014).
*Cở sở trí tuệ nhân tạo.*
Nhà xuất bản Khoa học và Kỹ thuật.

📄 Nguyen, T. (2018).
Artificial intelligence slides.
Technical report, HCMC University of Sciences.

📄 Russell, S. and Norvig, P. (2016).
*Artificial intelligence: a modern approach.*
Pearson Education Limited.