# LEARNING PROBLEM

Bùi Tiến Lên

2023

# Contents

Learning
Components

A Simple
Learning Model
Hypothesis Set
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

## 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 |
| $\boldsymbol{X}\boldsymbol{Y}$ | matrix multiplication |
| $\boldsymbol{X}^{-1}$ | inverse |

**Learning Components**

**A Simple Learning Model**
Hypothesis Set
Learning Algorithm

**Feasibility Of Learning**
Probability to the rescue

**Risk and Empirical Risk**
Loss function
Empirical risk
Regularizer

# 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?

**Learning Components**

**A Simple Learning Model**
Hypothesis Set
Learning Algorithm

**Feasibility Of Learning**
Probability to the rescue

**Risk and Emprical Risk**
Loss function
Empirical risk
Regularizer

## Problem Statement

Formalization

- Input: **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**
Hypothesis Set
Learning Algorithm

**Feasibility Of
Learning**
Probability to the
rescue

**Risk and
Emprical Risk**
Loss function
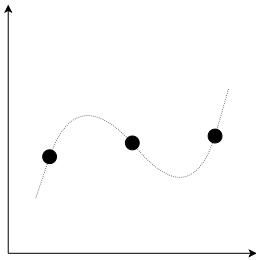Empirical risk
Regularizer

# Inductive Bias

### Theorem 1 (No Free Lunch Theorems)

*An unbiased learner can never generalize.*

### Concept 1

**An inductive bias** of a learner is the set of assumptions a learner uses to predict results given inputs it has not yet encountered.

- **Consider**: arbitrarily wiggly functions or random truth tables.



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

**Learning Components**

A Simple
Learning Model
Hypothesis Set
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
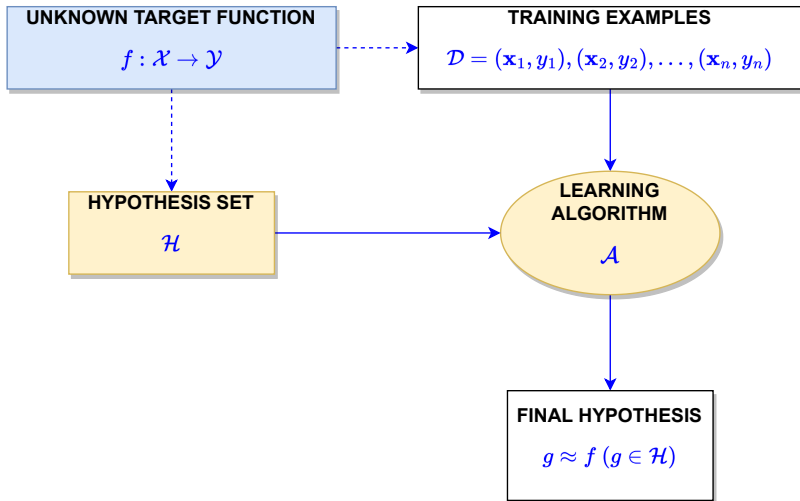Empirical risk
Regularizer

# Inductive Bias (cont.)

### 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**.
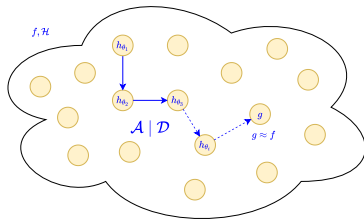
**Learning Components**

A Simple Learning Model
Hypothesis Set
Learning Algorithm

Feasibility Of Learning
Probability to the rescue

Risk and Empirical Risk
Loss function
Empirical risk
Regularizer

# Components of Learning

**Learning
Components**

**A Simple
Learning Model**
Hypothesis Set
Learning Algorithm

**Feasibility Of
Learning**
Probability to the
rescue

**Risk and
Emprical Risk**
Loss function
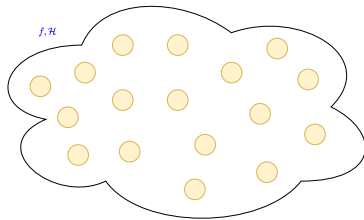Empirical risk
Regularizer

## Learning Model

The two components are referred as the
**learning model**

- The **hypothesis set** $\mathcal{H}$ is a set of
  functions that is potentially similar
  to $f$

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



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

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

**Learning Components**

A Simple
Learning Model
Hypothesis Set
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
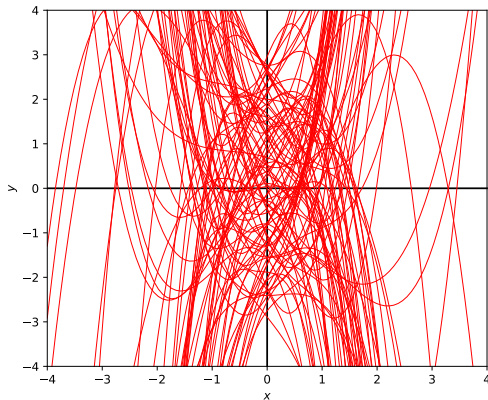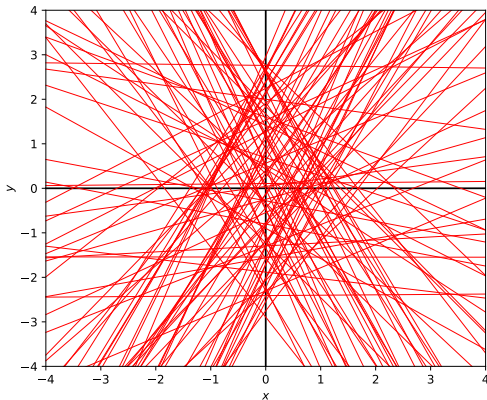Regularizer

# What is hypothesis set

### Concept 2

**Hypothesis set** is a set of potential functions, models or solutions

- Hypothesis set can be **finite**. For example
  - {*guilty*, *not guilty*}
  - {accept, *reject*}
  - {happy, sad}
  - $\{1, 2, 3, 4, 5, 6\}$

Learning
Components

A Simple
Learning Model
Hypothesis Set
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
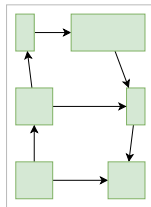Empirical risk
Regularizer

# What is hypothesis set (cont.)

- Hypothesis set can be **infinite**. For example, sets of functions $y = \theta_0 + \theta_1 x$ and $y = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$

**Learning Components**

**A Simple Learning Model**
Hypothesis Set
Learning Algorithm

**Feasibility Of Learning**
Probability to the rescue

**Risk and Emprical Risk**
Loss function
Empirical risk
Regularizer

# Parameter representations

- Each element of hypothesis set often indexed by **parameters** or **weights** ($\theta$ or $w$)
- Two basic representations for parameters: **factored**, and **structured**
  1. Factored: a paramater set consists of a vector of attribute values; values can be boolean, real-valued, or one of a fixed set of symbols.
  2. Structured: a paramater set includes objects, each of which may have attributes of its own as well as relationships to other objects.

# A Simple Learning Model

- Hypothesis Set
- Learning Algorithm

Learning
Components

A Simple
Learning Model

**Hypothesis Set**
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

# A Simple Hypothesis Set

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

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

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

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

$$h(\boldsymbol{x}) = h_{\boldsymbol{w}, threshold}(\boldsymbol{x}) = \text{sign}\left(\left(\sum_{i=1}^{d} w_i x_i\right) - threshold\right) \tag{2}$$

Learning
Components

A Simple
Learning Model

**Hypothesis Set**
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Empirical Risk
Loss function
Empirical risk
Regularizer

# A Simple Hypothesis Set (cont.)

- Set $w_0 = -threshold$

$$h(x) = h_{\mathbf{w}}(x) = \text{sign}\left(\left(\sum_{i=1}^{d} w_i x_i\right) + w_0\right) \tag{3}$$

- Introduce an artificial coordinate $x_0 = 1$

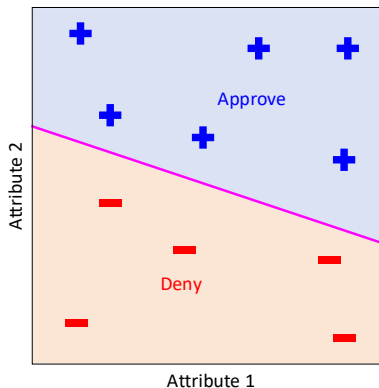$$h(x) = h_{\mathbf{w}}(x) = \text{sign}\left(\sum_{i=0}^{d} w_i x_i\right) \tag{4}$$

- In vector form, the perceptron implements

$$h(x) = h_{\mathbf{w}}(x) = \text{sign}\left(\mathbf{w}^\mathsf{T} \mathbf{x}\right) \tag{5}$$

Learning
Components

A Simple
Learning Model
Hypothesis Set
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

## 2D Model Visualization

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

Learning
Components

A Simple
Learning Model
Hypothesis Set
**Learning Algorithm**

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

## A Simple Learning Algorithm

- **The performance measure**: *the error rate*
- We uses the simple learning algorithm (**perceptron learning algorithm** - **PLA**) to find $w$

$$\arg \min_{w} E(h_{w}(x), y \mid \mathcal{D}) \tag{6}$$

Learning
Components

A Simple
Learning Model
Hypothesis Set
**Learning Algorithm**

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

## A Simple Learning Algorithm (cont.)

- Given the training set

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

1. **Init $\mathbf{w}$**

2. **Repeat until satisfied**

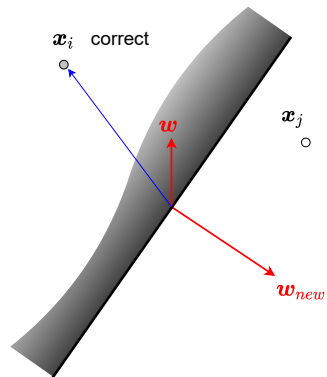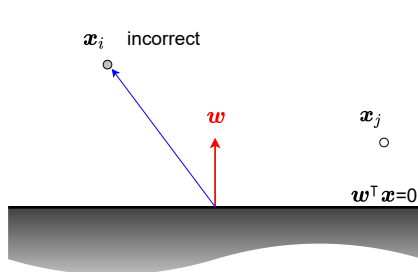   - At iteration $t = 1, 2, 3, ...$, pick a *misclassified* point $(\mathbf{x}_i, y_i)$

   $$sign(\mathbf{w}^\mathsf{T}\mathbf{x}_i) \neq y_i \tag{7}$$

   - and update the weight vector

   $$\mathbf{w} \leftarrow \mathbf{w} + y_i\mathbf{x}_i \tag{8}$$

Learning
Components

A Simple
Learning Model
Hypothesis Set
**Learning Algorithm**

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

# A Simple Explanation

Learning
Components

A Simple
Learning Model
Hypothesis Set
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

## Is It Learning Algorithm?

Learning
Components

A Simple
Learning Model
Hypothesis Set
**Learning Algorithm**

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

# A Learning Puzzle



y = -1

y = +1

y = ?

# Feasibility Of Learning

- Probability to the rescue

Learning
Components

A Simple
Learning Model
Hypothesis Set
Learning Algorithm

**Feasibility Of
Learning**
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

# Feasibility Of Learning

The **feasibility of learning** is thus split into two questions:

1. Can we make **the performance** good enough?
   - run our learning algorithm on the actual data $\mathcal{D}$ and see how good we can get.

2. Can we make sure that **the performance** inside of $\mathcal{D}$ is close enough to **the performance** outside of $\mathcal{D}$?
   - probability theory

Learning
Components

A Simple
Learning Model
Hypothesis Set
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

# 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
Hypothesis Set
Learning Algorithm

Feasibility Of Learning
**Probability to the rescue**

Risk and Empirical Risk
Loss function
Empirical risk
Regularizer

# 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**
Hypothesis Set
Learning Algorithm

**Feasibility Of Learning**
**Probability to the rescue**

**Risk and Empirical Risk**
Loss function
Empirical risk
Regularizer

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

- In a big sample (large $N$), $\nu$ is probably close $\mu$ (within $\epsilon$)
- Formally,

$$P[|\nu - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N} \text{ for any } \epsilon > 0 \tag{9}$$

   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
Hypothesis Set
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

## Connection to Learning

| Bin problem |
|---|
| The unknown is a number $\mu$ |
| a marble ● |
| 🟢 |
| 🔴 |

| Learning problem |
|---|
| The unknown is a function $f : \mathcal{X} \to \mathcal{Y}$ |
| a point $\boldsymbol{x} \in \mathcal{X}$ |
| hypothesis got it right $h(x) = f(x)$ |
| hypothesis got it wrong $h(x) \neq f(x)$ |

Learning
Components

A Simple
Learning Model
Hypothesis Set
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Empirical Risk
Loss function
Empirical risk
Regularizer

## Connection to Learning (cont.)

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

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

$$= \frac{1}{N} \sum_{n=1}^{N} \mathbb{I}(h(\boldsymbol{x}_n) \neq f(\boldsymbol{x}_n))$$

where $\mathbb{I}(...) = 1$ if the statement is **true**, and $\mathbb{I}(...) = 0$ if the statement is **false**

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

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

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

Learning
Components

A Simple
Learning Model
Hypothesis Set
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

## Connection to Learning (cont.)

- The Hoeffding inequality becomes:

$$P[|E_{in}(h) - E_{out}(h)| > \epsilon] \leq 2e^{-2\epsilon^2 N} \text{ for any } \epsilon > 0 \qquad (10)$$

> *In a big sample $\mathcal{D}$, **the performance** inside of $\mathcal{D}$ is close enough to **the performance** outside of $\mathcal{D}$*

# Risk and Emprical Risk

- Loss function
- Empirical risk
- Regularizer

Learning
Components

A Simple
Learning Model
Hypothesis Set
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

## Loss function

### Concept 3

Given a hypothesis $\hat{y} = h(\boldsymbol{x}) \in \mathcal{H}$, a non-negative real-valued **loss function** $\ell(\hat{y}, y)$ which measures how different the prediction $\hat{y}$ of a hypothesis is from the true outcome $y$.

Learning
Components

A Simple
Learning Model
Hypothesis Set
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

## Loss Functions for Binary Classification

- Zero-one loss

$$\mathbb{I}(h(\boldsymbol{x}) \neq y) \tag{11}$$

- Log loss (logistic regression)

$$\log(1 + e^{-h(\boldsymbol{x})y}) \tag{12}$$

- Exponential loss (AdaBoost)

$$e^{-h(\boldsymbol{x})y} \tag{13}$$

**Learning Components**

**A Simple Learning Model**
Hypothesis Set
Learning Algorithm

**Feasibility Of Learning**
Probability to the rescue

**Risk and Empirical Risk**
**Loss function**
Empirical risk
Regularizer

## Loss Functions for Regression

- Squared loss

$$(h(\boldsymbol{x}) - y)^2 \tag{14}$$

- Absolute loss

$$|h(\boldsymbol{x}) - y| \tag{15}$$

Learning
Components

A Simple
Learning Model
Hypothesis Set
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

# Risk

### Concept 4

The **risk** $E$ associated with hypothesis $h(\boldsymbol{x})$ is defined as the expectation of the loss function

$$E(h) = \mathbb{E}[\ell(h(\boldsymbol{x}), y)] = \int \ell(h(\boldsymbol{x}), y) dp(\boldsymbol{x}, y) \tag{16}$$

Learning
Components

A Simple
Learning Model
Hypothesis Set
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

# Empirical Risk

## Concept 5

The **empirical risk** $\hat{E}$ is the average of the loss function on the training set
$\mathcal{D} = \{(\boldsymbol{x}_1, y_1), (\boldsymbol{x}_2, y_2), ...(\boldsymbol{x}_N, y_N)\}$

$$\hat{E} = \frac{1}{N} \sum_{i=1}^{N} \ell(h_{\boldsymbol{w}}(\boldsymbol{x}_i), y_i) \tag{17}$$

## Theorem 2

*The empirial risk is unbiased estimate of the risk*

Learning
Components

A Simple
Learning Model
Hypothesis Set
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
**Empirical risk**
Regularizer

# Empirical Risk (cont.)

### Concept 6

Empirical risk of hypothesis $h_{\boldsymbol{w}}(x)$ with a loss function $\ell$ and a regularizer *reg*

$$\hat{E} = \frac{1}{N} \sum_{i=1}^{N} \underbrace{\ell(h_{\boldsymbol{w}}(\boldsymbol{x}_i), y_i)}_{Loss} + \underbrace{\lambda reg(\boldsymbol{w})}_{Regularizer} \tag{18}$$

Learning
Components

A Simple
Learning Model
Hypothesis Set
Learning Algorithm

Feasibility Of
Learning
Probability to the
rescue

Risk and
Emprical Risk
Loss function
Empirical risk
Regularizer

# The empirical risk minimization principle

## Principle

The learning algorithm should choose a hypothesis $h_{\mathbf{w}}$ which minimizes the empirical risk

$$h_{\mathbf{w}} = \arg \min_{h_{\mathbf{w}} \in \mathcal{H}} \hat{E}(h_{\mathbf{w}} \mid \mathcal{D}) \tag{19}$$

# Regularizers

### Theorem 3

*For each $\lambda \geq 0$, there exists $B \geq 0$. such that the two formulations are equivalent,*

$$\arg\min_{\mathbf{w}} \sum_{i=1}^{N} \ell(h_{\mathbf{w}}(\mathbf{x}_i), y_i) + \lambda reg(\mathbf{w}) \tag{20}$$

$$\arg\min_{\mathbf{w}} \sum_{i=1}^{N} \ell(h_{\mathbf{w}}(\mathbf{x}_i), y_i) \text{ subject to } reg(\mathbf{w}) \leq B \tag{21}$$

**Learning Components**

**A Simple Learning Model**
Hypothesis Set
Learning Algorithm

**Feasibility Of Learning**
Probability to the rescue

**Risk and Emprical Risk**
Loss function
Empirical risk
**Regularizer**

# Regularizers (cont.)

- $L_2$-regularization

$$reg(\boldsymbol{w}) = \boldsymbol{w}^\top \boldsymbol{w} = \|\boldsymbol{w}\|_2^2 \tag{22}$$

- $L_1$-regularization

$$reg(\boldsymbol{w}) = \|\boldsymbol{w}\|_1 \tag{23}$$

## 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.

📄 Russell, S. and Norvig, P. (2021).
*Artificial intelligence: a modern approach*.
Pearson Education Limited.