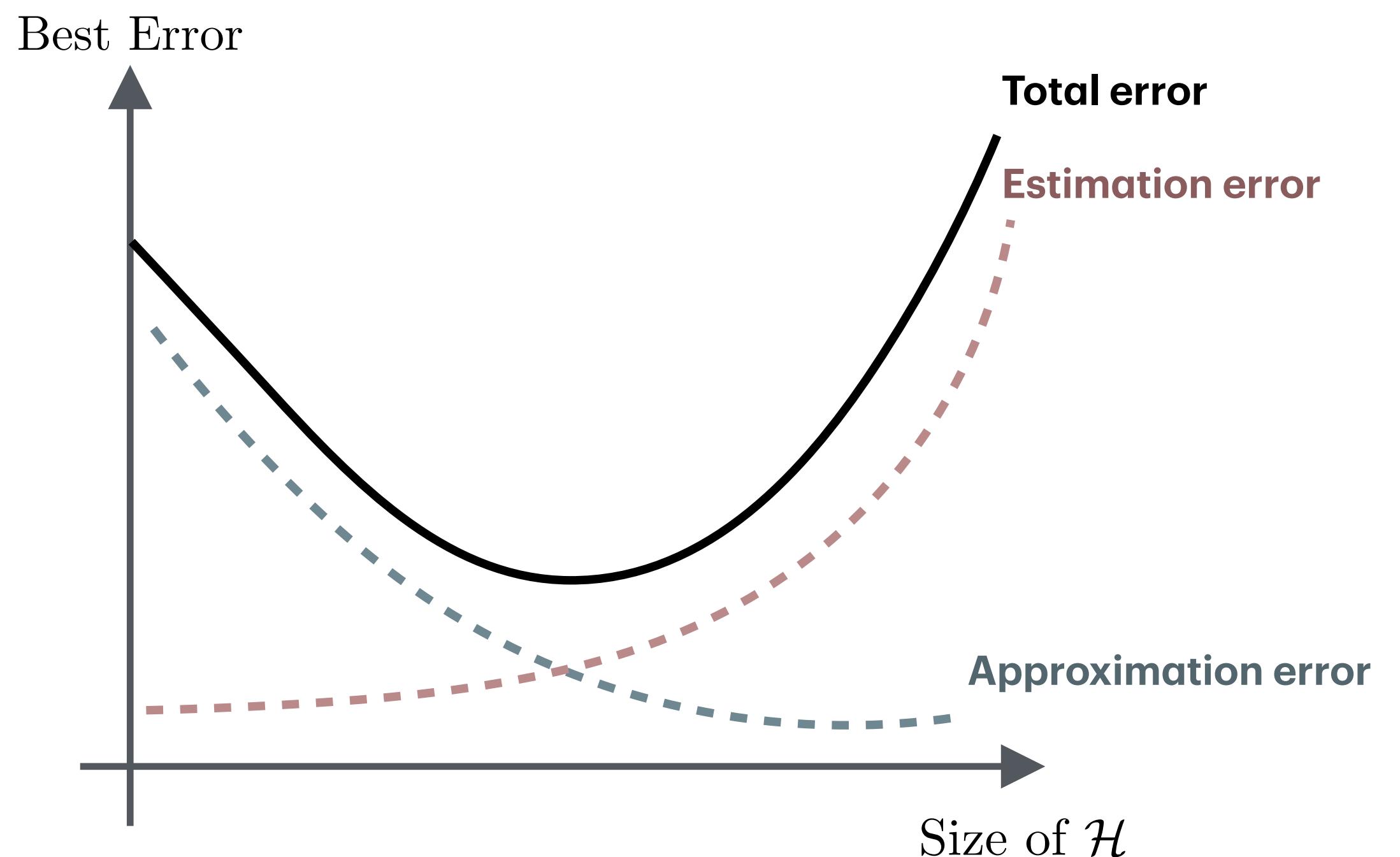
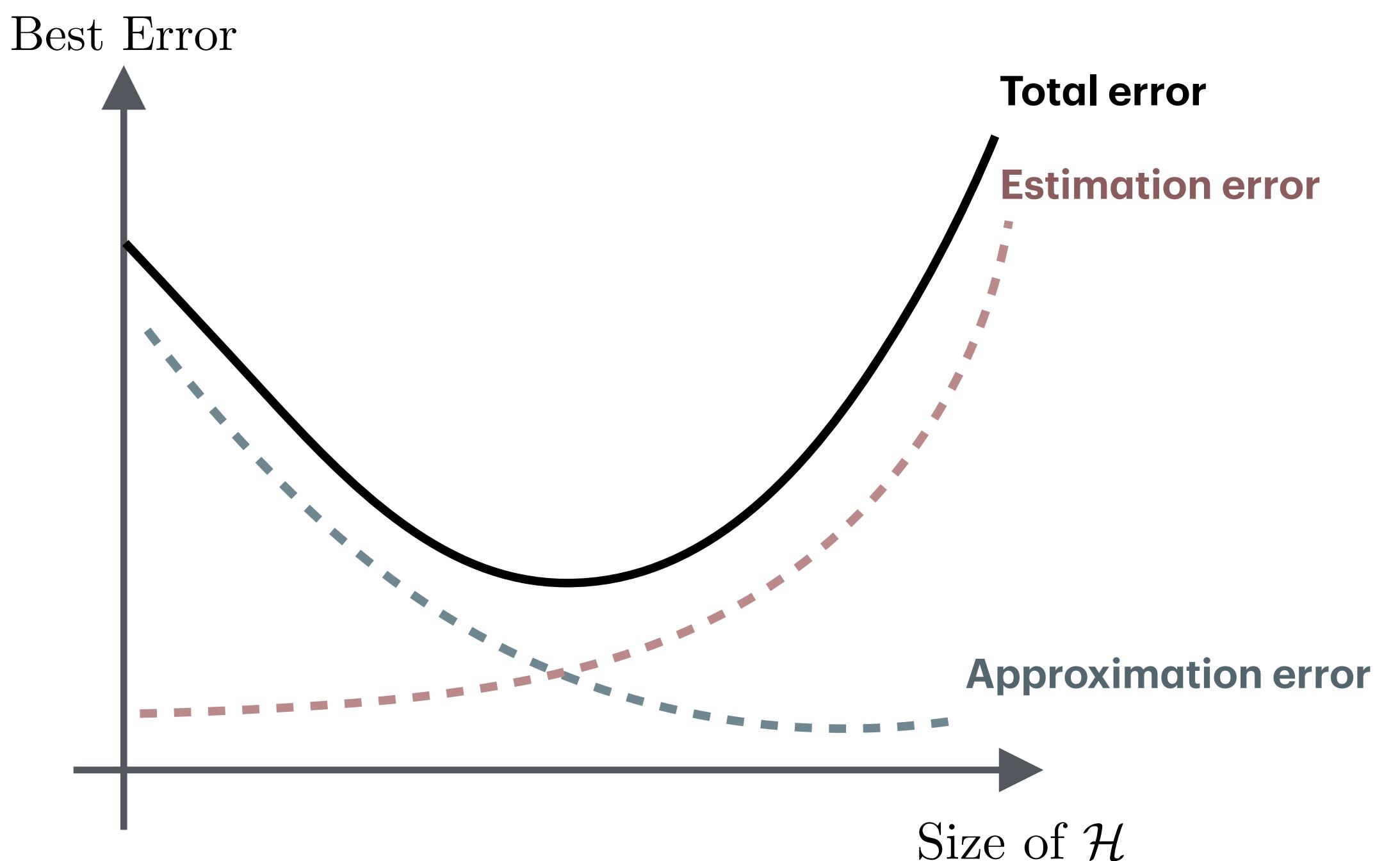


# **LESSON 1**

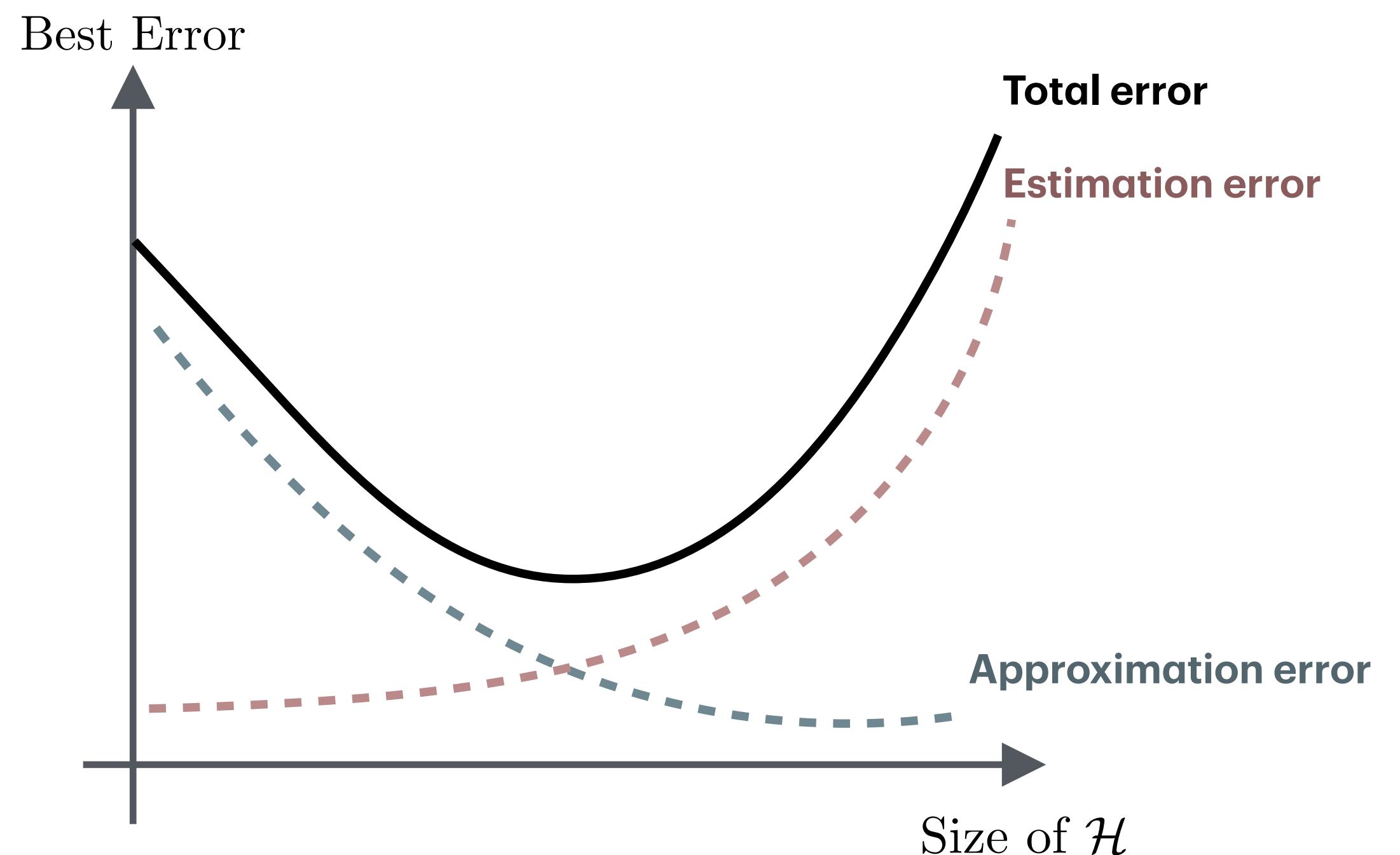
## **INTRO TO STATISTICAL LEARNING**



1. Supervised Learning Setting
2. Estimation vs Approximation
3. Maximal inequalities
4. Rademacher Complexity



## 1. Supervised Learning Setting



# **The supervised learning setup**

- General setup

# The supervised learning setup

## ■ General setup

**Supervised learning (SL)** is a paradigm in machine learning where input objects and a desired output value train a model. The training data is processed, building a function that maps new data to expected output values.

# The supervised learning setup

## ■ General setup

**Supervised learning (SL)** is a paradigm in machine learning where input objects and a desired output value train a model. The training data is processed, building a function that maps new data to expected output values.

### Definitions

Formally, a supervised learning is a problem of the form

$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

where

# The supervised learning setup

## ■ General setup

**Supervised learning (SL)** is a paradigm in machine learning where input objects and a desired output value train a model. The training data is processed, building a function that maps new data to expected output values.

### Definitions

Formally, a supervised learning is a problem of the form

$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

where  $\mathcal{H}$  is a set of functions called the **hypothesis class**

# The supervised learning setup

## ■ General setup

Supervised learning (**SL**) is a paradigm in machine learning where input objects and a desired output value train a model. The training data is processed, building a function that maps new data to expected output values.

### Definitions

Formally, a supervised learning is a problem of the form

$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

where  $\mathcal{H}$  is a set of functions called the **hypothesis class**

$f \in \mathcal{H}$  is called an **estimator**

# The supervised learning setup

## ■ General setup

Supervised learning (**SL**) is a paradigm in machine learning where input objects and a desired output value train a model. The training data is processed, building a function that maps new data to expected output values.

### Definitions

Formally, a supervised learning is a problem of the form

$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

where  $\mathcal{H}$  is a set of functions called the **hypothesis class**

$f \in \mathcal{H}$  is called an **estimator**

$\mathcal{D}$  is a **data distribution** from which we can sample **input-output pairs**  $x, y \in \mathcal{X} \times \mathcal{Y}$

# The supervised learning setup

## ■ General setup

Supervised learning (**SL**) is a paradigm in machine learning where input objects and a desired output value train a model. The training data is processed, building a function that maps new data to expected output values.

### Definitions

Formally, a supervised learning is a problem of the form

$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

where  $\mathcal{H}$  is a set of functions called the **hypothesis class**

$f \in \mathcal{H}$  is called an **estimator**

$\mathcal{D}$  is a **data distribution** from which we can sample **input-output pairs**  $x, y \in \mathcal{X} \times \mathcal{Y}$

$\mathcal{L} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  is an **error function** aimed at penalizing the discrepancy between the ground truth  $y$  and the estimate  $f(x)$ .

# The supervised learning setup

## ■ General setup

Supervised learning (**SL**) is a paradigm in machine learning where input objects and a desired output value train a model. The training data is processed, building a function that maps new data to expected output values.

### Definitions

Formally, a supervised learning is a problem of the form

$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

where  $\mathcal{H}$  is a set of functions called the **hypothesis class**

$f \in \mathcal{H}$  is called an **estimator**

$\mathcal{D}$  is a **data distribution** from which we can sample **input-output pairs**  $x, y \in \mathcal{X} \times \mathcal{Y}$

$\mathcal{L} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  is an **error function** aimed at penalizing the discrepancy between the ground truth  $y$  and the estimate  $f(x)$ .

$f : \mathcal{X} \rightarrow \mathcal{Y}$  is called an **estimator**.

# The supervised learning setup

## ■ General setup

Supervised learning (**SL**) is a paradigm in machine learning where input objects and a desired output value train a model. The training data is processed, building a function that maps new data to expected output values.

### Definitions

Formally, a supervised learning is a problem of the form

$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

where  $\mathcal{H}$  is a set of functions called the **hypothesis class**

$f \in \mathcal{H}$  is called an **estimator**

$\mathcal{D}$  is a **data distribution** from which we can sample **input-output pairs**  $x, y \in \mathcal{X} \times \mathcal{Y}$

$\mathcal{L} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  is an **error function** aimed at penalizing the discrepancy between the ground truth  $y$  and the estimate  $f(x)$ .

$f : \mathcal{X} \rightarrow \mathcal{Y}$  is called an **estimator**.

### Remarks :

- When  $\mathcal{Y}$  is finite, this above problem is called a **classification problem**.

# The supervised learning setup

## ■ General setup

Supervised learning (**SL**) is a paradigm in machine learning where input objects and a desired output value train a model. The training data is processed, building a function that maps new data to expected output values.

### Definitions

Formally, a supervised learning is a problem of the form

$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

where  $\mathcal{H}$  is a set of functions called the **hypothesis class**

$f \in \mathcal{H}$  is called an **estimator**

$\mathcal{D}$  is a **data distribution** from which we can sample **input-output pairs**  $x, y \in \mathcal{X} \times \mathcal{Y}$

$\mathcal{L} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  is an **error function** aimed at penalizing the discrepancy between the ground truth  $y$  and the estimate  $f(x)$ .

$f : \mathcal{X} \rightarrow \mathcal{Y}$  is called an **estimator**.

### Remarks :

- When  $\mathcal{Y}$  is finite, this above problem is called a **classification problem**.
- When  $\mathcal{Y}$  is **not** finite, this above problem is called a **regression problem**.

# The supervised learning setup

## ■ General setup

Supervised learning (**SL**) is a paradigm in machine learning where input objects and a desired output value train a model. The training data is processed, building a function that maps new data to expected output values.

### Definitions

Formally, a supervised learning is a problem of the form

$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

where  $\mathcal{H}$  is a set of functions called the **hypothesis class**

$f \in \mathcal{H}$  is called an **estimator**

$\mathcal{D}$  is a **data distribution** from which we can sample **input-output pairs**  $x, y \in \mathcal{X} \times \mathcal{Y}$

$\mathcal{L} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  is an **error function** aimed at penalizing the discrepancy between the ground truth  $y$  and the estimate  $f(x)$ .

$f : \mathcal{X} \rightarrow \mathcal{Y}$  is called an **estimator**.

## ■ Empirical Risk Minimization

In practice, one often cannot access the whole distribution  $\mathcal{D}$ , but only a finite sample of it,  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ , often called the **training set**.

### Remarks :

- When  $\mathcal{Y}$  is finite, this above problem is called a **classification problem**.
- When  $\mathcal{Y}$  is **not** finite, this above problem is called a **regression problem**.

# The supervised learning setup

## ■ General setup

Supervised learning (**SL**) is a paradigm in machine learning where input objects and a desired output value train a model. The training data is processed, building a function that maps new data to expected output values.

### Definitions

Formally, a supervised learning is a problem of the form

$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

where  $\mathcal{H}$  is a set of functions called the **hypothesis class**

$f \in \mathcal{H}$  is called an **estimator**

$\mathcal{D}$  is a **data distribution** from which we can sample **input-output pairs**  $x, y \in \mathcal{X} \times \mathcal{Y}$

$\mathcal{L} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  is an **error function** aimed at penalizing the discrepancy between the ground truth  $y$  and the estimate  $f(x)$ .

$f : \mathcal{X} \rightarrow \mathcal{Y}$  is called an **estimator**.

## ■ Empirical Risk Minimization

In practice, one often cannot access the whole distribution  $\mathcal{D}$ , but only a finite sample of it,  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ , often called the **training set**.

### Definition

**Empirical minimization** consists in replacing the general supervised learning problem by a minimization over the training set. Formally, we solve

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, f(x_i))$$

### Remarks :

- When  $\mathcal{Y}$  is finite, this above problem is called a **classification problem**.
- When  $\mathcal{Y}$  is **not** finite, this above problem is called a **regression problem**.

# The supervised learning setup

## ■ General setup

Supervised learning (**SL**) is a paradigm in machine learning where input objects and a desired output value train a model. The training data is processed, building a function that maps new data to expected output values.

### Definitions

Formally, a supervised learning is a problem of the form

$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

where  $\mathcal{H}$  is a set of functions called the **hypothesis class**

$f \in \mathcal{H}$  is called an **estimator**

$\mathcal{D}$  is a **data distribution** from which we can sample **input-output pairs**  $x, y \in \mathcal{X} \times \mathcal{Y}$

$\mathcal{L} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  is an **error function** aimed at penalizing the discrepancy between the ground truth  $y$  and the estimate  $f(x)$ .

$f : \mathcal{X} \rightarrow \mathcal{Y}$  is called an **estimator**.

### Remarks :

- When  $\mathcal{Y}$  is finite, this above problem is called a **classification problem**.
- When  $\mathcal{Y}$  is **not** finite, this above problem is called a **regression problem**.

## ■ Empirical Risk Minimization

In practice, one often cannot access the whole distribution  $\mathcal{D}$ , but only a finite sample of it,  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ , often called the **training set**.

### Definition

**Empirical minimization** consists in replacing the general supervised learning problem by a minimization over the training set. Formally, we solve

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, f(x_i))$$

### Remarks :

- The hypothesis class is typically defined as a parameterized family of functions :  
$$\mathcal{H} = \{x \mapsto f(w, x), w \in \mathbb{R}^d\}$$

# The supervised learning setup

## ■ General setup

Supervised learning (**SL**) is a paradigm in machine learning where input objects and a desired output value train a model. The training data is processed, building a function that maps new data to expected output values.

### Definitions

Formally, a supervised learning is a problem of the form

$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

where  $\mathcal{H}$  is a set of functions called the **hypothesis class**

$f \in \mathcal{H}$  is called an **estimator**

$\mathcal{D}$  is a **data distribution** from which we can sample **input-output pairs**  $x, y \in \mathcal{X} \times \mathcal{Y}$

$\mathcal{L} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  is an **error function** aimed at penalizing the discrepancy between the ground truth  $y$  and the estimate  $f(x)$ .

$f : \mathcal{X} \rightarrow \mathcal{Y}$  is called an **estimator**.

### Remarks :

- When  $\mathcal{Y}$  is finite, this above problem is called a **classification problem**.
- When  $\mathcal{Y}$  is **not** finite, this above problem is called a **regression problem**.

## ■ Empirical Risk Minimization

In practice, one often cannot access the whole distribution  $\mathcal{D}$ , but only a finite sample of it,  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ , often called the **training set**.

### Definition

**Empirical minimization** consists in replacing the general supervised learning problem by a minimization over the training set. Formally, we solve

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, f(x_i))$$

### Remarks :

- The hypothesis class is typically defined as a parameterized family of functions :  
$$\mathcal{H} = \{x \mapsto f(w, x), w \in \mathbb{R}^d\}$$
- The above problem can thus be replaced by an optimization problem in  $\mathbb{R}^d$  :

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, f(w, x_i))$$

# The supervised learning setup

## ■ General setup

Supervised learning (**SL**) is a paradigm in machine learning where input objects and a desired output value train a model. The training data is processed, building a function that maps new data to expected output values.

### Definitions

Formally, a supervised learning is a problem of the form

$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

where  $\mathcal{H}$  is a set of functions called the **hypothesis class**

$f \in \mathcal{H}$  is called an **estimator**

$\mathcal{D}$  is a **data distribution** from which we can sample **input-output pairs**  $x, y \in \mathcal{X} \times \mathcal{Y}$

$\mathcal{L} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  is an **error function** aimed at penalizing the discrepancy between the ground truth  $y$  and the estimate  $f(x)$ .

$f : \mathcal{X} \rightarrow \mathcal{Y}$  is called an **estimator**.

### Remarks :

- When  $\mathcal{Y}$  is finite, this above problem is called a **classification problem**.
- When  $\mathcal{Y}$  is **not** finite, this above problem is called a **regression problem**.

## ■ Empirical Risk Minimization

In practice, one often cannot access the whole distribution  $\mathcal{D}$ , but only a finite sample of it,  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ , often called the **training set**.

### Definition

**Empirical minimization** consists in replacing the general supervised learning problem by a minimization over the training set. Formally, we solve

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, f(x_i))$$

### Remarks :

- The hypothesis class is typically defined as a parameterized family of functions :  
$$\mathcal{H} = \{x \mapsto f(w, x), w \in \mathbb{R}^d\}$$
- The above problem can thus be replaced by an optimization problem in  $\mathbb{R}^d$  :

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, f(w, x_i))$$

- Hence, we can use the gradient descent, with respect to the hyperparameter  $w$ ,

$$w_{t+1} = w_t - \alpha \frac{\partial}{\partial w} \left[ \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, f(w, x_i)) \right] (w_t)$$

to learn the optimal parameter  $w^*$ , and thus the optimal predictor  $x \mapsto f(w^*, x)$ .

## **2 classical examples**

- Linear regression

## 2 classical examples

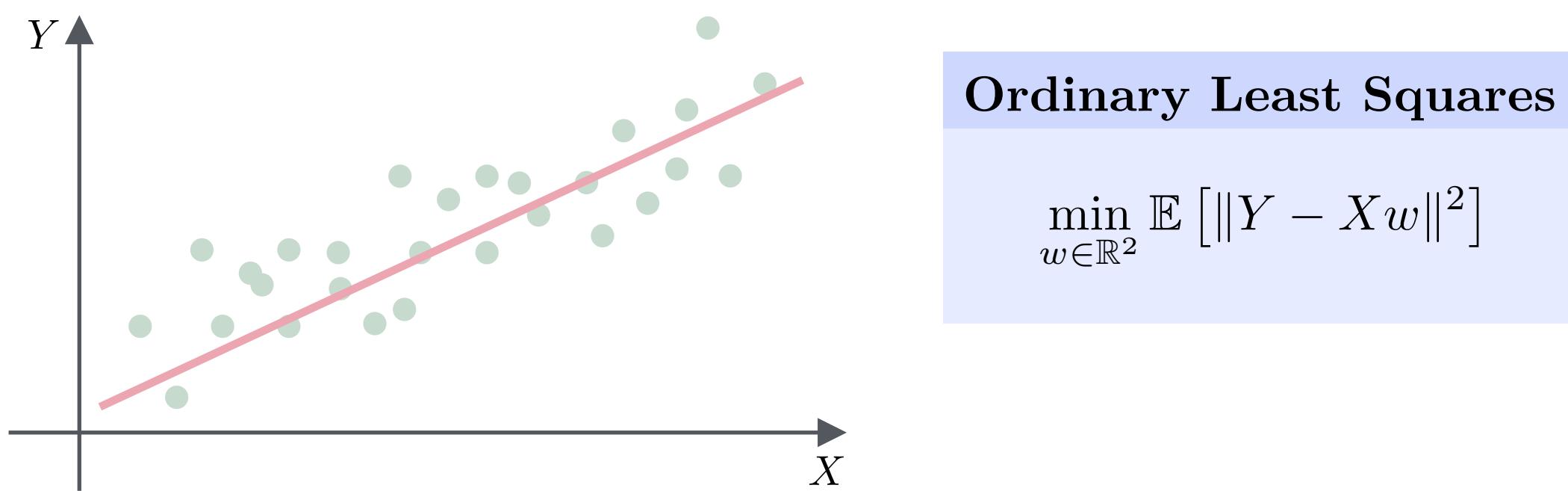
- Linear regression

Linear regression seeks to estimate a linear relationship between a **continuous** scalar response  $Y$ , a one or more explanatory variables captured in a vector  $X$ .

## 2 classical examples

### ■ Linear regression

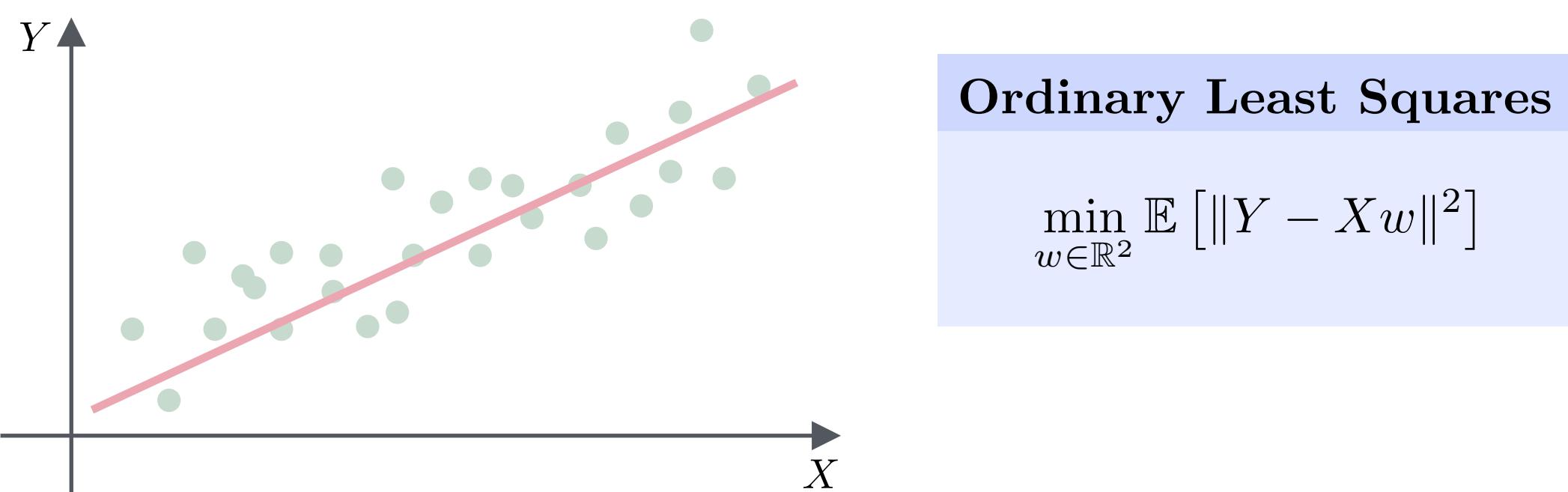
Linear regression seeks to estimate a linear relationship between a **continuous** scalar response  $Y$ , a one or more explanatory variables captured in a vector  $X$ .



## 2 classical examples

### ■ Linear regression

Linear regression seeks to estimate a linear relationship between a **continuous** scalar response  $Y$ , a one or more explanatory variables captured in a vector  $X$ .

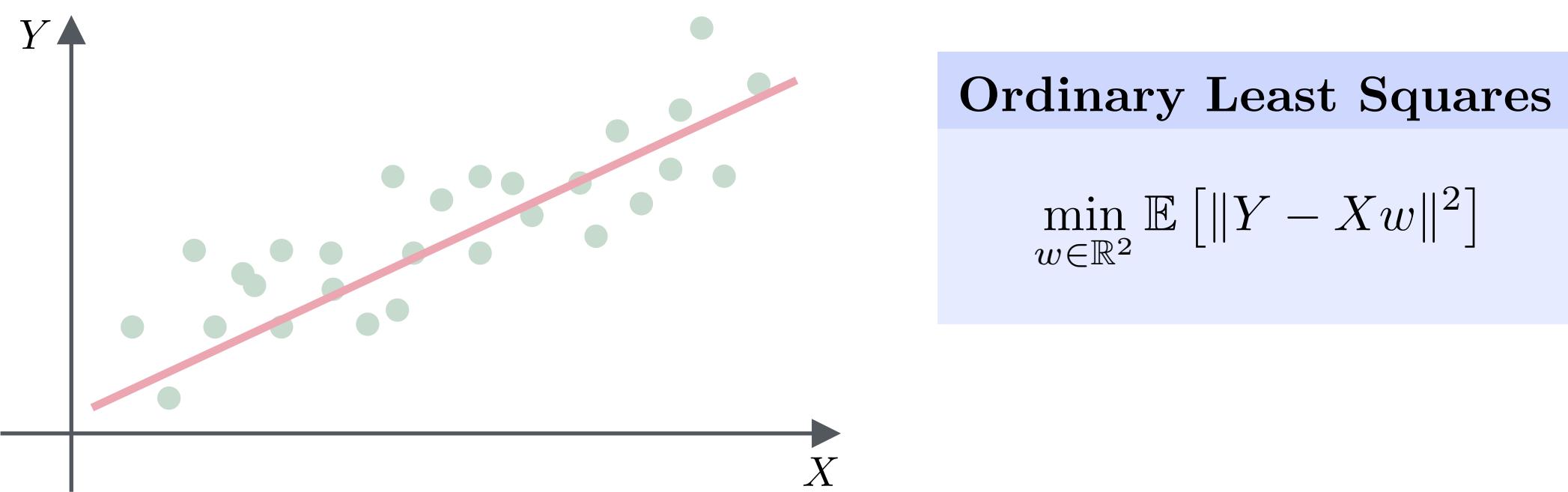


$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

## 2 classical examples

### ■ Linear regression

Linear regression seeks to estimate a linear relationship between a **continuous** scalar response  $Y$ , a one or more explanatory variables captured in a vector  $X$ .



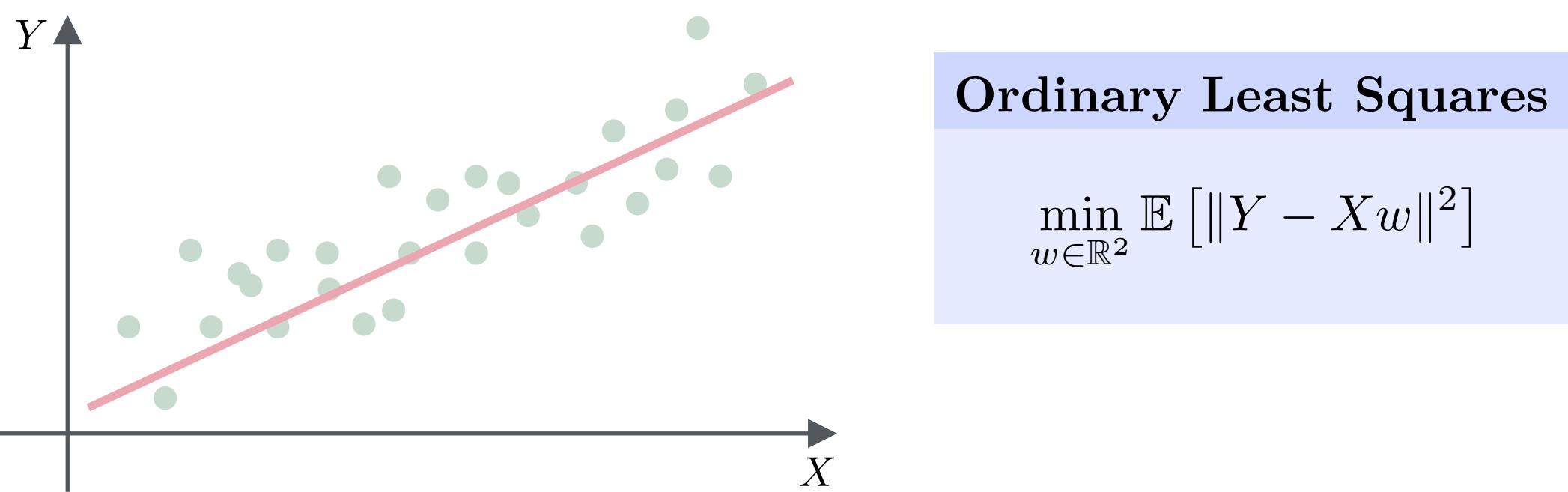
$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

$$\mathcal{H} = \mathbb{R}^2$$

## 2 classical examples

### ■ Linear regression

Linear regression seeks to estimate a linear relationship between a **continuous** scalar response  $Y$ , a one or more explanatory variables captured in a vector  $X$ .



$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

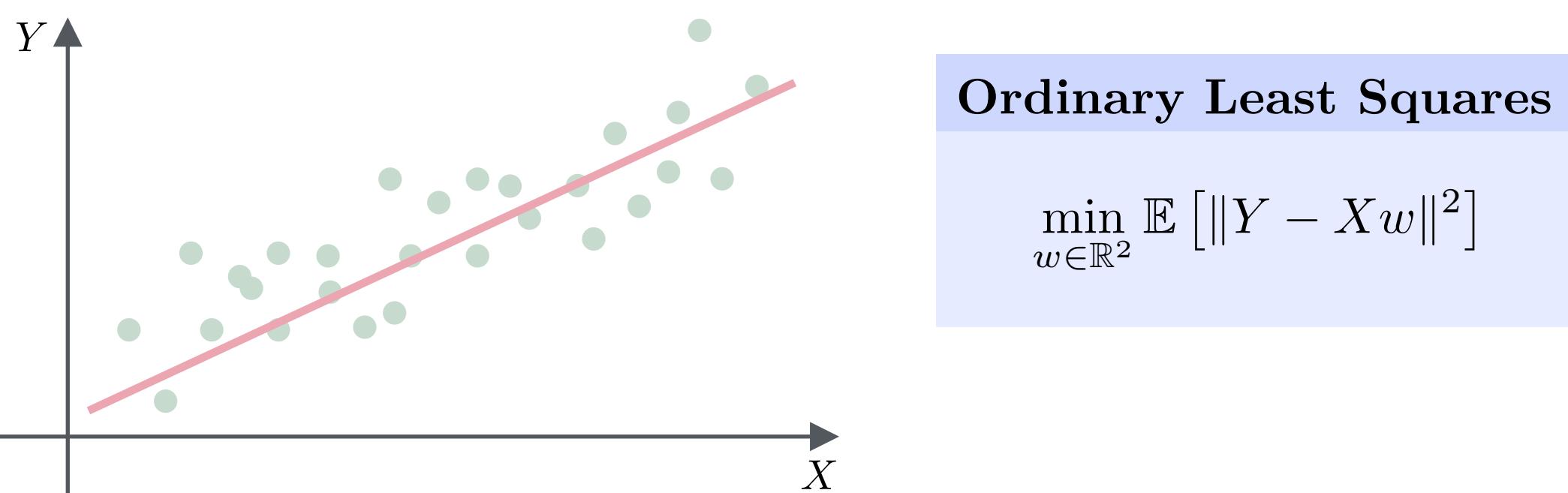
└─────────────────  $\mathcal{L}(y, \hat{y}) = (y - \hat{y})^2$

$$\mathcal{H} = \mathbb{R}^2$$

## 2 classical examples

### ■ Linear regression

Linear regression seeks to estimate a linear relationship between a **continuous** scalar response  $Y$ , a one or more explanatory variables captured in a vector  $X$ .



$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

$\mathcal{L}(y, \hat{y}) = (y - \hat{y})^2$

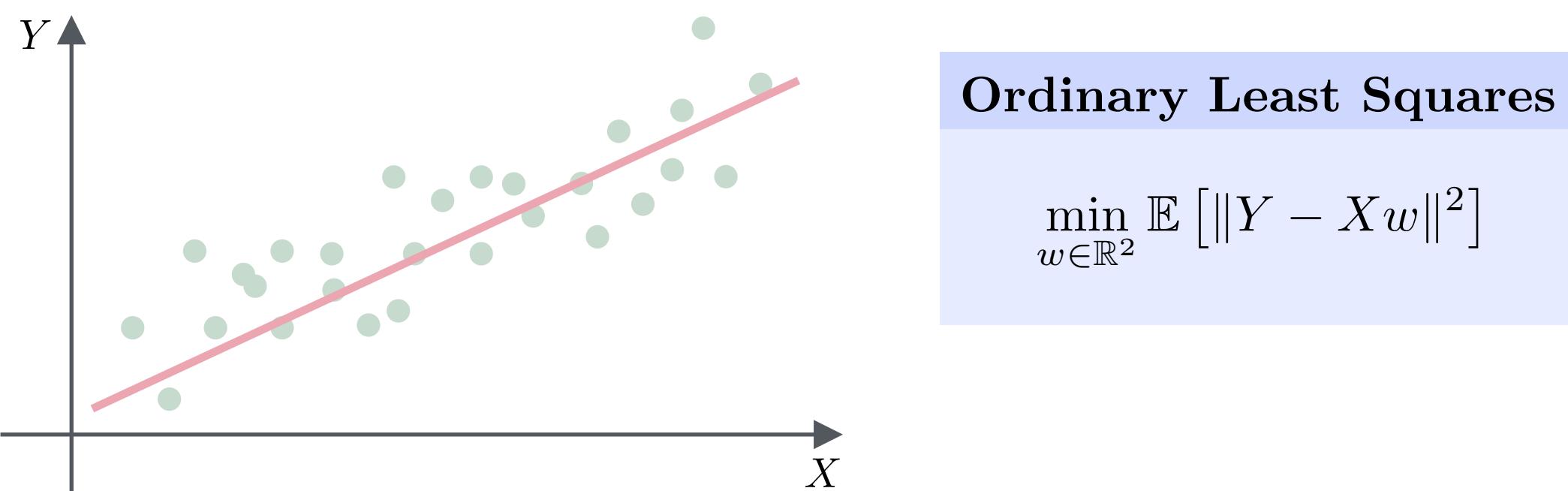
$f(x) = w^\top x$

$\mathcal{H} = \mathbb{R}^2$

## 2 classical examples

### ■ Linear regression

Linear regression seeks to estimate a linear relationship between a **continuous** scalar response  $Y$ , a one or more explanatory variables captured in a vector  $X$ .



$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

$\mathcal{L}(y, \hat{y}) = (y - \hat{y})^2$

$f(x) = w^\top x$

$\mathcal{H} = \mathbb{R}^2$

### ■ Logistic regression

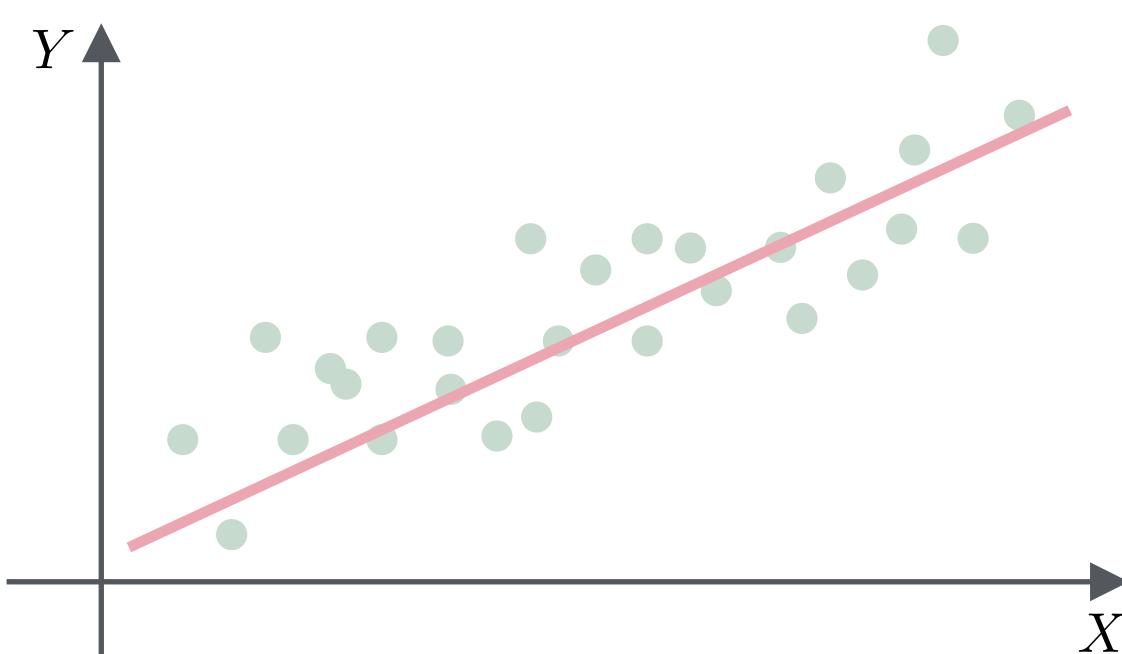
## 2 classical examples

### ■ Linear regression

Linear regression seeks to estimate a linear relationship between a **continuous** scalar response  $Y$ , a one or more explanatory variables captured in a vector  $X$ .

### ■ Logistic regression

Logistic regression seeks to a linear relationship between a discrete response  $Y$  taking values in some finite space  $\mathcal{Y}$ , and one or more explanatory variables captured in a vector  $X$ .



#### Ordinary Least Squares

$$\min_{w \in \mathbb{R}^2} \mathbb{E} [\|Y - Xw\|^2]$$

$$\mathcal{L}(y, \hat{y}) = (y - \hat{y})^2$$

$$\min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))]$$

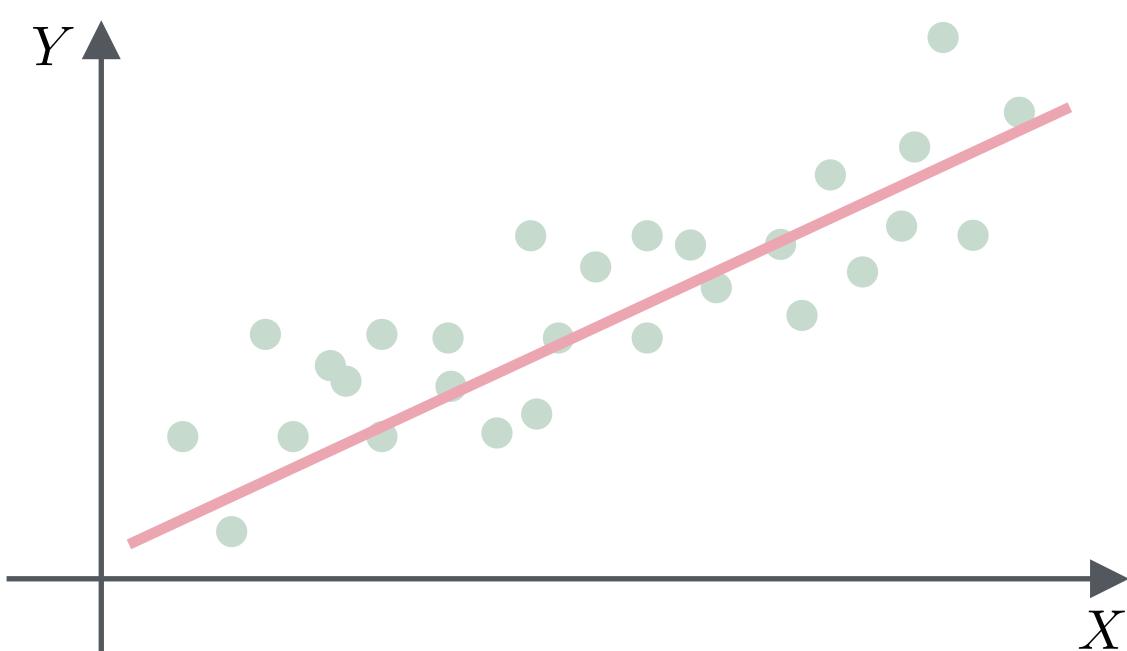
$$f(x) = w^\top x$$

$$\mathcal{H} = \mathbb{R}^2$$

## 2 classical examples

### ■ Linear regression

Linear regression seeks to estimate a linear relationship between a **continuous** scalar response  $Y$ , a one or more explanatory variables captured in a vector  $X$ .

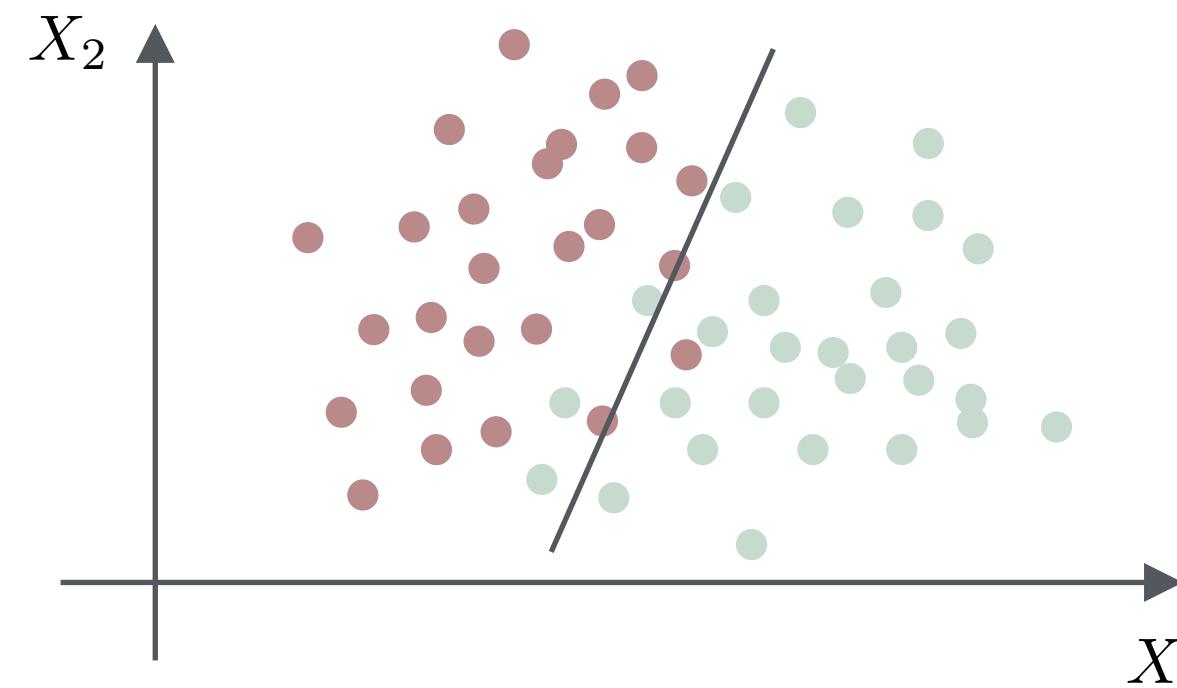


#### Ordinary Least Squares

$$\min_{w \in \mathbb{R}^2} \mathbb{E} [\|Y - Xw\|^2]$$

### ■ Logistic regression

Logistic regression seeks to a linear relationship between a discrete response  $Y$  taking values in some finite space  $\mathcal{Y}$ , and one or more explanatory variables captured in a vector  $X$ .



#### Logistic Regression

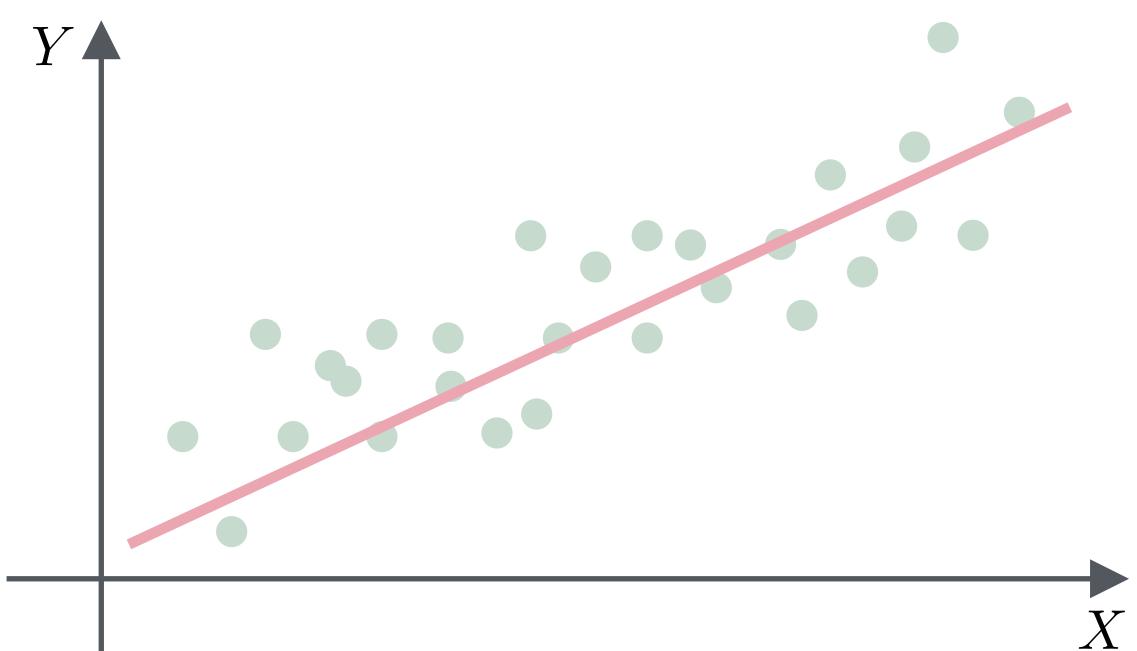
$$\min_{w \in \mathbb{R}^3} \mathbb{E} [\log(1 + \exp(-y w^T X))]$$

$$\begin{aligned} \min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))] & \quad \boxed{\mathcal{L}(y, \hat{y}) = (y - \hat{y})^2} \\ & \quad \boxed{f(x) = w^\top x} \\ & \quad \boxed{\mathcal{H} = \mathbb{R}^2} \end{aligned}$$

## 2 classical examples

### ■ Linear regression

Linear regression seeks to estimate a linear relationship between a **continuous** scalar response  $Y$ , a one or more explanatory variables captured in a vector  $X$ .

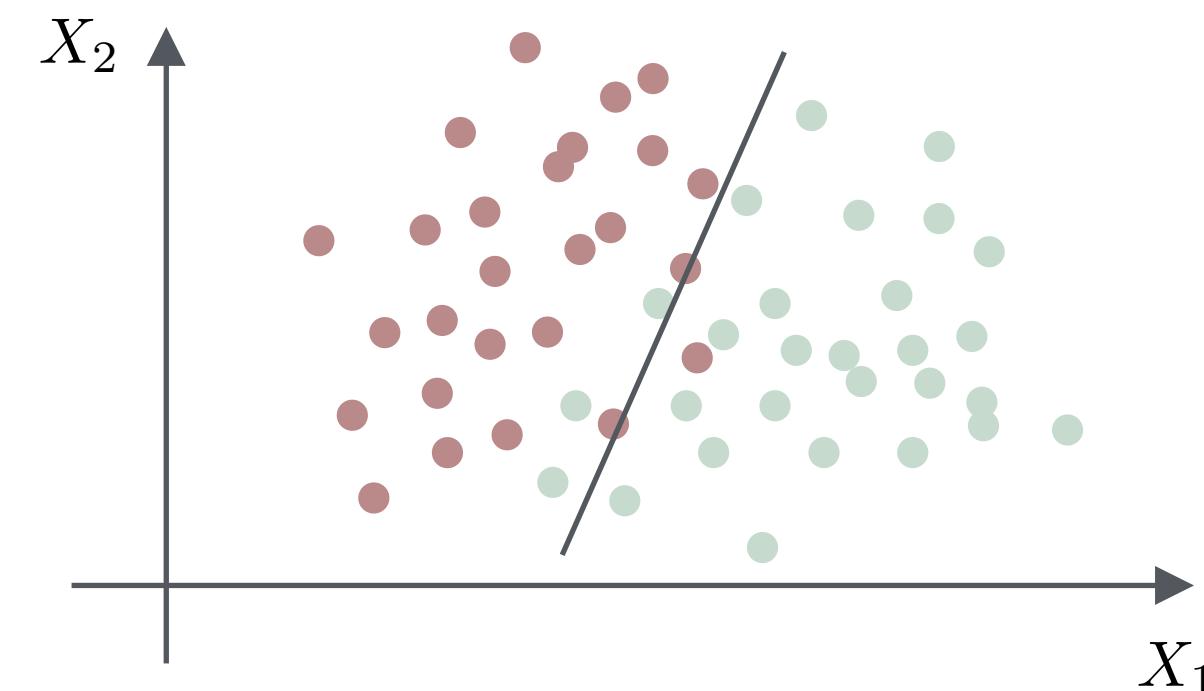


#### Ordinary Least Squares

$$\min_{w \in \mathbb{R}^2} \mathbb{E} [\|Y - Xw\|^2]$$

### ■ Logistic regression

Logistic regression seeks to a linear relationship between a discrete response  $Y$  taking values in some finite space  $\mathcal{Y}$ , and one or more explanatory variables captured in a vector  $X$ .



#### Logistic Regression

$$\min_{w \in \mathbb{R}^3} \mathbb{E} [\log(1 + \exp(-y w^T X))]$$

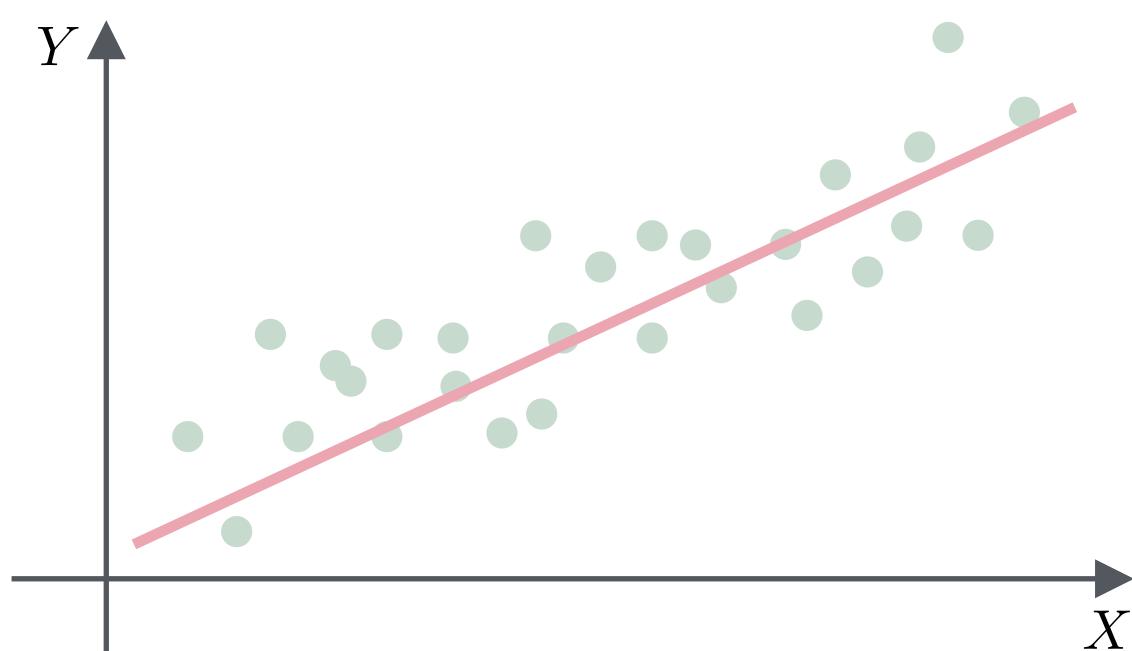
$$\begin{aligned} & \mathcal{L}(y, \hat{y}) = (y - \hat{y})^2 \\ & \min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))] \\ & f(x) = w^\top x \\ & \mathcal{H} = \mathbb{R}^2 \end{aligned}$$

$$\begin{aligned} & \min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))] \\ & \mathcal{H} = \mathbb{R}^3 \end{aligned}$$

## 2 classical examples

### ■ Linear regression

Linear regression seeks to estimate a linear relationship between a **continuous** scalar response  $Y$ , a one or more explanatory variables captured in a vector  $X$ .

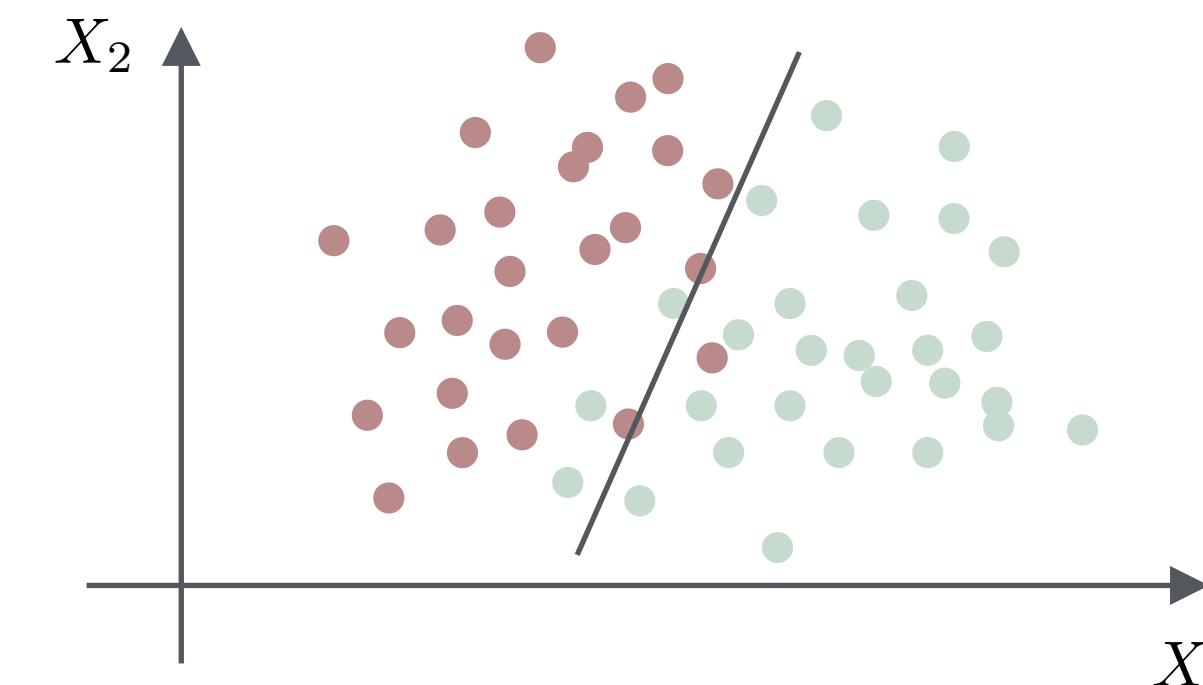


#### Ordinary Least Squares

$$\min_{w \in \mathbb{R}^2} \mathbb{E} [\|Y - Xw\|^2]$$

### ■ Logistic regression

Logistic regression seeks to a linear relationship between a discrete response  $Y$  taking values in some finite space  $\mathcal{Y}$ , and one or more explanatory variables captured in a vector  $X$ .



#### Logistic Regression

$$\min_{w \in \mathbb{R}^3} \mathbb{E} [\log(1 + \exp(-y w^T X))]$$

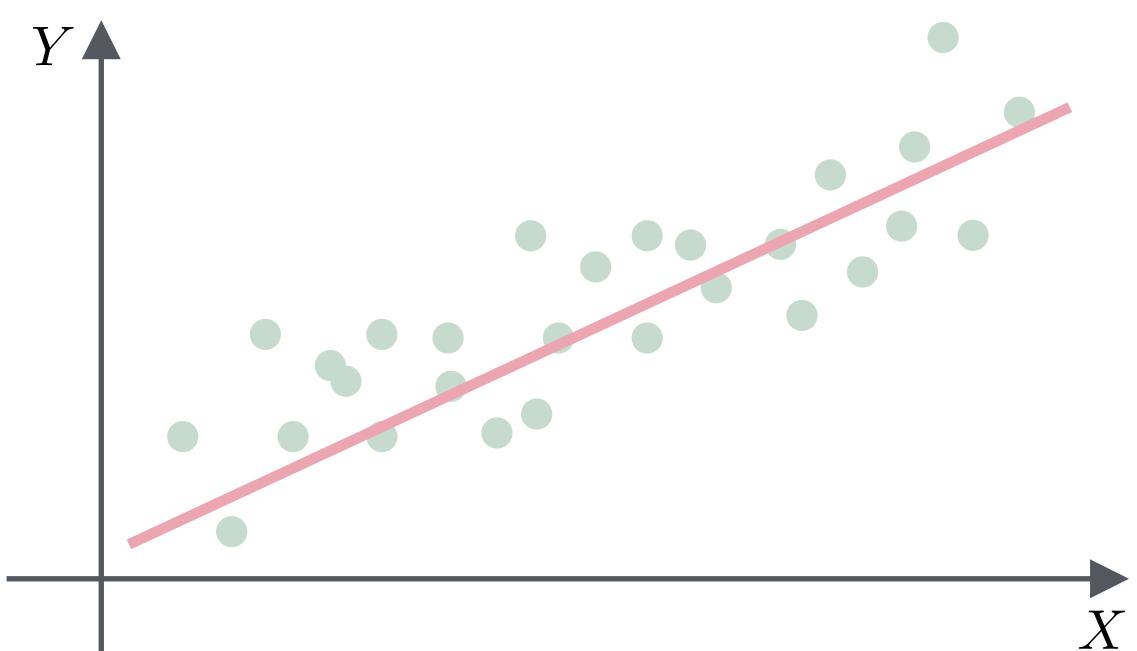
$$\begin{aligned} \min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))] & \quad \mathcal{L}(y, \hat{y}) = (y - \hat{y})^2 \\ f(x) &= w^\top x \\ \mathcal{H} &= \mathbb{R}^2 \end{aligned}$$

$$\begin{aligned} \min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))] & \quad \mathcal{L}(y, \hat{y}) = \log(1 + \exp(-y \hat{y})) \\ \mathcal{H} &= \mathbb{R}^3 \end{aligned}$$

## 2 classical examples

### ■ Linear regression

Linear regression seeks to estimate a linear relationship between a **continuous** scalar response  $Y$ , a one or more explanatory variables captured in a vector  $X$ .

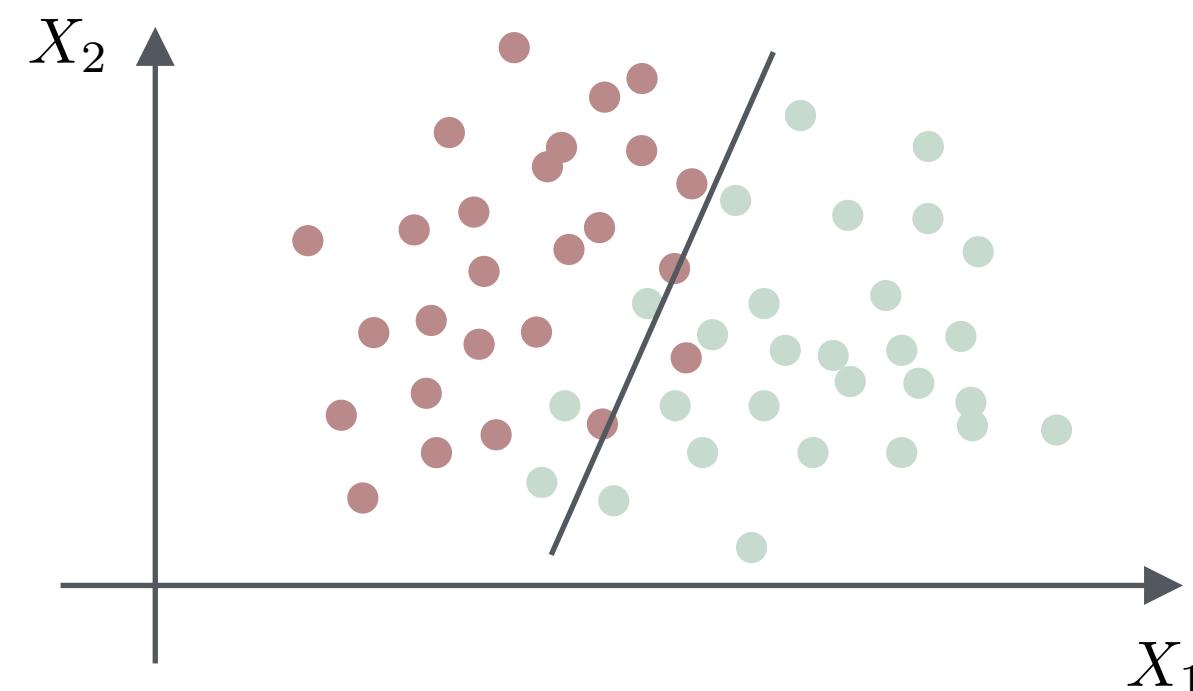


#### Ordinary Least Squares

$$\min_{w \in \mathbb{R}^2} \mathbb{E} [\|Y - Xw\|^2]$$

### ■ Logistic regression

Logistic regression seeks to a linear relationship between a discrete response  $Y$  taking values in some finite space  $\mathcal{Y}$ , and one or more explanatory variables captured in a vector  $X$ .



#### Logistic Regression

$$\min_{w \in \mathbb{R}^3} \mathbb{E} [\log(1 + \exp(-y w^T X))]$$

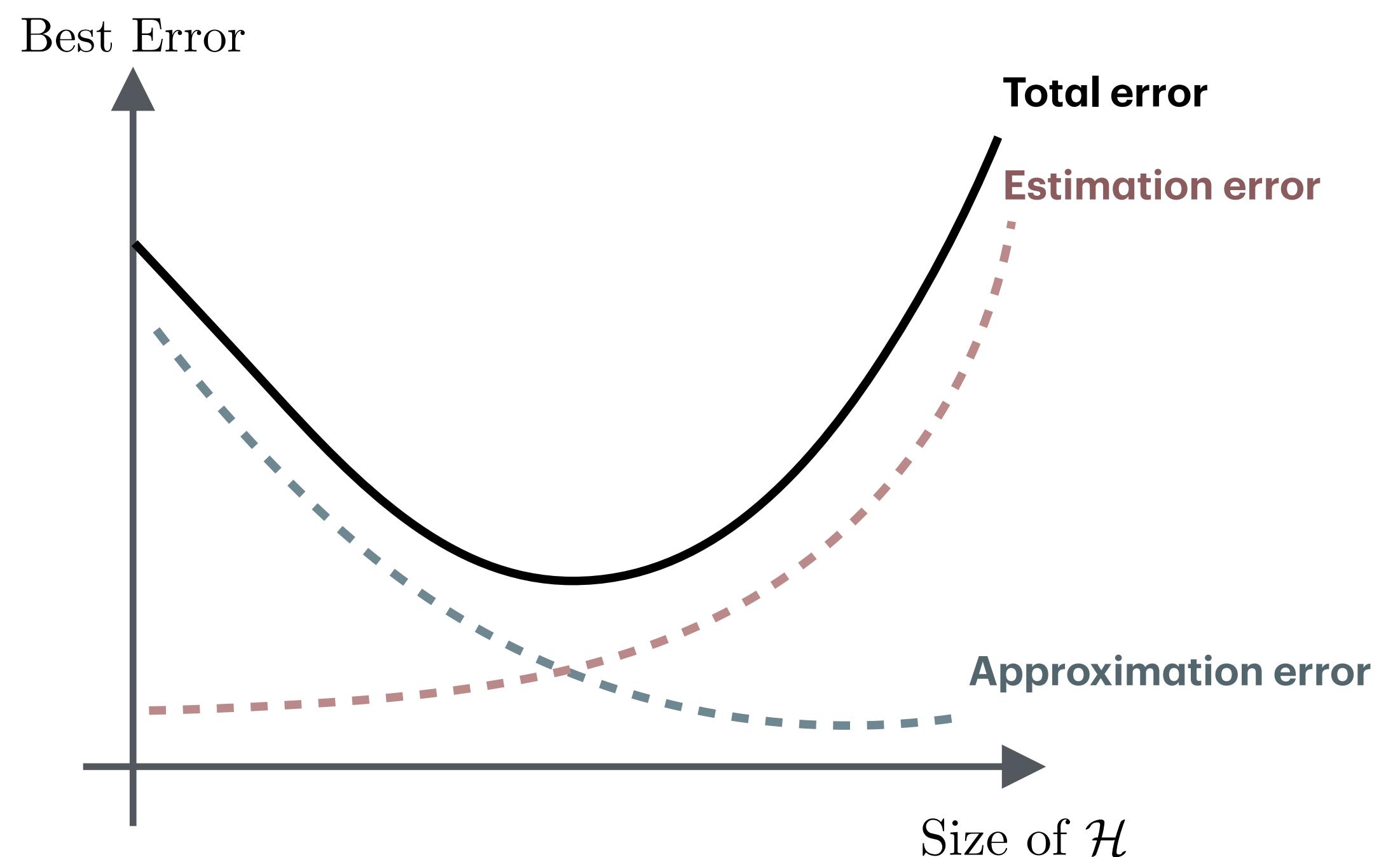
$$\begin{aligned} & \mathcal{L}(y, \hat{y}) = (y - \hat{y})^2 \\ \min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))] & \quad f(x) = w^\top x \\ & \mathcal{H} = \mathbb{R}^2 \end{aligned}$$

$$\begin{aligned} & \mathcal{L}(y, \hat{y}) = \log(1 + \exp(-y \hat{y})) \\ \min_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(y, f(x))] & \quad f(x) = w^\top x \\ & \mathcal{H} = \mathbb{R}^3 \end{aligned}$$

# **Implementation**



1. Supervised Learning Setting
2. Estimation vs Approximation



# **The supervised learning setup**

- Excess Risk

# The supervised learning setup

## ■ Excess Risk

### Definition

For any estimator  $f \in \mathcal{H}$ , we define the **excess risk** of  $f$ , denoted  $r(f)$  as

$$r(f) = \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)].$$

# The supervised learning setup

## ■ Excess Risk

### Definition

For any estimator  $f \in \mathcal{H}$ , we define the **excess risk** of  $f$ , denoted  $r(f)$  as

$$r(f) = \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)].$$

### Remarks :

- Thus, supervised learning consists in solving  $\inf_{f \in \mathcal{H}} r(f)$ .

# The supervised learning setup

## ■ Excess Risk

### Definition

For any estimator  $f \in \mathcal{H}$ , we define the **excess risk** of  $f$ , denoted  $r(f)$  as

$$r(f) = \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)].$$

### Remarks :

- Thus, supervised learning consists in solving  $\inf_{f \in \mathcal{H}} r(f)$ .
- Let's introduce as well  $f^*$ , and  $f^{**}$  as

$$f^* \in \operatorname{argmin}_{f \in \mathcal{H}} r(f) \quad f^{**} \in \operatorname{argmin}_{f \in \mathcal{F}} r(f)$$

where  $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y}\}$ .

# The supervised learning setup

## ■ Excess Risk

### Definition

For any estimator  $f \in \mathcal{H}$ , we define the **excess risk** of  $f$ , denoted  $r(f)$  as

$$r(f) = \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)].$$

### Remarks :

- Thus, supervised learning consists in solving  $\inf_{f \in \mathcal{H}} r(f)$ .
- Let's introduce as well  $f^*$ , and  $f^{**}$  as

$$f^* \in \operatorname{argmin}_{f \in \mathcal{H}} r(f) \quad f^{**} \in \operatorname{argmin}_{f \in \mathcal{F}} r(f)$$

where  $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y}\}$ .

- Clearly, we have  $r(f^{**}) \leq r(f^*)$ .

# The supervised learning setup

## ■ Excess Risk

### Definition

For any estimator  $f \in \mathcal{H}$ , we define the **excess risk** of  $f$ , denoted  $r(f)$  as

$$r(f) = \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)].$$

### Proposition : Bayes decision rule

We have, for all  $x \in \mathcal{X}$ ,

$$f^{**}(x) = \operatorname{argmin}_{\hat{y} \in \mathcal{Y}} \mathbb{E} [\ell(\hat{y}, Y) | X = x]$$

### Remarks :

- Thus, supervised learning consists in solving  $\inf_{f \in \mathcal{H}} r(f)$ .
- Let's introduce as well  $f^*$ , and  $f^{**}$  as

$$f^* \in \operatorname{argmin}_{f \in \mathcal{H}} r(f) \quad f^{**} \in \operatorname{argmin}_{f \in \mathcal{F}} r(f)$$

where  $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y}\}$ .

- Clearly, we have  $r(f^{**}) \leq r(f^*)$ .

# The supervised learning setup

## ■ Excess Risk

### Definition

For any estimator  $f \in \mathcal{H}$ , we define the **excess risk** of  $f$ , denoted  $r(f)$  as

$$r(f) = \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)].$$

### Proposition : Bayes decision rule

We have, for all  $x \in \mathcal{X}$ ,

$$f^{**}(x) = \operatorname{argmin}_{\hat{y} \in \mathcal{Y}} \mathbb{E} [\ell(\hat{y}, Y) | X = x]$$

**Proof :** on the black board.

### Remarks :

- Thus, supervised learning consists in solving  $\inf_{f \in \mathcal{H}} r(f)$ .
- Let's introduce as well  $f^*$ , and  $f^{**}$  as

$$f^* \in \operatorname{argmin}_{f \in \mathcal{H}} r(f) \quad f^{**} \in \operatorname{argmin}_{f \in \mathcal{F}} r(f)$$

where  $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y}\}$ .

- Clearly, we have  $r(f^{**}) \leq r(f^*)$ .

# The supervised learning setup

## ■ Excess Risk

### Definition

For any estimator  $f \in \mathcal{H}$ , we define the **excess risk** of  $f$ , denoted  $r(f)$  as

$$r(f) = \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)].$$

### Proposition : Bayes decision rule

We have, for all  $x \in \mathcal{X}$ ,

$$f^{**}(x) = \operatorname{argmin}_{\hat{y} \in \mathcal{Y}} \mathbb{E} [\ell(\hat{y}, Y) | X = x]$$

**Proof :** on the black board.

### Remarks :

- $f^{**}$  is often called **bayes estimator**.

### Remarks :

- Thus, supervised learning consists in solving  $\inf_{f \in \mathcal{H}} r(f)$ .
- Let's introduce as well  $f^*$ , and  $f^{**}$  as

$$f^* \in \operatorname{argmin}_{f \in \mathcal{H}} r(f) \quad f^{**} \in \operatorname{argmin}_{f \in \mathcal{F}} r(f)$$

where  $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y}\}$ .

- Clearly, we have  $r(f^{**}) \leq r(f^*)$ .

# The supervised learning setup

## ■ Excess Risk

### Definition

For any estimator  $f \in \mathcal{H}$ , we define the **excess risk** of  $f$ , denoted  $r(f)$  as

$$r(f) = \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)].$$

### Remarks :

- Thus, supervised learning consists in solving  $\inf_{f \in \mathcal{H}} r(f)$ .
- Let's introduce as well  $f^*$ , and  $f^{**}$  as

$$f^* \in \operatorname{argmin}_{f \in \mathcal{H}} r(f) \quad f^{**} \in \operatorname{argmin}_{f \in \mathcal{F}} r(f)$$

where  $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y}\}$ .

- Clearly, we have  $r(f^{**}) \leq r(f^*)$ .

### Proposition : Bayes decision rule

We have, for all  $x \in \mathcal{X}$ ,

$$f^{**}(x) = \operatorname{argmin}_{\hat{y} \in \mathcal{Y}} \mathbb{E} [\ell(\hat{y}, Y) | X = x]$$

**Proof :** on the black board.

### Remarks :

- $f^{**}$  is often called **bayes estimator**.
- Computing  $f^{**}$  requires the knowledge of the conditional distribution  $\mathbb{P}_{Y|X}$ . This is often impossible to obtain.

# The supervised learning setup

## ■ Excess Risk

### Definition

For any estimator  $f \in \mathcal{H}$ , we define the **excess risk** of  $f$ , denoted  $r(f)$  as

$$r(f) = \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)].$$

### Remarks :

- Thus, supervised learning consists in solving  $\inf_{f \in \mathcal{H}} r(f)$ .
- Let's introduce as well  $f^*$ , and  $f^{**}$  as

$$f^* \in \operatorname{argmin}_{f \in \mathcal{H}} r(f) \quad f^{**} \in \operatorname{argmin}_{f \in \mathcal{F}} r(f)$$

where  $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y}\}$ .

- Clearly, we have  $r(f^{**}) \leq r(f^*)$ .

### Proposition : Bayes decision rule

We have, for all  $x \in \mathcal{X}$ ,

$$f^{**}(x) = \operatorname{argmin}_{\hat{y} \in \mathcal{Y}} \mathbb{E} [\ell(\hat{y}, Y) | X = x]$$

**Proof :** on the black board.

### Remarks :

- $f^{**}$  is often called **bayes estimator**.
- Computing  $f^{**}$  requires the knowledge of the conditional distribution  $\mathbb{P}_{Y|X}$ . This is often impossible to obtain.
- In general, we suffer from an **estimation/approximation** tradeoff :

$$r(f) - \inf_{g \in \mathcal{F}} r(g) = r(f) - \inf_{h \in \mathcal{H}} r(h) + \inf_{h \in \mathcal{H}} r(h) - \inf_{g \in \mathcal{F}} r(g)$$

# The supervised learning setup

## ■ Excess Risk

### Definition

For any estimator  $f \in \mathcal{H}$ , we define the **excess risk** of  $f$ , denoted  $r(f)$  as

$$r(f) = \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)].$$

### Remarks :

- Thus, supervised learning consists in solving  $\inf_{f \in \mathcal{H}} r(f)$ .
- Let's introduce as well  $f^*$ , and  $f^{**}$  as

$$f^* \in \operatorname{argmin}_{f \in \mathcal{H}} r(f) \quad f^{**} \in \operatorname{argmin}_{f \in \mathcal{F}} r(f)$$

where  $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y}\}$ .

- Clearly, we have  $r(f^{**}) \leq r(f^*)$ .

### Proposition : Bayes decision rule

We have, for all  $x \in \mathcal{X}$ ,

$$f^{**}(x) = \operatorname{argmin}_{\hat{y} \in \mathcal{Y}} \mathbb{E} [\ell(\hat{y}, Y) | X = x]$$

**Proof :** on the black board.

### Remarks :

- $f^{**}$  is often called **bayes estimator**.
- Computing  $f^{**}$  requires the knowledge of the conditional distribution  $\mathbb{P}_{Y|X}$ . This is often impossible to obtain.
- In general, we suffer from an **estimation/approximation** tradeoff :

$$r(f) - \inf_{g \in \mathcal{F}} r(g) = r(f) - \inf_{h \in \mathcal{H}} r(h) + \inf_{h \in \mathcal{H}} r(h) - \inf_{g \in \mathcal{F}} r(g)$$

Estimation error

# The supervised learning setup

## ■ Excess Risk

### Definition

For any estimator  $f \in \mathcal{H}$ , we define the **excess risk** of  $f$ , denoted  $r(f)$  as

$$r(f) = \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)].$$

### Remarks :

- Thus, supervised learning consists in solving  $\inf_{f \in \mathcal{H}} r(f)$ .
- Let's introduce as well  $f^*$ , and  $f^{**}$  as

$$f^* \in \operatorname{argmin}_{f \in \mathcal{H}} r(f) \quad f^{**} \in \operatorname{argmin}_{f \in \mathcal{F}} r(f)$$

where  $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y}\}$ .

- Clearly, we have  $r(f^{**}) \leq r(f^*)$ .

### Proposition : Bayes decision rule

We have, for all  $x \in \mathcal{X}$ ,

$$f^{**}(x) = \operatorname{argmin}_{\hat{y} \in \mathcal{Y}} \mathbb{E} [\ell(\hat{y}, Y) | X = x]$$

**Proof :** on the black board.

### Remarks :

- $f^{**}$  is often called **bayes estimator**.
- Computing  $f^{**}$  requires the knowledge of the conditional distribution  $\mathbb{P}_{Y|X}$ . This is often impossible to obtain.
- In general, we suffer from an **estimation/approximation** tradeoff :

$$r(f) - \inf_{g \in \mathcal{F}} r(g) = r(f) - \inf_{h \in \mathcal{H}} r(h) + \inf_{h \in \mathcal{H}} r(h) - \inf_{g \in \mathcal{F}} r(g)$$

Estimation error
Approximation error

# The supervised learning setup

## ■ Excess Risk

### Definition

For any estimator  $f \in \mathcal{H}$ , we define the **excess risk** of  $f$ , denoted  $r(f)$  as

$$r(f) = \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)].$$

### Remarks :

- Thus, supervised learning consists in solving  $\inf_{f \in \mathcal{H}} r(f)$ .
- Let's introduce as well  $f^*$ , and  $f^{**}$  as

$$f^* \in \operatorname{argmin}_{f \in \mathcal{H}} r(f) \quad f^{**} \in \operatorname{argmin}_{f \in \mathcal{F}} r(f)$$

where  $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y}\}$ .

- Clearly, we have  $r(f^{**}) \leq r(f^*)$ .

### Proposition : Bayes decision rule

We have, for all  $x \in \mathcal{X}$ ,

$$f^{**}(x) = \operatorname{argmin}_{\hat{y} \in \mathcal{Y}} \mathbb{E} [\ell(\hat{y}, Y) | X = x]$$

**Proof :** on the black board.

### Remarks :

- $f^{**}$  is often called **bayes estimator**.
- Computing  $f^{**}$  requires the knowledge of the conditional distribution  $\mathbb{P}_{Y|X}$ . This is often impossible to obtain.
- In general, we suffer from an **estimation/approximation** tradeoff :

$$r(f) - \inf_{g \in \mathcal{F}} r(g) = r(f) - \inf_{h \in \mathcal{H}} r(h) + \inf_{h \in \mathcal{H}} r(h) - \inf_{g \in \mathcal{F}} r(g)$$

Estimation error
Approximation error

- And we need to work within the ERM framework, which leads us to work with the **empirical excess risk**  $R(f)$ , defined as

$$R(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(X_i), Y_i)$$

# The supervised learning setup

## ■ Excess Risk

### Definition

For any estimator  $f \in \mathcal{H}$ , we define the **excess risk** of  $f$ , denoted  $r(f)$  as

$$r(f) = \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)].$$

### Remarks :

- Thus, supervised learning consists in solving  $\inf_{f \in \mathcal{H}} r(f)$ .
- Let's introduce as well  $f^*$ , and  $f^{**}$  as

$$f^* \in \operatorname{argmin}_{f \in \mathcal{H}} r(f) \quad f^{**} \in \operatorname{argmin}_{f \in \mathcal{F}} r(f)$$

where  $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y}\}$ .

- Clearly, we have  $r(f^{**}) \leq r(f^*)$ .

### Proposition : Bayes decision rule

We have, for all  $x \in \mathcal{X}$ ,

$$f^{**}(x) = \operatorname{argmin}_{\hat{y} \in \mathcal{Y}} \mathbb{E} [\ell(\hat{y}, Y) | X = x]$$

**Proof :** on the black board.

### Remarks :

- $f^{**}$  is often called **bayes estimator**.
- Computing  $f^{**}$  requires the knowledge of the conditional distribution  $\mathbb{P}_{Y|X}$ . This is often impossible to obtain.
- In general, we suffer from an **estimation/approximation** tradeoff :

$$r(f) - \inf_{g \in \mathcal{F}} r(g) = r(f) - \inf_{h \in \mathcal{H}} r(h) + \inf_{h \in \mathcal{H}} r(h) - \inf_{g \in \mathcal{F}} r(g)$$

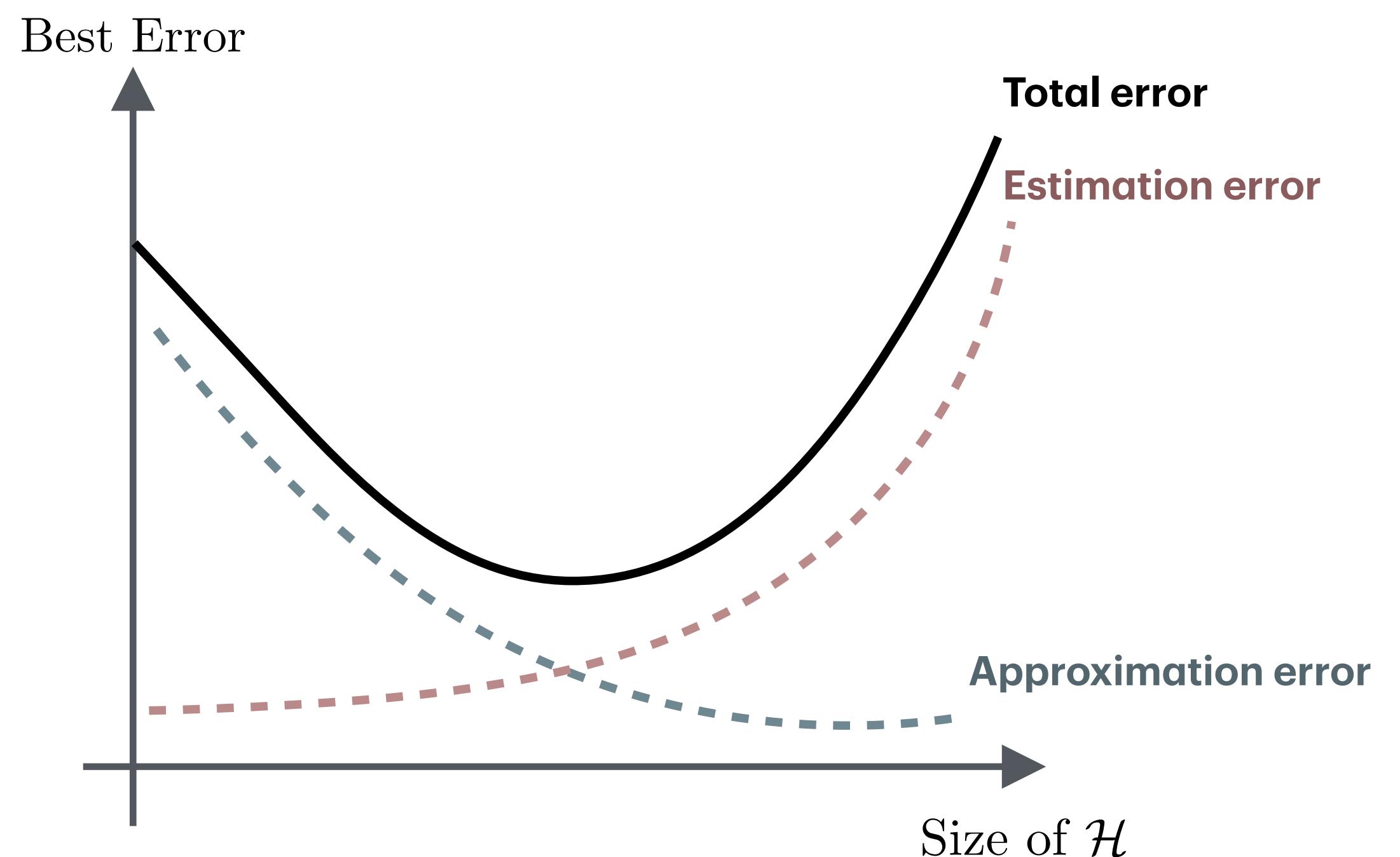
Estimation error
Approximation error

- And we need to work within the ERM framework, which leads us to work with the **empirical excess risk**  $R(f)$ , defined as

$$R(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(X_i), Y_i)$$

Note that here, the  $(X_i, Y_i)$  are **random samples**

1. Supervised Learning Setting
2. Estimation vs Approximation
3. Maximal inequalities



# Excess risk decompositon and maximal inequalities

## ■ Excess Risk Decomposition

### Proposition : Excess Risk Decompositon

For any  $f \in \mathcal{H}$ , we have

$$r(f) - r(f^*) \leq R(f) - R(\tilde{f}^*) + \sup_{f \in \mathcal{H}} \{R(f) - r(f)\} + \sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$$

$$\text{where } \tilde{f}^* = \operatorname{argmin}_{f \in \mathcal{H}} R(f) = \operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \ell(f(x_i), y_i).$$

# Excess risk decompositon and maximal inequalities

## ■ Excess Risk Decomposition

### Proposition : Excess Risk Decompositon

For any  $f \in \mathcal{H}$ , we have

$$r(f) - r(f^*) \leq R(f) - R(\tilde{f}^*) + \sup_{f \in \mathcal{H}} \{R(f) - r(f)\} + \sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$$

$$\text{where } \tilde{f}^* = \operatorname{argmin}_{f \in \mathcal{H}} R(f) = \operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \ell(f(x_i), y_i).$$

### Remarks :

- $R(f) - R(\tilde{f}^*)$  may be reduced through optimization methods like gradient descent.

# Excess risk decompositon and maximal inequalities

## ■ Excess Risk Decomposition

### Proposition : Excess Risk Decompositon

For any  $f \in \mathcal{H}$ , we have

$$r(f) - r(f^*) \leq R(f) - R(\tilde{f}^*) + \sup_{f \in \mathcal{H}} \{R(f) - r(f)\} + \sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$$

$$\text{where } \tilde{f}^* = \operatorname{argmin}_{f \in \mathcal{H}} R(f) = \operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \ell(f(x_i), y_i).$$

### Remarks :

- $R(f) - R(\tilde{f}^*)$  may be reduced through optimization methods like gradient descent.
- In this chapter, we will focus on the next terms,  
 $\sup_{f \in \mathcal{H}} \{R(f) - r(f)\}$ , and  $\sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$ .

# Excess risk decompositon and maximal inequalities

## ■ Excess Risk Decomposition

### Proposition : Excess Risk Decompositon

For any  $f \in \mathcal{H}$ , we have

$$r(f) - r(f^*) \leq R(f) - R(\tilde{f}^*) + \sup_{f \in \mathcal{H}} \{R(f) - r(f)\} + \sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$$

where  $\tilde{f}^* = \underset{f \in \mathcal{H}}{\operatorname{argmin}} R(f) = \underset{f \in \mathcal{H}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^N \ell(f(x_i), y_i)$ .

### Remarks :

- $R(f) - R(\tilde{f}^*)$  may be reduced through optimization methods like gradient descent.
- In this chapter, we will focus on the next terms,  $\sup_{f \in \mathcal{H}} \{R(f) - r(f)\}$ , and  $\sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$ .
- Our goal is to derive **maximal inequalities**, i.e. inequalities of the form

$$\mathbb{E} \left[ \sup_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \frac{\psi(\mathcal{H})}{n^\alpha}$$

where  $\psi$ , and  $\alpha$  are to determine.

# Excess risk decompositon and maximal inequalities

## ■ Excess Risk Decomposition

### Proposition : Excess Risk Decompositon

For any  $f \in \mathcal{H}$ , we have

$$r(f) - r(f^*) \leq R(f) - R(\tilde{f}^*) + \sup_{f \in \mathcal{H}} \{R(f) - r(f)\} + \sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$$

where  $\tilde{f}^* = \underset{f \in \mathcal{H}}{\operatorname{argmin}} R(f) = \underset{f \in \mathcal{H}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^N \ell(f(x_i), y_i)$ .

## ■ The case $\mathcal{H}$ finite.

### ■ Let us first consider the case $\mathcal{H}$ finite.

#### Remarks :

- $R(f) - R(\tilde{f}^*)$  may be reduced through optimization methods like gradient descent.
- In this chapter, we will focus on the next terms,  $\sup_{f \in \mathcal{H}} \{R(f) - r(f)\}$ , and  $\sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$ .
- Our goal is to derive **maximal inequalities**, i.e. inequalities of the form

$$\mathbb{E} \left[ \sup_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \frac{\psi(\mathcal{H})}{n^\alpha}$$

where  $\psi$ , and  $\alpha$  are to determine.

# Excess risk decompositon and maximal inequalities

## ■ Excess Risk Decomposition

### Proposition : Excess Risk Decompositon

For any  $f \in \mathcal{H}$ , we have

$$r(f) - r(f^*) \leq R(f) - R(\tilde{f}^*) + \sup_{f \in \mathcal{H}} \{R(f) - r(f)\} + \sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$$

where  $\tilde{f}^* = \underset{f \in \mathcal{H}}{\operatorname{argmin}} R(f) = \underset{f \in \mathcal{H}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^N \ell(f(x_i), y_i)$ .

## ■ The case $\mathcal{H}$ finite.

- Let us first consider the case  $\mathcal{H}$  finite.
- As a consequence of the law of large number, we know that  $R(f) \xrightarrow{n \rightarrow \infty} r(f)$  almost surely, for all  $f \in \mathcal{H}$ .

### Remarks :

- $R(f) - R(\tilde{f}^*)$  may be reduced through optimization methods like gradient descent.
- In this chapter, we will focus on the next terms,  $\sup_{f \in \mathcal{H}} \{R(f) - r(f)\}$ , and  $\sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$ .
- Our goal is to derive **maximal inequalities**, i.e. inequalities of the form

$$\mathbb{E} \left[ \sup_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \frac{\psi(\mathcal{H})}{n^\alpha}$$

where  $\psi$ , and  $\alpha$  are to determine.

# Excess risk decompositon and maximal inequalities

## ■ Excess Risk Decomposition

### Proposition : Excess Risk Decompositon

For any  $f \in \mathcal{H}$ , we have

$$r(f) - r(f^*) \leq R(f) - R(\tilde{f}^*) + \sup_{f \in \mathcal{H}} \{R(f) - r(f)\} + \sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$$

where  $\tilde{f}^* = \operatorname{argmin}_{f \in \mathcal{H}} R(f) = \operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \ell(f(x_i), y_i)$ .

## ■ The case $\mathcal{H}$ finite.

- Let us first consider the case  $\mathcal{H}$  finite.
- As a consequence of the law of large number, we know that  $R(f) \xrightarrow{n \rightarrow \infty} r(f)$  almost surely, for all  $f \in \mathcal{H}$ .
- Here we seek to establish a **non-asymptotic** bound that holds uniformly over  $\mathcal{H}$ .

### Remarks :

- $R(f) - R(\tilde{f}^*)$  may be reduced through optimization methods like gradient descent.
- In this chapter, we will focus on the next terms,  $\sup_{f \in \mathcal{H}} \{R(f) - r(f)\}$ , and  $\sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$ .
- Our goal is to derive **maximal inequalities**, i.e. inequalities of the form

$$\mathbb{E} \left[ \sup_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \frac{\psi(\mathcal{H})}{n^\alpha}$$

where  $\psi$ , and  $\alpha$  are to determine.

# Excess risk decompositon and maximal inequalities

## ■ Excess Risk Decomposition

### Proposition : Excess Risk Decompositon

For any  $f \in \mathcal{H}$ , we have

$$r(f) - r(f^*) \leq R(f) - R(\tilde{f}^*) + \sup_{f \in \mathcal{H}} \{R(f) - r(f)\} + \sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$$

where  $\tilde{f}^* = \operatorname{argmin}_{f \in \mathcal{H}} R(f) = \operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \ell(f(x_i), y_i)$ .

### Remarks :

- $R(f) - R(\tilde{f}^*)$  may be reduced through optimization methods like gradient descent.
- In this chapter, we will focus on the next terms,  $\sup_{f \in \mathcal{H}} \{R(f) - r(f)\}$ , and  $\sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$ .
- Our goal is to derive **maximal inequalities**, i.e. inequalities of the form

$$\mathbb{E} \left[ \sup_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \frac{\psi(\mathcal{H})}{n^\alpha}$$

where  $\psi$ , and  $\alpha$  are to determine.

## ■ The case $\mathcal{H}$ finite.

- Let us first consider the case  $\mathcal{H}$  finite.
- As a consequence of the law of large number, we know that  $R(f) \xrightarrow{n \rightarrow \infty} r(f)$  almost surely, for all  $f \in \mathcal{H}$ .
- Here we seek to establish a **non-asymptotic** bound that holds uniformly over  $\mathcal{H}$ .

### Proposition : Hoeffding's inequality

Let  $X$  be a real-valued random variable such that  $a \leq X - \mathbb{E}[X] \leq b$ . Then for any  $\lambda \in \mathbb{R}$ , we have

$$\mathbb{E} \left[ e^{\lambda(X - \mathbb{E}[X])} \right] \leq e^{\lambda^2 \frac{(b-a)^2}{8}}.$$

# Excess risk decompositon and maximal inequalities

## ■ Excess Risk Decomposition

### Proposition : Excess Risk Decompositon

For any  $f \in \mathcal{H}$ , we have

$$r(f) - r(f^*) \leq R(f) - R(\tilde{f}^*) + \sup_{f \in \mathcal{H}} \{R(f) - r(f)\} + \sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$$

where  $\tilde{f}^* = \operatorname{argmin}_{f \in \mathcal{H}} R(f) = \operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \ell(f(x_i), y_i)$ .

### Remarks :

- $R(f) - R(\tilde{f}^*)$  may be reduced through optimization methods like gradient descent.
- In this chapter, we will focus on the next terms,  $\sup_{f \in \mathcal{H}} \{R(f) - r(f)\}$ , and  $\sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$ .
- Our goal is to derive **maximal inequalities**, i.e. inequalities of the form

$$\mathbb{E} \left[ \sup_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \frac{\psi(\mathcal{H})}{n^\alpha}$$

where  $\psi$ , and  $\alpha$  are to determine.

## ■ The case $\mathcal{H}$ finite.

- Let us first consider the case  $\mathcal{H}$  finite.
- As a consequence of the law of large number, we know that  $R(f) \xrightarrow{n \rightarrow \infty} r(f)$  almost surely, for all  $f \in \mathcal{H}$ .
- Here we seek to establish a **non-asymptotic** bound that holds uniformly over  $\mathcal{H}$ .

### Proposition : Hoeffding's inequality

Let  $X$  be a real-valued random variable such that  $a \leq X - \mathbb{E}[X] \leq b$ . Then for any  $\lambda \in \mathbb{R}$ , we have

$$\mathbb{E} \left[ e^{\lambda(X - \mathbb{E}[X])} \right] \leq e^{\lambda^2 \frac{(b-a)^2}{8}}.$$

**Proof :** on the black board.

# Excess risk decompositon and maximal inequalities

## ■ Excess Risk Decomposition

### Proposition : Excess Risk Decompositon

For any  $f \in \mathcal{H}$ , we have

$$r(f) - r(f^*) \leq R(f) - R(\tilde{f}^*) + \sup_{f \in \mathcal{H}} \{R(f) - r(f)\} + \sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$$

where  $\tilde{f}^* = \operatorname{argmin}_{f \in \mathcal{H}} R(f) = \operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \ell(f(x_i), y_i)$ .

### Remarks :

- $R(f) - R(\tilde{f}^*)$  may be reduced through optimization methods like gradient descent.
- In this chapter, we will focus on the next terms,  $\sup_{f \in \mathcal{H}} \{R(f) - r(f)\}$ , and  $\sup_{f \in \mathcal{H}} \{r(f) - R(f)\}$ .
- Our goal is to derive **maximal inequalities**, i.e. inequalities of the form

$$\mathbb{E} \left[ \sup_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \frac{\psi(\mathcal{H})}{n^\alpha}$$

where  $\psi$ , and  $\alpha$  are to determine.

## ■ The case $\mathcal{H}$ finite.

- Let us first consider the case  $\mathcal{H}$  finite.
- As a consequence of the law of large number, we know that  $R(f) \xrightarrow{n \rightarrow \infty} r(f)$  almost surely, for all  $f \in \mathcal{H}$ .
- Here we seek to establish a **non-asymptotic** bound that holds uniformly over  $\mathcal{H}$ .

### Proposition : Hoeffding's inequality

Let  $X$  be a real-valued random variable such that  $a \leq X - \mathbb{E}[X] \leq b$ . Then for any  $\lambda \in \mathbb{R}$ , we have

$$\mathbb{E} \left[ e^{\lambda(X - \mathbb{E}[X])} \right] \leq e^{\lambda^2 \frac{(b-a)^2}{8}}.$$

**Proof :** on the black board.

### Remarks :

- Hoeffding's inequality is a classical example of concentration inequality.