

# Statistical Learning

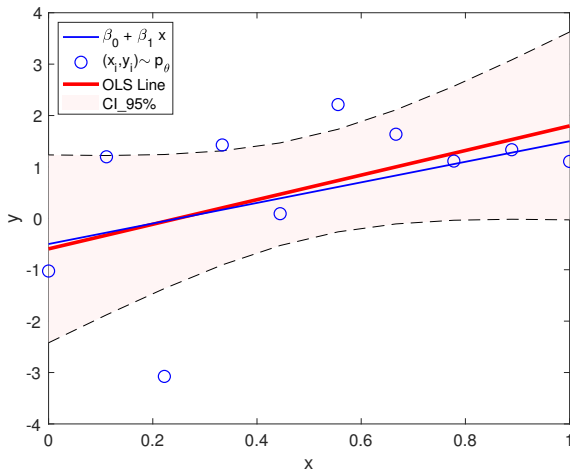
Master Spécialisé Intelligence Artificielle de Confiance (IAC)  
@ Centrale Supélec en partenariat avec l'IRT SystemX  
2024/2025.

FAÏCEL CHAMROUKHI



 [chamroukhi.com](http://chamroukhi.com)

## ■ Regression



- The data are represented by a random pair  $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$  where  $\mathbf{X}$  is a vector of descriptors for some variable of interest  $Y$
  - The objective is **Prediction**, i.e. to seek for a prediction function  $h : \mathcal{X} \rightarrow \mathcal{Y}$  for which  $\hat{y} = h(\mathbf{x})$  is a good approximation of the true output  $y$
  - In a **regression** problem : typically  $\mathbf{X} \in \mathbb{R}^p$ ,  $Y \in \mathcal{Y} = \mathbb{R}^d$
- ↪ We will mainly focus on parametric probabilistic models of the form

$$Y = h(X) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr.  $P(Y|X, h)$  can be computed in terms of  $P_\theta(Y - h(X))$ .  
 $P(X, Y|h) = P(Y|X, h)P(X|h)$

- Data : a random sample  $(\mathbf{X}_i, Y_i)_{i=1}^n$  with observed values  $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- **Data-Scientist's role** : given the **data**, choose a **prediction function**  $h$  that attempts to “minimize” the prediction error for of all possible data (**risk**)  $R(h)$ , under a chosen **loss** function  $\ell$  measuring the error of predicting  $Y$  by  $h(X)$ .
  - ↪ minimize the **empirical** (data- $\mathcal{D}_n$ -driven) **risk**  $R_n(h)$
  - ↪ Minimizing  $R_n(h)$  may require an optimization **algorithm**  $\mathcal{A}$
- Data-Scientist's “**Toolbox**” : {Data, hypothesis, loss, risk, algorithm}

- The data are represented by a random pair  $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$  where  $\mathbf{X}$  is a vector of descriptors for some variable of interest  $Y$
  - The objective is **Prediction**, i.e. to seek for a prediction function  $h : \mathcal{X} \rightarrow \mathcal{Y}$  for which  $\hat{y} = h(\mathbf{x})$  is a good approximation of the true output  $y$
  - In a **regression** problem : typically  $\mathbf{X} \in \mathbb{R}^p$ ,  $Y \in \mathcal{Y} = \mathbb{R}^d$
- ↪ We will mainly focus on parametric probabilistic models of the form

$$Y = h(X) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr.  $P(Y|X, h)$  can be computed in terms of  $P_\theta(Y - h(X))$ .  
 $P(X, Y|h) = P(Y|X, h)P(X|h)$

- Data : a random sample  $(\mathbf{X}_i, Y_i)_{i=1}^n$  with observed values  $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- **Data-Scientist's role** : given the **data**, choose a **prediction function**  $h$  that attempts to “minimize” the prediction error for of all possible data (**risk**)  $R(h)$ , under a chosen **loss** function  $\ell$  measuring the error of predicting  $Y$  by  $h(X)$ .
  - ↪ minimize the **empirical** (data- $\mathcal{D}_n$ -driven) **risk**  $R_n(h)$
  - ↪ Minimizing  $R_n(h)$  may require an optimization **algorithm**  $\mathcal{A}$
- Data-Scientist's “**Toolbox**” : {Data, hypothesis, loss, risk, algorithm}

- The data are represented by a random pair  $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$  where  $\mathbf{X}$  is a vector of descriptors for some variable of interest  $Y$
- The objective is **Prediction**, i.e. to seek for a prediction function  $h : \mathcal{X} \rightarrow \mathcal{Y}$  for which  $\hat{y} = h(\mathbf{x})$  is a good approximation of the true output  $y$
- In a **regression** problem : typically  $\mathbf{X} \in \mathbb{R}^p$ ,  $Y \in \mathcal{Y} = \mathbb{R}^d$

↪ We will mainly focus on parametric probabilistic models of the form

$$Y = h(X) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr.  $P(Y|X, h)$  can be computed in terms of  $P_\theta(Y - h(X))$ .  
 $P(X, Y|h) = P(Y|X, h)P(X|h)$

- Data : a random sample  $(\mathbf{X}_i, Y_i)_{i=1}^n$  with observed values  $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- **Data-Scientist's role** : given the **data**, choose a **prediction function**  $h$  that attempts to “minimize” the prediction error for of all possible data (**risk**)  $R(h)$ , under a chosen **loss** function  $\ell$  measuring the error of predicting  $Y$  by  $h(X)$ .

↪ minimize the **empirical** (data- $\mathcal{D}_n$ -driven) **risk**  $R_n(h)$

↪ Minimizing  $R_n(h)$  may require an optimization **algorithm**  $\mathcal{A}$

- Data-Scientist's “**Toolbox**” : {Data, hypothesis, loss, risk, algorithm}

- The data are represented by a random pair  $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$  where  $\mathbf{X}$  is a vector of descriptors for some variable of interest  $Y$
  - The objective is **Prediction**, i.e. to seek for a prediction function  $h : \mathcal{X} \rightarrow \mathcal{Y}$  for which  $\hat{y} = h(\mathbf{x})$  is a good approximation of the true output  $y$
  - In a **regression** problem : typically  $\mathbf{X} \in \mathbb{R}^p$ ,  $Y \in \mathcal{Y} = \mathbb{R}^d$
- ↪ We will mainly focus on parametric probabilistic models of the form

$$Y = h(\mathbf{X}) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr.  $P(Y|\mathbf{X}, h)$  can be computed in terms of  $P_\theta(Y - h(\mathbf{X}))$ .  
 $P(\mathbf{X}, Y|h) = P(Y|\mathbf{X}, h)P(\mathbf{X}|h)$

- Data : a random sample  $(\mathbf{X}_i, Y_i)_{i=1}^n$  with observed values  $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- Data-Scientist's role : given the **data**, choose a **prediction function**  $h$  that attempts to "minimize" the prediction error for of all possible data (**risk**)  $R(h)$ , under a chosen **loss** function  $\ell$  measuring the error of predicting  $Y$  by  $h(\mathbf{X})$ .
  - ↪ minimize the **empirical** (data- $\mathcal{D}_n$ -driven) **risk**  $R_n(h)$
  - ↪ Minimizing  $R_n(h)$  may require an optimization **algorithm**  $\mathcal{A}$
- Data-Scientist's "**Toolbox**" : {Data, hypothesis, loss, risk, algorithm}

- The data are represented by a random pair  $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$  where  $\mathbf{X}$  is a vector of descriptors for some variable of interest  $Y$
  - The objective is **Prediction**, i.e. to seek for a prediction function  $h : \mathcal{X} \rightarrow \mathcal{Y}$  for which  $\hat{y} = h(\mathbf{x})$  is a good approximation of the true output  $y$
  - In a **regression** problem : typically  $\mathbf{X} \in \mathbb{R}^p$ ,  $Y \in \mathcal{Y} = \mathbb{R}^d$
- ↪ We will mainly focus on parametric probabilistic models of the form

$$Y = h(\mathbf{X}) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr.  $P(Y|\mathbf{X}, h)$  can be computed in terms of  $P_\theta(Y - h(\mathbf{X}))$ .  
 $P(\mathbf{X}, Y|h) = P(Y|\mathbf{X}, h)P(\mathbf{X}|h)$

- **Data** : a random sample  $(\mathbf{X}_i, Y_i)_{i=1}^n$  with observed values  $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- **Data-Scientist's role** : given the **data**, choose a **prediction function**  $h$  that attempts to “minimize” the prediction error for of all possible data (**risk**)  $R(h)$ , under a chosen **loss** function  $\ell$  measuring the error of predicting  $Y$  by  $h(\mathbf{X})$ .
  - ↪ minimize the **empirical** (data- $\mathcal{D}_n$ -driven) **risk**  $R_n(h)$
  - ↪ Minimizing  $R_n(h)$  may require an optimization **algorithm**  $\mathcal{A}$
- Data-Scientist's “**Toolbox**” : {Data, hypothesis, loss, risk, algorithm}

- The data are represented by a random pair  $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$  where  $\mathbf{X}$  is a vector of descriptors for some variable of interest  $Y$
  - The objective is **Prediction**, i.e. to seek for a prediction function  $h : \mathcal{X} \rightarrow \mathcal{Y}$  for which  $\hat{y} = h(\mathbf{x})$  is a good approximation of the true output  $y$
  - In a **regression** problem : typically  $\mathbf{X} \in \mathbb{R}^p$ ,  $Y \in \mathcal{Y} = \mathbb{R}^d$
- ↪ We will mainly focus on parametric probabilistic models of the form

$$Y = h(\mathbf{X}) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr.  $P(Y|\mathbf{X}, h)$  can be computed in terms of  $P_\theta(Y - h(\mathbf{X}))$ .  
 $P(\mathbf{X}, Y|h) = P(Y|\mathbf{X}, h)P(\mathbf{X}|h)$

- Data : a random sample  $(\mathbf{X}_i, Y_i)_{i=1}^n$  with observed values  $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- **Data-Scientist's role** : given the **data**, choose a **prediction function**  $h$  that attempts to “minimize” the prediction error for of all possible data (**risk**)  $R(h)$ , under a chosen **loss** function  $\ell$  measuring the error of predicting  $Y$  by  $h(\mathbf{X})$ .
  - ↪ minimize the **empirical** (data- $\mathcal{D}_n$ -driven) **risk**  $R_n(h)$
  - ↪ Minimizing  $R_n(h)$  may require an optimization **algorithm**  $\mathcal{A}$
- Data-Scientist's “**Toolbox**” : {Data, hypothesis, loss, risk, algorithm}



- The data are represented by a random pair  $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$  where  $\mathbf{X}$  is a vector of descriptors for some variable of interest  $Y$
  - The objective is **Prediction**, i.e. to seek for a prediction function  $h : \mathcal{X} \rightarrow \mathcal{Y}$  for which  $\hat{y} = h(\mathbf{x})$  is a good approximation of the true output  $y$
  - In a **regression** problem : typically  $\mathbf{X} \in \mathbb{R}^p$ ,  $Y \in \mathcal{Y} = \mathbb{R}^d$
- ↪ We will mainly focus on parametric probabilistic models of the form

$$Y = h(\mathbf{X}) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr.  $P(Y|\mathbf{X}, h)$  can be computed in terms of  $P_\theta(Y - h(\mathbf{X}))$ .  
 $P(\mathbf{X}, Y|h) = P(Y|\mathbf{X}, h)P(\mathbf{X}|h)$

- Data : a random sample  $(\mathbf{X}_i, Y_i)_{i=1}^n$  with observed values  $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- **Data-Scientist's role** : given the **data**, choose a **prediction function**  $h$  that attempts to “minimize” the prediction error for of all possible data (**risk**)  $R(h)$ , under a chosen **loss** function  $\ell$  measuring the error of predicting  $Y$  by  $h(\mathbf{X})$ .
  - ↪ minimize the **empirical** (data- $\mathcal{D}_n$ -driven) **risk**  $R_n(h)$
  - ↪ Minimizing  $R_n(h)$  may require an optimization **algorithm**  $\mathcal{A}$
- Data-Scientist's “**Toolbox**” : {Data, hypothesis, loss, risk, algorithm}

- The data are represented by a random pair  $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$  where  $\mathbf{X}$  is a vector of descriptors for some variable of interest  $Y$
  - The objective is **Prediction**, i.e. to seek for a prediction function  $h : \mathcal{X} \rightarrow \mathcal{Y}$  for which  $\hat{y} = h(\mathbf{x})$  is a good approximation of the true output  $y$
  - In a **regression** problem : typically  $\mathbf{X} \in \mathbb{R}^p$ ,  $Y \in \mathcal{Y} = \mathbb{R}^d$
- ↪ We will mainly focus on parametric probabilistic models of the form

$$Y = h(\mathbf{X}) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr.  $P(Y|\mathbf{X}, h)$  can be computed in terms of  $P_\theta(Y - h(\mathbf{X}))$ .  
 $P(\mathbf{X}, Y|h) = P(Y|\mathbf{X}, h)P(\mathbf{X}|h)$

- Data : a random sample  $(\mathbf{X}_i, Y_i)_{i=1}^n$  with observed values  $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- **Data-Scientist's role** : given the **data**, choose a **prediction function**  $h$  that attempts to “minimize” the prediction error for of all possible data (**risk**)  $R(h)$ , under a chosen **loss** function  $\ell$  measuring the error of predicting  $Y$  by  $h(\mathbf{X})$ .
  - ↪ minimize the **empirical** (data- $\mathcal{D}_n$ -driven) **risk**  $R_n(h)$
  - ↪ Minimizing  $R_n(h)$  may require an optimization **algorithm**  $\mathcal{A}$
- Data-Scientist's “**Toolbox**” : {Data, hypothesis, loss, risk, algorithm}

- The data are represented by a random pair  $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$  where  $\mathbf{X}$  is a vector of descriptors for some variable of interest  $Y$
  - The objective is **Prediction**, i.e. to seek for a prediction function  $h : \mathcal{X} \rightarrow \mathcal{Y}$  for which  $\hat{y} = h(\mathbf{x})$  is a good approximation of the true output  $y$
  - In a **regression** problem : typically  $\mathbf{X} \in \mathbb{R}^p$ ,  $Y \in \mathcal{Y} = \mathbb{R}^d$
- ↪ We will mainly focus on parametric probabilistic models of the form

$$Y = h(\mathbf{X}) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr.  $P(Y|\mathbf{X}, h)$  can be computed in terms of  $P_\theta(Y - h(\mathbf{X}))$ .  
 $P(\mathbf{X}, Y|h) = P(Y|\mathbf{X}, h)P(\mathbf{X}|h)$

- Data : a random sample  $(\mathbf{X}_i, Y_i)_{i=1}^n$  with observed values  $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- **Data-Scientist's role** : given the **data**, choose a **prediction function**  $h$  that attempts to “minimize” the prediction error for of all possible data (**risk**)  $R(h)$ , under a chosen **loss** function  $\ell$  measuring the error of predicting  $Y$  by  $h(\mathbf{X})$ .
  - ↪ minimize the **empirical** (data- $\mathcal{D}_n$ -driven) **risk**  $R_n(h)$
  - ↪ Minimizing  $R_n(h)$  may require an optimization **algorithm**  $\mathcal{A}$
- Data-Scientist's “**Toolbox**” : {Data, hypothesis, loss, risk, algorithm}

## Objective

Regression models the relationship between two variables  $X$  and  $Y$

- Temperature ( $Y$ ) of some water source, given the air temperate ( $x$ )
- Price ( $Y$ ) of an apartment given its surface ( $x_1$ ) and number of rooms ( $x_2$ )

Vocabulary :

- The  $x$ 's are called inputs/predictors/covariates/features/descriptors/exogenous/Explanatory/independent variables
- The  $y$ 's are called output/outcome/response/endogenous/variable of interest/Explained/dependent variable
- **Simple** regression :  $x \in \mathbb{R}$
- **Univariate** regression :  $y \in \mathbb{R}$
- **Multiple** regression :  $x \in \mathbb{R}^p$
- **Multivariate** regression :  $y \in \mathbb{R}^d$
- **Functional** regression : when  $x$  and/or  $y$  are functional data

## Objective

Regression models the relationship between two variables  $X$  and  $Y$

- Temperature ( $Y$ ) of some water source, given the air temperate ( $x$ )
- Price ( $Y$ ) of an apartment given its surface ( $x_1$ ) and number of rooms ( $x_2$ )

Vocabulary :

- The  $x$ 's are called inputs/predictors/covariates/features/descriptors/exogenous/Explanatory/independent variables
- The  $y$ 's are called output/outcome/response/endogenous/variable of interest/Explained/dependent variable
- **Simple** regression :  $x \in \mathbb{R}$
- **Multiple** regression :  $x \in \mathbb{R}^p$
- **Univariate** regression :  $y \in \mathbb{R}$
- **Multivariate** regression :  $y \in \mathbb{R}^d$
- **Functional** regression : when  $x$  and/or  $y$  are functional data

- Consider the random pair  $(X, Y)$  where  $X \in \mathcal{X} \subset \mathbb{R}^p$  is the predictor and  $Y \in \mathcal{Y} \subset \mathbb{R}^d$  is the response.
- A regression model can be phrased as

$$Y = f(X) + \varepsilon$$

where  $f : \mathcal{X} \rightarrow \mathcal{Y}$  is the **regression function** (parametric or not, linear or not..)  
 $\varepsilon$  is a random variable : noise/residual/error

- Standard hypotheses : The error terms  $\varepsilon$  are

(i) centered :  $\mathbb{E}(\varepsilon_i) = 0$  (for all  $i$ )

(ii) uncorrelated with the covariates :  $\mathbb{E}(\varepsilon_i X_j) = 0$  (for all  $i, j$ )

(iii) homoskedastic, with limited variance :  $\mathbb{E}(\varepsilon_i^2 | X_i) = \sigma^2 < \infty$  (for all  $i$ )

(iv) uncorrelated with each other :  $\mathbb{E}[\varepsilon_i \varepsilon_j | X] = 0$  (for all  $i \neq j$ )

(v) identically distributed :  $\varepsilon_i \underset{\text{id}}{\sim} p$  (for all  $i$ )

The covariates  $X$  can be deterministic (fixed design) or random

The errors  $\varepsilon_i$  can be supposed normal for statistical inference

- Consider the random pair  $(X, Y)$  where  $X \in \mathcal{X} \subset \mathbb{R}^p$  is the predictor and  $Y \in \mathcal{Y} \subset \mathbb{R}^d$  is the response.
- A regression model can be phrased as

$$Y = f(X) + \varepsilon$$

where  $f : \mathcal{X} \rightarrow \mathcal{Y}$  is the **regression function** (parametric or not, linear or not..)  
 $\varepsilon$  is a random variable : noise/residual/error

- Standard hypotheses : The error terms  $\varepsilon$  are
  - (i) centered :  $\mathbb{E}(\varepsilon_i) = 0$  (for all  $i$ )
  - (ii) uncorrelated with the covariates :  $\mathbb{E}(\varepsilon_i X_j) = 0$  (for all  $i, j$ )
  - (iii) homoskedastic, with limited variance :  $\mathbb{E}(\varepsilon_i^2 | X_i) = \sigma^2 < \infty$  (for all  $i$ )
  - (iv) uncorrelated with each other :  $\mathbb{E}[\varepsilon_i \varepsilon_j | X] = 0$  (for all  $i \neq j$ )
  - (v) identically distributed :  $\varepsilon_i \underset{\text{id}}{\sim} p$  (for all  $i$ )

The covariates  $X$  can be deterministic (fixed design) or random

The errors  $\varepsilon_i$  can be supposed normal for statistical inference

## Def. Regression function

$$h: \mathcal{X} \subset \mathbb{R}^p \rightarrow \mathcal{Y} \subset \mathbb{R}^d$$
$$x \mapsto h(x)$$

is a regression function, parametric or not, linear or not, ...

Example : Linear prediction functions : Consider  $\mathcal{H} = \{h(x) = \langle x, \theta \rangle = \theta^T x\}$ , the set of linear functions in  $X$  of the form  $h(x; \theta) = \mathbb{E}_\theta[Y|X] = \beta_0 + \beta^T X$  and  $\theta = (\beta_0, \beta^T)^T$ .

$$h: \mathbb{R}^p \rightarrow \mathbb{R}$$
$$x \mapsto \langle x, \theta \rangle = \theta^T x$$

The predicted values of  $Y_i$ 's for new covariates  $X_i = x_i$ s correspond to

$$\hat{y}_i = h(x_i)$$

Example : Linear prediction functions (cont.) :  $\hat{y}_i = \langle x_i, \theta \rangle = \theta^T x_i$

Q : How good we are in prediction on a particular pair  $(x, y)$  ?



## Def. Regression function

$$h: \mathcal{X} \subset \mathbb{R}^p \rightarrow \mathcal{Y} \subset \mathbb{R}^d$$
$$x \mapsto h(x)$$

is a regression function, parametric or not, linear or not, ...

Example : Linear prediction functions : Consider  $\mathcal{H} = \{h(x) = \langle x, \theta \rangle = \theta^T x\}$ , the set of linear functions in  $X$  of the form  $h(x; \theta) = \mathbb{E}_\theta[Y|X] = \beta_0 + \beta^T X$  and  $\theta = (\beta_0, \beta^T)^T$ .

$$h: \mathbb{R}^p \rightarrow \mathbb{R}$$
$$x \mapsto \langle x, \theta \rangle = \theta^T x$$

The predicted values of  $Y_i$ 's for new covariates  $X_i = x_i$ s correspond to

$$\hat{y}_i = h(x_i)$$

Example : Linear prediction functions (cont.) :  $\hat{y}_i = \langle x_i, \theta \rangle = \theta^T x_i$

Q : How good we are in prediction on a particular pair  $(x, y)$  ?

## Def. Regression function

$$h: \mathcal{X} \subset \mathbb{R}^p \rightarrow \mathcal{Y} \subset \mathbb{R}^d$$
$$x \mapsto h(x)$$

is a regression function, parametric or not, linear or not, ...

Example : Linear prediction functions : Consider  $\mathcal{H} = \{h(x) = \langle x, \theta \rangle = \theta^T x\}$ , the set of linear functions in  $X$  of the form  $h(x; \theta) = \mathbb{E}_\theta[Y|X] = \beta_0 + \beta^T X$  and  $\theta = (\beta_0, \beta^T)^T$ .

$$h: \mathbb{R}^p \rightarrow \mathbb{R}$$
$$x \mapsto \langle x, \theta \rangle = \theta^T x$$

The **predicted** values of  $Y_i$ 's for new covariates  $X_i = x_i$ s correspond to

$$\hat{y}_i = h(x_i)$$

Example : Linear prediction functions (cont.) :  $\hat{y}_i = \langle x_i, \theta \rangle = \theta^T x_i$

Q : How good we are in prediction on a particular pair  $(x, y)$  ?

## Def. Regression function

$$h: \mathcal{X} \subset \mathbb{R}^p \rightarrow \mathcal{Y} \subset \mathbb{R}^d$$
$$x \mapsto h(x)$$

is a regression function, parametric or not, linear or not, ...

Example : Linear prediction functions : Consider  $\mathcal{H} = \{h(x) = \langle x, \theta \rangle = \theta^T x\}$ , the set of linear functions in  $X$  of the form  $h(x; \theta) = \mathbb{E}_\theta[Y|X] = \beta_0 + \beta^T X$  and  $\theta = (\beta_0, \beta^T)^T$ .

$$h: \mathbb{R}^p \rightarrow \mathbb{R}$$
$$x \mapsto \langle x, \theta \rangle = \theta^T x$$

The **predicted** values of  $Y_i$ 's for new covariates  $X_i = x_i$ s correspond to

$$\hat{y}_i = h(x_i)$$

Example : Linear prediction functions (cont.) :  $\hat{y}_i = \langle x_i, \theta \rangle = \theta^T x_i$

**Q** : How good we are in prediction on a particular pair  $(x, y)$  ?

## Def. Loss function

$$\begin{aligned} \ell: \mathcal{Y} \times \mathcal{Y} &\rightarrow \mathbb{R} \\ (y, h(x)) &\mapsto \ell(y, h(x)) \end{aligned}$$

It measures how good we are on a particular  $(x, y)$  pair.

(We assume that the distribution of the test data is the same as for the training data).

### Examples of loss functions in regression

- Square ( $\ell_2$ )-loss :  $\ell_2(y, h(x)) = (y - h(x))^2$
- Absolute ( $\ell_1$ )-loss :  $\ell_1(y, h(x)) = |y - h(x)|$
- Huber loss :  $\ell_\delta(y, h(x)) = \begin{cases} \frac{1}{2}(y - h(x))^2 & \text{if } |y - h(x)| \leq \delta, \\ \delta (|y - h(x)| - \frac{1}{2}\delta), & \text{otherwise.} \end{cases}$

## Def. Loss function

$$\begin{aligned} \ell: \mathcal{Y} \times \mathcal{Y} &\rightarrow \mathbb{R} \\ (y, h(x)) &\mapsto \ell(y, h(x)) \end{aligned}$$

It measures how good we are on a particular  $(x, y)$  pair.

(We assume that the distribution of the test data is the same as for the training data).

### Examples of loss functions in regression

- Square ( $\ell_2$ )-loss :  $\ell_2(y, h(x)) = (y - h(x))^2$
- Absolute ( $\ell_1$ )-loss :  $\ell_1(y, h(x)) = |y - h(x)|$
- Huber loss :  $\ell_\delta(y, h(x)) = \begin{cases} \frac{1}{2}(y - h(x))^2 & \text{if } |y - h(x)| \leq \delta, \\ \delta (|y - h(x)| - \frac{1}{2}\delta), & \text{otherwise.} \end{cases}$

## Def. Loss function

$$\begin{aligned} \ell: \mathcal{Y} \times \mathcal{Y} &\rightarrow \mathbb{R} \\ (y, h(x)) &\mapsto \ell(y, h(x)) \end{aligned}$$

It measures how good we are on a particular  $(x, y)$  pair.

(We assume that the distribution of the test data is the same as for the training data).

### Examples of loss functions in regression

- Square ( $\ell_2$ )-loss :  $\ell_2(y, h(x)) = (y - h(x))^2$
- Absolute ( $\ell_1$ )-loss :  $\ell_1(y, h(x)) = |y - h(x)|$
- Huber loss :  $\ell_\delta(y, h(x)) = \begin{cases} \frac{1}{2}(y - h(x))^2 & \text{if } |y - h(x)| \leq \delta, \\ \delta (|y - h(x)| - \frac{1}{2}\delta), & \text{otherwise.} \end{cases}$

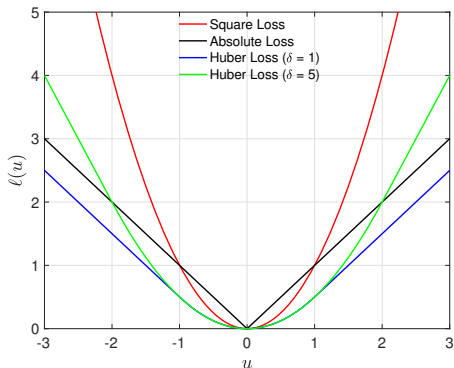


FIGURE – Some loss functions in regression : curves of  $\ell(u)$  for  $u = y - h(x)$ ;  $y \in \mathbb{R}$ .

- Square loss :  $\ell_2(u) = (u)^2$
- Absolute loss :  $\ell_1(u) = |u|$
- Huber loss :  $\ell_\delta(u) = \begin{cases} \frac{1}{2}(u)^2 & \text{if } |u| \leq \delta, \\ \delta(|u| - \frac{1}{2}\delta), & \text{otherwise.} \end{cases}$

- **Risk** : the *Expected loss* :

$$R(h) = \mathbb{E}_P[\ell(Y, h(X))] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, h(x)) dP(x, y)$$

- ↪ the error of approximating  $Y$  by model/hypothesis  $h(X)$  as measured by a chosen loss function  $\ell(Y, h(X))$  given the pair  $(X, Y)$  with (unknown) joint distribution  $P$ ,
- ↪ prediction error : measures the generalization performance of the function  $h$ .
- **Squared Risk** : Under the squared loss  $\ell(y, h(x)) = (y - h(x))^2$  :

$$R(h) = \mathbb{E}_P[(Y - h(X))^2] = \int_{\mathcal{X} \times \mathcal{Y}} (y - h(x))^2 dP(x, y).$$

↪ This is the most used risk in regression

Q : what is the best function  $h$  ? or equivalently, when the risk  $R(h)$  is optimal ?



- Risk : the *Expected loss* :

$$R(h) = \mathbb{E}_P[\ell(Y, h(X))] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, h(x)) dP(x, y)$$

- ↔ the error of approximating  $Y$  by model/hypothesis  $h(X)$  as measured by a chosen loss function  $\ell(Y, h(X))$  given the pair  $(X, Y)$  with (unknown) joint distribution  $P$ ,
- ↔ prediction error : measures the generalization performance of the function  $h$ .

- **Squared Risk** : Under the squared loss  $\ell(y, h(x)) = (y - h(x))^2$  :

$$R(h) = \mathbb{E}_P[(Y - h(X))^2] = \int_{\mathcal{X} \times \mathcal{Y}} (y - h(x))^2 dP(x, y).$$

↔ This is the most used risk in regression

Q : what is the best function  $h$ ? or equivalently, when the risk  $R(h)$  is optimal?

- Risk : the *Expected loss* :

$$R(h) = \mathbb{E}_P[\ell(Y, h(X))] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, h(x)) dP(x, y)$$

- ↔ the error of approximating  $Y$  by model/hypothesis  $h(X)$  as measured by a chosen loss function  $\ell(Y, h(X))$  given the pair  $(X, Y)$  with (unknown) joint distribution  $P$ ,
- ↔ prediction error : measures the generalization performance of the function  $h$ .

- **Squared Risk** : Under the squared loss  $\ell(y, h(x)) = (y - h(x))^2$  :

$$R(h) = \mathbb{E}_P[(Y - h(X))^2] = \int_{\mathcal{X} \times \mathcal{Y}} (y - h(x))^2 dP(x, y).$$

- ↔ This is the most used risk in regression

**Q** : what is the best function  $h$  ? or equivalently, when the risk  $R(h)$  is optimal ?

## Def. regression function

$$h: \mathcal{X} \rightarrow \mathcal{Y}$$
$$x \mapsto h(x)$$

## Theorem (Bayes regression function under the square loss)

Consider the  $\ell_2$ -loss,  $\ell(Y, h(X)) = (h(X) - Y)^2$ , then, the Bayes rule minimizing the corresponding regression risk (the best prediction function) of a regression function  $h(x)$

$$R(h) = \mathbb{E}_X((Y - h(X))^2 | X = x)$$

is given by the conditional expectation

$$h^*(x) = \mathbb{E}(Y | X = x).$$

## Def. regression function

$$h: \mathcal{X} \rightarrow \mathcal{Y}$$
$$x \mapsto h(x)$$

## Theorem (Bayes regression function under the square loss)

Consider the  $\ell_2$ -loss,  $\ell(Y, h(X)) = (h(X) - Y)^2$ , then, the Bayes rule minimizing the corresponding regression risk (the best prediction function) of a regression function  $h(x)$

$$R(h) = \mathbb{E}_X((Y - h(X))^2 | X = x)$$

is given by the conditional expectation

$$h^*(x) = \mathbb{E}(Y | X = x).$$

## Proof : Bayes predictor under the squared loss.

Under the square loss,  $\ell(y, h(x)) = (h(x) - y)^2$ , the best prediction function is

$$h^*(x) = \mathbb{E}[Y|X = x]$$

For  $X = x$ , consider the conditional risk :

$\mathbb{E}[\ell(Y, h(X))|X = x] = \mathbb{E}_{Y|X=x}[(h(X) - Y)^2|X = x] = \int_{\mathcal{Y}} (h(x) - y)^2 p(y|x) dy$ ,  
optimizing the Risk by differentiating w.r.t  $h(x)$  and setting the derivative to 0 :

$$\begin{aligned} \frac{\partial R(h(x))}{\partial h(x)} &= 2 \int (h(x) - y)p(y|x) dy = 2[h(x) \int p(y|x) dy - \int yp(y|x) dy] \\ &= 2(h(x) - \mathbb{E}[Y|X = x]) \end{aligned}$$

which is zero at  $h(x)^* = \mathbb{E}[Y|X = x]$ . Then

$$h^*(x) = \arg \min_{h(x) \in \mathcal{Y}} \mathbb{E}[\ell(Y, h(X))|X = x] = \mathbb{E}[Y|X = x]$$



**Goal :** estimate  $h^*$ , knowing only the data sample  $D_n = (X_i, Y_i)_{i=1}^n$  and loss  $\ell$ .

- Then *Expected loss*  $R(h)$  depends on the joint distribution  $P$  of the pair  $(X, Y)$ . In real situations  $P$  is unknown, as we only have a sample  $D_n = (X_i, Y_i)_{1 \leq i \leq n}$ ,
- ↔ We attempt to minimize the **Empirical Risk**  $R_n(h) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, h(X_i))$  to estimate  $h^*$  (within a family  $\mathcal{H}$ ):

$$\hat{h}_n \in \arg \min_{h \in \mathcal{H}} R_n(h).$$

## MSE and Ordinary Least Squares (OLS) :

- Squared error : Under the squared loss (the standard in regression) :  $\ell_2(y, h(x)) = (y - h(x))^2$ , the empirical risk  $R_n(h)$  is the empirical square loss<sup>1</sup>

$$R_n(h) = \frac{1}{n} \sum_{i=1}^n \|Y_i - h(X_i)\|_2^2$$

- ERM :  $\hat{h}_n = \arg \min_{h \in \mathcal{H}} R_n(h)$  is the **Ordinary Least Squares Estimator** of  $h$
- **Liner regression** : Consider  $\mathcal{H} = \{h_\theta(x) = \alpha + \beta^T x\}$ , the set of linear functions in  $x$

1. also called the Mean Squared Error (MSE), or the mean Residual Squared Sum (RSS) when the ML problem is phrased as an error model  $Y = h(X) + \epsilon$ ,  $\epsilon \sim p$

- Then *Expected loss*  $R(h)$  depends on the joint distribution  $P$  of the pair  $(X, Y)$ . In real situations  $P$  is unknown, as we only have a sample  $D_n = (X_i, Y_i)_{1 \leq i \leq n}$ .
- ⇒ We attempt to minimize the **Empirical Risk**  $R_n(h) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, h(X_i))$  to estimate  $h^*$  (within a family  $\mathcal{H}$ ):

$$\hat{h}_n \in \arg \min_{h \in \mathcal{H}} R_n(h).$$

## MSE and Ordinary Least Squares (OLS) :

- Squared error : Under the squared loss (the standard in regression) :  
 $\ell_2(y, h(x)) = (y - h(x))^2$ , the empirical risk  $R_n(h)$  is the empirical square loss<sup>1</sup>

$$R_n(h) = \frac{1}{n} \sum_{i=1}^n \|Y_i - h(X_i)\|_2^2$$

- ERM :  $\hat{h}_n = \arg \min_{h \in \mathcal{H}} R_n(h)$  is the **Ordinary Least Squares Estimator** of  $h$
- **Liner regression** : Consider  $\mathcal{H} = \{h_\theta(x) = \alpha + \beta^T x\}$ , the set of linear functions in  $x$

1. also called the Mean Squared Error (MSE), or the mean Residual Squared Sum (RSS) when the ML problem is phrased as an error model  $Y = h(X) + \epsilon$ ,  $\epsilon \sim p$

- Then *Expected loss*  $R(h)$  depends on the joint distribution  $P$  of the pair  $(X, Y)$ . In real situations  $P$  is unknown, as we only have a sample  $D_n = (X_i, Y_i)_{1 \leq i \leq n}$ .
- ⇨ We attempt to minimize the **Empirical Risk**  $R_n(h) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, h(X_i))$  to estimate  $h^*$  (within a family  $\mathcal{H}$ ):

$$\hat{h}_n \in \arg \min_{h \in \mathcal{H}} R_n(h).$$

## MSE and Ordinary Least Squares (OLS) :

- Squared error : Under the squared loss (the standard in regression) :  
 $\ell_2(y, h(x)) = (y - h(x))^2$ , the empirical risk  $R_n(h)$  is the empirical square loss<sup>1</sup>

$$R_n(h) = \frac{1}{n} \sum_{i=1}^n \|Y_i - h(X_i)\|_2^2$$

- ERM :  $\hat{h}_n = \arg \min_{h \in \mathcal{H}} R_n(h)$  is the **Ordinary Least Squares Estimator** of  $h$
- **Liner regression** : Consider  $\mathcal{H} = \{h_\theta(x) = \alpha + \beta^T x\}$ , the set of linear functions in  $x$

1. also called the Mean Squared Error (MSE), or the mean Residual Squared Sum (RSS) when the ML problem is phrased as an error model  $Y = h(X) + \epsilon$ ,  $\epsilon \sim p$



## Simple Linear Regression

- We model the pair  $(X, Y)$  where the predictor  $X \in \mathbb{R}$  and the response  $Y \in \mathbb{R}$
- An observed  $Y$  given a single scalar predictor  $x$ , is said to satisfy the simple linear regression model when

$$h(x) = \mathbb{E}[Y|X = x] = \int_{\mathcal{Y}} yp(y|x)dy = \beta_0 + \beta_1 x$$

i.e., equivalently

$$\begin{aligned} Y &= \beta_0 + \beta_1 X + \epsilon, \\ \mathbb{E}[\epsilon|X] &= 0 \text{ and } \mathbb{V}[\epsilon|X] = \sigma^2 \end{aligned}$$

$\beta_0$  (the **intercept**) and  $\beta_1$  (the **slope**) : unknown **regression coefficients**  
 $\sigma^2$  an unknown **noise variance**

## Bayes Risk

$$R^* = R(\theta^*) = \mathbb{E}[(Y - h^*(X))^2] = \mathbb{E}_X \mathbb{E}_{Y|X}[(Y - h^*(X))^2|X] = \mathbb{E}_X \mathbb{E}[\epsilon^2|X] = \sigma^2$$

- We model the pair  $(X, Y)$  where the predictor  $X \in \mathbb{R}$  and the response  $Y \in \mathbb{R}$
- An observed  $Y$  given a single scalar predictor  $x$ , is said to satisfy the simple linear regression model when

$$h(x) = \mathbb{E}[Y|X = x] = \int_{\mathcal{Y}} yp(y|x)dy = \beta_0 + \beta_1x$$

i.e., equivalently

$$\begin{aligned} Y &= \beta_0 + \beta_1 X + \epsilon, \\ \mathbb{E}[\epsilon|X] &= 0 \text{ and } \mathbb{V}[\epsilon|X] = \sigma^2 \end{aligned}$$

$\beta_0$  (the **intercept**) and  $\beta_1$  (the **slope**) : unknown **regression coefficients**  
 $\sigma^2$  an unknown **noise variance**

## Bayes Risk

$$R^* = R(\theta^*) = \mathbb{E}[(Y - h^*(X))^2] = \mathbb{E}_X \mathbb{E}_{Y|X}[(Y - h^*(X))^2|X] = \mathbb{E}_X \mathbb{E}[\epsilon^2|X] = \sigma^2$$

- We model the pair  $(X, Y)$  where the predictor  $X \in \mathbb{R}$  and the response  $Y \in \mathbb{R}$
- An observed  $Y$  given a single scalar predictor  $x$ , is said to satisfy the simple linear regression model when

$$h(x) = \mathbb{E}[Y|X = x] = \int_{\mathcal{Y}} yp(y|x)dy = \beta_0 + \beta_1x$$

i.e., equivalently

$$\begin{aligned} Y &= \beta_0 + \beta_1X + \epsilon, \\ \mathbb{E}[\epsilon|X] &= 0 \text{ and } \mathbb{V}[\epsilon|X] = \sigma^2 \end{aligned}$$

$\beta_0$  (the **intercept**) and  $\beta_1$  (the **slope**) : unknown **regression coefficients**  
 $\sigma^2$  an unknown **noise variance**

## Bayes Risk

$$R^* = R(\theta^*) = \mathbb{E}[(Y - h^*(X))^2] = \mathbb{E}_X \mathbb{E}_{Y|X}[(Y - h^*(X))^2|X] = \mathbb{E}_X \mathbb{E}[\epsilon^2|X] = \sigma^2$$

- We model the pair  $(X, Y)$  where the predictor  $X \in \mathbb{R}$  and the response  $Y \in \mathbb{R}$
- An observed  $Y$  given a single scalar predictor  $x$ , is said to satisfy the simple linear regression model when

$$h(x) = \mathbb{E}[Y|X = x] = \int_{\mathcal{Y}} yp(y|x)dy = \beta_0 + \beta_1x$$

i.e., equivalently

$$\begin{aligned} Y &= \beta_0 + \beta_1X + \epsilon, \\ \mathbb{E}[\epsilon|X] &= 0 \text{ and } \mathbb{V}[\epsilon|X] = \sigma^2 \end{aligned}$$

$\beta_0$  (the **intercept**) and  $\beta_1$  (the **slope**) : unknown **regression coefficients**  
 $\sigma^2$  an unknown **noise variance**

## Bayes Risk

$$R^* = R(\theta^*) = \mathbb{E}[(Y - h^*(X))^2] = \mathbb{E}_X \mathbb{E}_{Y|X}[(Y - h^*(X))^2|X] = \mathbb{E}_X \mathbb{E}[\epsilon^2|X] = \sigma^2$$

# Ordinary Least Squares (OLS) for SLR

Let  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  and  $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$  be the empirical (sample) means.

## Theorem (Ordinary Least Squares (OLS) for SLR)

If  $(\hat{\beta}_0^{OLS}, \hat{\beta}_1^{OLS})$  are the OLS Estimators of  $(\beta_0, \beta_1)$ , then

$$\begin{aligned}\hat{\beta}_0^{OLS} &= \bar{Y} - \hat{\beta}_1 \bar{X}, \\ \hat{\beta}_1^{OLS} &= \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}.\end{aligned}$$

- We have  $\hat{\beta}_1^{OLS} = S_{XY} / S_X^2$  where  $S_{XY} = n^{-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$  is the sample covariance and  $S_X^2 = n^{-1} \sum_{i=1}^n (X_i - \bar{X})^2$  is the sample variance.
- An estimator  $\hat{\sigma}^2$  of the variance  $\sigma^2$  can be taken as the empirical variance

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{h}(X_i))^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - (\hat{\beta}_0^{OLS} + \hat{\beta}_1^{OLS} X_i))^2$$

We'll see its construction later in connection with Gaussian regression, as the Maximum-Likelihood Estimator (MLE)

# Ordinary Least Squares (OLS) for SLR

Let  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  and  $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$  be the empirical (sample) means.

## Theorem (Ordinary Least Squares (OLS) for SLR)

If  $(\hat{\beta}_0^{OLS}, \hat{\beta}_1^{OLS})$  are the OLS Estimators of  $(\beta_0, \beta_1)$ , then

$$\hat{\beta}_0^{OLS} = \bar{Y} - \hat{\beta}_1 \bar{X},$$

$$\hat{\beta}_1^{OLS} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}.$$

- We have  $\hat{\beta}_1^{OLS} = S_{XY} / S_X^2$  where  $S_{XY} = n^{-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$  is the sample covariance and  $S_X^2 = n^{-1} \sum_{i=1}^n (X_i - \bar{X})^2$  is the sample variance.
- An estimator  $\hat{\sigma}^2$  of the variance  $\sigma^2$  can be taken as the empirical variance

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{h}(X_i))^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - (\hat{\beta}_0^{OLS} + \hat{\beta}_1^{OLS} X_i))^2$$

We'll see its construction later in connection with Gaussian regression, as the Maximum-Likelihood Estimator (MLE)

## Proof of the OLS for SLR.

The OLS estimates  $(\hat{\beta}_0, \hat{\beta}_1)$  are the ERM, i.e minimizing the residual squared sum (RSS)

$$Q(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2, \text{ i.e. } (\hat{\beta}_0, \hat{\beta}_1) = \arg \min_{(\beta_0, \beta_1) \in \mathbb{R}^2} Q(\beta_0, \beta_1).$$

F.O.C : Deriving  $Q$  w.r.t  $(\beta_0, \beta_1)$  we get

$$\frac{\partial Q}{\partial \beta_0} = \frac{\partial \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))$$

$$\frac{\partial Q}{\partial \beta_1} = \frac{\partial \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2}{\partial \beta_1} = -2 \sum_{i=1}^n x_i (y_i - (\beta_0 + \beta_1 x_i)).$$

and setting to zero we obtain

$$\begin{aligned} \sum_{i=1}^n y_i - \sum_{i=1}^n \hat{\beta}_0 - \hat{\beta}_1 \sum_{i=1}^n x_i &= 0 \\ \sum_{i=1}^n x_i y_i - \hat{\beta}_0 \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= 0. \end{aligned}$$

which gives the *normal equations* :

$$\begin{aligned} n\hat{\beta}_0 &= \sum_{i=1}^n y_i - \hat{\beta}_1 \sum_{i=1}^n x_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \end{aligned}$$





# Ordinary Least Squares (OLS) for SLR

## Proof of the OLS for SLR (cont.)

The first normal equation one gives

$$\begin{aligned}\widehat{\beta}_0 &= \sum_{i=1}^n y_i - \widehat{\beta}_1 \sum_{i=1}^n x_i \\ \widehat{\beta}_0 &= \frac{1}{n} \sum_{i=1}^n y_i - \widehat{\beta}_1 \left( \frac{1}{n} \sum_{i=1}^n x_i \right) = \bar{y} - \widehat{\beta}_1 \bar{x}.\end{aligned}$$

The second gives

$$\begin{aligned}\widehat{\beta}_0 \sum_{i=1}^n x_i + \widehat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \\ \sum_{i=1}^n x_i \bar{y} - \widehat{\beta}_1 \sum_{i=1}^n x_i \bar{x} + \widehat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i,\end{aligned}$$

we finally obtain :  $\widehat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \bar{y}}{\sum_{i=1}^n x_i^2 - \sum_{i=1}^n x_i \bar{x}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}.$

S.O.C :

$$\det \begin{pmatrix} \frac{\partial^2 Q}{\partial \beta_0^2} & \frac{\partial^2 Q}{\partial \beta_0 \partial \beta_1} \\ \frac{\partial^2 Q}{\partial \beta_0 \partial \beta_1} & \frac{\partial^2 Q}{\partial \beta_1^2} \end{pmatrix} = \det \begin{pmatrix} 2n & 2 \sum_i x_i \\ 2 \sum_i x_i & 2 \sum_i x_i^2 \end{pmatrix} = 4n \sum_i (x_i - \bar{x})^2 > 0.$$

This determinant is zero if all the  $x_i$ 's take the same value. At least two distinct  $x_i$ 's are necessary to estimate the coefficients  $(\beta_0, \beta_1)$  (to fit the line).

# Statistical Properties of the OLS Estimator

Let  $S_X^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$  be the empirical variance of the  $n$  covariates  $X^n = (X_1, \dots, X_n)$ .

## Theorem : Linearity, Unbiasedness and Variance of the OLS

The OLS Estimators  $(\hat{\beta}_0, \hat{\beta}_1)$  of  $(\beta_0, \beta_1)$  are **linear** in  $Y_i$  and **unbiased**, with

$$\mathbb{V}(\hat{\beta}_0 | X^n) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right) = \frac{\sigma^2}{n} \left( 1 + \frac{\bar{X}^2}{S_X^2} \right)$$

$$\mathbb{V}(\hat{\beta}_1 | X^n) = \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sigma^2}{n S_X^2}$$

$$\text{Cov}(\hat{\beta}_0, \hat{\beta}_1 | X^n) = -\frac{\bar{X} \sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} = -\frac{\bar{X} \sigma^2}{n S_X^2}$$

- Estimates of these statistics are obtained by replacing the variance  $\sigma^2$  by its estimator  $\hat{\sigma}^2$  (eg., the corrected MLE). The estimated standard errors  $\hat{s}_e$  of  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are given by

$$\hat{s}_e(\hat{\beta}_0 | X^n) = \frac{\hat{\sigma}}{\sqrt{n}} \sqrt{\left( 1 + \frac{\bar{X}^2}{S_X^2} \right)} = \frac{\hat{\sigma}}{\sqrt{n} S_X} \sqrt{\frac{\sum_{i=1}^n X_i^2}{n}}$$

$$\hat{s}_e(\hat{\beta}_1 | X^n) = \frac{\hat{\sigma}}{\sqrt{n} S_X}; \quad \widehat{\text{Cov}}(\hat{\beta}_0, \hat{\beta}_1 | X^n) = -\frac{\bar{X} \hat{\sigma}^2}{n S_X^2}$$

# Statistical Properties of the OLS Estimator

Let  $S_X^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$  be the empirical variance of the  $n$  covariates  $X^n = (X_1, \dots, X_n)$ .

## Theorem : Linearity, Unbiasedness and Variance of the OLS

The OLS Estimators  $(\hat{\beta}_0, \hat{\beta}_1)$  of  $(\beta_0, \beta_1)$  are **linear** in  $Y_i$  and **unbiased**, with

$$\mathbb{V}(\hat{\beta}_0|X^n) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right) = \frac{\sigma^2}{n} \left( 1 + \frac{\bar{X}^2}{S_X^2} \right)$$

$$\mathbb{V}(\hat{\beta}_1|X^n) = \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sigma^2}{nS_X^2}$$

$$\text{Cov}(\hat{\beta}_0, \hat{\beta}_1|X^n) = -\frac{\bar{X}\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} = -\frac{\bar{X}\sigma^2}{nS_X^2}$$

- Estimates of these statistics are obtained by replacing the variance  $\sigma^2$  by its estimator  $\hat{\sigma}^2$  (eg., the corrected MLE). The estimated standard errors  $\hat{\text{se}}$  of  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are given by

$$\hat{\text{se}}(\hat{\beta}_0|X^n) = \frac{\hat{\sigma}}{\sqrt{n}} \sqrt{\left( 1 + \frac{\bar{X}^2}{S_X^2} \right)} = \frac{\hat{\sigma}}{\sqrt{n}S_X} \sqrt{\frac{\sum_{i=1}^n X_i^2}{n}}$$

$$\hat{\text{se}}(\hat{\beta}_1|X^n) = \frac{\hat{\sigma}}{\sqrt{n}S_X}; \quad \widehat{\text{Cov}}(\hat{\beta}_0, \hat{\beta}_1|X^n) = -\frac{\bar{X}\hat{\sigma}^2}{nS_X^2}$$

# Statistical Properties of the OLS Estimator

Let  $S_X^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$  be the empirical variance of the  $n$  covariates  $X^n = (X_1, \dots, X_n)$ .

## Theorem : Linearity, Unbiasedness and Variance of the OLS

The OLS Estimators  $(\hat{\beta}_0, \hat{\beta}_1)$  of  $(\beta_0, \beta_1)$  are **linear** in  $Y_i$  and **unbiased**, with

$$\mathbb{V}(\hat{\beta}_0 | X^n) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right) = \frac{\sigma^2}{n} \left( 1 + \frac{\bar{X}^2}{S_X^2} \right) = \frac{\sigma^2}{n S_X^2} \left( \frac{1}{n} \sum_{i=1}^n X_i^2 \right)$$

$$\mathbb{V}(\hat{\beta}_1 | X^n) = \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sigma^2}{n S_X^2}$$

$$\text{Cov}(\hat{\beta}_0, \hat{\beta}_1 | X^n) = -\frac{\bar{X} \sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} = -\frac{\bar{X} \sigma^2}{n S_X^2}$$

That is,  $\text{Cov} \left( (\hat{\beta}_0, \hat{\beta}_1)^T \right) = \frac{\sigma^2}{n S_X^2} \begin{pmatrix} \frac{1}{n} \sum_{i=1}^n X_i^2 & -\bar{X} \\ -\bar{X} & 1 \end{pmatrix}$

- Estimates of these statistics are obtained by replacing the variance  $\sigma^2$  by its estimator  $\hat{\sigma}^2$  (eg., the corrected MLE). The estimated standard errors  $\hat{\text{se}}$  of  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are given by

$$\hat{\text{se}}(\hat{\beta}_0 | X^n) = \frac{\hat{\sigma}}{\sqrt{n S_X^2}} \sqrt{\frac{1}{n} \sum_{i=1}^n X_i^2}; \quad \hat{\text{se}}(\hat{\beta}_1 | X^n) = \frac{\hat{\sigma}}{\sqrt{n S_X^2}}; \quad \widehat{\text{Cov}}(\hat{\beta}_0, \hat{\beta}_1 | X^n) = -\frac{\bar{X} \hat{\sigma}^2}{n S_X^2}$$

# Statistical Properties of the OLS Estimator

Let  $S_X^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$  be the empirical variance of the  $n$  covariates  $X^n = (X_1, \dots, X_n)$ .

## Theorem : Linearity, Unbiasedness and Variance of the OLS

The OLS Estimators  $(\hat{\beta}_0, \hat{\beta}_1)$  of  $(\beta_0, \beta_1)$  are **linear** in  $Y_i$  and **unbiased**, with

$$\mathbb{V}(\hat{\beta}_0 | X^n) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right) = \frac{\sigma^2}{n} \left( 1 + \frac{\bar{X}^2}{S_X^2} \right) = \frac{\sigma^2}{n S_X^2} \left( \frac{1}{n} \sum_{i=1}^n X_i^2 \right)$$

$$\mathbb{V}(\hat{\beta}_1 | X^n) = \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sigma^2}{n S_X^2}$$

$$\text{Cov}(\hat{\beta}_0, \hat{\beta}_1 | X^n) = -\frac{\bar{X} \sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} = -\frac{\bar{X} \sigma^2}{n S_X^2}$$

That is,  $\text{Cov} \left( (\hat{\beta}_0, \hat{\beta}_1)^T \right) = \frac{\sigma^2}{n S_X^2} \begin{pmatrix} \frac{1}{n} \sum_{i=1}^n X_i^2 & -\bar{X} \\ -\bar{X} & 1 \end{pmatrix}$

- Estimates of these statistics  $\sigma^2$  are obtained by replacing the variance  $\sigma^2$  by its estimator  $\hat{\sigma}^2$  (eg., the corrected MLE). The estimated standard errors  $\hat{\text{se}}$  of  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are given by

$$\hat{\text{se}}(\hat{\beta}_0 | X^n) = \frac{\hat{\sigma}}{\sqrt{n} S_X} \sqrt{\frac{1}{n} \sum_{i=1}^n X_i^2}; \quad \hat{\text{se}}(\hat{\beta}_1 | X^n) = \frac{\hat{\sigma}}{\sqrt{n} S_X}; \quad \widehat{\text{Cov}}(\hat{\beta}_0, \hat{\beta}_1 | X^n) = -\frac{\bar{X} \hat{\sigma}^2}{n S_X^2}$$

## Proof : Linearity of the OLS.

To simplify notation, let  $w_i = \frac{X_i - \bar{X}}{\sum_{i=1}^n (X_i - \bar{X})^2}$ .

Then we have :

$$\begin{aligned}\sum_{i=1}^n w_i &= \frac{\sum_{i=1}^n X_i}{\sum_{i=1}^n (X_i - \bar{X})^2} - \frac{\sum_{i=1}^n \bar{X}}{\sum_{i=1}^n (X_i - \bar{X})^2} \\ &= \frac{n\bar{X}}{\sum_{i=1}^n (X_i - \bar{X})^2} - \frac{n\bar{X}}{\sum_{i=1}^n (X_i - \bar{X})^2} \\ &= 0\end{aligned}$$

We can then write :

$$\begin{aligned}\hat{\beta}_1 &= \sum_{i=1}^n w_i (Y_i - \bar{Y}) = \sum_{i=1}^n w_i Y_i - \bar{Y} \sum_{i=1}^n w_i = \sum_{i=1}^n w_i Y_i \\ \hat{\beta}_0 &= \bar{Y} - \hat{\beta}_1 \bar{X} = \frac{1}{n} \sum_{i=1}^n Y_i - \bar{X} \sum_{i=1}^n w_i Y_i = \sum_{i=1}^n \left( \frac{1}{n} - \bar{X} w_i \right) Y_i\end{aligned}$$

who are **linear** in  $Y$ .



## Proof : Unbiasdness of the OLS.

We have  $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$ , with  $\mathbb{E}[\epsilon_i|X] = 0$  and  $\mathbb{V}[\epsilon_i|X] = \sigma^2$ , then  $\mathbb{E}[Y_i|X_i] = \beta_0 + \beta_1 X_i$   
Knowing that  $\hat{\beta}_1 = \sum_{i=1}^n w_i Y_i$  and  $\hat{\beta}_0 = \sum_{i=1}^n \left(\frac{1}{n} - \bar{X} w_i\right) Y_i$ , we can then write

$$\begin{aligned}\mathbb{E}[\hat{\beta}_1|X^n] &= \sum_{i=1}^n w_i (\beta_0 + \beta_1 x_i) = \beta_0 \sum_{i=1}^n w_i + \beta_1 \sum_{i=1}^n w_i x_i \\ &= \beta_1 \sum_{i=1}^n w_i x_i - \bar{x} \beta_1 \sum_{i=1}^n w_i = \beta_1 \sum_{i=1}^n w_i (x_i - \bar{x}) = \beta_1 \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = \beta_1\end{aligned}$$

$$\begin{aligned}\mathbb{E}[\hat{\beta}_0|X^n] &= \sum_{i=1}^n \left(\frac{1}{n} - \bar{x} w_i\right) (\beta_0 + \beta_1 x_i) = \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i) - \bar{x} \sum_{i=1}^n w_i (\beta_0 + \beta_1 x_i) \\ &= \beta_0 + \beta_1 \bar{x} - \underbrace{\bar{x} \beta_0 \sum_{i=1}^n w_i}_0 - \underbrace{\beta_1 \bar{x} \sum_{i=1}^n w_i x_i}_1 = \beta_0\end{aligned}$$

We used the fact that :

$$\sum_{i=1}^n w_i x_i = \sum_{i=1}^n w_i x_i - \bar{x} \sum_{i=1}^n w_i = \sum_{i=1}^n w_i (x_i - \bar{x}) = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = 1$$

## Proof : Variance of the OLS.

$$w_i = \frac{X_i - \bar{X}}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

$$\hat{\beta}_1 = \sum_{i=1}^n w_i Y_i$$

$$\begin{aligned}\text{Var}(\hat{\beta}_1 | X^n) &= \mathbb{E}[\hat{\beta}_1^2 | X^n] - 0 \\ &= \mathbb{E} \left[ \sum_{i=1}^n w_i Y_i \sum_{j=1}^n w_j Y_j \middle| X^n \right] \\ &= \sum_{i=1}^n \sum_{j=1}^n w_i w_j \mathbb{E} \left[ Y_i Y_j \middle| X^n \right] = \sum_{i=1}^n \sum_{j=1}^n w_i w_j \sigma^2 \mathbf{1}_{i=j} \\ &= \sum_{i=1}^n w_i^2 \sigma^2 = \sum_{i=1}^n \frac{(X_i - \bar{X})^2}{(\sum_{i=1}^n (X_i - \bar{X})^2)^2} \sigma^2 \\ &= \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \\ &= \frac{\sigma^2}{n S_X^2}\end{aligned}$$





## Proof : Variance of the OLS.

$$\begin{aligned}\hat{\beta}_0 &= \sum_{i=1}^n \left(\frac{1}{n} - \bar{x}w_i\right) Y_i \\ \text{Var}(\hat{\beta}_0|X^n) &= \mathbb{E}[\hat{\beta}_0^2|X^n] - 0 \\ &= \mathbb{E}\left[\sum_{i=1}^n \left(\frac{1}{n} - \bar{x}w_i\right) Y_i \sum_{j=1}^n \left(\frac{1}{n} - \bar{x}w_j\right) Y_j\right] \\ &= \sum_{i=1}^n \sum_{j=1}^n \left(\frac{1}{n} - \bar{x}w_i\right) \left(\frac{1}{n} - \bar{x}w_j\right) \mathbb{E}[Y_i Y_j|X^n] \\ &= \sum_{i=1}^n \sum_{j=1}^n \left(\frac{1}{n} - \bar{x}w_i\right) \left(\frac{1}{n} - \bar{x}w_j\right) \sigma^2 \mathbf{1}_{i=j} \\ &= \sigma^2 \sum_{i=1}^n \left(\frac{1}{n} - \bar{x}w_i\right)^2 = \sigma^2 \left(\sum_{i=1}^n \frac{1}{n^2} - \frac{2\bar{x}}{n} \sum_{i=1}^n w_i + \bar{x}^2 \sum_{i=1}^n w_i^2\right) \\ &= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right) \\ &= \frac{\sigma^2}{n} \left(1 + \frac{\bar{x}}{S_X^2}\right)\end{aligned}$$

Since  $S_X^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}$ , then

$$\text{Var}(\hat{\beta}_0|X^n) = \frac{\sigma^2}{nS_X^2} \left(\frac{1}{n} \sum_{i=1}^n x_i^2\right)$$



## Proof : Variance of the OLS.

$$\begin{aligned}\text{Cov}(\hat{\beta}_0, \hat{\beta}_1 | X^n) &= \text{Cov} \left( \sum_{i=1}^n \left( \frac{1}{n} - \bar{x}w_i \right) Y_i, \sum_{j=1}^n w_j Y_j \middle| X^n \right) \\ &= \sum_{i=1}^n \sum_{j=1}^n \left( \frac{1}{n} - \bar{x}w_i \right) w_j \text{Cov} (Y_i, Y_j | X^n) \\ &= \sum_{i=1}^n \sum_{j=1}^n \left( \frac{1}{n} - \bar{x}w_i \right) w_j \sigma^2 \mathbf{1}_{i=j} \\ &= \sigma^2 \sum_{i=1}^n \left( \frac{1}{n} - \bar{x}w_i \right) w_i = \sigma^2 \left( \frac{1}{n} \sum_{i=1}^n w_i - \bar{x} \sum_{i=1}^n w_i^2 \right) \\ &= - \frac{\bar{x} \sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ &= - \frac{\bar{x} \sigma^2}{n S_X^2}\end{aligned}$$



## Theorem (Gauss-Markov)

The OLS estimator is the unique linear unbiased estimator with minimum variance :  
The OLS estimator is the Best Linear Unbiased Estimator (BLUE)

We give the proof of this result in the multiple regression part.

## Theorem (Consistency)

The OLS estimators  $\widehat{\beta}_0$  and  $\widehat{\beta}_1$  are consistent.

### Proof : Consistency.

Let  $(\widehat{\beta}_j^{(n)})_{n \in \mathbb{N}^*}$ ,  $j \in \{0, 1\}$ , be an estimator sequence derived from increasing sample sizes. By Bienaymé-Tchebychev ineq.

$$\forall \lambda > 0, \mathbb{P} \left( \left| \widehat{\beta}_j^{(n)} - \mathbb{E}[\widehat{\beta}_j^{(n)}] \right| \geq \lambda \right) \leq \frac{\mathbb{V} \left( \widehat{\beta}_j^{(n)} \right)}{\lambda^2}.$$

Then  $\mathbb{P} \left( \left| \widehat{\beta}_1^{(n)} - \beta_1 \right| \geq \lambda \right) \leq \frac{\sigma^2}{\lambda^2 S_X^2}$  and  $\mathbb{P} \left( \left| \widehat{\beta}_0^{(n)} - \beta_0 \right| \geq \lambda \right) \leq \frac{\sigma^2}{\lambda^2} \left( 1 + \frac{\bar{x}}{S_X^2} \right)$ .

Thus  $0 \leq \lim_{n \rightarrow \infty} \mathbb{P} \left( \left| \widehat{\beta}_1^{(n)} - \beta_1 \right| \geq \lambda \right) \leq \frac{\sigma^2}{\lambda^2 S_X^2} \lim_{n \rightarrow \infty} \frac{1}{n} = 0$ ,

and  $0 \leq \lim_{n \rightarrow \infty} \mathbb{P} \left( \left| \widehat{\beta}_0^{(n)} - \beta_0 \right| \geq \lambda \right) \leq \frac{\sigma^2 \left( 1 + \frac{\bar{x}}{S_X^2} \right)}{\lambda^2} \lim_{n \rightarrow \infty} \frac{1}{n} = 0$ ,

That is  $\lim_{n \rightarrow \infty} \mathbb{P} \left( \left| \widehat{\beta}_j^{(n)} - \beta_j \right| \geq \lambda \right) = 0$ ; Then the sequence  $(\widehat{\beta}_j^{(n)})_{n \in \mathbb{N}^*}$  converges in probability to  $\beta_j$  ( $\text{plim}_{n \rightarrow \infty} \widehat{\beta}_j^{(n)} = \beta_j$ ): the OLS Estimators  $\widehat{\beta}_j$  are consistent.  $\square$

## Theorem (Consistency)

The OLS estimators  $\widehat{\beta}_0$  and  $\widehat{\beta}_1$  are consistent.

## Proof : Consistency.

Let  $(\widehat{\beta}_j^{(n)})_{n \in \mathbb{N}^*}$ ,  $j \in \{0, 1\}$ , be an estimator sequence derived from increasing sample sizes. By Bienaymé-Tchebychev ineq.

$$\forall \lambda > 0, \mathbb{P} \left( \left| \widehat{\beta}_j^{(n)} - \mathbb{E}[\widehat{\beta}_j^{(n)}] \right| \geq \lambda \right) \leq \frac{\mathbb{V} \left( \widehat{\beta}_j^{(n)} \right)}{\lambda^2}.$$

Then  $\mathbb{P} \left( \left| \widehat{\beta}_1^{(n)} - \beta_1 \right| \geq \lambda \right) \leq \frac{\sigma^2}{\lambda^2 S_X^2}$  and  $\mathbb{P} \left( \left| \widehat{\beta}_0^{(n)} - \beta_0 \right| \geq \lambda \right) \leq \frac{\sigma^2}{\lambda^2} \left( 1 + \frac{\bar{x}}{S_X^2} \right)$ .

Thus  $0 \leq \lim_{n \rightarrow \infty} \mathbb{P} \left( \left| \widehat{\beta}_1^{(n)} - \beta_1 \right| \geq \lambda \right) \leq \frac{\sigma^2}{\lambda^2 S_X^2} \lim_{n \rightarrow \infty} \frac{1}{n} = 0$ ,

and  $0 \leq \lim_{n \rightarrow \infty} \mathbb{P} \left( \left| \widehat{\beta}_0^{(n)} - \beta_0 \right| \geq \lambda \right) \leq \frac{\sigma^2 \left( 1 + \frac{\bar{x}}{S_X^2} \right)}{\lambda^2} \lim_{n \rightarrow \infty} \frac{1}{n} = 0$ ,

That is  $\lim_{n \rightarrow \infty} \mathbb{P} \left( \left| \widehat{\beta}_j^{(n)} - \beta_j \right| \geq \lambda \right) = 0$ ; Then the sequence  $(\widehat{\beta}_j^{(n)})_{n \in \mathbb{N}^*}$  converges in probability to  $\beta_j$  ( $\text{plim}_{n \rightarrow \infty} \widehat{\beta}_j^{(n)} = \beta_j$ ): the OLS Estimators  $\widehat{\beta}_j$  are consistent.  $\square$

## Theorem (Asymptotic normality)

The OLS estimators  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are asymptotically normal.

### Proof : Asymptotic normality.

By the Central Limit Theorem (CLT), the sequence  $(Z_n)_{n \in \mathbb{N}^*}$  such that

$Z_n = \frac{\hat{\beta}_j^{(n)} - \mathbb{E}[\hat{\beta}_j^{(n)}]}{\sqrt{\text{v}(\hat{\beta}_j^{(n)})}}$  converges to a standard normal random variable.

Then the limit distribution of  $\frac{\hat{\beta}_j^{(n)} - \beta_j}{\sqrt{\text{v}(\hat{\beta}_j^{(n)})}}$  is  $\mathcal{N}(0, 1)$  :

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\sigma^2}{nS_X^2}}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1) \text{ and } \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\frac{\sigma^2}{nS_X^2} \left(\frac{1}{n} \sum_{i=1}^n X_i^2\right)}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1).$$

Equivalently :

$$\sqrt{n}(\hat{\beta}_1 - \beta_1) \xrightarrow{\mathcal{L}} \mathcal{N}\left(0, \frac{\sigma^2}{S_X^2}\right) \text{ and } \sqrt{n}(\hat{\beta}_0 - \beta_0) \xrightarrow{\mathcal{L}} \mathcal{N}\left(0, \frac{\sigma^2}{S_X^2} \left(\frac{1}{n} \sum_{i=1}^n X_i^2\right)\right). \quad \square$$

## Theorem (Asymptotic normality)

The OLS estimators  $\widehat{\beta}_0$  and  $\widehat{\beta}_1$  are asymptotically normal.

## Proof : Asymptotic normality.

By the Central Limit Theorem (CLT), the sequence  $(Z_n)_{n \in \mathbb{N}^*}$  such that

$$Z_n = \frac{\widehat{\beta}_j^{(n)} - \mathbb{E}[\widehat{\beta}_j^{(n)}]}{\sqrt{\text{V}(\widehat{\beta}_j^{(n)})}}$$
 converges to a standard normal random variable.

Then the limit distribution of  $\frac{\widehat{\beta}_j^{(n)} - \beta_j}{\sqrt{\text{V}(\widehat{\beta}_j^{(n)})}}$  is  $\mathcal{N}(0, 1)$  :

$$\frac{\widehat{\beta}_1 - \beta_1}{\sqrt{\frac{\sigma^2}{nS_X^2}}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1) \text{ and } \frac{\widehat{\beta}_0 - \beta_0}{\sqrt{\frac{\sigma^2}{nS_X^2} \left(\frac{1}{n} \sum_{i=1}^n X_i^2\right)}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1).$$

Equivalently :

$$\sqrt{n}(\widehat{\beta}_1 - \beta_1) \xrightarrow{\mathcal{L}} \mathcal{N}\left(0, \frac{\sigma^2}{S_X^2}\right) \text{ and } \sqrt{n}(\widehat{\beta}_0 - \beta_0) \xrightarrow{\mathcal{L}} \mathcal{N}\left(0, \frac{\sigma^2}{S_X^2} \left(\frac{1}{n} \sum_{i=1}^n X_i^2\right)\right). \quad \square$$

## Theorem (Efficiency)

The OLS estimators  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are efficient (achieve minimum variance (CRLB)).

## Proof : Efficiency.

Consider the SLR with normal errors  $\epsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$ . Then  $Y_i | x_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\beta_0 + \beta_1 x_i, \sigma^2)$ . Then, the Fisher information matrix can be defined as

$$\mathcal{I}_n(\beta_0, \beta_1) = -\mathbb{E} \left[ \left( \frac{\partial^2 L(\beta_0, \beta_1)}{\partial \beta_i \partial \beta_j} \right)_{i,j=0,1} \right]$$

with  $L(\beta_0, \beta_1)$  is the conditional log-likelihood function as given by

$$\begin{aligned} L(\beta_0, \beta_1) &= \log p(y_1, \dots, y_n | x_1, \dots, x_n; \beta_0, \beta_1) = \log \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left( \frac{y_i - (\beta_0 + \beta_1 x_i)}{\sigma} \right)^2} \\ &= \sum_{i=1}^n \log \frac{1}{\sigma \sqrt{2\pi}} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2. \end{aligned}$$

□



## Theorem (Efficiency)

The OLS estimators  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are efficient (achieve minimum variance (CRLB)).

## Proof : Efficiency.

Consider the SLR with normal errors  $\epsilon_i \underset{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$ . Then  $Y_i|x_i \underset{\text{iid}}{\sim} \mathcal{N}(\beta_0 + \beta_1 x_i, \sigma^2)$ . Then, the Fisher information matrix can be defined as

$$\mathcal{I}_n(\beta_0, \beta_1) = -\mathbb{E} \left[ \left( \frac{\partial^2 L(\beta_0, \beta_1)}{\partial \beta_i \partial \beta_j} \right)_{i,j=0,1} \right]$$

with  $L(\beta_0, \beta_1)$  is the conditional log-likelihood function as given by

$$\begin{aligned} L(\beta_0, \beta_1) &= \log p(y_1, \dots, y_n | x_1, \dots, x_n; \beta_0, \beta_1) = \log \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left( \frac{y_i - (\beta_0 + \beta_1 x_i)}{\sigma} \right)^2} \\ &= \sum_{i=1}^n \log \frac{1}{\sigma \sqrt{2\pi}} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2. \end{aligned}$$



## Proof : Efficiency (cont.)

$$\frac{\partial L(\beta_0, \beta_1)}{\partial \beta_0} = \frac{\partial}{\partial \beta_0} \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \right\} = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))$$

$$\frac{\partial L(\beta_0, \beta_1)}{\partial \beta_1} = \frac{\partial}{\partial \beta_1} \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \right\} = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i)) x_i$$

$$\frac{\partial^2 L(\beta_0, \beta_1)}{\partial \beta_0^2} = \frac{\partial}{\partial \beta_0} \left\{ \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i)) \right\} = -\frac{n}{\sigma^2}$$

$$\frac{\partial^2 L(\beta_0, \beta_1)}{\partial \beta_1^2} = \frac{\partial}{\partial \beta_1} \left\{ \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i)) x_i \right\} = -\frac{1}{\sigma^2} \sum_{i=1}^n x_i^2$$

$$\frac{\partial^2 L(\beta_0, \beta_1)}{\partial \beta_0 \partial \beta_1} = \frac{\partial}{\partial \beta_0} \left\{ \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i)) x_i \right\} = -\frac{1}{\sigma^2} \sum_{i=1}^n x_i = -\frac{n}{\sigma^2} \bar{x}$$

Then the Fisher information is (fixed design here)

$$\mathcal{I}_n(\beta_0, \beta_1) = -\mathbb{E} \begin{pmatrix} -\frac{n}{\sigma^2} & -\frac{n}{\sigma^2} \bar{x} \\ -\frac{n}{\sigma^2} \bar{x} & -\frac{1}{\sigma^2} \sum_{i=1}^n x_i^2 \end{pmatrix} = \frac{n}{\sigma^2} \begin{pmatrix} 1 & \bar{x} \\ \bar{x} & \frac{1}{n} \sum_{i=1}^n x_i^2 \end{pmatrix}$$



## Proof : Efficiency (cont.)

The Cramer-Rao Lower Bound is given by the inverse of the Fisher information matrix

$$\mathcal{I}_n^{-1}(\beta_0, \beta_1) = \frac{\sigma^2}{n} \begin{pmatrix} 1 & \bar{x} \\ \bar{x} & \frac{1}{n} \sum_{i=1}^n x_i^2 \end{pmatrix}^{-1} = \frac{\sigma^2}{n} \times \frac{1}{\frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2} \begin{pmatrix} \frac{1}{n} \sum_{i=1}^n x_i^2 & -\bar{x} \\ -\bar{x} & 1 \end{pmatrix}$$

Since  $S_X^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$ , finally we get

$$\begin{aligned} \mathcal{I}_n^{-1}(\beta_0, \beta_1) &= \frac{\sigma^2}{nS_X^2} \begin{pmatrix} \frac{1}{n} \sum_{i=1}^n x_i^2 & -\bar{x} \\ -\bar{x} & 1 \end{pmatrix} \\ &= \text{Cov} \left( (\hat{\beta}_0, \hat{\beta}_1)^T \right) \end{aligned}$$

The OLS Estimators  $\hat{\beta}_0$  and  $\hat{\beta}_1$  then achieve the Cramer-Rao Lower Bound. □

Let  $se(\hat{\beta}_1) = \sqrt{\frac{\sigma^2}{nS_X^2}}$  and  $se(\hat{\beta}_0) = \sqrt{\frac{\sigma^2}{n} \left(1 + \frac{\bar{x}}{S_X^2}\right)}$ . We then have

$$Z_j = \frac{\hat{\beta}_j - \beta_j}{se(\hat{\beta}_j)} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1).$$

Since the variance  $\sigma^2$  is unknown, we use instead its best estimator : the corrected MLE

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - (\hat{\beta}_0^{\text{OLS}} + \hat{\beta}_1^{\text{OLS}} X_i))^2$$

We then use instead the statistic

$$T_j = \frac{\hat{\beta}_j - \beta_j}{\hat{se}(\hat{\beta}_j)}$$

where  $\hat{se}(\hat{\beta}_j)$  corresponds to replacing  $\sigma^2$  by  $\hat{\sigma}^2$  in  $se(\hat{\beta}_j)$ .

We know that

$$U = \frac{(n-2)\widehat{\sigma}^2}{\sigma^2} \sim \chi_{n-2}^2$$

Then we finally have

$$T_j = \frac{\widehat{\beta}_j - \beta_j}{\widehat{\text{se}}(\widehat{\beta}_j)} = \frac{Z_j}{\sqrt{\frac{U}{n-2}}} \sim \mathcal{T}_{n-2}$$

Let  $\mathcal{T}_{1-\frac{\alpha}{2}} = \mathbb{P}(T_j \leq \frac{\alpha}{2})$ , i.e the quantile of order  $\frac{\alpha}{2}$  of the Student's law with  $n-2$  degrees of freedom. An approximate  $1-\alpha$  confidence interval for  $\widehat{\beta}_j$  is then given by

$$\mathbb{P}(-\mathcal{T}_{n-2, \frac{\alpha}{2}} \leq T_j \leq \mathcal{T}_{n-2, \frac{\alpha}{2}}) = 1 - \alpha$$

which corresponds to

$$\mathbb{P}\left(\beta_j - \mathcal{T}_{n-2, \frac{\alpha}{2}} \widehat{\text{se}}(\widehat{\beta}_j) \leq \widehat{\beta}_j \leq \beta_j + \mathcal{T}_{n-2, \frac{\alpha}{2}} \widehat{\text{se}}(\widehat{\beta}_j)\right) = 1 - \alpha.$$

We finally obtain

$$\text{CI}_{1-\alpha}(\widehat{\beta}_j) = \left[ \beta_j - \mathcal{T}_{n-2, \frac{\alpha}{2}} \widehat{\text{se}}(\widehat{\beta}_j), \beta_j + \mathcal{T}_{n-2, \frac{\alpha}{2}} \widehat{\text{se}}(\widehat{\beta}_j) \right].$$

Finally :

$$\text{CI}_{1-\alpha}(\hat{\beta}_0) = \left[ \hat{\beta}_0 \pm \mathcal{T}_{1-\frac{\alpha}{2}}^{n-2} \sqrt{\sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)} \right]$$

$$\text{CI}_{1-\alpha}(\hat{\beta}_1) = \left[ \hat{\beta}_1 \pm \mathcal{T}_{1-\frac{\alpha}{2}}^{n-2} \sqrt{\frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \right]$$

## Confidence interval for the regression line

We can construct a confidence interval for  $\hat{h}(x_i) = \mathbb{E}[Y_i|X_i = x_i; \hat{\beta}_0, \hat{\beta}_1] = \hat{\beta}_0 + \hat{\beta}_1 x_i$ .  
Since we have

$$\begin{aligned}\text{Var}(\hat{\beta}_0 + \hat{\beta}_1 x_i) &= \text{Var}(\hat{\beta}_0) + \text{Var}(\hat{\beta}_1 x_i) + 2 \text{Cov}(\hat{\beta}_0, \hat{\beta}_1 x_i) \\ &= \text{Var}(\hat{\beta}_0) + x_i^2 \text{Var}(\hat{\beta}_1) + 2 \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) x_i \\ &= \frac{\sigma^2}{n} \left( 1 + \frac{\bar{x}^2}{S_X^2} \right) + x_i^2 \frac{\sigma^2}{n S_X^2} - 2 \frac{\bar{x} \sigma^2}{n S_X^2} x_i \\ &= \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_X^2} + \frac{x_i^2}{n S_X^2} - 2 \frac{\bar{x}}{n S_X^2} x_i \right) \\ &= \sigma^2 \left( \frac{1}{n} + \frac{(x_i - \bar{x})^2}{n S_X^2} \right) \\ &= \sigma^2 \left( \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2} \right)\end{aligned}$$

We then obtain  $\text{CI}_{1-\alpha}(\hat{\beta}_0 + \hat{\beta}_1 x_i) = \left[ \hat{\beta}_0 + \hat{\beta}_1 x_i \pm \mathcal{T}_{1-\frac{\alpha}{2}}^{n-2} \sqrt{\sigma^2 \left( \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2} \right)} \right]$

## Prediction interval

Given a new input  $x_*$ , the predicted output  $Y_*$  is given by

$$Y_* = \hat{h}(x_*) + \epsilon_* = \hat{\beta}_0 + \hat{\beta}_1 x_* + \epsilon_*$$

Since  $\epsilon_*$  is not observed, then independent from the training set, the variance of the predicted value is then

$$\begin{aligned}\text{Var}(Y_*) &= \text{Var}(\hat{\beta}_0 + \hat{\beta}_1 x_*) + \text{Var}(\epsilon_*) \\ &= \sigma^2 \left( \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2} \right) + \sigma^2 \\ &= \sigma^2 \left( \mathbf{1} + \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2} \right)\end{aligned}$$

and we have  $\text{CI}_{1-\alpha}(Y_*) = \left[ \hat{\beta}_0 + \hat{\beta}_1 x_* \pm \mathcal{T}_{1-\frac{\alpha}{2}}^{n-2} \sqrt{\sigma^2 \left( \mathbf{1} + \frac{1}{n} + \frac{(x_* - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2} \right)} \right]$ .

This one (on  $Y_*$ ) is larger compared to the previous one (on  $\mathbb{E}[Y_* | X = x_*]$ ).



## Regression with Gaussian errors

Let  $\mathcal{X} = \mathbb{R}$ ,  $y \in \mathbb{R}$  and  $h: \mathcal{X} \rightarrow \mathcal{Y}$  s.t.  $x \mapsto \beta_0 + \beta_1 x$ , and consider the following model

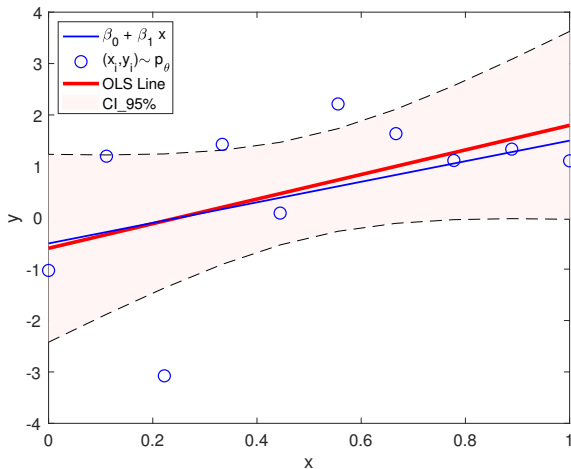
$$Y_i = h(X_i; \beta_0, \beta_1) + \varepsilon_i \quad \text{with} \quad \varepsilon_i | X \underset{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

- Empirical Risk : under the square loss  $R_n(\beta_0, \beta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$
- Empirical Risk Minimizer :  $(\hat{\beta}_0, \hat{\beta}_1)_n = \arg \min_{(\beta_0, \beta_1)} R_n(\beta_0, \beta_1)$
- Conditional Maximum Likelihood Risk

$$\text{Data model : } Y_i | X_i \underset{\text{iid}}{\sim} \mathcal{N}(\beta_0 + \beta_1 x_i, \sigma^2) : p_{\theta}(y_i | x_i) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left( \frac{y - (\beta_0 + \beta_1 x_i)}{\sigma} \right)^2}$$

$$\log L(\theta) = \sum_{i=1}^n \log p_{\theta}(y_i | x_i) = -\frac{1}{2\sigma^2} \underbrace{\sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2}_{R_n(\beta)} - \frac{n}{2} \log \sigma^2 - \frac{n}{2} \log(2\pi)$$

- Conditional MLE :  $=(\hat{\beta}_0^{(n)}, \hat{\beta}_1^{(n)}) = \arg \max_{(\beta_0, \beta_1)} \log L(\theta)$
- ↪ Then we have :  $\arg \min_{(\beta_0, \beta_1)} R_n(\beta_0, \beta_1) = \arg \max_{(\beta_0, \beta_1)} \log L(\theta)$ .
- For both we can take the sample variance  $\sigma^2$  :  $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - (\hat{\beta}_0, \hat{\beta}_1 x_i))^2$  which is the Maximum-Likelihood Estimator



$y_i = -\frac{1}{2} + 2x_i + \epsilon_i$ ,  $\epsilon_i \sim \mathcal{N}(0, 1)$  and  $x_i$ 's are  $n$  values from a uniform grid in  $[0, 1]$

**Practical work:** Coding session to implement from the scratch the confidence interval calculation for regression

$$y_i = -\frac{1}{2} + 2x_i + \epsilon_i,$$

$\epsilon_i \sim \mathcal{N}(0, 1)$  and  $x_i$ 's are values from a uniform grid in  $[0, 1]$

## Practical work session : Confidence intervals and prediction using linear regression

- Simulated data
- Appart data prediction

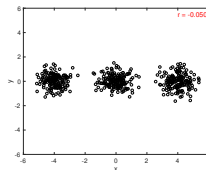
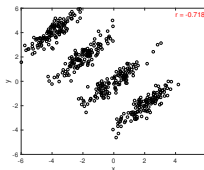
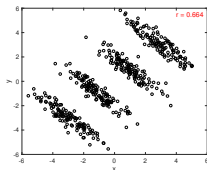
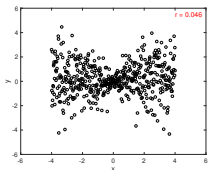
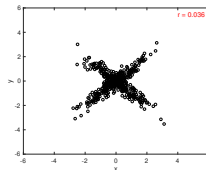
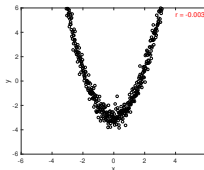
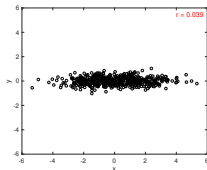
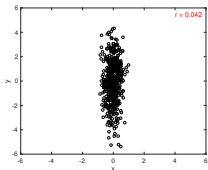
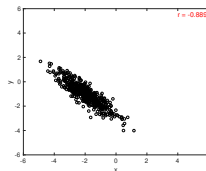
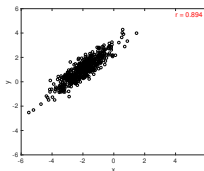
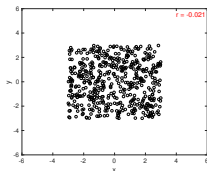
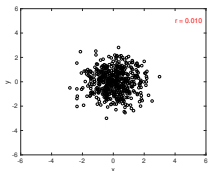
Python, R, and MATLAB code provided during the session and accessible here on the course's page :

# Goodness of Fit

# Correlation coefficient

Correlation coefficient  $\rho_{XY} = \frac{\text{Cov}(X,Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} = \frac{\sigma_{X,Y}}{\sigma_X\sigma_Y} \in [-1, 1]$ .

Sample Correlation coefficient :  $r_{XY} = \frac{S_{XY}}{S_X S_Y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$ .



Measures the **quality of fit**

- $SST = \sum_{i=1}^n (Y_i - \bar{Y})^2$  : sum of the squares of the deviations around  $\bar{Y}$  :  
a measure of the total variability for the  $n$  given observations,  $Y_i$ 's.
- $SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$  : The sum of the squares of the deviations around the  $\hat{Y}_i$ 's :  
A measure of the variability in  $Y$  that remains after the regression is fitted.
- $\frac{SSE}{SST} = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$  : proportion of the total variability unexplained by the fitted regression
- $R^2$  : proportion of the total variability accounted for by the regression :

$$R^2 = 1 - \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = 1 - \frac{SSE}{SST}$$

- For Simple Linear Regression, i.e.  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$  :  $R^2 = r^2$   
 $R^2$  is the correlation coefficient (squared)
- NB : In the general case,  $R^2$  is not a coefficient correlation (and should not be confused with).

- In simple linear regression :  $R^2 \simeq 1$  indicates that the empirical correlation coefficient between the response  $y$  and the predictor  $x$  is close to 1, so that a modeling by a line is satisfactory
- In general : High value of  $R^2$  indicates that the regression model is well fitted to the data. However, it is not an indicator of how good the prediction capability of the fitted model is. For example, a model with  $R^2 \approx 1$  will have high variance (and hence over-fits the data)

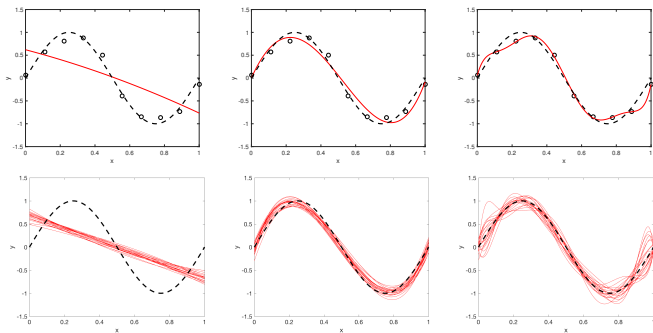


FIGURE – Sample ( $\circ$ ), True function ( $--$ ), realizations of the fitted prediction function ( $—$ )



