



CentraleSupélec

# Statistics and Learning

Lecturers (alphabetic order):

Julien Bect, Gilles Faÿ, Ziad Kobeissi, Laurent Le Brusquet,  
Vincent Lescarret, Arshak Minasyan, Arthur Tenenhaus<sup>†</sup> & Xujia Zhu

<sup>†</sup> Course coordinator

## Lecture 8/9

# Regularization and model selection

### Course objectives

- ▶ Introduction to regularization for regression and classification.
- ▶ Estimation of generalization error.
- ▶ Selection of hyperparameter values and model selection.

# Lecture outline

- 1 – Regularized regression (or classification): penalization
- 2 – Estimation of the risk (generalization error)
- 3 – Hyper-parameters, model selection
- 4 – Exercises and solutions
- 5 – Appendices

# Lecture outline

## 1 – Regularized regression (or classification): penalization

1.1 – Limitations of “ordinary least squares”

1.2 – Ridge regression

1.3 – LASSO regression

## 2 – Estimation of the risk (generalization error)

## 3 – Hyper-parameters, model selection

## 4 – Exercises and solutions

## 5 – Appendices

# Lecture outline

## 1 – Regularized regression (or classification): penalization

### 1.1 – Limitations of “ordinary least squares”

### 1.2 – Ridge regression

### 1.3 – LASSO regression

## 2 – Estimation of the risk (generalization error)

## 3 – Hyper-parameters, model selection

## 4 – Exercises and solutions

## 5 – Appendices

# Limitations of “ordinary least squares”

Recall that  $\underline{X}$  has size **#individuals**  $\times$  **#variables** ( $n \times (p + 1)$ ).

Critical situations for (ordinary) linear regression:

- ▶ when  $\underline{X}^T \underline{X}$  is singular
- ▶ or poorly conditioned

## Typical cases

- 1 when the number of variables is large
- 2 when there are strong correlations between explanatory variables

# Limitations of “ordinary least squares”

Recall that  $\underline{X}$  has size  $\# \text{individuals} \times \# \text{variables}$  ( $n \times (p + 1)$ ).

**Critical situations for (ordinary) linear regression:**

- ▶ when  $\underline{X}^T \underline{X}$  is singular
- ▶ or poorly conditioned

## Typical cases

- 1 when the number of variables is large
- 2 when there are strong correlations between explanatory variables

# Limitations of “ordinary least squares”

Recall that  $\underline{X}$  has size  $\# \text{individuals} \times \# \text{variables}$  ( $n \times (p + 1)$ ).

**Critical situations for (ordinary) linear regression:**

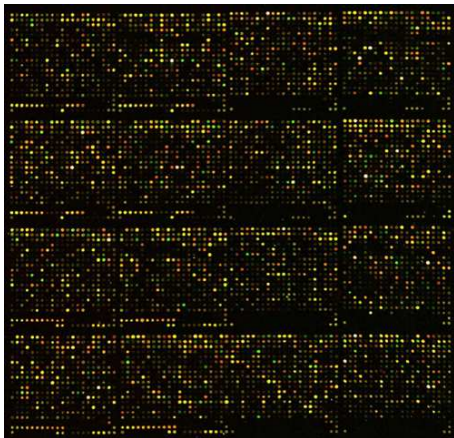
- ▶ when  $\underline{X}^\top \underline{X}$  is singular
- ▶ or poorly conditioned

## Typical cases

- 1 when the number of variables is large
- 2 when there are strong correlations between explanatory variables



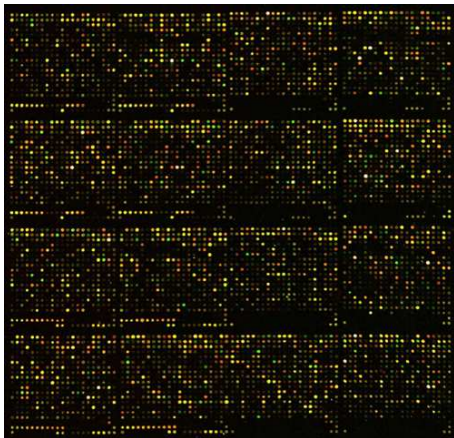
Example:  $p \gg n$



*Subset of a microarray for transcriptome analysis,  
 $p \approx 25000$  for one patient*

Typically,  $n \approx 10$  or  $100$ !

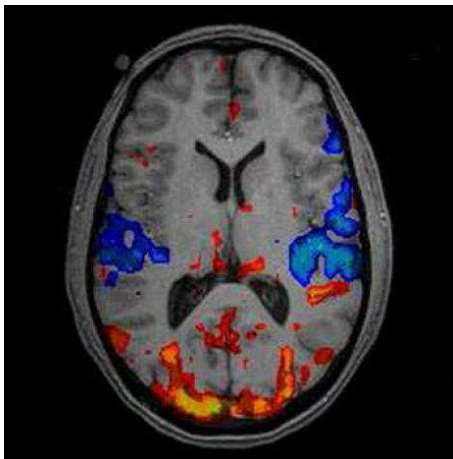
Example:  $p \gg n$



*Subset of a microarray for transcriptome analysis,  
 $p \approx 25000$  for one patient*

Typically,  $n \approx 10$  or  $100$ !

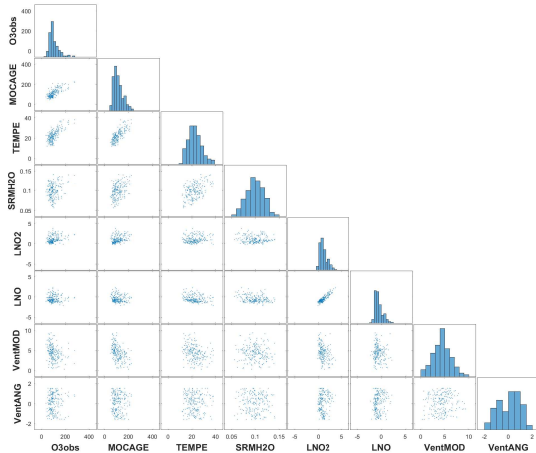
Example:  $p \gg n$



*Functional Magnetic Resonance Imaging (fMRI), with approximately,  $p \approx 300000$  voxels*

Typically,  $n \approx 10$  or  $100$ !

# Example: strong correlation between explanatory variables



“Ozone” example → correlation between variables NO and NO2

## Example: strong correlation... (cont'd)

Vector  $\hat{\beta}$  obtained by OLS regression:

$\beta_0$	MOCAGE	TEMPE	RMH2O	NO2	NO	VentMOD	VentANG
103.4	28.9	22.5	-3.2	-34.4	37.9	1.4	2.6

Observations:

- ▶ The negative coefficient associated to NO2 is surprising  
    ▮ hazardous interpretation of the coefficients
- ▶ The least influential variables (small coefficients) could perhaps be removed from the model?

# One possible solution: penalized regression

A **penalty** term is added to the empirical risk:

$$\hat{\beta} = \arg \min_{\beta} \underbrace{\|\underline{Y} - \underline{X}\beta\|^2}_{\text{data "fidelity"}} + \underbrace{\lambda}_{\text{hyperparameter}} \underbrace{\Omega(\beta)}_{\text{penalty}}. \quad (\star)$$

Expected benefits of penalization:

- ▶ make the solution of  $(\star)$  unique,
- ▶ take prior information into account (this is related to the Bayesian approach),
- ▶ avoid over-fitting when the family of predictor functions is “large” (for linear models:  $p \gg n$ ),
- ▶ make it easier to interpret the resulting model.

# One possible solution: penalized regression

A **penalty** term is added to the empirical risk:

$$\hat{\beta} = \arg \min_{\beta} \underbrace{\|\underline{Y} - \underline{X}\beta\|^2}_{\text{data "fidelity"}} + \underbrace{\lambda}_{\text{hyperparameter}} \underbrace{\Omega(\beta)}_{\text{penalty}}. \quad (\star)$$

## Expected benefits of penalization:

- ▶ make the solution of  $(\star)$  **unique**,
- ▶ take **prior information** into account  
(this is related to the Bayesian approach),
- ▶ **avoid over-fitting** when the family of predictor functions is  
"large" (for linear models:  $p \gg n$ ),
- ▶ make it **easier to interpret** the resulting model.

# Lecture outline

## 1 – Regularized regression (or classification): penalization

1.1 – Limitations of “ordinary least squares”

1.2 – Ridge regression

1.3 – LASSO regression

## 2 – Estimation of the risk (generalization error)

## 3 – Hyper-parameters, model selection

## 4 – Exercises and solutions

## 5 – Appendices



# Ridge regression

## Penalty

$$\Omega(\beta) = \|\beta\|^2$$

$$\hat{\beta}^{\text{RIDGE}} = \arg \min_{\beta} \|\underline{Y} - \underline{X}\beta\|^2 + \lambda \|\beta\|^2$$

It can be proved that ( $\Rightarrow$  see PC) :

$$\hat{\beta}^{\text{RIDGE}} = \left( \underline{X}^{\top} \underline{X} + \lambda I_{p+1} \right)^{-1} \underline{X}^{\top} \underline{Y}.$$

$\Rightarrow$  When  $\lambda \nearrow$ , the conditioning of  $(\underline{X}^{\top} \underline{X} + \lambda I_{p+1})$  improves.

Remark:  $\hat{\beta}^{\text{RIDGE}}$  has a Bayesian interpretation ( $\Rightarrow$  see PC).

# Ridge regression

## Penalty

$$\Omega(\beta) = \|\beta\|^2$$

$$\hat{\beta}^{\text{RIDGE}} = \arg \min_{\beta} \|\underline{Y} - \underline{X}\beta\|^2 + \lambda \|\beta\|^2$$

It can be proved that ( $\Rightarrow$  see PC) :

$$\hat{\beta}^{\text{RIDGE}} = \left( \underline{X}^{\top} \underline{X} + \lambda I_{p+1} \right)^{-1} \underline{X}^{\top} \underline{Y}.$$

$\Rightarrow$  When  $\lambda \nearrow$ , the **conditioning** of  $(\underline{X}^{\top} \underline{X} + \lambda I_{p+1})$  **improves**.

Remark:  $\hat{\beta}^{\text{RIDGE}}$  has a Bayesian interpretation ( $\Rightarrow$  see PC).

# Ridge regression

## Penalty

$$\Omega(\beta) = \|\beta\|^2$$

$$\hat{\beta}^{\text{RIDGE}} = \arg \min_{\beta} \|\underline{Y} - \underline{X}\beta\|^2 + \lambda \|\beta\|^2$$

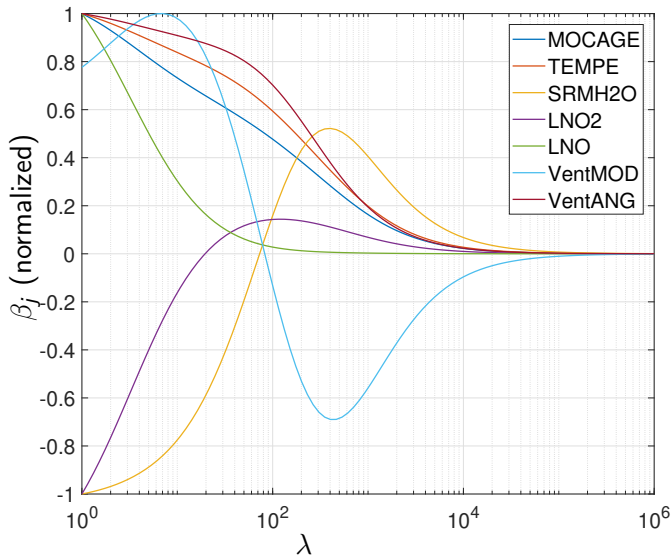
It can be proved that ( $\Rightarrow$  see PC) :

$$\hat{\beta}^{\text{RIDGE}} = \left( \underline{X}^{\top} \underline{X} + \lambda I_{p+1} \right)^{-1} \underline{X}^{\top} \underline{Y}.$$

$\Rightarrow$  When  $\lambda \nearrow$ , the conditioning of  $(\underline{X}^{\top} \underline{X} + \lambda I_{p+1})$  improves.

Remark:  $\hat{\beta}^{\text{RIDGE}}$  has a Bayesian interpretation ( $\Rightarrow$  see PC).

# “Ozone” example: Evolution of $\hat{\beta}^{RIDGE}$ as a function of $\lambda$



# Lecture outline

## 1 – Regularized regression (or classification): penalization

1.1 – Limitations of “ordinary least squares”

1.2 – Ridge regression

1.3 – LASSO regression

## 2 – Estimation of the risk (generalization error)

## 3 – Hyper-parameters, model selection

## 4 – Exercises and solutions

## 5 – Appendices

# LASSO regression

## Penalty

$$\Omega(\beta) = \|\beta\|_1 = \sum_{j=1}^n |\beta_j|$$

$$\hat{\beta}^{\text{LASSO}} = \arg \min_{\beta} \|\underline{Y} - \underline{X}\beta\|^2 + \lambda \|\beta\|_1 \quad (\star)$$

## Minimization of the criterion

- ▶ no explicit solution for  $\hat{\beta}^{\text{LASSO}}$  (except in some cases,  
    ⇒ exercise 1)
- ⇒ dedicated algorithms

# LASSO regression

## Penalty

$$\Omega(\beta) = \|\beta\|_1 = \sum_{j=1}^n |\beta_j|$$

$$\hat{\beta}^{\text{LASSO}} = \arg \min_{\beta} \|\underline{Y} - \underline{X}\beta\|^2 + \lambda \|\beta\|_1 \quad (\star)$$

## Minimization of the criterion

- ▶ **no explicit solution** for  $\hat{\beta}^{\text{LASSO}}$  (except in some cases,  
    ▶ **exercice 1**)
- ▶ dedicated algorithms

# LASSO regression: reformulation

$$\hat{\beta}^{\text{LASSO}} = \arg \min_{\beta} \|\underline{Y} - \underline{X}\beta\|^2 + \lambda \|\beta\|_1 \quad (\star)$$

- ▶ Let  $\hat{\beta}^{\text{OLS}}$  denote the OLS estimator of  $\beta$ :

$$\hat{\beta}^{\text{LASSO}} = \hat{\beta}^{\text{OLS}} \quad \text{for } \lambda = 0$$

- ▶ Since  $\|\underline{Y} - \underline{X}\beta\|^2 = \|\underline{X}(\beta - \hat{\beta}^{\text{OLS}})\|^2 + c$ , we have:

$$\hat{\beta}^{\text{LASSO}} = \arg \min_{\beta} \|\underline{X}(\beta - \hat{\beta}^{\text{OLS}})\|^2 + \lambda \|\beta\|_1$$

- ▶ Reformulation with a constraint: it can be proved that there exists  $c_{\lambda} \in \mathbb{R}^+$  such that

$$\hat{\beta}^{\text{LASSO}} = \arg \min_{\beta} \|\underline{X}(\beta - \hat{\beta}^{\text{OLS}})\|^2 \quad \text{such that } \|\beta\|_1 \leq c_{\lambda}$$

(and similarly for  $\hat{\beta}^{\text{RIDGE}}$ )



## LASSO regression: reformulation

$$\hat{\beta}^{\text{LASSO}} = \arg \min_{\beta} \|\underline{Y} - \underline{X}\beta\|^2 + \lambda \|\beta\|_1 \quad (\star)$$

- ▶ Let  $\hat{\beta}^{\text{OLS}}$  denote the OLS estimator of  $\beta$ :

$$\hat{\beta}^{\text{LASSO}} = \hat{\beta}^{\text{OLS}} \quad \text{for } \lambda = 0$$

- ▶ Since  $\|\underline{Y} - \underline{X}\beta\|^2 = \|\underline{X}(\beta - \hat{\beta}^{\text{OLS}})\|^2 + c$ , we have:

$$\hat{\beta}^{\text{LASSO}} = \arg \min_{\beta} \|\underline{X}(\beta - \hat{\beta}^{\text{OLS}})\|^2 + \lambda \|\beta\|_1$$

- ▶ Reformulation with a constraint: it can be proved that there exists  $c_{\lambda} \in \mathbb{R}^+$  such that

$$\hat{\beta}^{\text{LASSO}} = \arg \min_{\beta} \|\underline{X}(\beta - \hat{\beta}^{\text{OLS}})\|^2 \text{ such that } \|\beta\|_1 \leq c_{\lambda}$$

(and similarly for  $\hat{\beta}^{\text{RIDGE}}$ )

## LASSO regression: reformulation

$$\hat{\beta}^{\text{LASSO}} = \arg \min_{\beta} \|\underline{Y} - \underline{X}\beta\|^2 + \lambda \|\beta\|_1 \quad (\star)$$

- ▶ Let  $\hat{\beta}^{\text{OLS}}$  denote the OLS estimator of  $\beta$ :

$$\hat{\beta}^{\text{LASSO}} = \hat{\beta}^{\text{OLS}} \quad \text{for } \lambda = 0$$

- ▶ Since  $\|\underline{Y} - \underline{X}\beta\|^2 = \|\underline{X}(\beta - \hat{\beta}^{\text{OLS}})\|^2 + c$ , we have:

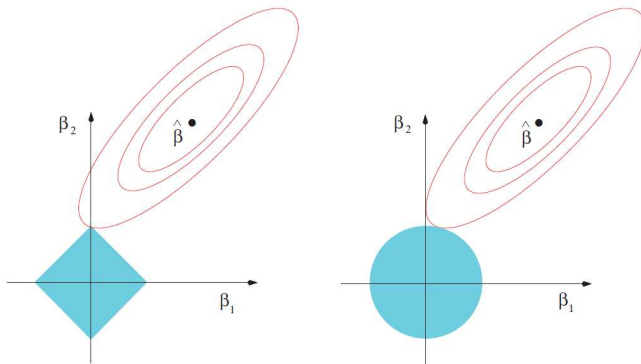
$$\hat{\beta}^{\text{LASSO}} = \arg \min_{\beta} \|\underline{X}(\beta - \hat{\beta}^{\text{OLS}})\|^2 + \lambda \|\beta\|_1$$

- ▶ Reformulation with a **constraint**: it can be proved that there exists  $c_{\lambda} \in \mathbb{R}^+$  such that

$$\hat{\beta}^{\text{LASSO}} = \arg \min_{\beta} \|\underline{X}(\beta - \hat{\beta}^{\text{OLS}})\|^2 \text{ such that } \|\beta\|_1 \leq c_{\lambda}$$

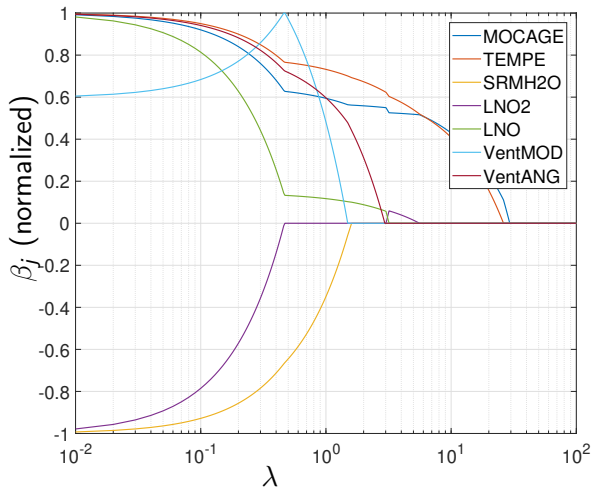
(and similarly for  $\hat{\beta}^{\text{RIDGE}}$ )

# LASSO regression: intuitive interpretation



**FIGURE 3.11.** *Estimation picture for the lasso (left) and ridge regression (right). Shown are contours of the error and constraint functions. The solid blue areas are the constraint regions  $|\beta_1| + |\beta_2| \leq t$  and  $\beta_1^2 + \beta_2^2 \leq t^2$ , respectively, while the red ellipses are the contours of the least squares error function.*

## “Ozone” example: $\hat{\beta}^{\text{LASSO}}$ versus $\lambda$



When  $\lambda \nearrow$ , the number of coefficients equal to zero  $\nearrow$

## “Ozone” example: $\hat{\beta}^{\text{LASSO}}$ for several $\lambda$

**With  $\lambda = 0$  (OLS)**

$\beta_0$	MOCAGE	TEMPE	RMH2O	NO2	NO	VentMOD	VentANG
103.4	28.9	22.5	-3.2	-34.4	37.9	1.4	2.6

⇒ The coefficient for NO2 may seem surprising

**With  $\lambda = 0.5$**

$\beta_0$	MOCAGE	TEMPE	RMH2O	NO2	NO	VentMOD	VentANG
103.4	18.1	17.2	-2.1	0	4.9	2.2	1.9

⇒ One of the two correlated variables is discarded,  
makes it easier to interpret the coefficients

**With  $\lambda = 3$**

$\beta_0$	MOCAGE	TEMPE	RMH2O	NO2	NO	VentMOD	VentANG
103.4	15.9	14.1	0	0	2.2	0	0

⇒ The remaining variables are progressively discarded

Choice of the hyper-parameter  $\lambda$  ?

## “Ozone” example: $\hat{\beta}^{\text{LASSO}}$ for several $\lambda$

With  $\lambda = 0$  (OLS)

$\beta_0$	MOCAGE	TEMPE	RMH2O	NO2	NO	VentMOD	VentANG
103.4	28.9	22.5	-3.2	-34.4	37.9	1.4	2.6

⇒ The coefficient for NO2 may seem surprising

With  $\lambda = 0.5$

$\beta_0$	MOCAGE	TEMPE	RMH2O	NO2	NO	VentMOD	VentANG
103.4	18.1	17.2	-2.1	0	4.9	2.2	1.9

⇒ One of the two correlated variables is discarded,  
makes it easier to **interpret the coefficients**

With  $\lambda = 3$

$\beta_0$	MOCAGE	TEMPE	RMH2O	NO2	NO	VentMOD	VentANG
103.4	15.9	14.1	0	0	2.2	0	0

⇒ The remaining variables are progressively discarded

Choice of the hyper-parameter  $\lambda$  ?

## “Ozone” example: $\hat{\beta}^{\text{LASSO}}$ for several $\lambda$

With  $\lambda = 0$  (OLS)

$\beta_0$	MOCAGE	TEMPE	RMH2O	NO2	NO	VentMOD	VentANG
103.4	28.9	22.5	-3.2	-34.4	37.9	1.4	2.6

⇒ The coefficient for NO2 may seem surprising

With  $\lambda = 0.5$

$\beta_0$	MOCAGE	TEMPE	RMH2O	NO2	NO	VentMOD	VentANG
103.4	18.1	17.2	-2.1	0	4.9	2.2	1.9

⇒ One of the two correlated variables is discarded,  
makes it easier to interpret the coefficients

With  $\lambda = 3$

$\beta_0$	MOCAGE	TEMPE	RMH2O	NO2	NO	VentMOD	VentANG
103.4	15.9	14.1	0	0	2.2	0	0

⇒ The remaining variables are progressively discarded

Choice of the hyper-parameter  $\lambda$  ?

## “Ozone” example: $\hat{\beta}^{\text{LASSO}}$ for several $\lambda$

With  $\lambda = 0$  (OLS)

$\beta_0$	MOCAGE	TEMPE	RMH2O	NO2	NO	VentMOD	VentANG
103.4	28.9	22.5	-3.2	-34.4	37.9	1.4	2.6

⇒ The coefficient for NO2 may seem surprising

With  $\lambda = 0.5$

$\beta_0$	MOCAGE	TEMPE	RMH2O	NO2	NO	VentMOD	VentANG
103.4	18.1	17.2	-2.1	0	4.9	2.2	1.9

⇒ One of the two correlated variables is discarded,  
makes it easier to interpret the coefficients

With  $\lambda = 3$

$\beta_0$	MOCAGE	TEMPE	RMH2O	NO2	NO	VentMOD	VentANG
103.4	15.9	14.1	0	0	2.2	0	0

⇒ The remaining variables are progressively discarded

Choice of the hyper-parameter  $\lambda$  ?



# Lecture outline

1 – Regularized regression (or classification): penalization

2 – Estimation of the risk (generalization error)

2.1 – Problem

2.2 – Zoom in on an illuminating special case

2.3 – Training set and test set

3 – Hyper-parameters, model selection

4 – Exercises and solutions

5 – Appendices

# Lecture outline

1 – Regularized regression (or classification): penalization

2 – Estimation of the risk (generalization error)

2.1 – Problem

2.2 – Zoom in on an illuminating special case

2.3 – Training set and test set

3 – Hyper-parameters, model selection

4 – Exercises and solutions

5 – Appendices

## Problem

Back to the **general setting** (regression/classification).

Let  $\hat{h}$  be a predictor  $\mathcal{X} \rightarrow \mathcal{Y}$  learned from data:

$$\hat{h}(x) = \hat{h}(x; (X_1, Y_1), \dots, (X_n, Y_n)) = \hat{h}(x; \underline{X}, \underline{Y}).$$

Recall that, given a loss function  $L$ , we define the risk, or generalization error :

$$\begin{aligned}\mathcal{R}(\hat{h}) &= \mathbb{E} \left( L(Y, \hat{h}(X)) \mid \underline{X}, \underline{Y} \right) \\ &= \iint_{\mathcal{X} \times \mathcal{Y}} L(y, \hat{h}(x)) P^{X,Y}(\mathrm{d}x, \mathrm{d}y).\end{aligned}$$

Examples.  $L(y, \tilde{y}) = (y - \tilde{y})^2$ ,  $L(y, \tilde{y}) = |y - \tilde{y}|$ ,  $L(y, \tilde{y}) = \mathbb{1}_{y \neq \tilde{y}}$ , ...

## Problem

How can we estimate this risk (which depends on  $P^{X,Y}$ ) ?

## Problem

Back to the general setting (regression/classification).

Let  $\hat{h}$  be a predictor  $\mathcal{X} \rightarrow \mathcal{Y}$  learned from data:

$$\hat{h}(x) = \hat{h}(x; (X_1, Y_1), \dots, (X_n, Y_n)) = \hat{h}(x; \underline{X}, \underline{Y}).$$

Recall that, given a loss function  $L$ , we define the **risk**, or **generalization error** :

$$\begin{aligned}\mathcal{R}(\hat{h}) &= \mathbb{E} \left( L(Y, \hat{h}(X)) \mid \underline{X}, \underline{Y} \right) \\ &= \iint_{\mathcal{X} \times \mathcal{Y}} L(y, \hat{h}(x)) P^{\underline{X}, \underline{Y}}(dx, dy).\end{aligned}$$

Examples.  $L(y, \tilde{y}) = (y - \tilde{y})^2$ ,  $L(y, \tilde{y}) = |y - \tilde{y}|$ ,  $L(y, \tilde{y}) = \mathbb{1}_{y \neq \tilde{y}}$ , ...

## Problem

How can we estimate this risk (which depends on  $P^{\underline{X}, \underline{Y}}$ ) ?

## Problem

Back to the general setting (regression/classification).

Let  $\hat{h}$  be a predictor  $\mathcal{X} \rightarrow \mathcal{Y}$  learned from data:

$$\hat{h}(x) = \hat{h}(x; (X_1, Y_1), \dots, (X_n, Y_n)) = \hat{h}(x; \underline{X}, \underline{Y}).$$

Recall that, given a loss function  $L$ , we define the risk, or generalization error :

$$\begin{aligned}\mathcal{R}(\hat{h}) &= \mathbb{E} \left( L(Y, \hat{h}(X)) \mid \underline{X}, \underline{Y} \right) \\ &= \iint_{\mathcal{X} \times \mathcal{Y}} L(y, \hat{h}(x)) P^{\underline{X}, \underline{Y}}(dx, dy).\end{aligned}$$

Examples.  $L(y, \tilde{y}) = (y - \tilde{y})^2$ ,  $L(y, \tilde{y}) = |y - \tilde{y}|$ ,  $L(y, \tilde{y}) = \mathbb{1}_{y \neq \tilde{y}}$ , ...

## Problem

How can we **estimate this risk** (which depends on  $P^{\underline{X}, \underline{Y}}$ ) ?

## Refresher: empirical risk

We call **empirical risk** the risk

$$\hat{\mathcal{R}}_n = \iint_{\mathcal{X} \times \mathcal{Y}} L(y, \hat{h}(x)) \hat{P}_n(dx, dy) = \frac{1}{n} \sum_{i=1}^n L(Y_i, \hat{h}(X_i))$$

computed with  $P^{X,Y}$  equal to  $\hat{P}_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i, Y_i}$ .

### Question

Is this empirical risk  $\hat{\mathcal{R}}_n$ , in general, a “good” estimator of the true risk  $\mathcal{R}(\hat{h})$  ?



the data is used twice !

**Intuition:** It is “risky” to estimate the risk from the error observed on the same data already used to construct  $\hat{h}$ ...

## Refresher: empirical risk

We call empirical risk the risk

$$\hat{\mathcal{R}}_n = \iint_{\mathcal{X} \times \mathcal{Y}} L(y, \hat{h}(x)) \hat{P}_n(dx, dy) = \frac{1}{n} \sum_{i=1}^n L(Y_i, \hat{h}(X_i))$$

computed with  $P^{X,Y}$  equal to  $\hat{P}_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i, Y_i}$ .

### Question

Is this empirical risk  $\hat{\mathcal{R}}_n$ , in general, a “good” estimator of the true risk  $\mathcal{R}(\hat{h})$  ?



the data is used twice !

**Intuition:** It is “risky” to estimate the risk from the error observed on the same data already used to construct  $\hat{h}$ ...

## Refresher: empirical risk

We call empirical risk the risk

$$\hat{\mathcal{R}}_n = \iint_{\mathcal{X} \times \mathcal{Y}} L(y, \hat{h}(x)) \hat{P}_n(dx, dy) = \frac{1}{n} \sum_{i=1}^n L(\mathbf{Y}_i, \hat{h}(\mathbf{X}_i))$$

computed with  $P^{\mathbf{X}, \mathbf{Y}}$  equal to  $\hat{P}_n = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{X}_i, \mathbf{Y}_i}$ .

### Question

Is this empirical risk  $\hat{\mathcal{R}}_n$ , in general, a “good” estimator of the true risk  $\mathcal{R}(\hat{h})$  ?



the data is used twice !

**Intuition:** It is “risky” to estimate the risk from the error observed on the same data already used to construct  $\hat{h}$ ...



# Lecture outline

1 – Regularized regression (or classification): penalization

2 – Estimation of the risk (generalization error)

2.1 – Problem

2.2 – Zoom in on an illuminating special case

2.3 – Training set and test set

3 – Hyper-parameters, model selection

4 – Exercises and solutions

5 – Appendices

## Zoom in on an illuminating special case

Consider the case of “ordinary” linear regression:

- ▶  $h(x) = \beta_0 + \beta_1 x^{(1)} + \dots + \beta_p x^{(p)},$
- ▶ quadratic loss:  $L(y, \tilde{y}) = (y - \tilde{y})^2,$
- ▶  $p + 1 \leq n$  and  $\underline{X}^\top \underline{X}$  an a.s. invertible  $(p + 1) \times (p + 1)$  matrix.

Empirical risk minimization :  $\hat{\beta} = (\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \underline{Y}.$

Remark: link between  $\hat{\mathcal{R}}_n$  and the coefficient  $R^2$  of determination:

$$\begin{aligned} R^2 &= 1 - \frac{\text{RSS}(\hat{\beta})}{\text{TSS}} = 1 - \frac{\sum_{i=1}^n (Y_i - \hat{\beta}^\top X_i)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \\ &= 1 - \frac{\hat{\mathcal{R}}_n}{\widehat{\text{var}}_n(Y)} \quad \text{with } \widehat{\text{var}}_n(Y) = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2. \end{aligned}$$

## Zoom in on an illuminating special case

Consider the case of “ordinary” linear regression:

- ▶  $h(x) = \beta_0 + \beta_1 x^{(1)} + \dots + \beta_p x^{(p)},$
- ▶ quadratic loss:  $L(y, \tilde{y}) = (y - \tilde{y})^2,$
- ▶  $p + 1 \leq n$  and  $\underline{X}^\top \underline{X}$  an a.s. invertible  $(p + 1) \times (p + 1)$  matrix.

Empirical risk minimization :  $\hat{\beta} = (\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \underline{Y}.$

Remark: link between  $\hat{\mathcal{R}}_n$  and the coefficient  $R^2$  of determination:

$$\begin{aligned} R^2 &= 1 - \frac{\text{RSS}(\hat{\beta})}{\text{TSS}} = 1 - \frac{\sum_{i=1}^n (Y_i - \hat{\beta}^\top X_i)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \\ &= 1 - \frac{\hat{\mathcal{R}}_n}{\widehat{\text{var}}_n(Y)} \quad \text{with } \widehat{\text{var}}_n(Y) = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2. \end{aligned}$$

## Zoom in on an illuminating special case

Consider the case of “ordinary” linear regression:

- ▶  $h(x) = \beta_0 + \beta_1 x^{(1)} + \dots + \beta_p x^{(p)},$
- ▶ quadratic loss:  $L(y, \tilde{y}) = (y - \tilde{y})^2,$
- ▶  $p + 1 \leq n$  and  $\underline{X}^\top \underline{X}$  an a.s. invertible  $(p + 1) \times (p + 1)$  matrix.

Empirical risk minimization :  $\hat{\beta} = (\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \underline{Y}.$

Remark: link between  $\hat{\mathcal{R}}_n$  and the coefficient  $R^2$  of determination:

$$\begin{aligned} R^2 &= 1 - \frac{\text{RSS}(\hat{\beta})}{\text{TSS}} = 1 - \frac{\sum_{i=1}^n (Y_i - \hat{\beta}^\top X_i)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \\ &= 1 - \frac{\hat{\mathcal{R}}_n}{\widehat{\text{var}}_n(Y)} \quad \text{with } \widehat{\text{var}}_n(Y) = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2. \end{aligned}$$

## Zoom on an illuminating special case (cont'd)

Consider the generalization error wrt responses only:

$$\tilde{\mathcal{R}}_n = \mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n \left( \tilde{Y}_i - \hat{\beta}^\top X_i \right)^2 \mid \underline{X}, \underline{Y} \right),$$

with, for all  $i$ ,  $\tilde{Y}_i$  and  $Y_i$  iid conditionally to  $\underline{X}$ .

### Proposition

Assume that the unknown distribution  $P^{X,Y}$  is such that  $Y_i = \beta^\top X_i + \varepsilon_i$ , with  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ , independent of  $X_i$ . Then

$$\mathbb{E} \left( \tilde{\mathcal{R}}_n \right) = \sigma^2 \left( 1 + \frac{p+1}{n} \right),$$

$$\mathbb{E} \left( \hat{\mathcal{R}}_n \right) = \sigma^2 \left( 1 - \frac{p+1}{n} \right).$$

## Zoom on an illuminating special case (cont'd)

Consider the generalization error wrt responses only:

$$\tilde{\mathcal{R}}_n = \mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n \left( \tilde{Y}_i - \hat{\beta}^\top X_i \right)^2 \mid \underline{X}, \underline{Y} \right),$$

with, for all  $i$ ,  $\tilde{Y}_i$  and  $Y_i$  iid conditionally to  $\underline{X}$ .

### Proposition

Assume that the unknown distribution  $P^{X,Y}$  is such that  $Y_i = \beta^\top X_i + \varepsilon_i$ , with  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ , independent of  $X_i$ . Then

$$\mathbb{E} \left( \tilde{\mathcal{R}}_n \right) = \sigma^2 \left( 1 + \frac{p+1}{n} \right),$$

$$\mathbb{E} \left( \hat{\mathcal{R}}_n \right) = \sigma^2 \left( 1 - \frac{p+1}{n} \right).$$

## Zoom on an illuminating special case (cont'd)

Consider the generalization error wrt responses only:

$$\tilde{\mathcal{R}}_n = \mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n \left( \tilde{Y}_i - \hat{\beta}^\top X_i \right)^2 \mid \underline{X}, \underline{Y} \right),$$

with, for all  $i$ ,  $\tilde{Y}_i$  and  $Y_i$  iid conditionally to  $\underline{X}$ .

### Proposition

Assume that the unknown distribution  $P^{X,Y}$  is such that  $Y_i = \beta^\top X_i + \varepsilon_i$ , with  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ , independent of  $X_i$ . Then

$$\mathbb{E} \left( \tilde{\mathcal{R}}_n \right) = \sigma^2 \left( 1 + \frac{p+1}{n} \right),$$

$$\mathbb{E} \left( \hat{\mathcal{R}}_n \right) = \sigma^2 \left( 1 - \frac{p+1}{n} \right).$$

## Zoom on an illuminating special case (cont'd)

**Interpretation.** On average, the empirical risk under-estimates the generalization error:

$$\mathbb{E} \left( \tilde{\mathcal{R}}_n - \hat{\mathcal{R}}_n \right) = 2 \frac{p+1}{n} \sigma^2 > 0.$$

Another way of looking at this result. Set

$$\eta = \frac{p+1}{n} = \frac{\text{number of coefficients}}{\text{sample size}}.$$

Then

$$\frac{\mathbb{E} \left( \tilde{\mathcal{R}}_n \right)}{\mathbb{E} \left( \hat{\mathcal{R}}_n \right)} = \frac{1+\eta}{1-\eta} \xrightarrow{\eta \rightarrow 1} +\infty.$$



## Zoom on an illuminating special case (cont'd)

**Interpretation.** On average, the empirical risk under-estimates the generalization error:

$$\mathbb{E} \left( \tilde{\mathcal{R}}_n - \hat{\mathcal{R}}_n \right) = 2 \frac{p+1}{n} \sigma^2 > 0.$$

Another way of looking at this result. Set

$$\eta = \frac{p+1}{n} = \frac{\text{number of coefficients}}{\text{sample size}}.$$

Then

$$\frac{\mathbb{E} \left( \tilde{\mathcal{R}}_n \right)}{\mathbb{E} \left( \hat{\mathcal{R}}_n \right)} = \frac{1 + \eta}{1 - \eta} \xrightarrow[\eta \rightarrow 1]{} +\infty.$$

## Zoom on an illuminating special case (cont'd)

**Proof.** Let us compute first  $\mathbb{E}(\tilde{\mathcal{R}}_n \mid \underline{X})$  with (reminder)

$$\tilde{\mathcal{R}}_n = \mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n \left( \tilde{Y}_i - \hat{\beta}^\top X_i \right)^2 \mid \underline{X}, \underline{Y} \right).$$

We have  $\mathbb{E}(\tilde{Y}_i \mid \underline{X}) = \mathbb{E}(\hat{\beta}^\top X_i \mid \underline{X}) = \beta^\top X_i$ , therefore

$$\begin{aligned} \mathbb{E}(\tilde{\mathcal{R}}_n \mid \underline{X}) &= \frac{1}{n} \sum_{i=1}^n \text{var}(\tilde{Y}_i - \hat{\beta}^\top X_i \mid \underline{X}) \\ &= \frac{1}{n} \sum_{i=1}^n \left( \underbrace{\text{var}(\tilde{Y}_i \mid \underline{X})}_{=\sigma^2} + \underbrace{\text{var}(\hat{\beta}^\top X_i \mid \underline{X})}_{=0} \right). \end{aligned}$$

## Zoom on an illuminating special case (cont'd)

**Proof.** Let us compute first  $\mathbb{E}(\tilde{\mathcal{R}}_n \mid \underline{X})$  with (reminder)

$$\tilde{\mathcal{R}}_n = \mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n \left( \tilde{Y}_i - \hat{\beta}^\top X_i \right)^2 \mid \underline{X}, \underline{Y} \right).$$

We have  $\mathbb{E}(\tilde{Y}_i \mid \underline{X}) = \mathbb{E}(\hat{\beta}^\top X_i \mid \underline{X}) = \beta^\top X_i$ , therefore

$$\begin{aligned} \mathbb{E}(\tilde{\mathcal{R}}_n \mid \underline{X}) &= \frac{1}{n} \sum_{i=1}^n \text{var}(\tilde{Y}_i - \hat{\beta}^\top X_i \mid \underline{X}) \\ &= \frac{1}{n} \sum_{i=1}^n \left( \underbrace{\text{var}(\tilde{Y}_i \mid \underline{X})}_{=\sigma^2} + \underbrace{\text{var}(\hat{\beta}^\top X_i \mid \underline{X})}_{=0} \right). \end{aligned}$$

## Zoom on an illuminating special case (cont'd)

**Proof.** Let us compute first  $\mathbb{E}(\tilde{\mathcal{R}}_n \mid \underline{X})$  with (reminder)

$$\tilde{\mathcal{R}}_n = \mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n \left( \tilde{Y}_i - \hat{\beta}^\top X_i \right)^2 \mid \underline{X}, \underline{Y} \right).$$

We have  $\mathbb{E}(\tilde{Y}_i \mid \underline{X}) = \mathbb{E}(\hat{\beta}^\top X_i \mid \underline{X}) = \beta^\top X_i$ , therefore

$$\begin{aligned} \mathbb{E}(\tilde{\mathcal{R}}_n \mid \underline{X}) &= \frac{1}{n} \sum_{i=1}^n \text{var}(\tilde{Y}_i - \hat{\beta}^\top X_i \mid \underline{X}) \\ &= \frac{1}{n} \sum_{i=1}^n \left( \underbrace{\text{var}(\tilde{Y}_i \mid \underline{X})}_{=\sigma^2} + \underbrace{\text{var}(\hat{\beta}^\top X_i \mid \underline{X})}_{=\circledast} \right). \end{aligned}$$

## Zoom on an illuminating special case (cont'd)

We already know that  $\text{var}(\hat{\beta} \mid \underline{X}) = \sigma^2 (\underline{X}^\top \underline{X})^{-1}$ . Therefore:

$$\begin{aligned} \circledast &= \text{var}(\hat{\beta}^\top X_i \mid \underline{X}) \\ &= X_i^\top \text{var}(\hat{\beta} \mid \underline{X}) X_i \\ &= \sigma^2 X_i^\top (\underline{X}^\top \underline{X})^{-1} X_i \\ &= \sigma^2 \text{tr}\left((\underline{X}^\top \underline{X})^{-1} X_i X_i^\top\right). \end{aligned}$$

By noting that  $\underline{X}^\top \underline{X} = \sum_i X_i X_i^\top$ , we get:

$$\begin{aligned} \sum_i \text{var}(\hat{\beta}^\top X_i \mid \underline{X}) &= \sigma^2 \text{tr}\left((\underline{X}^\top \underline{X})^{-1} \sum_i X_i X_i^\top\right) \\ &= \sigma^2 \text{tr}(I_{p+1}) = \sigma^2 (p+1). \end{aligned}$$

## Zoom on an illuminating special case (cont'd)

We already know that  $\text{var}(\hat{\beta} \mid \underline{X}) = \sigma^2 (\underline{X}^\top \underline{X})^{-1}$ . Therefore:

$$\begin{aligned} \circledast &= \text{var}(\hat{\beta}^\top X_i \mid \underline{X}) \\ &= X_i^\top \text{var}(\hat{\beta} \mid \underline{X}) X_i \\ &= \sigma^2 X_i^\top (\underline{X}^\top \underline{X})^{-1} X_i \\ &= \sigma^2 \text{tr}\left((\underline{X}^\top \underline{X})^{-1} X_i X_i^\top\right). \end{aligned}$$

By noting that  $\underline{X}^\top \underline{X} = \sum_i X_i X_i^\top$ , we get:

$$\begin{aligned} \sum_i \text{var}(\hat{\beta}^\top X_i \mid \underline{X}) &= \sigma^2 \text{tr}\left((\underline{X}^\top \underline{X})^{-1} \sum_i X_i X_i^\top\right) \\ &= \sigma^2 \text{tr}(I_{p+1}) = \sigma^2 (p+1). \end{aligned}$$

## Zoom on an illuminating special case (cont'd)

We already know that  $\text{var}(\hat{\beta} \mid \underline{X}) = \sigma^2 (\underline{X}^\top \underline{X})^{-1}$ . Therefore:

$$\begin{aligned} \circledast &= \text{var}(\hat{\beta}^\top X_i \mid \underline{X}) \\ &= X_i^\top \text{var}(\hat{\beta} \mid \underline{X}) X_i \\ &= \sigma^2 X_i^\top (\underline{X}^\top \underline{X})^{-1} X_i \\ &= \sigma^2 \text{tr}\left((\underline{X}^\top \underline{X})^{-1} X_i X_i^\top\right). \end{aligned}$$

By noting that  $\underline{X}^\top \underline{X} = \sum_i X_i X_i^\top$ , we get:

$$\begin{aligned} \sum_i \text{var}(\hat{\beta}^\top X_i \mid \underline{X}) &= \sigma^2 \text{tr}\left((\underline{X}^\top \underline{X})^{-1} \sum_i X_i X_i^\top\right) \\ &= \sigma^2 \text{tr}(I_{p+1}) = \sigma^2 (p+1). \end{aligned}$$

## Zoom on an illuminating special case (cont'd)

We already know that  $\text{var}(\hat{\beta} \mid \underline{X}) = \sigma^2 (\underline{X}^\top \underline{X})^{-1}$ . Therefore:

$$\begin{aligned} \circledast &= \text{var}(\hat{\beta}^\top X_i \mid \underline{X}) \\ &= X_i^\top \text{var}(\hat{\beta} \mid \underline{X}) X_i \\ &= \sigma^2 X_i^\top (\underline{X}^\top \underline{X})^{-1} X_i \\ &= \sigma^2 \text{tr}\left((\underline{X}^\top \underline{X})^{-1} X_i X_i^\top\right). \end{aligned}$$

By noting that  $\underline{X}^\top \underline{X} = \sum_i X_i X_i^\top$ , we get:

$$\begin{aligned} \sum_i \text{var}(\hat{\beta}^\top X_i \mid \underline{X}) &= \sigma^2 \text{tr}\left((\underline{X}^\top \underline{X})^{-1} \sum_i X_i X_i^\top\right) \\ &= \sigma^2 \text{tr}(I_{p+1}) = \sigma^2 (p+1). \end{aligned}$$



## Zoom on an illuminating special case (cont'd)

Thus, we have:

$$\begin{aligned}\mathbb{E}\left(\tilde{\mathcal{R}}_n \mid \underline{X}\right) &= \frac{1}{n} \sum_{i=1}^n \left( \underbrace{\text{var}\left(\tilde{Y}_i \mid \underline{X}\right)}_{=\sigma^2} + \underbrace{\text{var}\left(\hat{\beta}^\top X_i \mid \underline{X}\right)}_{=0} \right) \\ &= \sigma^2 + \sigma^2 \frac{p+1}{n} = \sigma^2 \left(1 + \frac{p+1}{n}\right).\end{aligned}$$

Hence the result:  $\mathbb{E}\left(\tilde{\mathcal{R}}_n\right) = \sigma^2 \left(1 + \frac{p+1}{n}\right)$ .

Exercise: prove the second inequality, i.e.,

$$\mathbb{E}\left(\hat{\mathcal{R}}_n\right) = \sigma^2 \left(1 - \frac{p+1}{n}\right).$$

⇒ see PC



## Zoom on an illuminating special case (cont'd)

Thus, we have:

$$\begin{aligned}\mathbb{E}\left(\tilde{\mathcal{R}}_n \mid \underline{X}\right) &= \frac{1}{n} \sum_{i=1}^n \left( \underbrace{\text{var}\left(\tilde{Y}_i \mid \underline{X}\right)}_{=\sigma^2} + \underbrace{\text{var}\left(\hat{\beta}^\top X_i \mid \underline{X}\right)}_{=0} \right) \\ &= \sigma^2 + \sigma^2 \frac{p+1}{n} = \sigma^2 \left(1 + \frac{p+1}{n}\right).\end{aligned}$$

Hence the result:  $\mathbb{E}\left(\tilde{\mathcal{R}}_n\right) = \sigma^2 \left(1 + \frac{p+1}{n}\right)$ .

Exercise: prove the second inequality, i.e.,

$$\mathbb{E}\left(\hat{\mathcal{R}}_n\right) = \sigma^2 \left(1 - \frac{p+1}{n}\right).$$

⇒ see PC



# Lecture outline

1 – Regularized regression (or classification): penalization

2 – Estimation of the risk (generalization error)

2.1 – Problem

2.2 – Zoom in on an illuminating special case

2.3 – Training set and test set

3 – Hyper-parameters, model selection

4 – Exercises and solutions

5 – Appendices

# Training set and test set

**Conclusion/extrapolation.** The empirical risk is in general

- ▶ a **negatively biased estimator** of the risk,
- ▶ with a **bias that is increasing when  $p \nearrow$** .

Solution: split the data in two sets

- ▶ training data: used to construct  $\hat{h}$ ,
- ▶ test data: used to estimate the generalization error.

Example:

**training**  
(e.g., 80%)

**test**  
(20%)

# Training set and test set

**Conclusion/extrapolation.** The empirical risk is in general

- ▶ a negatively biased estimator of the risk,
- ▶ with a bias that is increasing when  $p \nearrow$ .

Solution: split the data in two sets

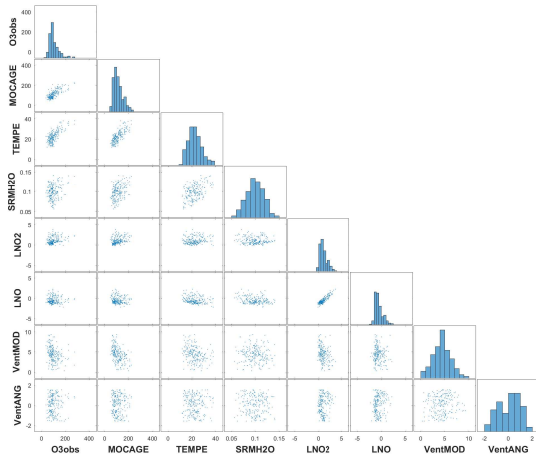
- ▶ **training** data: used to construct  $\hat{h}$ ,
- ▶ **test** data: used to estimate the generalization error.

Example:

**training**  
(e.g., 80%)

**test**  
(20%)

## Example “Ozone” (cont'd from lecture #6)



Goal: predict the ozone concentration on day  $t + 1$   
from data available on day  $t$

## “Ozone” example: 70/30

All 7 explanatory variables and their 21 interactions are used.

Results from 10 random splits, 70% / 30%:

$R^2$	$\hat{\mathcal{R}}_n$	$\hat{\mathcal{R}}_n^{\text{test}}$
0.77185	345.1	573.32
0.76831	371.41	496.03
0.77292	343.96	608.62
0.76093	350.53	606.14
0.78584	345.45	669.66
0.75459	399.9	476.61
0.71367	343.72	643.72
0.77689	377.32	524.74
0.8176	317.83	695.86
0.79784	373.18	554.25

# Lecture outline

1 – Regularized regression (or classification): penalization

2 – Estimation of the risk (generalization error)

3 – Hyper-parameters, model selection

3.1 – Problem

3.2 – Cross validation

3.3 – AIC criterion

4 – Exercises and solutions

5 – Appendices



# Lecture outline

1 – Regularized regression (or classification): penalization

2 – Estimation of the risk (generalization error)

3 – Hyper-parameters, model selection

3.1 – Problem

3.2 – Cross validation

3.3 – AIC criterion

4 – Exercises and solutions

5 – Appendices

## Problem #1: choosing a “good” family $\mathcal{H}$

**Example.** Selection of  $k$  variables among  $p$ . Let  $J \subset \{1, \dots, p\}$ :

$$h(x) = \beta_0 + \sum_{j \in J} \beta_j x^{(j)}.$$

⇒ Defines a family  $\mathcal{H}_J$  with  $k_J = \text{card}(J) + 1$  parameters.

**Example.** Expansion in a basis, truncated at order  $J$  :

$$h(x) = \sum_{k=0}^J \beta_j \psi_j(x).$$

⇒ Defines a family  $\mathcal{H}_J$  with  $k_J = J + 1$  parameters.

⇒ complement

**Problem: model selection**

How to choose the family  $\mathcal{H}_J$  (and, in particular, its “size”  $k_J$ ) ?

**Remark:** replace  $h(x)$  with  $\ln \frac{h(x)}{1-h(x)}$  for logistic regression.

## Problem #1: choosing a “good” family $\mathcal{H}$

**Example.** Selection of  $k$  variables among  $p$ . Let  $J \subset \{1, \dots, p\}$ :

$$h(x) = \beta_0 + \sum_{j \in J} \beta_j x^{(j)}.$$

⇒ Defines a family  $\mathcal{H}_J$  with  $k_J = \text{card}(J) + 1$  parameters.

**Example.** Expansion in a basis, truncated at order  $J$  :

$$h(x) = \sum_{k=0}^J \beta_k \psi_k(x).$$

⇒ Defines a family  $\mathcal{H}_J$  with  $k_J = J + 1$  parameters.

⇒ complement

Problem: model selection

How to choose the family  $\mathcal{H}_J$  (and, in particular, its “size”  $k_J$ ) ?

Remark: replace  $h(x)$  with  $\ln \frac{h(x)}{1-h(x)}$  for logistic regression.

## Problem #1: choosing a “good” family $\mathcal{H}$

**Example.** Selection of  $k$  variables among  $p$ . Let  $J \subset \{1, \dots, p\}$ :

$$h(x) = \beta_0 + \sum_{j \in J} \beta_j x^{(j)}.$$

⇒ Defines a family  $\mathcal{H}_J$  with  $k_J = \text{card}(J) + 1$  parameters.

**Example.** Expansion in a basis, truncated at order  $J$  :

$$h(x) = \sum_{k=0}^J \beta_k \psi_k(x).$$

⇒ Defines a family  $\mathcal{H}_J$  with  $k_J = J + 1$  parameters.

⇒ complement

### Problem: model selection

How to choose the family  $\mathcal{H}_J$  (and, in particular, its “size”  $k_J$ ) ?

Remark: replace  $h(x)$  with  $\ln \frac{h(x)}{1-h(x)}$  for logistic regression.

## Problem #2: choosing a regularization hyper-parameter

Most methods require some “tuning”...

- ▶ Ridge/LASSO regression:  $\hat{\beta} = \operatorname{argmin} \hat{\mathcal{R}}_{n,\lambda}^{\text{pen}}$ , with

$$\hat{\mathcal{R}}_{n,\lambda}^{\text{pen}}(\beta) = \hat{\mathcal{R}}_n(\beta) + \lambda \sum_j |\beta_j|^q, \quad q \in \{1, 2\},$$

- ▶ Choosing the number  $k$  of neighbors in a  $k$ -NN model:

$$h(x) = \frac{1}{k} \sum_{i \in \mathcal{V}_{n,k}(x)} y_i,$$

with  $\mathcal{V}_{n,k}(x)$  the indices of the  $k$  nearest neighbors of  $x$ .

### Problem: calibration

How to “tune” the values of such hyperparameters ?

## Problem #2: choosing a regularization hyper-parameter

Most methods require some “tuning”...

- ▶ Ridge/LASSO regression:  $\hat{\beta} = \operatorname{argmin} \hat{\mathcal{R}}_{n,\lambda}^{\text{pen}}$ , with

$$\hat{\mathcal{R}}_{n,\lambda}^{\text{pen}}(\beta) = \hat{\mathcal{R}}_n(\beta) + \lambda \sum_j |\beta_j|^q, \quad q \in \{1, 2\},$$

- ▶ Choosing the number  $k$  of neighbors in a  $k$ -NN model:

$$h(x) = \frac{1}{k} \sum_{i \in \mathcal{V}_{n,k}(x)} y_i,$$

with  $\mathcal{V}_{n,k}(x)$  the indices of the  $k$  nearest neighbors of  $x$ .

Problem: calibration

How to “tune” the values of such hyperparameters ?

## Problem #2: choosing a regularization hyper-parameter

Most methods require some “tuning”...

- ▶ Ridge/LASSO regression:  $\hat{\beta} = \operatorname{argmin} \hat{\mathcal{R}}_{n,\lambda}^{\text{pen}}$ , with

$$\hat{\mathcal{R}}_{n,\lambda}^{\text{pen}}(\beta) = \hat{\mathcal{R}}_n(\beta) + \lambda \sum_j |\beta_j|^q, \quad q \in \{1, 2\},$$

- ▶ Choosing the number  $k$  of neighbors in a  $k$ -NN model:

$$h(x) = \frac{1}{k} \sum_{i \in \mathcal{V}_{n,k}(x)} y_i,$$

with  $\mathcal{V}_{n,k}(x)$  the indices of the  $k$  nearest neighbors of  $x$ .

**Problem: calibration**

How to “tune” the values of such hyperparameters ?

# Over-fitting: beware!

## Idea

Choose the family  $\mathcal{H}_J$ , or the hyperparameter  $\lambda$ , in order to **minimize** (an estimation of) the **generalization error**.

 again, the empirical risk  $\hat{\mathcal{R}}_n$ , estimated on the training data, is not appropriate !

**Example.** Polynomial regression with  $x \in \mathbb{R}$ , degree  $\leq J$ :

$$h(x) = \beta_0 + \beta_1 x + \dots + \beta_J x^J,$$

with  $J = 2, 5, 8, 11$ .


Recall that, in linear regression, the empirical risk has a downward bias proportional to the number of parameters in the model.



# Over-fitting: beware!

## Idea

Choose the family  $\mathcal{H}_J$ , or the hyperparameter  $\lambda$ , in order to minimize (an estimation of) the generalization error.

 again, the empirical risk  $\hat{\mathcal{R}}_n$ , estimated on the training data, is not appropriate !

Example. Polynomial regression with  $x \in \mathbb{R}$ , degree  $\leq J$ :

$$h(x) = \beta_0 + \beta_1 x + \dots + \beta_J x^J,$$


with  $J = 2, 5, 8, 11$ .

Recall that, in linear regression, the empirical risk has a downward bias proportional to the number of parameters in the model.

# Over-fitting: beware!

## Idea

Choose the family  $\mathcal{H}_J$ , or the hyperparameter  $\lambda$ , in order to minimize (an estimation of) the generalization error.

 again, the empirical risk  $\hat{\mathcal{R}}_n$ , estimated on the training data, is not appropriate !

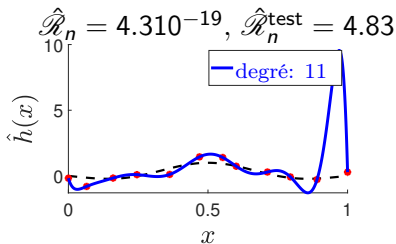
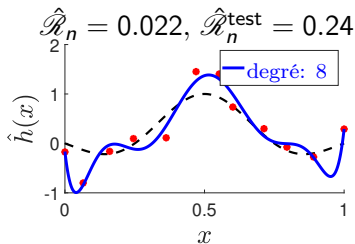
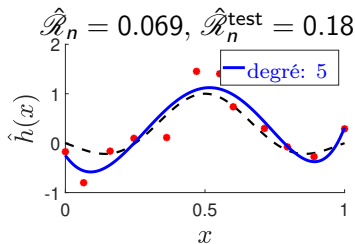
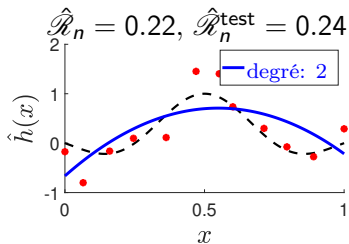
**Example.** Polynomial regression with  $x \in \mathbb{R}$ , degree  $\leq J$ :

$$h(x) = \beta_0 + \beta_1 x + \dots + \beta_J x^J,$$

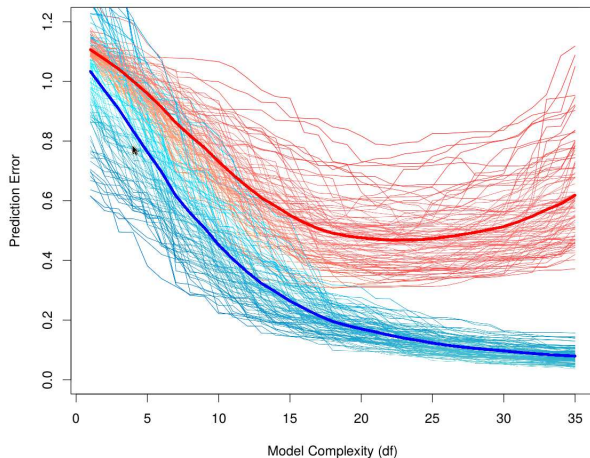
with  $J = 2, 5, 8, 11$ .

Recall that, in linear regression, the empirical risk has a downward bias proportional to the number of parameters in the model.

## Example: polynomial regression



# Understanding over-fitting: simulations



Blue: empirical risk  $\hat{\mathcal{R}}_n$  / Red: error on the test set

Figure from Hastie, Tibshirani & Friedman (2017).  
*The Elements of Statistical Learning* (12th edition), Springer.

## Let's recapitulate. . .

**Problem.** We want to estimate the error to choose  $\mathcal{H}$  or  $\lambda$  but. . .

- ▶ it should be done neither on the **training data**  
( $\Rightarrow$  **over-fitting** problem),
- ▶ nor on the test data  
( $\Rightarrow$  bias in the final estimation of the generalization error).

## Let's recapitulate...

**Problem.** We want to estimate the error to choose  $\mathcal{H}$  or  $\lambda$  but...

- ▶ it should be done neither on the training data  
( $\Rightarrow$  over-fitting problem),
- ▶ nor on the **test data**  
( $\Rightarrow$  **bias** in the final estimation of the generalization error).



# Lecture outline

1 – Regularized regression (or classification): penalization

2 – Estimation of the risk (generalization error)

3 – Hyper-parameters, model selection

3.1 – Problem

3.2 – Cross validation

3.3 – AIC criterion

4 – Exercises and solutions

5 – Appendices

# Solution: validation set

Idea: split the data in three sets

- ▶ **training** data: construct  $\hat{h}$  with given  $\mathcal{H}/\lambda$ ,
- ▶ **validation** set: choose  $\mathcal{H}$ ,  $\lambda$ , etc.
- ▶ **test** data: estimate the generalization error.

Simple validation (hold-out)

**training**  
(e.g., 60%)

**validation**  
(e.g., 20%)

**test**  
(e.g., 20%)



## Solution: validation set

Idea: split the data in three sets

- ▶ training data: construct  $\hat{h}$  with given  $\mathcal{H}/\lambda$ ,
- ▶ validation set: choose  $\mathcal{H}$ ,  $\lambda$ , etc.
- ▶ test data: estimate the generalization error.

Simple validation (hold-out)

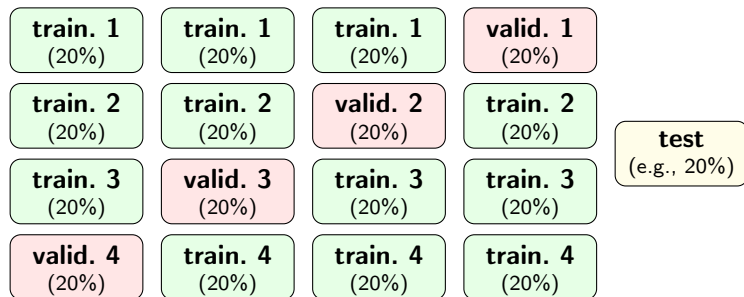
**training**  
(e.g., 60%)

**validation**  
(e.g., 20%)

**test**  
(e.g., 20%)

# Better validation: the cross validation method

*k*-fold cross-validation, here with  $k = 4$ :



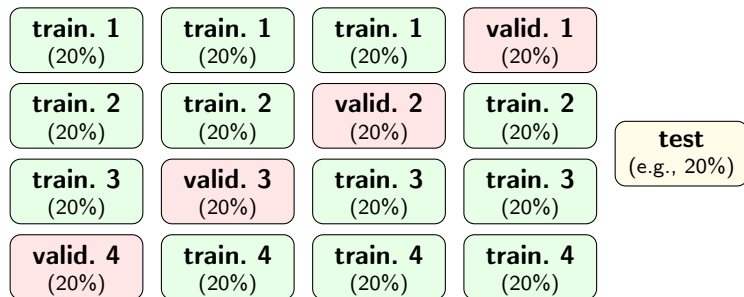
➡ the error is averaged over the  $k$  validation sets.

Special case: leave-one-out cross validation

▶  $k = n$  blocks (of size  $n/k = 1$ ).

# Better validation: the cross validation method

$k$ -fold cross-validation, here with  $k = 4$ :



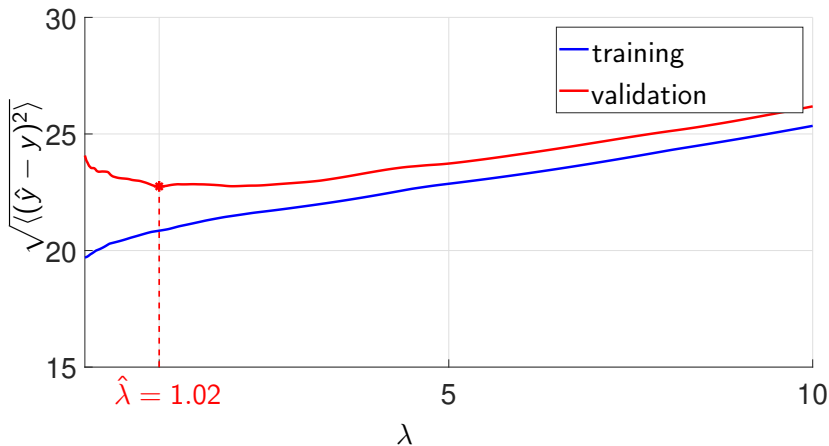
➡ the error is averaged over the  $k$  validation sets.

Special case: **leave-one-out** cross validation

►  $k = n$  blocks (of size  $n/k = 1$ ).

## “Ozone” example: LASSO / choice of $\lambda$

- ▶ Predictor: LASSO regression using all variables and their interactions
- ▶  $\hat{\lambda}$  obtained by CV (LOO)



## “Ozone” example: interactions

- ▶ We add variables of the form  $X^{(j)}X^{(j')}$  and  $X^{(j)}X^{(j')}X^{(j'')}$ .
- ▶ LASSO regression ( $L^1$  penalty).
- ▶ Hyper-parameter  $\lambda$  estimated through 10-fold CV.

model	$X^{(j)}$	$X^{(j)} X^{(j')}$	$X^{(j)} X^{(j')} X^{(j'')}$
total number of variables	7	35	119
number of selected variables ( $\beta_j \neq 0$ )	4	9	8
$\sqrt{MSE}$ CV (10-fold)	49.1	41.5	33.0
selected variables	MOCAGE TEMPE NO VentANG	MOCAGE TEMPE NO2 <b>MOCAGE · TEMPE</b> <b>TEMPE<sup>2</sup></b> TEMPE · MH2O TEMPE · NO2 NO2 · VentANG VentANG · VentANG	MOCAGE TEMPE NO2 MOCAGE · TEMPE TEMPE <sup>2</sup> TEMPE · RMH2O <b>TEMPE<sup>2</sup> · MOCAGE</b> VentANG <sup>2</sup> · TEMPE

# Lecture outline

1 – Regularized regression (or classification): penalization

2 – Estimation of the risk (generalization error)

3 – Hyper-parameters, model selection

3.1 – Problem

3.2 – Cross validation

3.3 – AIC criterion

4 – Exercises and solutions

5 – Appendices

## Another approach to model selection: the AIC criterion

Assumption: parametric statistical models  $\mathcal{M}_j$  for  $P^{Y|X}$ .

Denote by  $\hat{\theta}_j^{\text{MLE}}$  the MLE of  $\theta$  in model  $\mathcal{M}_j$ .

Then the AIC criterion can also be used for model selection:

$$\hat{j} = \operatorname{argmin} \operatorname{AIC}(j), \quad \operatorname{AIC}(j) = -2 \ln \mathcal{L} \left( \hat{\theta}_j^{\text{MLE}}; \underline{X}, \underline{Y} \right) + 2k_j,$$

with  $k_j$  the number of parameters in model  $\mathcal{M}_j$ .

⇒ see PC for a partial justification (OLS linear regression)

## Another approach to model selection: the AIC criterion

Assumption: parametric statistical models  $\mathcal{M}_j$  for  $P^{Y|X}$ .

Denote by  $\hat{\theta}_j^{\text{MLE}}$  the MLE of  $\theta$  in model  $\mathcal{M}_j$ .

Then the AIC criterion can also be used for model selection:

$$\hat{j} = \operatorname{argmin} \operatorname{AIC}(j), \quad \operatorname{AIC}(j) = -2 \ln \mathcal{L} \left( \hat{\theta}_j^{\text{MLE}}; \underline{X}, \underline{Y} \right) + 2k_j,$$

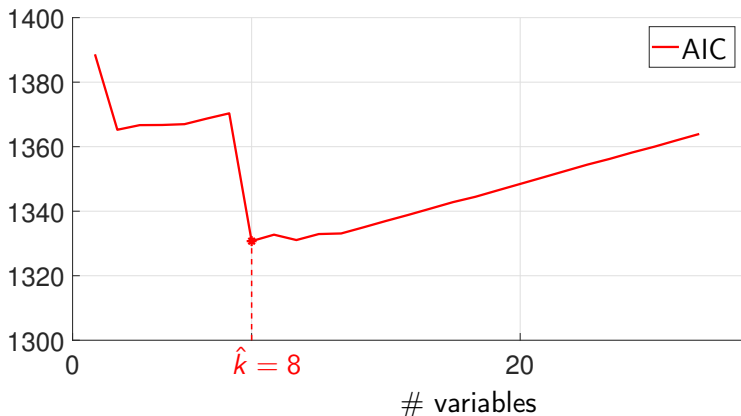
with  $k_j$  the number of parameters in model  $\mathcal{M}_j$ .

⇒ see PC for a partial justification (OLS linear regression)



## “Ozone” example: AIC

- Predictor obtained by the ordinary least squares method, on an increasing number of variables  
(linear terms first, then interactions)



# Lecture outline

1 – Regularized regression (or classification): penalization

2 – Estimation of the risk (generalization error)

3 – Hyper-parameters, model selection

4 – Exercises and solutions

4.1 – Questions

4.2 – Solutions

5 – Appendices

# Lecture outline

1 – Regularized regression (or classification): penalization

2 – Estimation of the risk (generalization error)

3 – Hyper-parameters, model selection

4 – Exercises and solutions

4.1 – Questions

4.2 – Solutions

5 – Appendices

## Exercise 1 (Penalized regression)

Let  $X_1, \dots, X_n$  represent the examples, taking values in  $\mathbb{R}^p$ , and  $Y_1, \dots, Y_n$  be the labels, taking values in  $\mathbb{R}$ . The relationship between  $Y_i$  and  $X_i$  is given by:

$$Y_i = \beta_1 X_i^{(1)} + \dots + \beta_p X_i^{(p)} + \varepsilon_i,$$

where  $\beta$  is the parameter vector to be estimated, and  $\varepsilon_i$  is a random variable following  $\mathcal{N}(0, \sigma^2)$ , independent of  $X_i$ .

We aim to estimate  $\beta$  by minimizing a criterion of the form

$$\frac{1}{2} \sum_{i=1}^n \left( Y_i - \beta^\top X_i \right)^2 + \lambda \mathcal{P}(\beta) \tag{1}$$

where  $\mathcal{P}$  is a penalty term, and  $\lambda \geq 0$  is a hyper-parameter.

We denote  $X = [X_1 \dots X_n]^\top$ , the  $n \times p$  matrix containing the observations. **We are considering the case where  $X^\top X = I_p$ .**

### Question

- 1 Give the expression of the estimator when  $\lambda = 0$ . Denote this estimator  $\hat{\beta}$ .
- 2 We consider a penalty of the form  $\mathcal{P}(\beta) = \|\beta\|_2^2$ . Give the expression of this estimator, denoted  $\hat{\beta}^R$ , and deduce that there exists a constant  $c_{1,\lambda}$  (to be specified) such that  $\hat{\beta}_j^R = c_{1,\lambda} \hat{\beta}_j$ ,  $j = 1, \dots, p$ .

## Question

- ③ We consider a penalty of the form  $\mathcal{P}(\beta) = \|\beta\|_1$ .  
To begin with, demonstrate that the minimum on  $\mathbb{R}$  of the function

$$f : \alpha \mapsto \frac{1}{2}(x - \alpha)^2 + \lambda |\alpha|$$

is achieved at  $\alpha^* = \text{sign}(x) \max(0, |x| - \lambda)$ .

- ④ Deduce the solution of the optimization problem (1) for  $\mathcal{P}(\beta) = \|\beta\|_1$ , which will be expressed in terms of  $\hat{\beta}$ . Denote this estimator  $\hat{\beta}^L$ .

# Lecture outline

1 – Regularized regression (or classification): penalization

2 – Estimation of the risk (generalization error)

3 – Hyper-parameters, model selection

4 – Exercises and solutions

4.1 – Questions

4.2 – Solutions

5 – Appendices

- ① We recognize the least squares criterion, and we have:

$$\hat{\beta} = (X^T X)^{-1} X^T Y = X^T Y$$

- ② This corresponds to ridge regression..

$$\begin{aligned}\hat{\beta}^R &= (X^T X + 2\lambda I)^{-1} X^T Y \\ &= (1 + 2\lambda)^{-1} \hat{\beta}\end{aligned}$$

Therefore  $\hat{\beta}_j^R = (1 + 2\lambda)^{-1} \hat{\beta}_j$ .



- ③ The function  $f$  is not differentiable, but it is differentiable at every point  $\alpha \neq 0$  and continuous at  $\alpha = 0$ . Thus, we can determine its minimum by analyzing its variations using the sign of the derivative, as if it were differentiable everywhere. The derivative at every  $\alpha \neq 0$  is given by

$$f'(\alpha) = \begin{cases} \alpha - x + \lambda & \text{si } \alpha > 0, \\ \alpha - x - \lambda & \text{si } \alpha < 0, \end{cases}$$

hence

$$f'(\alpha) > 0 \quad \Leftrightarrow \quad (\alpha > x - \lambda \text{ et } \alpha > 0) \text{ ou } (\alpha > x + \lambda \text{ et } \alpha < 0). \quad (2)$$

- ③ Let's consider, for example,  $x > 0$ . Then, the second case in the right-hand side of (2) is impossible, and we're left with:

$$f'(\alpha) > 0 \Leftrightarrow \alpha > x - \lambda \text{ et } \alpha > 0 \Leftrightarrow \alpha > \max(0, x - \lambda). \quad (3)$$

Similarly, still assuming  $x > 0$ ,

$$\begin{aligned} f'(\alpha) < 0 &\Leftrightarrow (\alpha < x - \lambda \text{ et } \alpha > 0) \text{ ou } (\alpha < x + \lambda \text{ et } \alpha < 0) \\ &\Leftrightarrow (0 < \alpha < \max(0, x - \lambda)) \text{ ou } (\alpha < 0) \\ &\Leftrightarrow (\alpha < \max(0, x - \lambda)) \text{ et } (\alpha \neq 0). \end{aligned}$$

Thus,  $f$  strictly decreases to the left of  $\max(0, x - \lambda)$ , and strictly increases to the right, which concludes the case  $x > 0$ . The case  $x < 0$  follows similarly.

- 4 Here, we'll manipulate the initial optimization problem to reduce it to the optimization problem from the previous question.:

$$\begin{aligned}\hat{\beta}^L &= \operatorname{argmin}_{\beta} \frac{1}{2} \|Y - X\beta\|^2 + \lambda \|\beta\|_1 \\ &= \operatorname{argmin}_{\beta} \frac{1}{2} \|Y - X\hat{\beta} + X\hat{\beta} - X\beta\|^2 + \lambda \|\beta\|_1 \\ &= \operatorname{argmin}_{\beta} \frac{1}{2} \left\{ \|Y - X\hat{\beta}\|^2 + \|X\hat{\beta} - X\beta\|^2 \right\} + \lambda \|\beta\|_1\end{aligned}$$

The cross product vanishes because the residual  $(Y - X\hat{\beta})$  is, by construction, orthogonal to any linear combination of columns of  $X$ , thus  $(Y - X\hat{\beta})^\top (X\hat{\beta} - X\beta) = 0$ .

- 4 Since the first term is independent of  $\beta$ , we have:

$$\begin{aligned}\hat{\beta}^L &= \operatorname{argmin}_{\beta} \frac{1}{2} \|X\hat{\beta} - X\beta\|^2 + \lambda \|\beta\|_1 \\ &= \operatorname{argmin}_{\beta} \frac{1}{2} (\hat{\beta} - \beta)^\top X^\top X (\hat{\beta} - \beta) + \lambda \|\beta\|_1 \\ &= \operatorname{argmin}_{\beta} \frac{1}{2} (\hat{\beta} - \beta)^\top (\hat{\beta} - \beta) + \lambda \|\beta\|_1 \\ &= \operatorname{argmin}_{\beta} \frac{1}{2} \sum_{j=1}^p (\hat{\beta}_j - \beta_j)^2 + \lambda |\beta_j|\end{aligned}$$

The problem is separable and, from the previous question, we have:

$$\hat{\beta}_j^L = \operatorname{sign}(\hat{\beta}_j) \max(0, |\hat{\beta}_j| - \lambda)$$

# Lecture outline

1 – Regularized regression (or classification): penalization

2 – Estimation of the risk (generalization error)

3 – Hyper-parameters, model selection

4 – Exercises and solutions

5 – Appendices

5.1 – Model building: feature engineering

# Lecture outline

1 – Regularized regression (or classification): penalization

2 – Estimation of the risk (generalization error)

3 – Hyper-parameters, model selection

4 – Exercises and solutions

5 – Appendices

5.1 – Model building: feature engineering

# Non-linearities in linear models...

If the empirical risk  $\hat{\mathcal{R}}(\hat{h})$  is high, several possible causes:

- ▶ **noise**: intrinsic difficulty in predicting  $Y$ 
  - ▮ irreducible **statistical error**.
- ▶ non-linearity of the optimal predictor wrt the  $X^{(j)}$ 's
  - ▮ reducible approximation error.

**Possible workaround:**  $x^{(1)}, \dots, x^{(p)} \mapsto \tilde{x}^{(1)}, \dots, \tilde{x}^{(q)}$

- ▶ with  $\tilde{x}^{(j)}$  function of  $x^{(1)}, \dots, x^{(p)}$ .
- ▶ The model is still linear with respect to  $\beta$ .

## Non-linearities in linear models...

If the empirical risk  $\hat{\mathcal{R}}(\hat{h})$  is high, several possible causes:

- ▶ noise: intrinsic difficulty in predicting  $Y$ 
  - ▮ irreducible statistical error.
- ▶ **non-linearity** of the optimal predictor wrt the  $X^{(j)}$ 's
  - ▮ reducible **approximation error**.

Possible workaround:  $x^{(1)}, \dots, x^{(p)} \mapsto \tilde{x}^{(1)}, \dots, \tilde{x}^{(q)}$

- ▶ with  $\tilde{x}^{(j)}$  function of  $x^{(1)}, \dots, x^{(p)}$ .
- ▶ The model is still linear with respect to  $\beta$ .



## Non-linearities in linear models...

If the empirical risk  $\hat{\mathcal{R}}(\hat{h})$  is high, several possible causes:

- ▶ noise: intrinsic difficulty in predicting  $Y$ 
  - ▮ irreducible statistical error.
- ▶ non-linearity of the optimal predictor wrt the  $X^{(j)}$ 's
  - ▮ reducible approximation error.

**Possible workaround:**  $x^{(1)}, \dots, x^{(p)} \mapsto \tilde{x}^{(1)}, \dots, \tilde{x}^{(q)}$

- ▶ with  $\tilde{x}^{(j)}$  function of  $x^{(1)}, \dots, x^{(p)}$ .
- ▶ The model is still **linear with respect to  $\beta$** .

# Examples

A few examples:

- ▶ **scalar transformations**:  $\ln(x^{(j)})$ ,  $\sqrt{x^{(j)}}$ ,  $(x^{(j)})^k \dots$
- ▶ **interactions** (here, of order two):  $x^{(j)}x^{(k)}$ ,  $j \neq k$ ,
- ▶ higher-order interactions,
- ▶ (truncated) expansion in a basis...



if  $q \gg p$ , **risk of over-fitting**.

Remarks: *feature engineering*

- ▶ Proposing new relevant variables
  - ▢ domain expertise (or model selection...?)
- ▶ The same principle can be used to *reduce* dimension
  - ▢ features extraction.

# Examples

A few examples:

- ▶ scalar transformations:  $\ln(x^{(j)})$ ,  $\sqrt{x^{(j)}}$ ,  $(x^{(j)})^k \dots$
- ▶ interactions (here, of order two):  $x^{(j)}x^{(k)}$ ,  $j \neq k$ ,
- ▶ higher-order interactions,
- ▶ (truncated) expansion in a basis...



if  $q \gg p$ , risk of over-fitting.

Remarks: *feature engineering*

- ▶ Proposing new relevant variables
  - ▢ *domain expertise* (or model selection...?)
- ▶ The same principle can be used to *reduce* dimension
  - ▢ *features extraction*.

# Expansion in a basis

## Principle

Let  $\{\psi_m\}_{m>0}$  be a function basis of  $L^2(\mathcal{X})^\dagger$ .

Consider  $\tilde{X}^{(m)} = \psi_m(X)$ ,  $m = 1, \dots, M$

⇒ truncated expansion in the basis  $\{\psi_m\}$ .

Examples of bases (preferably orthogonal):

- ▶ polynomial bases,
- ▶ wavelet bases,
- ▶ Fourier bases. . .

<sup>†</sup> or any other function space assumed to contain the optimal predictor  $h^*$ .

# Expansion in a basis

## Principle

Let  $\{\psi_m\}_{m>0}$  be a function basis of  $L^2(\mathcal{X})^\dagger$ .

Consider  $\tilde{X}^{(m)} = \psi_m(X)$ ,  $m = 1, \dots, M$

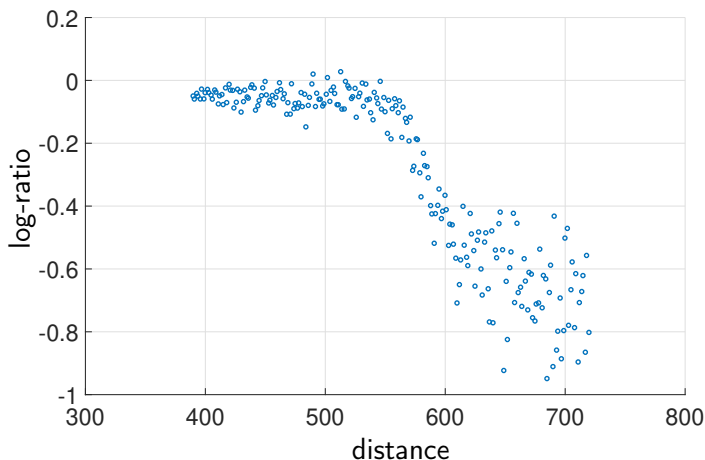
⇒ truncated expansion in the basis  $\{\psi_m\}$ .

**Examples of bases** (preferably orthogonal):

- ▶ polynomial bases,
- ▶ wavelet bases,
- ▶ Fourier bases. . .

<sup>†</sup> or any other function space assumed to contain the optimal predictor  $h^*$ .

## Example: LIDAR data



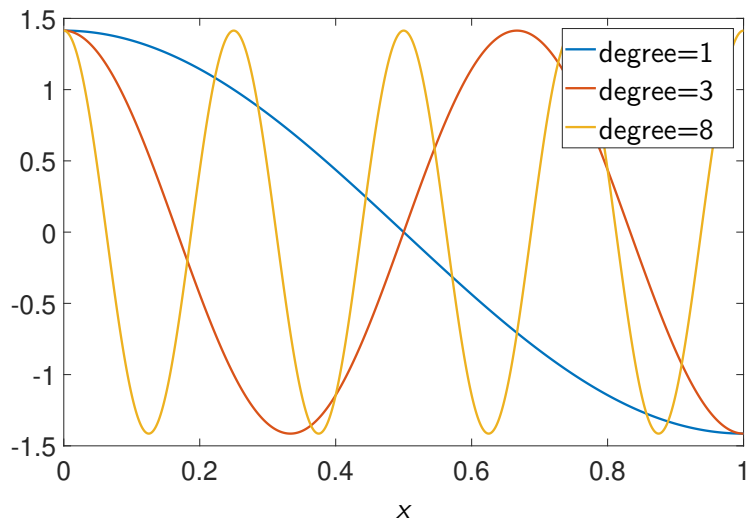
x-axis: distance travelled before the light is reflected back to its source

y-axis: logarithm of the ratio of received light from two laser sources

Data obtained from <http://matt-wand.utsacademics.info/webspr/lidar.html>

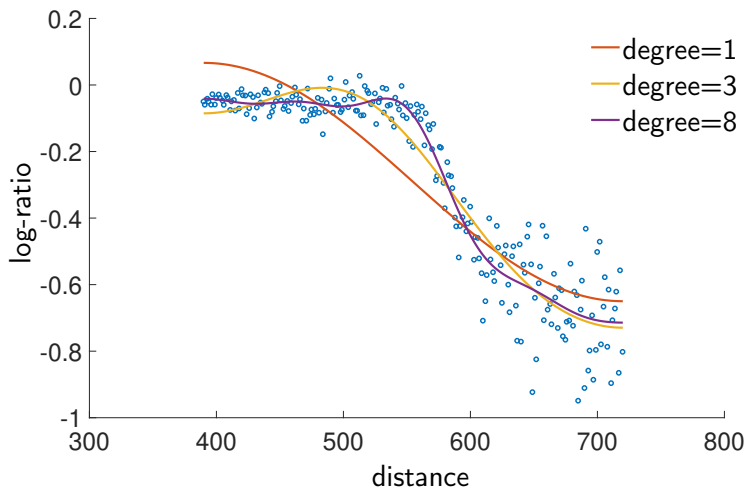
LIDAR: Light Detection And Ranging

## Basis of orthogonal cosines (basis of $L^2([0, 1])$ )



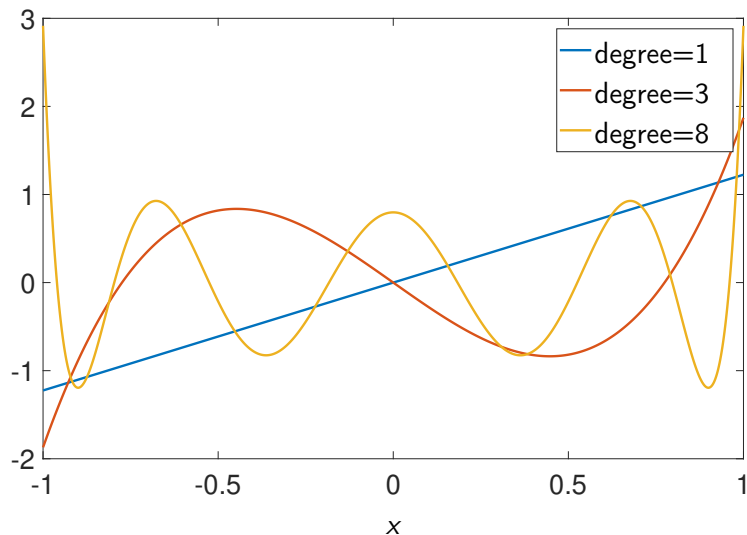
## Example: LIDAR data (cont'd)

Quadratic loss + basis of cosines



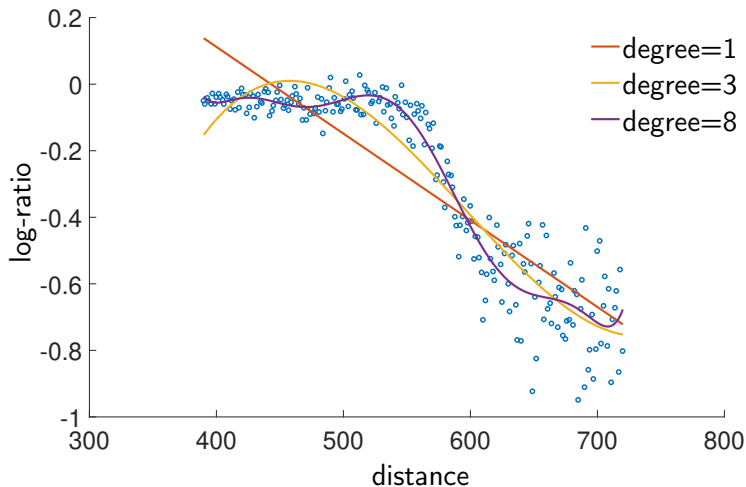


## Legendre polynomials (orthonormal basis of $L^2([-1, 1])$ )

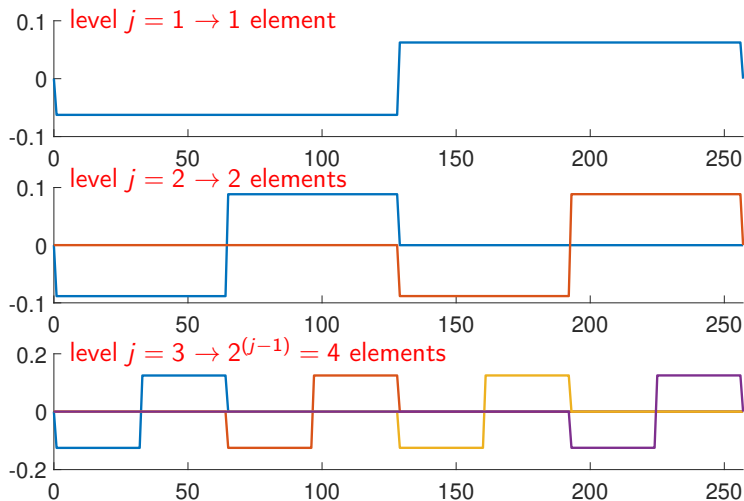


## Example: LIDAR data (cont'd)

Quadratic loss + Legendre polynomials

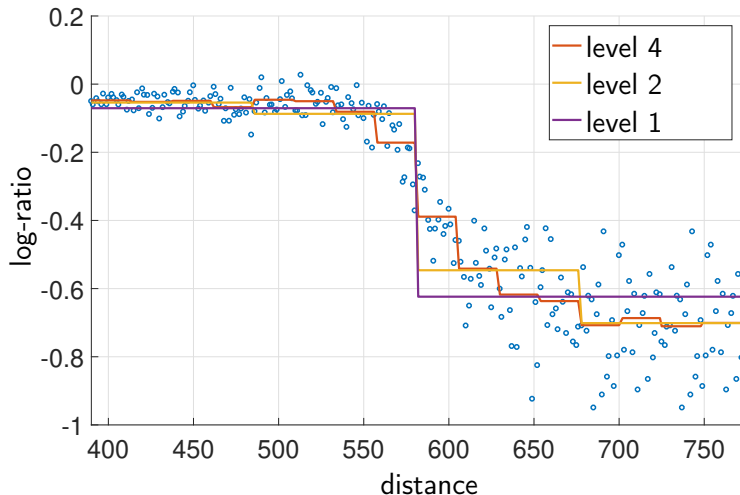


# Haar wavelet basis



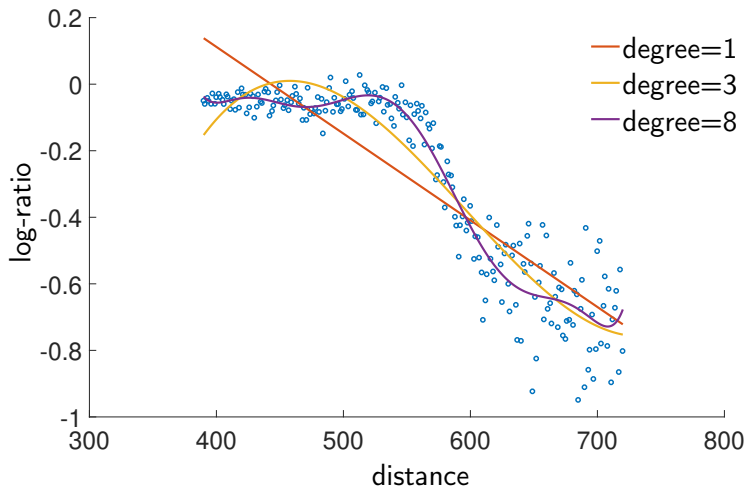
## Example: LIDAR data (cont'd)

Quadratic loss + Haar wavelets



## Example: LIDAR data (cont'd)

Quadratic loss + Legendre polynomials



## Example: LIDAR data (cont'd)

Model selection

