

Sélection de modèles ou d'estimateurs: une introduction

Sylvain Arlot^{1,2,3}

¹Université Paris-Saclay

²Inria Saclay, équipe-projet Celeste

³Institut Universitaire de France

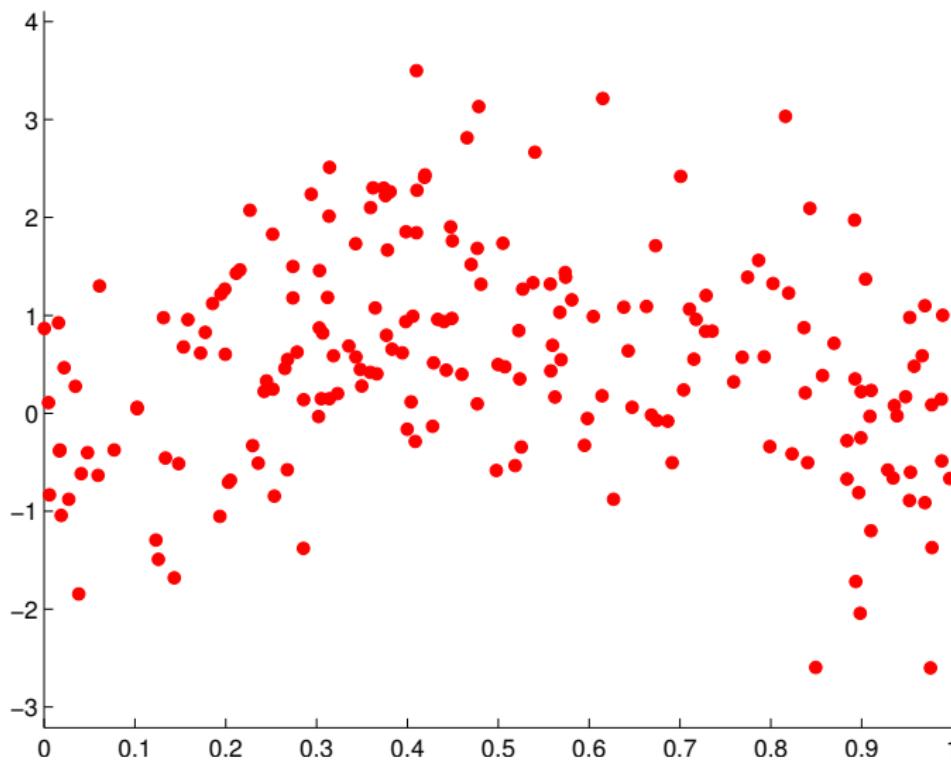
Séminaire Modal'X
23 Septembre 2021

Référence principale: arXiv:1901.07277

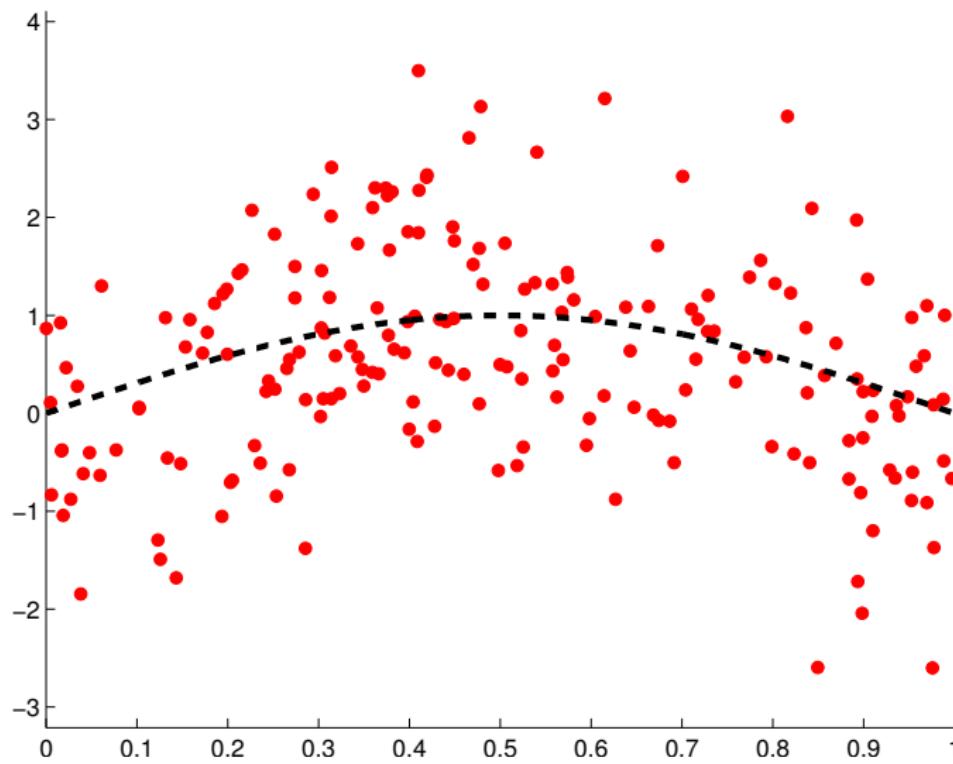
Outline

- 1 Introduction
 - 2 Fixed-design regression
 - 3 Model selection in fixed-design regression
 - 4 General prediction framework
 - 5 Estimator selection

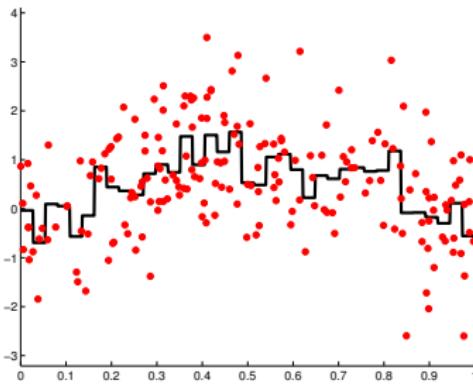
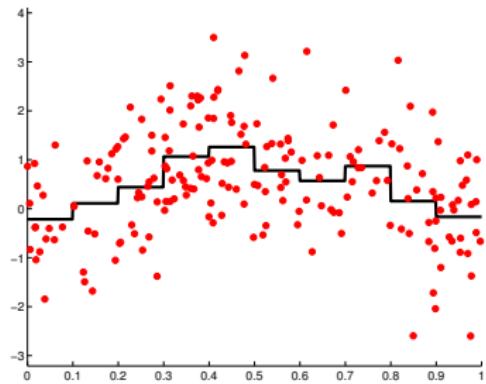
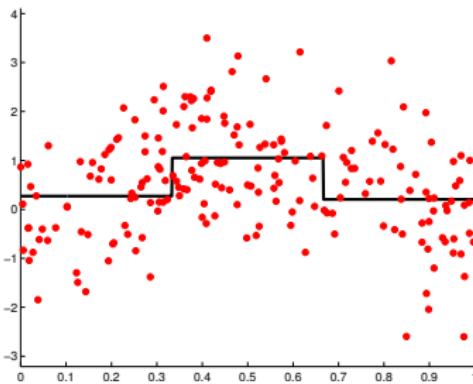
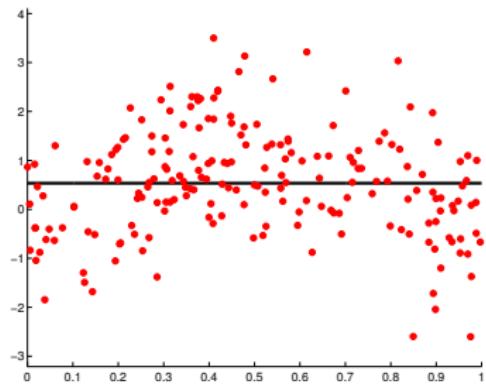
Regression: data $(X_1, Y_1), \dots, (X_n, Y_n)$



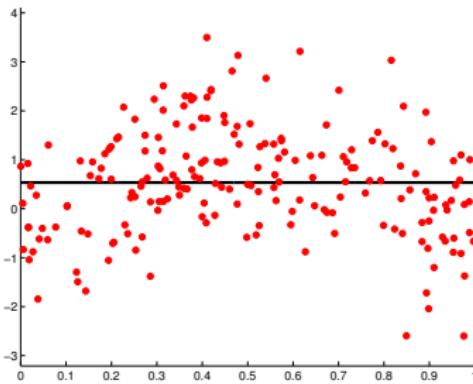
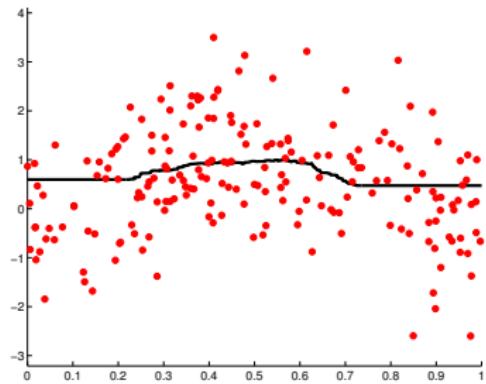
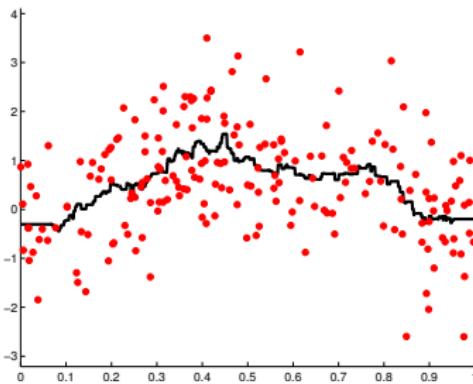
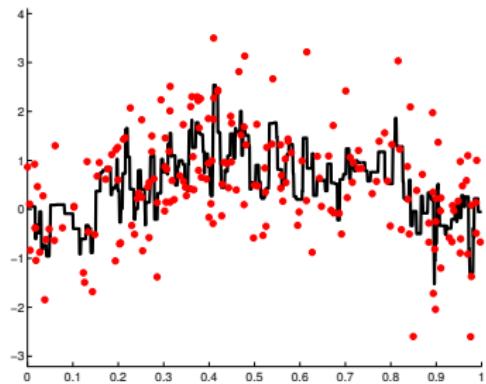
Goal: find the signal (denoising)



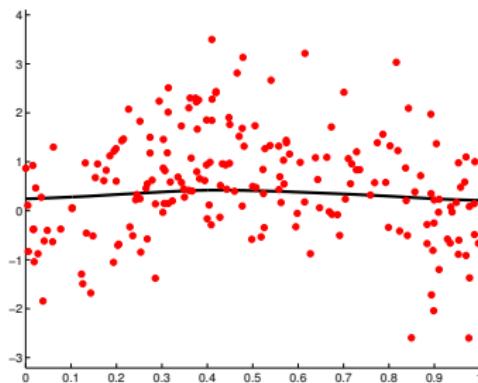
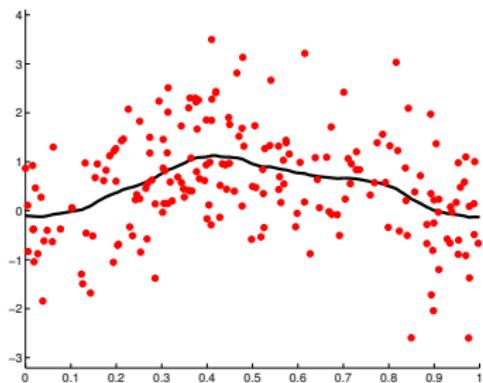
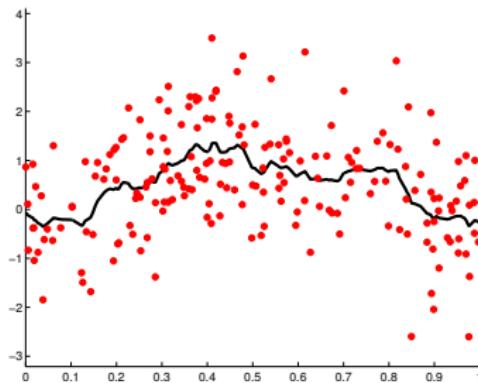
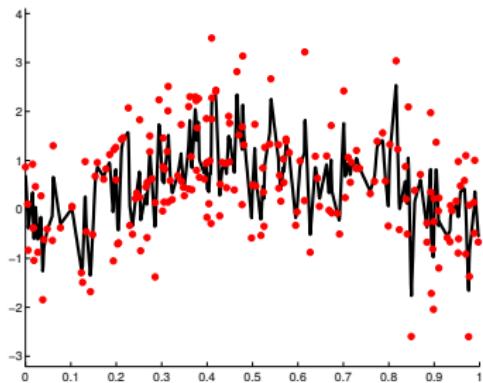
Model selection: regular regressograms, choose D ?



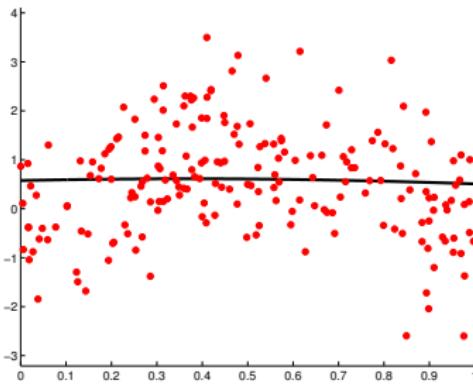
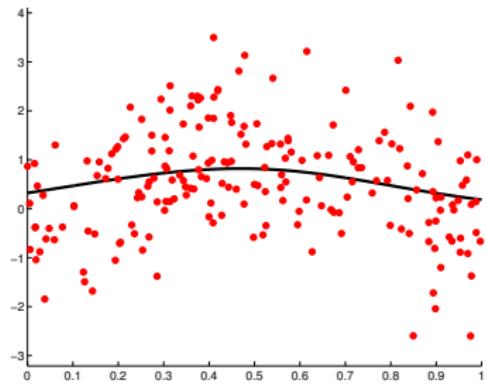
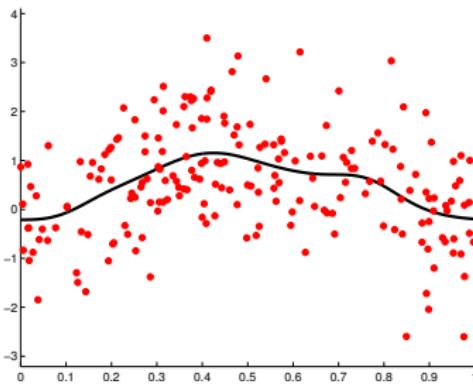
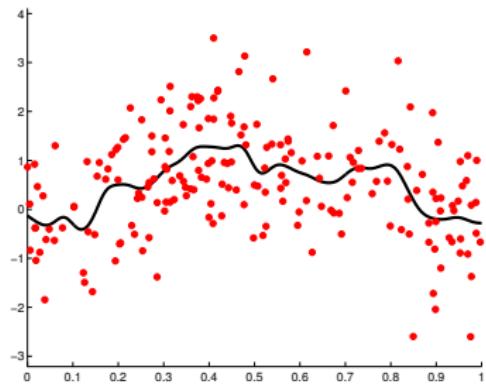
Estimator selection: k nearest neighbours, choose k ?



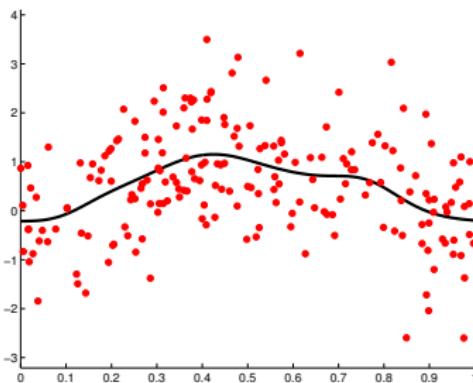
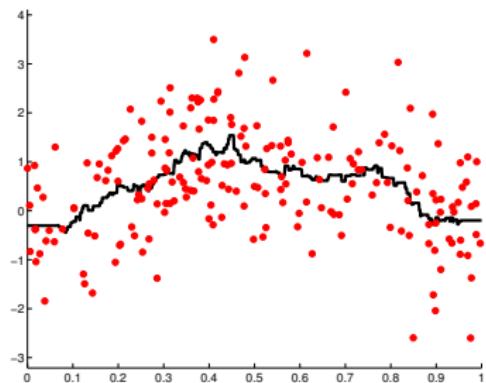
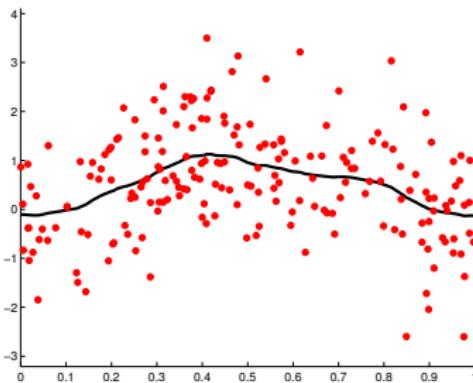
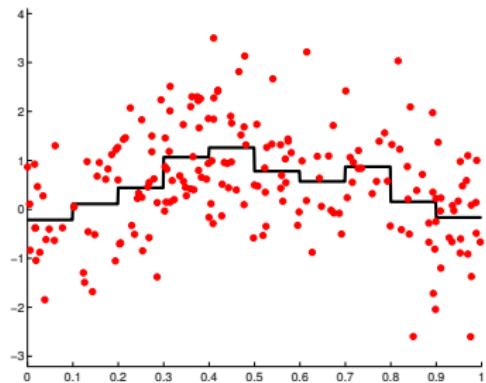
Estimator selection: kernel ridge, choose λ ?



Nadaraya-Watson, choose bandwidth h ?



Regressogram, ridge, k -NN or Nadaraya-Watson?



Outline

- 1 Introduction
- 2 Fixed-design regression
- 3 Model selection in fixed-design regression
- 4 General prediction framework
- 5 Estimator selection

Statistical framework: regression, least-squares risk

- Observations:

$$Y_i = f^*(x_i) + \varepsilon_i \in \mathbb{R}$$

with $\varepsilon_1, \dots, \varepsilon_n$ i.i.d. $\sim \mathcal{N}(0, \sigma^2)$, f^* and σ^2 unknown

- Fixed design: $x_i \in \mathcal{X}$ deterministic
- Notation: $Y = F + \varepsilon$ with

$$Y = (Y_i)_{1 \leqslant i \leqslant n}, \quad F = (f^*(x_i))_{1 \leqslant i \leqslant n}, \quad \varepsilon = (\varepsilon_i)_{1 \leqslant i \leqslant n} \in \mathbb{R}^n$$

- Least-squares risk of a predictor $t \in \mathbb{R}^n$ (" $t_i = t(x_i)$ "):

$$\frac{1}{n} \|t - F\|^2 = \frac{1}{n} \sum_{i=1}^n (t_i - F_i)^2$$

⇒ Estimator $\hat{F}(Y) \in \mathbb{R}^n$?

Least-squares estimators

- Natural idea: minimize an estimator of the risk $\frac{1}{n} \|t - F\|^2$

Least-squares estimators

- Natural idea: minimize an estimator of the risk $\frac{1}{n} \|t - F\|^2$
- Least-squares criterion:

$$\frac{1}{n} \|t - Y\|^2 = \frac{1}{n} \sum_{i=1}^n (t_i - Y_i)^2$$

$$\forall t \in \mathbb{R}^n , \quad \mathbb{E} \left[\frac{1}{n} \|t - Y\|^2 \right] = \frac{1}{n} \|t - F\|^2 + \frac{1}{n} \mathbb{E} [\|\varepsilon\|^2]$$

Least-squares estimators

- Natural idea: minimize an estimator of the risk $\frac{1}{n} \|t - F\|^2$
- Least-squares criterion:

$$\frac{1}{n} \|t - Y\|^2 = \frac{1}{n} \sum_{i=1}^n (t_i - Y_i)^2$$

$$\forall t \in \mathbb{R}^n, \quad \mathbb{E} \left[\frac{1}{n} \|t - Y\|^2 \right] = \frac{1}{n} \|t - F\|^2 + \frac{1}{n} \mathbb{E} [\|\varepsilon\|^2]$$

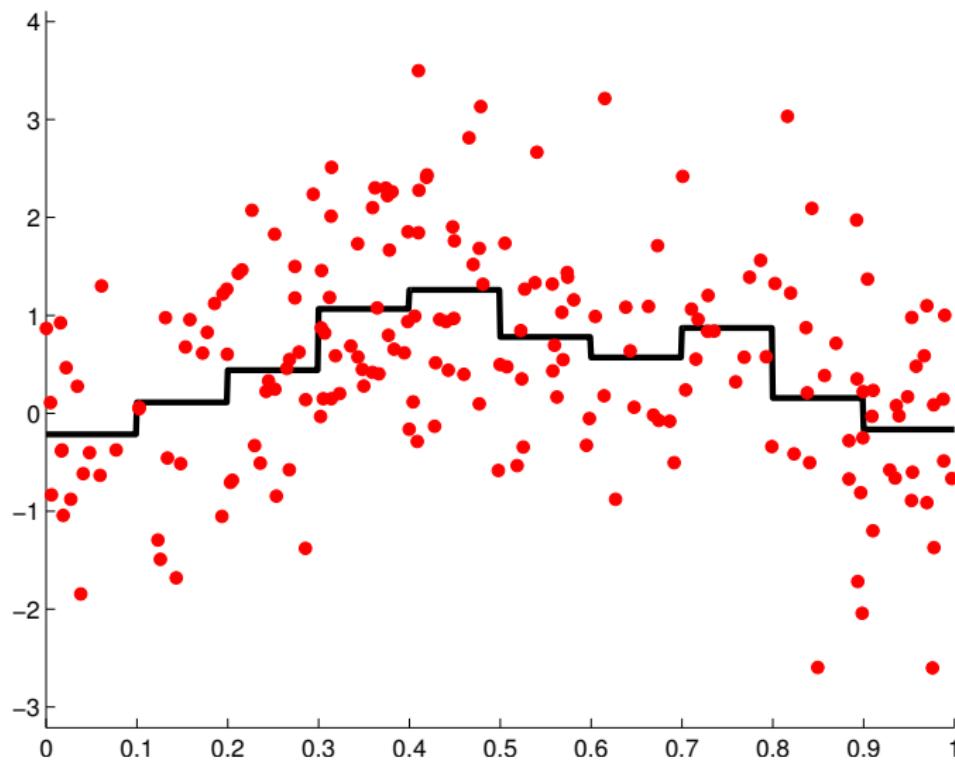
- Model: $S \subset \mathbb{R}^n \Rightarrow$ **Least-squares estimator** on S :

$$\hat{F}_S \in \operatorname{argmin}_{t \in S} \left\{ \frac{1}{n} \|t - Y\|^2 \right\} = \operatorname{argmin}_{t \in S} \left\{ \frac{1}{n} \sum_{i=1}^n (t_i - Y_i)^2 \right\}$$

so that

$$\hat{F}_S = \Pi_S(Y) \quad (\text{orthogonal projection})$$

Estimators: example: regressogram



Model examples

- **histograms** on some partition m of \mathcal{X}

$$\Rightarrow S_m = \text{vect}\{(\mathbb{1}_{x_i \in \lambda})_{1 \leq i \leq n} / \lambda \in m\}$$

\Rightarrow the least-squares estimator (regressogram) can be written

$$\hat{F}_{S_m}(x_i) = \sum_{\lambda \in m} \hat{\beta}_\lambda \mathbb{1}_{x_i \in \lambda} \quad \hat{\beta}_\lambda = \frac{1}{\text{Card } \{x_i \in \lambda\}} \sum_{x_i \in \lambda} Y_i$$

- **variable selection:** $x_i = (x_i^{(1)}, \dots, x_i^{(p)}) \in \mathbb{R}^p$ gathers p variables that can (linearly) explain Y_i

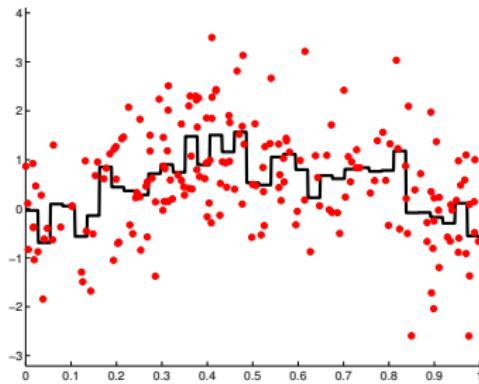
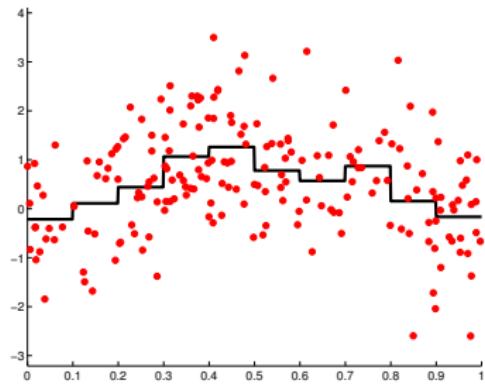
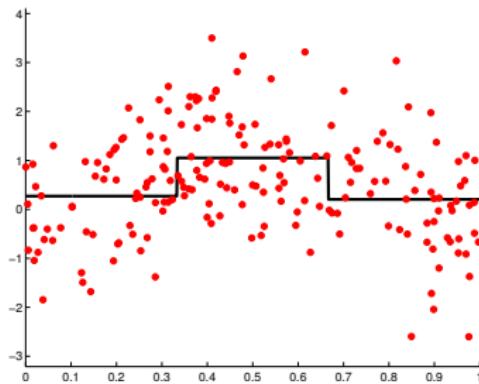
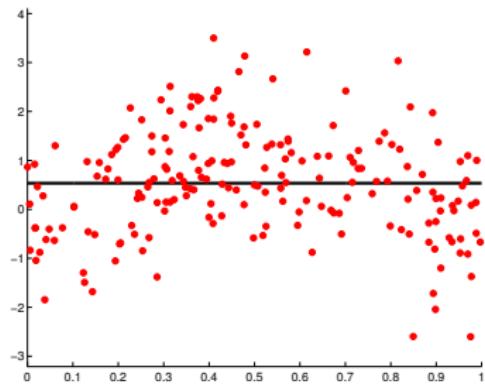
$$\forall m \subset \{1, \dots, p\} , \quad S_m = \text{vect} \left\{ x^{(j)} \text{ s.t. } j \in m \right\}$$

- S_m subspace generated by a subset of an orthogonal basis of \mathbb{R}^n (**Fourier, wavelets, ...**)

Outline

- 1 Introduction
- 2 Fixed-design regression
- 3 Model selection in fixed-design regression
- 4 General prediction framework
- 5 Estimator selection

Model selection: regular regressograms, choose D ?



Model selection

- Model collection $(S_m)_{m \in \mathcal{M}} \Rightarrow (\hat{F}_m)_{m \in \mathcal{M}} \Rightarrow \hat{m}(Y)$?

$$\hat{F}_m = \Pi_m Y = \Pi_{S_m} Y$$

- Goal: minimize the risk, i.e.,
Oracle inequality (in expectation or with a large probability):

$$\frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 \leq C \inf_{m \in \mathcal{M}} \left\{ \frac{1}{n} \left\| \hat{F}_m - F \right\|^2 \right\} + R_n$$

- Other possible goal: identify the “true” model

Approximation and estimation error

Fixed model S_m , linear subspace of \mathbb{R}^n , dimension $D_m = \dim(S_m)$.

- Approximation error of S_m :

$$\inf_{t \in S_m} \frac{1}{n} \|t - F\|^2 = \frac{1}{n} \|F_m - F\|^2 = \frac{1}{n} \|(I_n - \Pi_m)F\|^2$$

where $F_m = \Pi_m F$ orthogonal projection of F onto S_m

- $$\bullet \quad \widehat{F}_m \in S_m \Rightarrow \frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 \geq \frac{1}{n} \| F_m - F \|^2$$

Approximation and estimation error

Fixed model S_m , linear subspace of \mathbb{R}^n , dimension $D_m = \dim(S_m)$.

- Approximation error of S_m :

$$\inf_{t \in S_m} \frac{1}{n} \|t - F\|^2 = \frac{1}{n} \|F_m - F\|^2 = \frac{1}{n} \|(I_n - \Pi_m)F\|^2$$

where $F_m = \Pi_m F$ orthogonal projection of F onto S_m

- $\hat{F}_m \in S_m \Rightarrow \frac{1}{n} \left\| \hat{F}_m - F \right\|^2 \geq \frac{1}{n} \|F_m - F\|^2$
 - Estimation error of \hat{F}_m :

$$\frac{1}{n} \|\widehat{F}_m - F\|^2 - \frac{1}{n} \|F_m - F\|^2 \geq 0$$

Expectation of the estimation error

$$\begin{aligned}\|\hat{F}_m - F\|^2 &= \|\Pi_m(F + \varepsilon) - F\|^2 \\ &= \|\Pi_m F - F\|^2 + 2\underbrace{\langle \Pi_m F - F, \Pi_m \varepsilon \rangle}_{=0} + \|\Pi_m \varepsilon\|^2\end{aligned}$$

Expectation of the estimation error

$$\begin{aligned}\|\hat{F}_m - F\|^2 &= \|\Pi_m(F + \varepsilon) - F\|^2 \\ &= \|\Pi_m F - F\|^2 + 2\underbrace{\langle \Pi_m F - F, \Pi_m \varepsilon \rangle}_{=0} + \|\Pi_m \varepsilon\|^2\end{aligned}$$

⇒ Estimation error (of \hat{F}_m):

$$\frac{1}{n} \|\hat{F}_m - F\|^2 - \frac{1}{n} \|F_m - F\|^2 = \frac{1}{n} \|\Pi_m \varepsilon\|^2 = \frac{1}{n} \langle \Pi_m \varepsilon, \varepsilon \rangle$$

Expectation of the estimation error

$$\begin{aligned}\|\hat{F}_m - F\|^2 &= \|\Pi_m(F + \varepsilon) - F\|^2 \\ &= \|\Pi_m F - F\|^2 + 2\underbrace{\langle \Pi_m F - F, \Pi_m \varepsilon \rangle}_{=0} + \|\Pi_m \varepsilon\|^2\end{aligned}$$

⇒ Estimation error (of \hat{F}_m):

$$\frac{1}{n} \|\hat{F}_m - F\|^2 - \frac{1}{n} \|F_m - F\|^2 = \frac{1}{n} \|\Pi_m \varepsilon\|^2 = \frac{1}{n} \langle \Pi_m \varepsilon, \varepsilon \rangle$$

⇒ Expectation of the estimation error (of \hat{F}_m):

$$\frac{1}{n} \mathbb{E}[\langle \Pi_m \varepsilon, \varepsilon \rangle] = \frac{\sigma^2 \text{tr}(\Pi_m)}{n} = \frac{\sigma^2 D_m}{n}$$

Bias-variance trade-off

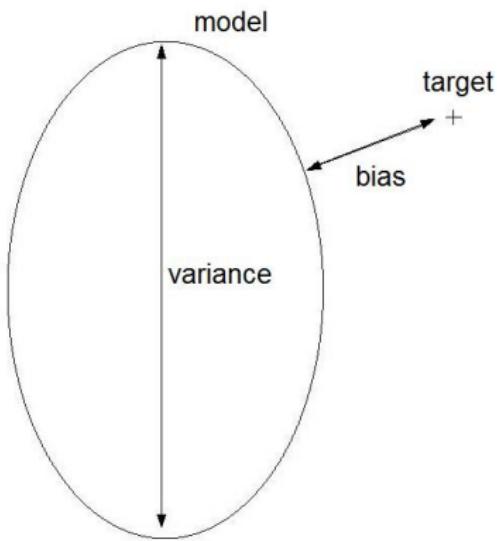
$$\mathbb{E} \left[\frac{1}{n} \|\hat{F}_m - F\|^2 \right] = \frac{1}{n} \|F_m - F\|^2 + \frac{\sigma^2 D_m}{n}$$

Approximation error / Bias:

$$\frac{1}{n} \|F_m - F\|^2 = \frac{1}{n} \|\Pi_m F - F\|^2$$

$\mathbb{E}[\text{Estimation error}]$ / Variance:

$$\frac{\sigma^2 D_m}{n}$$



Bias-variance trade-off

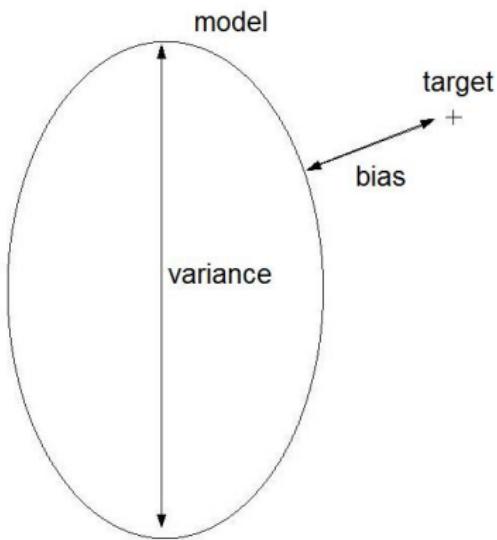
$$\mathbb{E} \left[\frac{1}{n} \|\hat{F}_m - F\|^2 \right] = \frac{1}{n} \|F_m - F\|^2 + \frac{\sigma^2 D_m}{n}$$

Approximation error / Bias:

$$\frac{1}{n} \|F_m - F\|^2 = \frac{1}{n} \|\Pi_m F - F\|^2$$

$\mathbb{E}[\text{Estimation error}]$ / Variance:

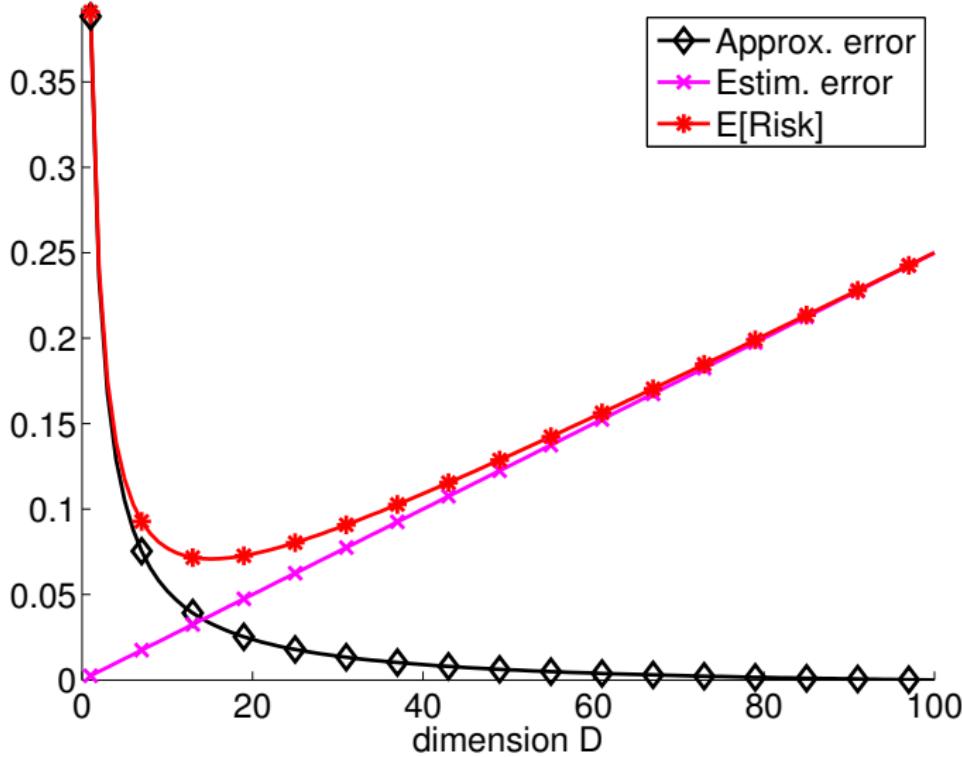
$$\frac{\sigma^2 D_m}{n}$$



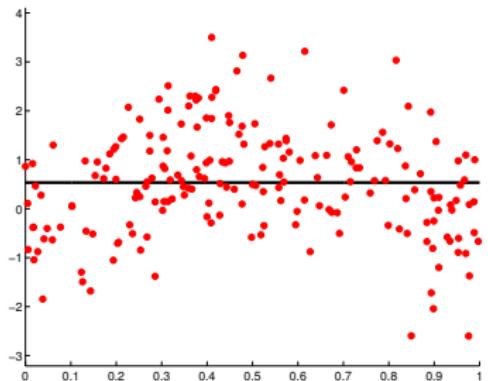
Bias-variance trade-off

↔ avoid **overfitting** and **underfitting**

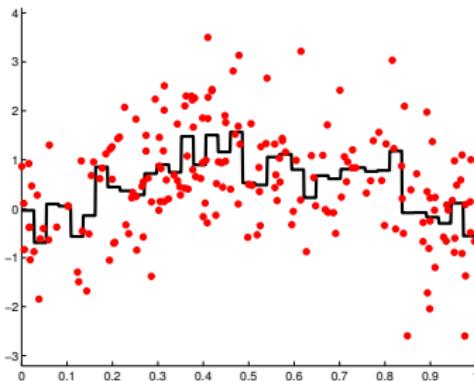
Bias-variance trade-off



Overfitting/underfitting: regressograms

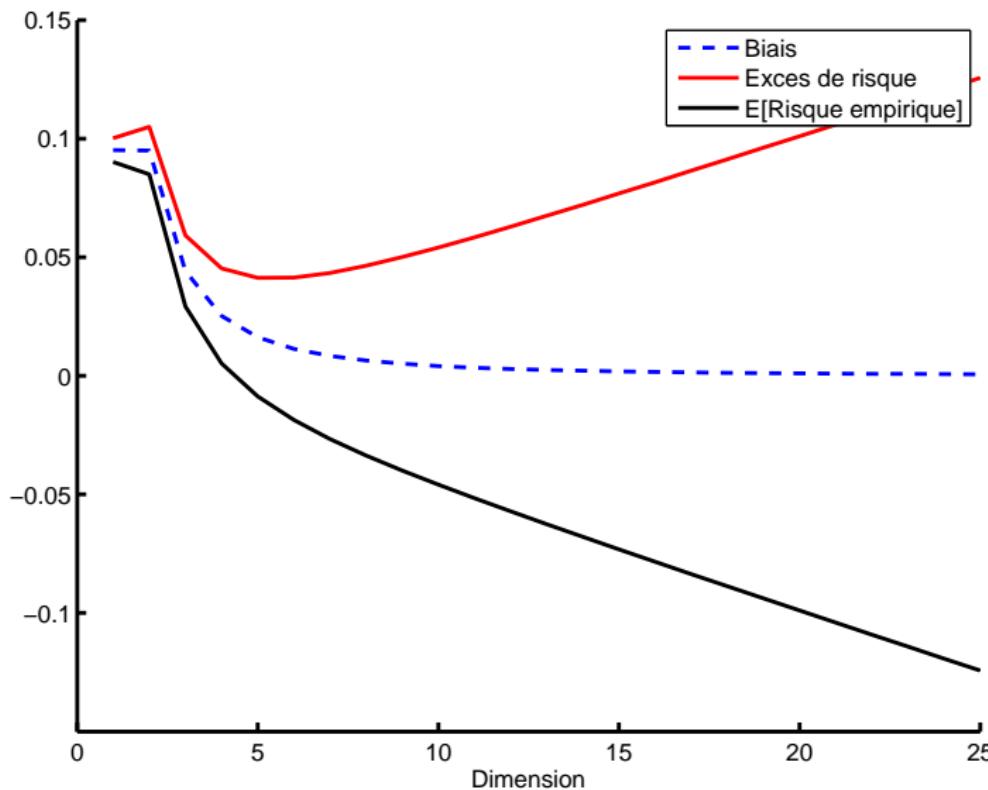


Underfitting
 $D = 1$ (too small)



Overfitting
 $D = 37$ (too large)

Why should the empirical risk be penalized?



Penalization

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \frac{1}{n} \|\hat{F}_m - Y\|^2 + \text{pen}(m) \right\}$$

Penalization

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \frac{1}{n} \|\hat{F}_m - Y\|^2 + \text{pen}(m) \right\}$$

- Ideal penalty:

$$\text{pen}_{\text{id}}(m) := \frac{1}{n} \|\hat{F}_m - F\|^2 - \frac{1}{n} \|\hat{F}_m - Y\|^2 = \text{Risk} - \text{Empirical risk}$$

- **Mallows' heuristic:** $\text{pen}(m) \approx \mathbb{E} [\text{pen}_{\text{id}}(m)]$
 \Rightarrow oracle inequality if $\text{Card}(\mathcal{M})$ not too large
(+ concentration inequalities)

Ideal penalty and its expectation

$$\begin{aligned}\|\hat{F}_m - Y\|^2 &= \|\hat{F}_m - F - \varepsilon\|^2 \\&= \|\hat{F}_m - F\|^2 + \|\varepsilon\|^2 - 2 \langle \hat{F}_m - F, \varepsilon \rangle \\&= \|\hat{F}_m - F\|^2 + \|\varepsilon\|^2 - 2 \langle \Pi_m F + \Pi_m \varepsilon - F, \varepsilon \rangle \\&= \|\hat{F}_m - F\|^2 + \|\varepsilon\|^2 - 2 \langle \Pi_m F - F, \varepsilon \rangle - 2 \langle \Pi_m \varepsilon, \varepsilon \rangle\end{aligned}$$

Ideal penalty and its expectation

$$\begin{aligned}\|\hat{F}_m - Y\|^2 &= \|\hat{F}_m - F - \varepsilon\|^2 \\&= \|\hat{F}_m - F\|^2 + \|\varepsilon\|^2 - 2 \langle \hat{F}_m - F, \varepsilon \rangle \\&= \|\hat{F}_m - F\|^2 + \|\varepsilon\|^2 - 2 \langle \Pi_m F + \Pi_m \varepsilon - F, \varepsilon \rangle \\&= \|\hat{F}_m - F\|^2 + \|\varepsilon\|^2 - 2 \langle \Pi_m F - F, \varepsilon \rangle - 2 \langle \Pi_m \varepsilon, \varepsilon \rangle\end{aligned}$$

⇒ Ideal penalty

$$\begin{aligned}\text{pen}_{\text{id}}(m) &= \frac{1}{n} \|\hat{F}_m - F\|^2 - \frac{1}{n} \|\hat{F}_m - Y\|^2 \\&= \frac{2}{n} \langle \Pi_m F - F, \varepsilon \rangle + \frac{2}{n} \langle \Pi_m \varepsilon, \varepsilon \rangle - \frac{1}{n} \|\varepsilon\|^2\end{aligned}$$

Ideal penalty and its expectation

⇒ Ideal penalty

$$\begin{aligned}\text{pen}_{\text{id}}(m) &= \frac{1}{n} \|\hat{F}_m - F\|^2 - \frac{1}{n} \|\hat{F}_m - Y\|^2 \\ &= \frac{2}{n} \langle \Pi_m F - F, \varepsilon \rangle + \frac{2}{n} \langle \Pi_m \varepsilon, \varepsilon \rangle - \frac{1}{n} \|\varepsilon\|^2\end{aligned}$$

⇒ Expectation of the ideal penalty

$$\begin{aligned}\mathbb{E}[\text{pen}_{\text{id}}(m)] &= \frac{2}{n} \mathbb{E}[\langle \Pi_m F - F, \varepsilon \rangle] + \frac{2}{n} \mathbb{E}[\langle \Pi_m \varepsilon, \varepsilon \rangle] - \frac{1}{n} \mathbb{E}[\|\varepsilon\|^2] \\ &= \underbrace{\frac{2\sigma^2 D_m}{n}}_{\text{Mallows' } C_p} - \sigma^2\end{aligned}$$

Towards theoretical guarantees: a key lemma

Lemma

Let $\text{crit} : \mathcal{M} \rightarrow \mathbb{R}$ be any function (possibly data-dependent).
On the event Ω on which, $\forall m, m' \in \mathcal{M}$

$$\begin{aligned} & \left[\text{crit}(m) - \frac{1}{n} \|\hat{F}_m - F\|^2 \right] - \left[\text{crit}(m') - \frac{1}{n} \|\hat{F}_{m'} - F\|^2 \right] \\ & \leq A(m) + B(m'), \end{aligned}$$

we have, $\forall \hat{m} \in \arg \min_{m \in \mathcal{M}} \{\text{crit}(m)\}$,

$$\frac{1}{n} \|\hat{F}_{\hat{m}} - F\|^2 - B(\hat{m}) \leq \inf_{m \in \mathcal{M}} \left\{ \frac{1}{n} \|\hat{F}_m - F\|^2 + A(m) \right\}$$

Proof of the key lemma

For any $m \in \mathcal{M}$, we have

$$\begin{aligned} & \frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 \\ &= \text{crit}(\hat{m}) + \frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 - \text{crit}(\hat{m}) \end{aligned}$$

Proof of the key lemma

For any $m \in \mathcal{M}$, we have

$$\begin{aligned} & \frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 \\ &= \text{crit}(\hat{m}) + \frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 - \text{crit}(\hat{m}) \\ &\leq \text{crit}(m) + \frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 - \text{crit}(\hat{m}) \end{aligned}$$

Proof of the key lemma

For any $m \in \mathcal{M}$, we have

$$\begin{aligned} & \frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 \\ &= \text{crit}(\hat{m}) + \frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 - \text{crit}(\hat{m}) \\ &\leq \text{crit}(m) + \frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 - \text{crit}(\hat{m}) \\ &= \frac{1}{n} \left\| \hat{F}_m - F \right\|^2 + \text{crit}(m) - \frac{1}{n} \left\| \hat{F}_m - F \right\|^2 + \frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 - \text{crit}(\hat{m}) \end{aligned}$$

Proof of the key lemma

For any $m \in \mathcal{M}$, we have

$$\begin{aligned} & \frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 \\ &= \text{crit}(\hat{m}) + \frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 - \text{crit}(\hat{m}) \\ &\leq \text{crit}(m) + \frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 - \text{crit}(\hat{m}) \\ &= \frac{1}{n} \left\| \hat{F}_m - F \right\|^2 + \text{crit}(m) - \frac{1}{n} \left\| \hat{F}_m - F \right\|^2 + \frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 - \text{crit}(\hat{m}) \\ &\leq \frac{1}{n} \left\| \hat{F}_m - F \right\|^2 + A(m) + B(\hat{m}) \end{aligned}$$

Proof of the key lemma

For any $m \in \mathcal{M}$, we have

$$\begin{aligned} & \frac{1}{n} \left\| \widehat{F}_{\widehat{m}} - F \right\|^2 \\ &= \text{crit}(\widehat{m}) + \frac{1}{n} \left\| \widehat{F}_{\widehat{m}} - F \right\|^2 - \text{crit}(\widehat{m}) \\ &\leq \text{crit}(m) + \frac{1}{n} \left\| \widehat{F}_{\widehat{m}} - F \right\|^2 - \text{crit}(\widehat{m}) \\ &= \frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 + \text{crit}(m) - \frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 + \frac{1}{n} \left\| \widehat{F}_{\widehat{m}} - F \right\|^2 - \text{crit}(\widehat{m}) \\ &\leq \frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 + A(m) + B(\widehat{m}) \end{aligned}$$

hence

$$\frac{1}{n} \left\| \widehat{F}_{\widehat{m}} - F \right\|^2 - B(\widehat{m}) \leq \inf_{m \in \mathcal{M}} \left\{ \frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 + A(m) \right\}. \quad \square$$

Key lemma (reformulated)

Lemma

Let $\text{pen} : \mathcal{M} \rightarrow \mathbb{R}$ be any penalty (possibly data-dependent).
On the event Ω on which, $\forall m, m' \in \mathcal{M}$

$$\begin{aligned} & [\text{pen}(m) - \text{pen}_{\text{id}}(m)] - [\text{pen}(m') - \text{pen}_{\text{id}}(m')] \\ & \leq A(m) + B(m'), \end{aligned}$$

we have $\forall \hat{m} \in \arg \min_{m \in \mathcal{M}} \left\{ \frac{1}{n} \|\hat{F}_m - Y\|^2 + \text{pen}(m) \right\}$

$$\frac{1}{n} \|\hat{F}_{\hat{m}} - F\|^2 - B(\hat{m}) \leq \inf_{m \in \mathcal{M}} \left\{ \frac{1}{n} \|\hat{F}_m - F\|^2 + A(m) \right\}$$

Proof: take $\text{crit}(m) = \frac{1}{n} \|\hat{F}_m - Y\|^2 + \text{pen}(m)$.

□

How to use the key lemma here?

- ➊ Take $\text{pen}(m) = \mathbb{E}[\text{pen}_{\text{id}}(m)]$ (up to a translation).
- ➋ Prove that $\forall m \in \mathcal{M}$, $\text{pen}_{\text{id}}(m)$ concentrates around its expectation.
- ➌ Union bound over $m \in \mathcal{M}$ (if $\text{Card}(\mathcal{M})$ “small”).
- ➍ Apply the key lemma with $A(m) \propto$ deviations of $\text{pen}_{\text{id}}(m)$.

Ideal penalty (reminder)

$$\begin{aligned}\text{pen}_{\text{id}}(m) &= \frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 - \frac{1}{n} \left\| \widehat{F}_m - Y \right\|^2 \\ &= \frac{2}{n} \langle \Pi_m F - F, \varepsilon \rangle + \frac{2}{n} \langle \Pi_m \varepsilon, \varepsilon \rangle - \frac{1}{n} \|\varepsilon\|^2\end{aligned}$$

Ideal penalty (reminder)

$$\begin{aligned}\text{pen}_{\text{id}}(m) &= \frac{1}{n} \|\hat{F}_m - F\|^2 - \frac{1}{n} \|\hat{F}_m - Y\|^2 \\ &= \frac{2}{n} \langle \Pi_m F - F, \varepsilon \rangle + \frac{2}{n} \langle \Pi_m \varepsilon, \varepsilon \rangle - \frac{1}{n} \|\varepsilon\|^2\end{aligned}$$

- Linear term:

$$\frac{2}{n} \langle \Pi_m F - F, \varepsilon \rangle \quad \Rightarrow \quad \text{mean} = 0, \quad \text{Gaussian distribution}$$

Ideal penalty (reminder)

$$\begin{aligned}\text{pen}_{\text{id}}(m) &= \frac{1}{n} \|\hat{F}_m - F\|^2 - \frac{1}{n} \|\hat{F}_m - Y\|^2 \\ &= \frac{2}{n} \langle \Pi_m F - F, \varepsilon \rangle + \frac{2}{n} \langle \Pi_m \varepsilon, \varepsilon \rangle - \frac{1}{n} \|\varepsilon\|^2\end{aligned}$$

- Linear term:

$$\frac{2}{n} \langle \Pi_m F - F, \varepsilon \rangle \quad \Rightarrow \quad \text{mean} = 0, \quad \text{Gaussian distribution}$$

- Quadratic term:

$$\frac{2}{n} \langle \Pi_m \varepsilon, \varepsilon \rangle \quad \Rightarrow \quad \text{mean} = \frac{2\sigma^2 D_m}{n}, \quad \chi^2 \text{ distribution}$$

Ideal penalty (reminder)

$$\begin{aligned}\text{pen}_{\text{id}}(m) &= \frac{1}{n} \|\hat{F}_m - F\|^2 - \frac{1}{n} \|\hat{F}_m - Y\|^2 \\ &= \frac{2}{n} \langle \Pi_m F - F, \varepsilon \rangle + \frac{2}{n} \langle \Pi_m \varepsilon, \varepsilon \rangle - \frac{1}{n} \|\varepsilon\|^2\end{aligned}$$

- Linear term:

$$\frac{2}{n} \langle \Pi_m F - F, \varepsilon \rangle \quad \Rightarrow \quad \text{mean} = 0, \quad \text{Gaussian distribution}$$

- Quadratic term:

$$\frac{2}{n} \langle \Pi_m \varepsilon, \varepsilon \rangle \quad \Rightarrow \quad \text{mean} = \frac{2\sigma^2 D_m}{n}, \quad \chi^2 \text{ distribution}$$

- Constant term:

$$\frac{1}{n} \|\varepsilon\|^2 \quad \Rightarrow \quad \text{can be discarded}$$

Concentration of the ideal penalty (1): linear term

Proposition (Gaussian concentration)

If $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ and $\alpha \in \mathbb{R}^n$, for every $x \geq 0$,

$$\mathbb{P}\left(|\langle \varepsilon, \alpha \rangle| \leq \sigma \sqrt{2x} \|\alpha\|\right) \geq 1 - 2e^{-x}.$$

Concentration of the ideal penalty (1): linear term

Proposition (Gaussian concentration)

If $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ and $\alpha \in \mathbb{R}^n$, for every $x \geq 0$,

$$\mathbb{P}\left(|\langle \varepsilon, \alpha \rangle| \leq \sigma\sqrt{2x} \|\alpha\|\right) \geq 1 - 2e^{-x}.$$

⇒ with probability $\geq 1 - 2e^{-x}$, for every $\theta > 0$,

$$\frac{2}{n} |\langle \Pi_m F - F, \varepsilon \rangle| \leq \frac{2\sigma\sqrt{2x}}{n} \|\Pi_m F - F\| \leq \frac{\theta}{n} \|\Pi_m F - F\|^2 + \frac{2x\sigma^2}{\theta n}$$

since $2ab \leq \theta a^2 + \theta^{-1} b^2 \quad \forall a, b \geq 0, \theta > 0$.

Concentration of the ideal penalty (1): linear term

Proposition (Gaussian concentration)

If $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ and $\alpha \in \mathbb{R}^n$, for every $x \geq 0$,

$$\mathbb{P}\left(|\langle \varepsilon, \alpha \rangle| \leq \sigma\sqrt{2x} \|\alpha\|\right) \geq 1 - 2e^{-x}.$$

⇒ with probability $\geq 1 - 2e^{-x}$, for every $\theta > 0$,

$$\frac{2}{n} |\langle \Pi_m F - F, \varepsilon \rangle| \leq \frac{2\sigma\sqrt{2x}}{n} \|\Pi_m F - F\| \leq \frac{\theta}{n} \|\Pi_m F - F\|^2 + \frac{2x\sigma^2}{\theta n}$$

since $2ab \leq \theta a^2 + \theta^{-1} b^2 \quad \forall a, b \geq 0, \theta > 0$.

- Can be generalized to **sub-Gaussian noise**, see arXiv:1901.07277 (Remark 1).

Concentration of the ideal penalty (2): quadratic term

Proposition (see A. & Bach 2011 (arXiv:0909.1884), Proposition 6)

If $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ and $M \in \mathcal{M}_n(\mathbb{R})$, for every $x \geq 0$,

$$\mathbb{P} \left(|\langle \varepsilon, M\varepsilon \rangle - \sigma^2 \operatorname{tr}(M)| \leqslant 2\sigma^2 \sqrt{x \operatorname{tr}(M^\top M)} + 2\sigma^2 \|M\| x \right) \geqslant 1 - 2e^{-x}.$$

Concentration of the ideal penalty (2): quadratic term

Proposition (see A. & Bach 2011 (arXiv:0909.1884), Proposition 6)

If $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ and $M \in \mathcal{M}_n(\mathbb{R})$, for every $x \geq 0$,

$$\mathbb{P} \left(|\langle \varepsilon, M\varepsilon \rangle - \sigma^2 \text{tr}(M)| \leqslant 2\sigma^2 \sqrt{x \text{tr}(M^\top M)} + 2\sigma^2 \|M\| x \right) \geqslant 1 - 2e^{-x}.$$

⇒ with probability $\geqslant 1 - 2e^{-x}$, for every $\theta > 0$,

$$\frac{2}{n} |\langle \Pi_m \varepsilon, \varepsilon \rangle - \sigma^2 D_m| \leqslant \frac{4\sigma^2}{n} \sqrt{x D_m} + \frac{4\sigma^2 x}{n}$$

Concentration of the ideal penalty (2): quadratic term

Proposition (see A. & Bach 2011 (arXiv:0909.1884), Proposition 6)

If $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ and $M \in \mathcal{M}_n(\mathbb{R})$, for every $x \geq 0$,

$$\mathbb{P} \left(|\langle \varepsilon, M\varepsilon \rangle - \sigma^2 \text{tr}(M)| \leqslant 2\sigma^2 \sqrt{x \text{tr}(M^\top M)} + 2\sigma^2 \|M\|_F x \right) \geqslant 1 - 2e^{-x}.$$

⇒ with probability $\geqslant 1 - 2e^{-x}$, for every $\theta > 0$,

$$\begin{aligned} \frac{2}{n} |\langle \Pi_m \varepsilon, \varepsilon \rangle - \sigma^2 D_m| &\leqslant \frac{4\sigma^2}{n} \sqrt{x D_m} + \frac{4\sigma^2 x}{n} \\ &\leqslant \frac{\theta \sigma^2 D_m}{n} + \left(4 + \frac{4}{\theta}\right) \frac{\sigma^2 x}{n}. \end{aligned}$$

Concentration of the ideal penalty (2): quadratic term

Proposition (see A. & Bach 2011 (arXiv:0909.1884), Proposition 6)

If $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ and $M \in \mathcal{M}_n(\mathbb{R})$, for every $x \geq 0$,

$$\mathbb{P} \left(|\langle \varepsilon, M\varepsilon \rangle - \sigma^2 \text{tr}(M)| \leqslant 2\sigma^2 \sqrt{x \text{tr}(M^\top M)} + 2\sigma^2 \|M\|_F x \right) \geqslant 1 - 2e^{-x}.$$

⇒ with probability $\geqslant 1 - 2e^{-x}$, for every $\theta > 0$,

$$\begin{aligned} \frac{2}{n} |\langle \Pi_m \varepsilon, \varepsilon \rangle - \sigma^2 D_m| &\leqslant \frac{4\sigma^2}{n} \sqrt{x D_m} + \frac{4\sigma^2 x}{n} \\ &\leqslant \frac{\theta \sigma^2 D_m}{n} + \left(4 + \frac{4}{\theta}\right) \frac{\sigma^2 x}{n}. \end{aligned}$$

- Can be generalized to **sub-Gaussian noise**, see arXiv:1901.07277 (Remark 1).

Concentration of the ideal penalty: summary

- With probability $\geq 1 - 4e^{-x}$, for every $\theta > 0$,

$$\begin{aligned} & \left| \text{pen}_{\text{id}}(m) + \frac{1}{n} \|\varepsilon\|^2 - \mathbb{E} \left[\text{pen}_{\text{id}}(m) + \frac{1}{n} \|\varepsilon\|^2 \right] \right| \\ & \leq \frac{\theta}{n} \|\Pi_m F - F\|^2 + \frac{2x\sigma^2}{\theta n} + \frac{\theta\sigma^2 D_m}{n} + \left(4 + \frac{4}{\theta}\right) \frac{\sigma^2 x}{n} \end{aligned}$$

Concentration of the ideal penalty: summary

- With probability $\geq 1 - 4e^{-x}$, for every $\theta > 0$,

$$\begin{aligned} & \left| \text{pen}_{\text{id}}(m) + \frac{1}{n} \|\varepsilon\|^2 - \mathbb{E} \left[\text{pen}_{\text{id}}(m) + \frac{1}{n} \|\varepsilon\|^2 \right] \right| \\ & \leq \frac{\theta}{n} \|\Pi_m F - F\|^2 + \frac{2x\sigma^2}{\theta n} + \frac{\theta\sigma^2 D_m}{n} + \left(4 + \frac{4}{\theta}\right) \frac{\sigma^2 x}{n} \\ & = \theta \mathbb{E} \left[\frac{1}{n} \|\widehat{F}_m - F\|^2 \right] + \left(6 + \frac{4}{\theta}\right) \frac{x\sigma^2}{n} \end{aligned}$$

Concentration of the ideal penalty: summary

- With probability $\geq 1 - 4e^{-x}$, for every $\theta \in (0, 1)$,

$$\begin{aligned} & \left| \text{pen}_{\text{id}}(m) + \frac{1}{n} \|\varepsilon\|^2 - \mathbb{E} \left[\text{pen}_{\text{id}}(m) + \frac{1}{n} \|\varepsilon\|^2 \right] \right| \\ & \leq \frac{\theta}{n} \|\Pi_m F - F\|^2 + \frac{2x\sigma^2}{\theta n} + \frac{\theta\sigma^2 D_m}{n} + \left(4 + \frac{4}{\theta}\right) \frac{\sigma^2 x}{n} \\ & = \theta \mathbb{E} \left[\frac{1}{n} \|\widehat{F}_m - F\|^2 \right] + \left(6 + \frac{4}{\theta}\right) \frac{x\sigma^2}{n} \\ & \leq \frac{\theta}{1-\theta} \frac{1}{n} \|\widehat{F}_m - F\|^2 + L(\theta) \frac{x\sigma^2}{n} =: A(m) =: B(m). \end{aligned}$$

Concentration of the ideal penalty: summary

- With probability $\geq 1 - 4e^{-x}$, for every $\theta \in (0, 1)$,

$$\begin{aligned} & \left| \text{pen}_{\text{id}}(m) + \frac{1}{n} \|\varepsilon\|^2 - \mathbb{E} \left[\text{pen}_{\text{id}}(m) + \frac{1}{n} \|\varepsilon\|^2 \right] \right| \\ & \leq \frac{\theta}{n} \|\Pi_m F - F\|^2 + \frac{2x\sigma^2}{\theta n} + \frac{\theta\sigma^2 D_m}{n} + \left(4 + \frac{4}{\theta}\right) \frac{\sigma^2 x}{n} \\ & = \theta \mathbb{E} \left[\frac{1}{n} \|\widehat{F}_m - F\|^2 \right] + \left(6 + \frac{4}{\theta}\right) \frac{x\sigma^2}{n} \\ & \leq \frac{\theta}{1-\theta} \frac{1}{n} \|\widehat{F}_m - F\|^2 + L(\theta) \frac{x\sigma^2}{n} =: A(m) =: B(m). \end{aligned}$$

⇒ end of step 2 (concentration of $\text{pen}_{\text{id}}(m)$)

Concentration of the ideal penalty: summary

- With probability $\geq 1 - 4e^{-x}$, for every $\theta \in (0, 1)$,

$$\begin{aligned}
& \left| \text{pen}_{\text{id}}(m) + \frac{1}{n} \|\varepsilon\|^2 - \mathbb{E} \left[\text{pen}_{\text{id}}(m) + \frac{1}{n} \|\varepsilon\|^2 \right] \right| \\
& \leq \frac{\theta}{n} \|\Pi_m F - F\|^2 + \frac{2x\sigma^2}{\theta n} + \frac{\theta\sigma^2 D_m}{n} + \left(4 + \frac{4}{\theta}\right) \frac{\sigma^2 x}{n} \\
& = \theta \mathbb{E} \left[\frac{1}{n} \|\widehat{F}_m - F\|^2 \right] + \left(6 + \frac{4}{\theta}\right) \frac{x\sigma^2}{n} \\
& \leq \frac{\theta}{1-\theta} \frac{1}{n} \|\widehat{F}_m - F\|^2 + L(\theta) \frac{x\sigma^2}{n} =: A(m) =: B(m).
\end{aligned}$$

- Step 3: **union bound** \Rightarrow with probability

$\geq 1 - 4 \text{Card}(\mathcal{M}) e^{-x}$, $\forall \theta \in (0, 1)$, $\forall m, m' \in \mathcal{M}$,

$$\begin{aligned}
& \left[\frac{2\sigma^2 D_m}{n} - \text{pen}_{\text{id}}(m) \right] - \left[\frac{2\sigma^2 D_{m'}}{n} - \text{pen}_{\text{id}}(m') \right] \\
& \leq A(m) + B(m').
\end{aligned}$$

Application of the key lemma (step 4)

- With probability $\geq 1 - 4 \text{Card}(\mathcal{M})e^{-x}$, $\forall \theta \in (0, 1)$,

$$\forall \hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \frac{1}{n} \left\| \hat{F}_m - Y \right\|^2 + \frac{2\sigma^2 D_m}{n} \right\},$$

$$\begin{aligned} & \frac{1 - 2\theta}{1 - \theta} \frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 - L(\theta) \frac{x\sigma^2}{n} \\ & \leq \frac{1}{1 - \theta} \inf_{m \in \mathcal{M}} \frac{1}{n} \left\| \hat{F}_m - F \right\|^2 + L(\theta) \frac{x\sigma^2}{n}, \end{aligned}$$

Application of the key lemma (step 4)

- With probability $\geq 1 - 4 \text{Card}(\mathcal{M})e^{-x}$, $\forall \theta \in (0, 1)$,

$$\forall \hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \frac{1}{n} \left\| \hat{F}_m - Y \right\|^2 + \frac{2\sigma^2 D_m}{n} \right\},$$

$$\begin{aligned} & (1 - 2\theta) \frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 \\ & \leq \inf_{m \in \mathcal{M}} \frac{1}{n} \left\| \hat{F}_m - F \right\|^2 + 2(1 - \theta)L(\theta) \frac{x\sigma^2}{n}, \end{aligned}$$

Application of the key lemma (step 4)

- With probability $\geq 1 - 4 \text{Card}(\mathcal{M})e^{-x}$, $\forall \theta \in (0, 1/2)$,

$$\forall \hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \frac{1}{n} \left\| \hat{F}_m - Y \right\|^2 + \frac{2\sigma^2 D_m}{n} \right\},$$

$$\begin{aligned} & \frac{1}{n} \left\| \hat{F}_{\hat{m}} - F \right\|^2 \\ & \leq \frac{1}{1-2\theta} \inf_{m \in \mathcal{M}} \frac{1}{n} \left\| \hat{F}_m - F \right\|^2 + \frac{2(1-\theta)L(\theta)x\sigma^2}{1-2\theta} \frac{x\sigma^2}{n}. \end{aligned}$$

Conclusion: oracle inequality for C_p

Theorem (Birgé & Massart 2007, reformulated = Theorem 1 in arXiv:1901.07277)

Assumptions: $\text{pen}(m) = \frac{2\sigma^2 D_m}{n}$, i.i.d. Gaussian noise.

Then, for every $x \geq 0$, with probability at least $1 - 4 \text{Card}(\mathcal{M}) e^{-x}$, for every $\theta \in (0, 1/6)$,

$$\frac{1}{n} \left\| \widehat{F}_{\widehat{m}} - F \right\|^2 \leq (1 + 3\theta) \underbrace{\inf_{m \in \mathcal{M}} \left\{ \frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 \right\}}_{\text{oracle risk}} + \frac{L'(\theta)\sigma^2 x}{n}.$$

Conclusion: oracle inequality for C_p

Theorem (Birgé & Massart 2007, reformulated = Theorem 1 in arXiv:1901.07277)

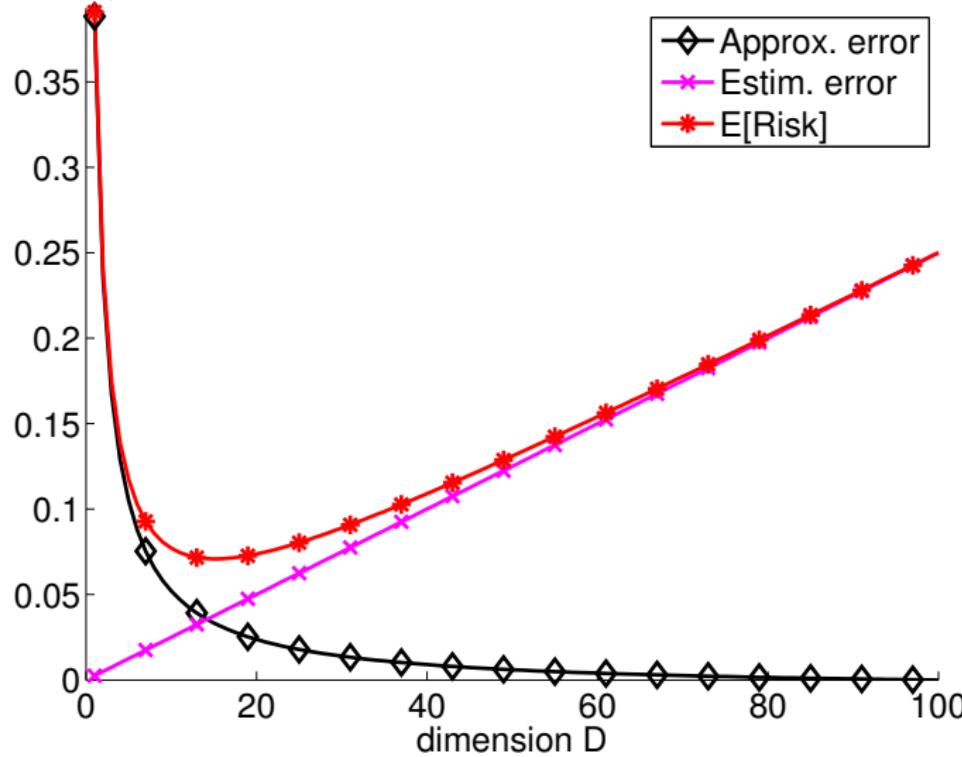
Assumptions: $\text{pen}(m) = \frac{2\sigma^2 D_m}{n}$, i.i.d. Gaussian noise.

Then, for every $x \geq 0$, with probability at least $1 - 4 \text{Card}(\mathcal{M}) e^{-x}$, for every $\theta \in (0, 1/6)$,

$$\frac{1}{n} \left\| \widehat{F}_{\widehat{m}} - F \right\|^2 \leq (1 + 3\theta) \underbrace{\inf_{m \in \mathcal{M}} \left\{ \frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 \right\}}_{\text{oracle risk}} + \frac{L'(\theta)\sigma^2 x}{n}.$$

Generalization (arXiv:1901.07277): sub-Gaussian noise.

Meaning of an oracle inequality



Outline

- 1 Introduction
- 2 Fixed-design regression
- 3 Model selection in fixed-design regression
- 4 General prediction framework
- 5 Estimator selection

General prediction setting

- Data: $D_n = (X_i, Y_i)_{1 \leq i \leq n} \in (\mathcal{X} \times \mathcal{Y})^n$ assumed i.i.d. $\sim P$

General prediction setting

- **Data:** $D_n = (X_i, Y_i)_{1 \leq i \leq n} \in (\mathcal{X} \times \mathcal{Y})^n$ assumed i.i.d. $\sim P$
- **Predictor:** $t : \mathcal{X} \rightarrow \mathcal{Y}$
new data $X_{n+1} \rightsquigarrow t(X_{n+1})$ “predicts” Y_{n+1}

General prediction setting

- **Data:** $D_n = (X_i, Y_i)_{1 \leq i \leq n} \in (\mathcal{X} \times \mathcal{Y})^n$ assumed i.i.d. $\sim P$
- **Predictor:** $t : \mathcal{X} \rightarrow \mathcal{Y}$
new data $X_{n+1} \rightsquigarrow t(X_{n+1})$ “predicts” Y_{n+1}
- **Risk (prediction error):** $\mathcal{R}(t) = \mathbb{E}[c(t(X), Y)]$ where
 $(X, Y) \sim P$
minimal for $t = f^*$ (Bayes predictor)
 \Rightarrow **Excess risk** $\ell(t, f^*) := \mathcal{R}(t) - \mathcal{R}(f^*) \geq 0.$

General prediction setting

- **Data:** $D_n = (X_i, Y_i)_{1 \leq i \leq n} \in (\mathcal{X} \times \mathcal{Y})^n$ assumed i.i.d. $\sim P$
- **Predictor:** $t : \mathcal{X} \rightarrow \mathcal{Y}$
new data $X_{n+1} \rightsquigarrow t(X_{n+1})$ “predicts” Y_{n+1}
- **Risk (prediction error):** $\mathcal{R}(t) = \mathbb{E}[c(t(X), Y)]$ where
 $(X, Y) \sim P$
minimal for $t = f^*$ (Bayes predictor)
 \Rightarrow **Excess risk** $\ell(t, f^*) := \mathcal{R}(t) - \mathcal{R}(f^*) \geq 0.$
- **Goal:** from D_n only, find t with $\mathcal{R}(t)$ minimal.

General prediction setting

- **Data:** $D_n = (X_i, Y_i)_{1 \leq i \leq n} \in (\mathcal{X} \times \mathcal{Y})^n$ assumed i.i.d. $\sim P$
- **Predictor:** $t : \mathcal{X} \rightarrow \mathcal{Y}$
new data $X_{n+1} \rightsquigarrow t(X_{n+1})$ “predicts” Y_{n+1}
- **Risk (prediction error):** $\mathcal{R}(t) = \mathbb{E}[c(t(X), Y)]$ where
 $(X, Y) \sim P$
minimal for $t = f^*$ (Bayes predictor)
 \Rightarrow **Excess risk** $\ell(t, f^*) := \mathcal{R}(t) - \mathcal{R}(f^*) \geq 0.$
- **Goal:** from D_n only, find t with $\mathcal{R}(t)$ minimal.
- More general setting possible, including density estimation with LS or KL risk.

Prediction setting: examples

- **Regression:** $\mathcal{Y} = \mathbb{R}$

- least squares: $c(y, y') = (y - y')^2$
 $\Rightarrow f^*(X) = \mathbb{E}[Y|X]$ and $\ell(t, f^*) = \mathbb{E}[(t(X) - f^*(X))^2]$
- L^p loss: $c(y, y') = |y - y'|^p$, $p \geq 1$
- Huber loss (robustness):

$$c(y, y') = \begin{cases} \frac{1}{2}(y - y')^2 & \text{if } |y - y'| \leq \delta \\ \delta(|y - y'| - \frac{\delta}{2}) & \text{otherwise} \end{cases}$$

Prediction setting: examples

- **Regression:** $\mathcal{Y} = \mathbb{R}$

- least squares: $c(y, y') = (y - y')^2$
 $\Rightarrow f^*(X) = \mathbb{E}[Y|X]$ and $\ell(t, f^*) = \mathbb{E}[(t(X) - f^*(X))^2]$
- L^p loss: $c(y, y') = |y - y'|^p$, $p \geq 1$
- Huber loss (robustness):

$$c(y, y') = \begin{cases} \frac{1}{2}(y - y')^2 & \text{if } |y - y'| \leq \delta \\ \delta(|y - y'| - \frac{\delta}{2}) & \text{otherwise} \end{cases}$$

- **Binary classification:** $\mathcal{Y} = \{0, 1\}$

- 0–1 loss: $c(y, y') = \mathbb{1}_{y \neq y'}$
 $\Rightarrow f^*(X) = \mathbb{1}_{\mathbb{E}[Y|X] \geq 1/2}$
- convex losses (hinge, logistic, exponential, ...)

- **Multi-class classification:** $\mathcal{Y} = \{0, \dots, M-1\}$

Link with fixed-design regression

- Random-design regression, least-squares loss:
 X_1, \dots, X_n i.i.d.
- Fixed-design regression:
 X_1, \dots, X_n deterministic

Link with fixed-design regression

- Random-design regression, least-squares loss:
 X_1, \dots, X_n i.i.d.
target: $f^*(X) = \mathbb{E}[Y|X]$
- Fixed-design regression:
 X_1, \dots, X_n deterministic
target $(f^*(X_i))_{1 \leq i \leq n} = (\mathbb{E}[Y_i|X_i])_{1 \leq i \leq n}$

Link with fixed-design regression

- Random-design regression, least-squares loss:

X_1, \dots, X_n i.i.d.

target: $f^*(X) = \mathbb{E}[Y|X]$

excess risk: $\ell(t, f^*) = \mathcal{R}(t) - \mathcal{R}(f^*) = \mathbb{E}[(t(X) - f^*(X))^2]$

- Fixed-design regression:

X_1, \dots, X_n deterministic

target $(f^*(X_i))_{1 \leqslant i \leqslant n} = (\mathbb{E}[Y_i|X_i])_{1 \leqslant i \leqslant n}$

excess risk: $\frac{1}{n} \sum_{i=1}^n (t(X_i) - f^*(X_i))^2$

↔ “ $X \sim \mathcal{U}(\{X_1, \dots, X_n\})$ ”

Outline

- 1 Introduction
- 2 Fixed-design regression
- 3 Model selection in fixed-design regression
- 4 General prediction framework
- 5 Estimator selection

Estimator selection

- Estimator collection $(\hat{f}_m)_{m \in \mathcal{M}} \Rightarrow \hat{m}(D_n)$?

Estimator selection

- Estimator collection $(\hat{f}_m)_{m \in \mathcal{M}} \Rightarrow \hat{m}(D_n)$?
- Examples:
 - model selection
 - parameter tuning (choosing k or the distance for k -NN, choice of a regularization parameter, choice of a kernel, etc.)
 - choice between different methods
 - ex.: k -NN vs. kernel ridge?

Estimator selection

- Estimator collection $(\hat{f}_m)_{m \in \mathcal{M}} \Rightarrow \hat{m}(D_n)$?
- Examples:
 - model selection
 - parameter tuning (choosing k or the distance for k -NN, choice of a regularization parameter, choice of a kernel, etc.)
 - choice between different methods
 - ex.: k -NN vs. kernel ridge?
- Goal: minimize the risk $\mathcal{R}(\hat{f}_{\hat{m}(D_n)}(D_n))$
- Other possible goal: identify the “best” estimator (or the “true” model)

Approximation / estimation error decomposition?

- No general decomposition of the risk between approximation and estimation error

Approximation / estimation error decomposition?

- No general decomposition of the risk between approximation and estimation error
- Sometimes possible:
 - empirical risk minimizer: $\hat{f}_m \in \operatorname{argmin}_{t \in S_m} \hat{\mathcal{R}}_n(t)$
 - “linear” estimators (fixed-design regression)
 - local averaging estimators

Approximation / estimation error decomposition?

- No general decomposition of the risk between approximation and estimation error
- Sometimes possible:
 - empirical risk minimizer: $\hat{f}_m \in \operatorname{argmin}_{t \in S_m} \hat{\mathcal{R}}_n(t)$
 - “linear” estimators (fixed-design regression)
 - local averaging estimators
- Always have to avoid overfitting and underfitting

Estimator selection: methods

- Classical approach:

$$\hat{m}(D_n) \in \operatorname{argmin}_{m \in \mathcal{M}} \{\text{crit}(m)\}$$

Estimator selection: methods

- Classical approach:

$$\hat{m}(D_n) \in \operatorname{argmin}_{m \in \mathcal{M}} \{\text{crit}(m)\}$$

- Examples:

- **penalization** (Mallows' C_p , AIC, BIC, structural risk minimization, ...)
- **cross-validation**
- FPE, GCV, ...

Estimator selection: methods

- Classical approach:

$$\hat{m}(D_n) \in \operatorname{argmin}_{m \in \mathcal{M}} \{\text{crit}(m)\}$$

- Examples:

- **penalization** (Mallows' C_p , AIC, BIC, structural risk minimization, ...)
- **cross-validation**
- FPE, GCV, ...

- How to choose crit?

Estimator selection: methods

- Classical approach:

$$\hat{m}(D_n) \in \operatorname{argmin}_{m \in \mathcal{M}} \{\text{crit}(m)\}$$

- Examples:

- **penalization** (Mallows' C_p , AIC, BIC, structural risk minimization, ...)
- **cross-validation**
- FPE, GCV, ...

- How to choose crit?

Idea: use the **key lemma**

The key lemma (revisited)

Lemma

Let $\text{crit} : \mathcal{M} \rightarrow \mathbb{R}$ be any function (possibly data-dependent).
On the event Ω on which, $\forall m, m' \in \mathcal{M}$

$$\begin{aligned} & [\text{crit}(m) - \mathcal{R}(\hat{f}_m)] - [\text{crit}(m') - \mathcal{R}(\hat{f}_{m'})] \\ & \leq A(m) + B(m'), \end{aligned}$$

we have $\forall \hat{m} \in \arg \min_{m \in \mathcal{M}} \{\text{crit}(m)\}$,

$$\ell(\hat{f}_{\hat{m}}, f^*) - B(\hat{m}) \leq \inf_{m \in \mathcal{M}} \left\{ \ell(\hat{f}_m, f^*) + A(m) \right\}$$

Proof of the key lemma

For any $m \in \mathcal{M}$, we have

$$\begin{aligned} & \mathcal{R}(\hat{f}_{\hat{m}}) \\ &= \text{crit}(\hat{m}) + \mathcal{R}(\hat{f}_{\hat{m}}) - \text{crit}(\hat{m}) \\ &\leq \text{crit}(m) + \mathcal{R}(\hat{f}_{\hat{m}}) - \text{crit}(\hat{m}) \\ &= \mathcal{R}(\hat{f}_m) + \text{crit}(m) - \mathcal{R}(\hat{f}_m) + \mathcal{R}(\hat{f}_{\hat{m}}) - \text{crit}(\hat{m}) \\ &\leq \mathcal{R}(\hat{f}_m) + A(m) + B(\hat{m}) \end{aligned}$$

hence

$$\mathcal{R}(\hat{f}_{\hat{m}}) - B(\hat{m}) \leq \inf_{m \in \mathcal{M}} \left\{ \mathcal{R}(\hat{f}_m) + A(m) \right\}$$

Proof of the key lemma

For any $m \in \mathcal{M}$, we have

$$\begin{aligned} & \mathcal{R}(\hat{f}_{\hat{m}}) \\ &= \text{crit}(\hat{m}) + \mathcal{R}(\hat{f}_{\hat{m}}) - \text{crit}(\hat{m}) \\ &\leq \text{crit}(m) + \mathcal{R}(\hat{f}_{\hat{m}}) - \text{crit}(\hat{m}) \\ &= \mathcal{R}(\hat{f}_m) + \text{crit}(m) - \mathcal{R}(\hat{f}_m) + \mathcal{R}(\hat{f}_{\hat{m}}) - \text{crit}(\hat{m}) \\ &\leq \mathcal{R}(\hat{f}_m) + A(m) + B(\hat{m}) \end{aligned}$$

hence

$$\mathcal{R}(\hat{f}_{\hat{m}}) - B(\hat{m}) \leq \inf_{m \in \mathcal{M}} \left\{ \mathcal{R}(\hat{f}_m) + A(m) \right\}$$

and $\ell(\hat{f}_{\hat{m}}, f^*) - B(\hat{m}) \leq \inf_{m \in \mathcal{M}} \left\{ \ell(\hat{f}_m, f^*) + A(m) \right\}$. □

Application 1: unbiased risk estimation

$$\forall m \in \mathcal{M}, \quad \mathbb{E}[\text{crit}(m; D_n)] \approx \mathbb{E}[\mathcal{R}(\hat{f}_m(D_n))]$$

- Examples: C_p , AIC, cross-validation, FPE, ...

Application 1: unbiased risk estimation

$$\forall m \in \mathcal{M}, \quad \mathbb{E}[\text{crit}(m; D_n)] \approx \mathbb{E}[\mathcal{R}(\hat{f}_m(D_n))]$$

- Examples: C_p , AIC, cross-validation, FPE, ...
- Concentration inequalities + \mathcal{M} “not too large”
⇒ with a large probability, $\forall m, m' \in \mathcal{M}$,

$$[\text{crit}(m) - \mathcal{R}(\hat{f}_m)] - [\text{crit}(m') - \mathcal{R}(\hat{f}_{m'})] \leq A(m) + B(m'),$$

$$\text{with } A(m) = B(m) \leq \epsilon_1 \ell(\hat{f}_m, f^*) + \epsilon_2.$$

Application 1: unbiased risk estimation

$$\forall m \in \mathcal{M}, \quad \mathbb{E}[\text{crit}(m; D_n)] \approx \mathbb{E}[\mathcal{R}(\hat{f}_m(D_n))]$$

- Examples: C_p , AIC, cross-validation, FPE, ...
- Concentration inequalities + \mathcal{M} “not too large”
⇒ with a large probability, $\forall m, m' \in \mathcal{M}$,

$$[\text{crit}(m) - \mathcal{R}(\hat{f}_m)] - [\text{crit}(m') - \mathcal{R}(\hat{f}_{m'})] \leq A(m) + B(m'),$$

$$\text{with } A(m) = B(m) \leq \epsilon_1 \ell(\hat{f}_m, f^*) + \epsilon_2.$$

- If $\epsilon_1 < 1$, by the key lemma

$$\ell(\hat{f}_m, f^*) \leq \frac{1 + \epsilon_1}{1 - \epsilon_1} \inf_{m \in \mathcal{M}} \{\ell(\hat{f}_m(D_n), f^*)\} + \frac{2\epsilon_2}{1 - \epsilon_1}$$

⇒ oracle inequality, first-order optimal if $\epsilon_1 \ll 1$ and $\epsilon_2 \ll \inf_{m \in \mathcal{M}} \{\ell(f^*, \hat{f}_m(D_n))\}$.

Application 2: upper bound on the risk

$$\forall m \in \mathcal{M}, \quad \text{crit}(m; D_n) \geq \mathcal{R}(\hat{f}_m(D_n)) \quad (\text{or } \ell(\hat{f}_m(D_n), f^*))$$

(with a large probability)

- **Examples:** BIC, structural risk minimization, ...

Application 2: upper bound on the risk

$$\forall m \in \mathcal{M}, \quad \text{crit}(m; D_n) \geq \mathcal{R}(\hat{f}_m(D_n)) \quad (\text{or } \ell(\hat{f}_m(D_n), f^*))$$

(with a large probability)

- **Examples:** BIC, structural risk minimization, ...
- Then, $\forall m, m' \in \mathcal{M}$,

$$[\text{crit}(m) - \mathcal{R}(\hat{f}_m)] - [\text{crit}(m') - \mathcal{R}(\hat{f}_{m'})] \leq A(m) + B(m'),$$

$$\text{with } A(m) = 0 \quad \text{and} \quad B(m) = \text{crit}(m) - \mathcal{R}(\hat{f}_m(D_n)).$$

Application 2: upper bound on the risk

$$\forall m \in \mathcal{M}, \quad \text{crit}(m; D_n) \geq \mathcal{R}(\hat{f}_m(D_n)) \quad (\text{or } \ell(\hat{f}_m(D_n), f^*))$$

(with a large probability)

- **Examples:** BIC, structural risk minimization, ...
- Then, $\forall m, m' \in \mathcal{M}$,

$$[\text{crit}(m) - \mathcal{R}(\hat{f}_m)] - [\text{crit}(m') - \mathcal{R}(\hat{f}_{m'})] \leq A(m) + B(m'),$$

$$\text{with } A(m) = 0 \text{ and } B(m) = \text{crit}(m) - \mathcal{R}(\hat{f}_m(D_n)).$$

- By the key lemma

$$\ell(\hat{f}_{\hat{m}}, f^*) \leq \inf_{m \in \mathcal{M}} \left\{ \ell(\hat{f}_m(D_n), f^*) + B(m) \right\}$$

⇒ **oracle inequality**, interesting if $B(m)$ small enough.

Estimator selection by penalization

$$\forall m \in \mathcal{M}, \quad \text{crit}(m; D_n) = \widehat{\mathcal{R}}_n(\widehat{f}_m(D_n)) + \text{pen}(m)$$

where $\widehat{\mathcal{R}}_n(t) := \frac{1}{n} \sum_{i=1}^n c(t(X_i), Y_i)$ (empirical risk)

Estimator selection by penalization

$$\forall m \in \mathcal{M}, \quad \text{crit}(m; D_n) = \widehat{\mathcal{R}}_n(\widehat{f}_m(D_n)) + \text{pen}(m)$$

where $\widehat{\mathcal{R}}_n(t) := \frac{1}{n} \sum_{i=1}^n c(t(X_i), Y_i)$ (empirical risk)

- Ideal penalty:

$$\text{pen}_{\text{id}}(m; D_n) := \mathcal{R}(\widehat{f}_m(D_n)) - \widehat{\mathcal{R}}_n(\widehat{f}_m(D_n)).$$

Estimator selection by penalization

$$\forall m \in \mathcal{M}, \quad \text{crit}(m; D_n) = \widehat{\mathcal{R}}_n(\widehat{f}_m(D_n)) + \text{pen}(m)$$

where $\widehat{\mathcal{R}}_n(t) := \frac{1}{n} \sum_{i=1}^n c(t(X_i), Y_i)$ (empirical risk)

- Ideal penalty:

$$\text{pen}_{\text{id}}(m; D_n) := \mathcal{R}(\widehat{f}_m(D_n)) - \widehat{\mathcal{R}}_n(\widehat{f}_m(D_n)).$$

- Unbiased risk estimation

$$\Leftrightarrow \forall m \in \mathcal{M}, \quad \mathbb{E}[\text{pen}(m; D_n)] \approx \mathbb{E}[\text{pen}_{\text{id}}(m; D_n)].$$

Estimator selection by penalization

$$\forall m \in \mathcal{M}, \quad \text{crit}(m; D_n) = \widehat{\mathcal{R}}_n(\widehat{f}_m(D_n)) + \text{pen}(m)$$

where $\widehat{\mathcal{R}}_n(t) := \frac{1}{n} \sum_{i=1}^n c(t(X_i), Y_i)$ (empirical risk)

- Ideal penalty:

$$\text{pen}_{\text{id}}(m; D_n) := \mathcal{R}(\widehat{f}_m(D_n)) - \widehat{\mathcal{R}}_n(\widehat{f}_m(D_n)).$$

- Unbiased risk estimation

$$\Leftrightarrow \forall m \in \mathcal{M}, \quad \mathbb{E}[\text{pen}(m; D_n)] \approx \mathbb{E}[\text{pen}_{\text{id}}(m; D_n)].$$

- Upper bound on the risk

$$\Leftrightarrow \forall m \in \mathcal{M}, \quad \text{pen}(m; D_n) \geq \text{pen}_{\text{id}}(m; D_n).$$