

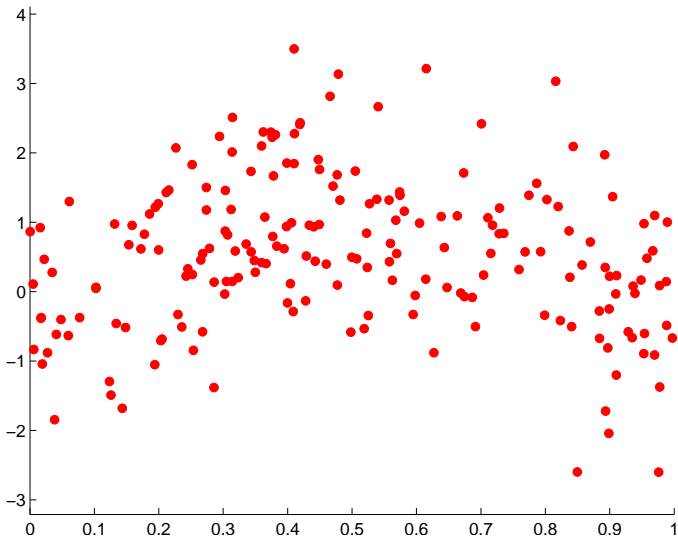
Data-driven calibration of linear estimators with minimal penalties, with an application to multi-task regression

Sylvain Arlot (joint works with F. Bach & M. Solnon)

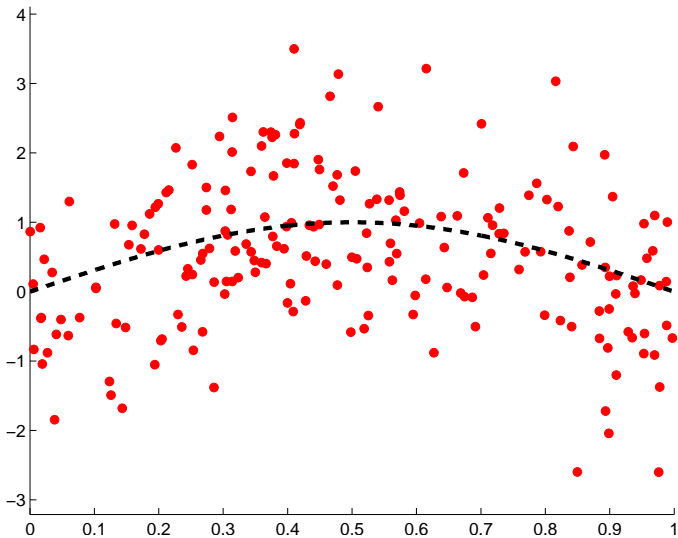
¹CNRS

²École Normale Supérieure (Paris), LIENS, Équipe SIERRA

Cambridge, November, 4th 2011

Regression: data $(x_1, Y_1), \dots, (x_n, Y_n)$ 

Goal: find the signal (denoising)



Statistical framework: regression, least-squares loss

- Observations: $Y = (Y_1, \dots, Y_n) \in \mathbb{R}^n$

$$Y_i = F_i + \varepsilon_i \quad (\text{e.g., } F_i = F(x_i))$$

with $Y_i \in \mathbb{R}$, $(\varepsilon_i)_{1 \leq i \leq n}$ i.i.d.

Statistical framework: regression, least-squares loss

- Observations: $Y = (Y_1, \dots, Y_n) \in \mathbb{R}^n$

$$Y_i = F_i + \varepsilon_i \quad (\text{e.g., } F_i = F(x_i))$$

with $Y_i \in \mathbb{R}$, $(\varepsilon_i)_{1 \leq i \leq n}$ i.i.d.

- Fixed design: x_i deterministic

Statistical framework: regression, least-squares loss

- Observations: $Y = (Y_1, \dots, Y_n) \in \mathbb{R}^n$

$$Y_i = F_i + \varepsilon_i \quad (\text{e.g., } F_i = F(x_i))$$

with $Y_i \in \mathbb{R}$, $(\varepsilon_i)_{1 \leq i \leq n}$ i.i.d.

- Fixed design: x_i deterministic
- **Least-squares loss** of a predictor $t \in \mathbb{R}^n$ (" $t_i = t(x_i)$ "):

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

Statistical framework: regression, least-squares loss

- Observations: $Y = (Y_1, \dots, Y_n) \in \mathbb{R}^n$

$$Y_i = F_i + \varepsilon_i \quad (\text{e.g., } F_i = F(x_i))$$

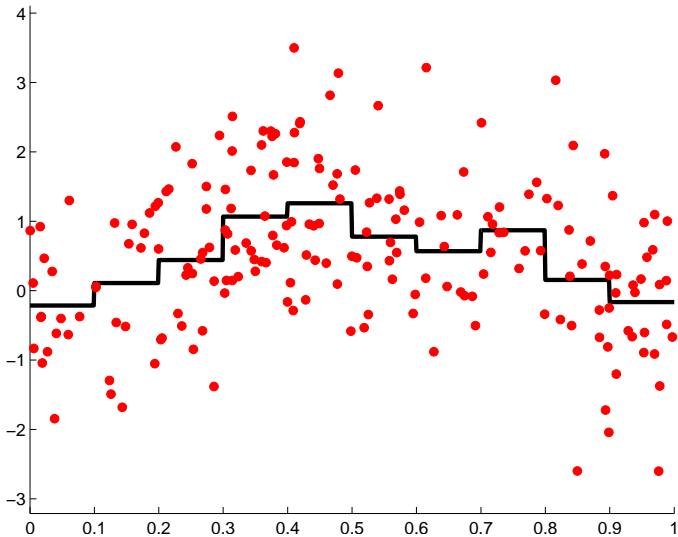
with $Y_i \in \mathbb{R}$, $(\varepsilon_i)_{1 \leq i \leq n}$ i.i.d.

- Fixed design: x_i deterministic
- **Least-squares loss** 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** $\hat{F}(Y) \in \mathbb{R}^n$?

Estimators: example: regressogram



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:**

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

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

Least-squares estimators

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

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

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

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

$$\widehat{F}_S \in \arg \min_{t \in S} \left\{ \frac{1}{n} \|t - Y\|^2 \right\} = \arg \min_{t \in S} \left\{ \frac{1}{n} \sum_{i=1}^n (t_i - Y_i)^2 \right\}$$

so that

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

Model examples

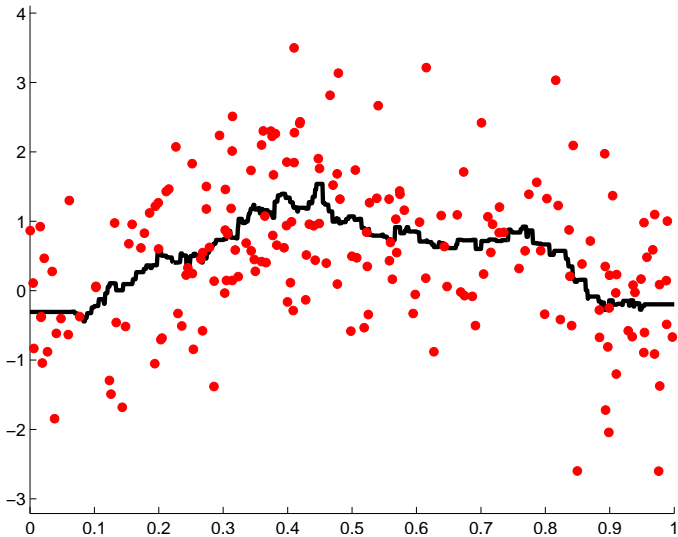
- **histograms** on some partition Λ of \mathcal{X}
⇒ the least-squares estimator (regressogram) can be written

$$\widehat{F}_m(x_i) = \sum_{\lambda \in \Lambda} \widehat{\beta}_\lambda \mathbb{1}_{x_i \in \lambda} \quad \widehat{\beta}_\lambda = \frac{1}{\text{Card}\{x_i \in \lambda\}} \sum_{x_i \in \lambda} Y_i$$

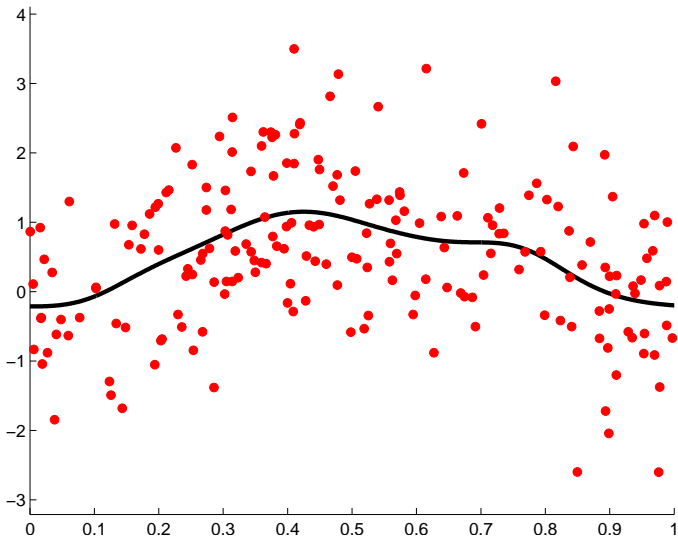
- subspace generated by a subset of an orthogonal basis of $L^2(\mu)$ (**Fourier, wavelets**, and so on)
- **variable selection**: $x_i = (x_i^{(1)}, \dots, x_i^{(p)}) \in \mathbb{R}^p$ gathers p variables that can (linearly) explain Y_i

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

k -nearest-neighbours estimator ($k = 20$)



Nadaraya-Watson estimator ($\sigma = 0.01$)



Linear estimators

- OLS: $\hat{F}_m = \Pi_{S_m} Y$ (projection onto S_m)
- (kernel) ridge regression, spline smoothing:

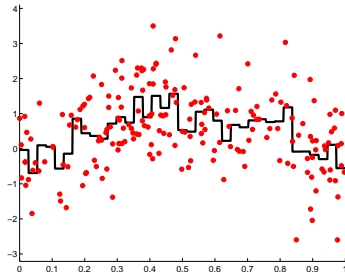
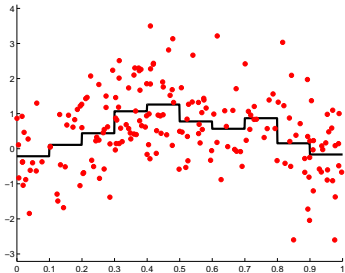
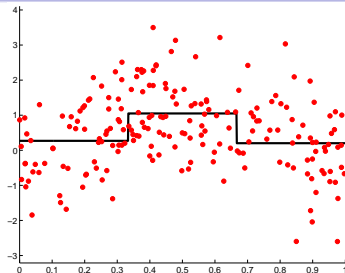
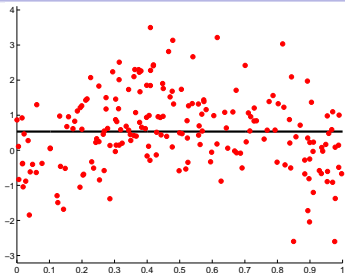
$$\hat{F}_i = \hat{f}(x_i) \quad \text{with} \quad \hat{f} \in \arg \min_{f \in \mathcal{F}_K} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - f(x_i))^2 + \lambda \|f\|_{\mathcal{F}_K}^2 \right\}$$

$$\Rightarrow \hat{F}_{\lambda, K} = K(K + \lambda I)^{-1} Y \quad \text{where} \quad K = (K(x_i, x_j))_{1 \leq i, j \leq n}$$

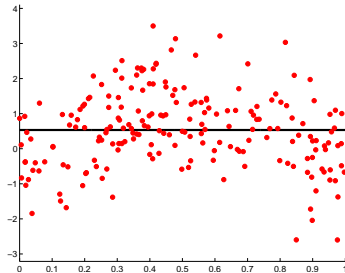
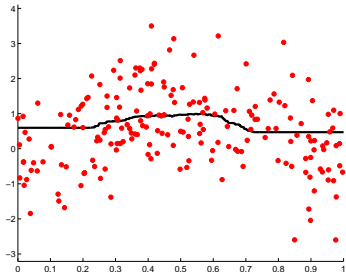
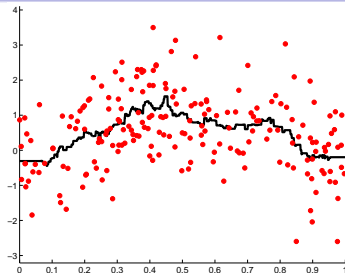
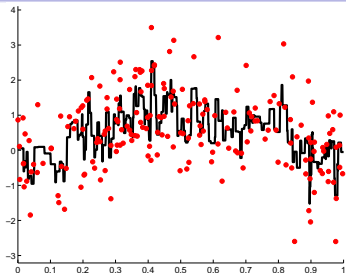
- k -nearest neighbours
- Nadaraya-Watson estimators

$$\hat{F} = AY \quad \text{where } A \text{ does not depend on } Y$$

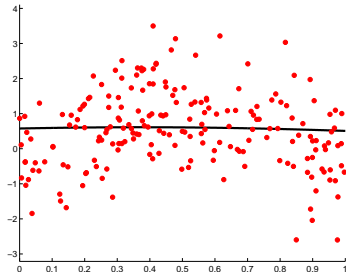
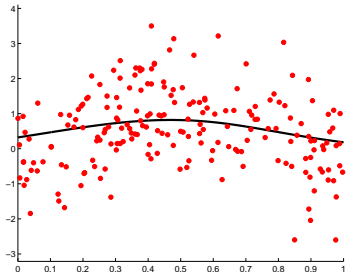
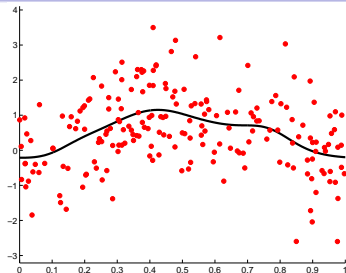
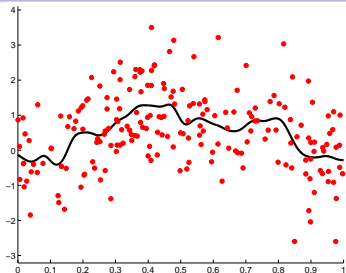
Estimator selection: regular regressograms



Estimator selection: k nearest neighbours



Estimator selection: Nadaraya-Watson



Estimator selection

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

Example: $\hat{F}_m = A_m Y$

Estimator selection

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

Example: $\hat{F}_m = A_m Y$

- Goal: minimize the risk, i.e.,

Oracle inequality (in expectation or with a large probability):

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

Estimator selection

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

Example: $\hat{F}_m = A_m Y$

- Goal: minimize the risk, i.e.,

Oracle inequality (in expectation or with a large probability):

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

- Examples:

- **model selection**
- **calibration** (choosing k or the distance for k -NN, choice of a regularization parameter, choice of a kernel, etc.)
- choice between **methods different in nature**
ex.: k -NN vs. smoothing splines?

Bias-variance trade-off

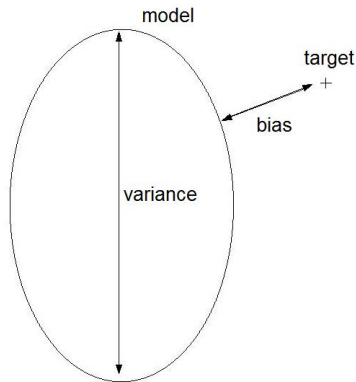
$$\mathbb{E} \left[\frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 \right] = \text{Bias} + \text{Variance}$$

Bias or Approximation error

$$\frac{1}{n} \left\| F_m - F \right\|^2 = \frac{1}{n} \left\| A_m F - F \right\|^2$$

Variance or Estimation error

$$\frac{\sigma^2 \text{tr} (A_m^\top A_m)}{n} \quad \text{OLS:} \quad \frac{\sigma^2 \dim(S_m)}{n}$$



Bias-variance trade-off

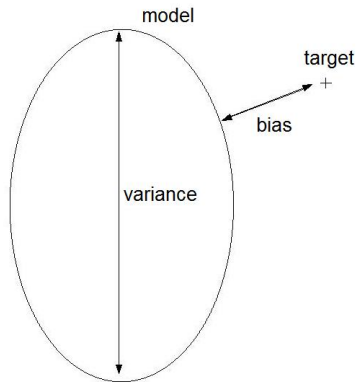
$$\mathbb{E} \left[\frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 \right] = \text{Bias} + \text{Variance}$$

Bias or Approximation error

$$\frac{1}{n} \left\| F_m - F \right\|^2 = \frac{1}{n} \left\| A_m F - F \right\|^2$$

Variance or Estimation error

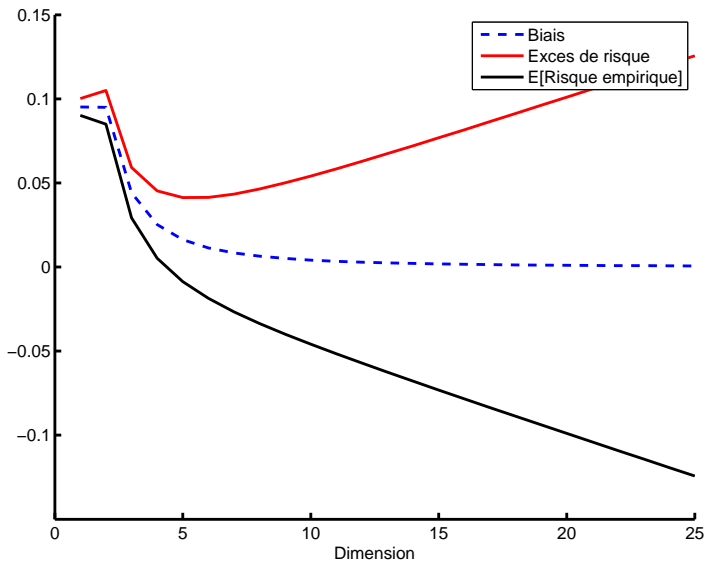
$$\frac{\sigma^2 \text{tr} (A_m^\top A_m)}{n} \quad \text{OLS:} \quad \frac{\sigma^2 \dim(S_m)}{n}$$



Bias-variance trade-off

⇔ avoid **overfitting** and **underfitting**

Why should the empirical risk be penalized?



Penalization

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

Penalization

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

- Ideal penalty:

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

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

Penalization

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

- Ideal penalty:

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

- 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)

\Rightarrow OLS: C_p : $2\sigma^2 D_m/n$ (Mallows, 1973)

\Rightarrow Linear estimators: C_L : $2\sigma^2 \text{tr}(A_m)/n$ (Mallows, 1973)

Oracle inequality

Theorem (Birgé & Massart 2007, A. & Bach 2009–2011)

Assumptions:

- $\text{pen}(m) = \frac{2C \text{tr}(A_m)}{n}$ with $|C\sigma^{-2} - 1| \leq L_0 \sqrt{\frac{\ln(n)}{n}}$

Then, with probability at least $1 - 3 \text{Card}(\mathcal{M})n^{-\delta}$, if $n \geq n_0$, for every $\eta \in (0, 1)$,

$$\frac{1}{n} \left\| \widehat{F}_{\widehat{m}} - F \right\|^2 \leq (1 + \eta) \inf_{m \in \mathcal{M}} \left\{ \frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 \right\} + K(\delta, L_0, L_1) \frac{\ln(n)\sigma^2}{\eta n}$$

Oracle inequality

Theorem (Birgé & Massart 2007, A. & Bach 2009–2011)

Assumptions:

- $\text{pen}(m) = \frac{2C \text{tr}(A_m)}{n}$ with $|C\sigma^{-2} - 1| \leq L_0 \sqrt{\frac{\ln(n)}{n}}$

Then, with probability at least $1 - 3 \text{Card}(\mathcal{M})n^{-\delta}$, if $n \geq n_0$, for every $\eta \in (0, 1)$,

$$\frac{1}{n} \left\| \widehat{F}_{\widehat{m}} - F \right\|^2 \leq (1 + \eta) \inf_{m \in \mathcal{M}} \left\{ \frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 \right\} + K(\delta, L_0, L_1) \frac{\ln(n)\sigma^2}{\eta n}$$

Note: ridge, $\lambda \in [0, +\infty] \Leftrightarrow \text{Card}(\mathcal{M}) \propto n$

Oracle inequality

Theorem (Birgé & Massart 2007, A. & Bach 2009–2011)

Assumptions:

- $\text{pen}(m) = \frac{2C \text{tr}(A_m)}{n}$ with $|C\sigma^{-2} - 1| \leq L_0 \sqrt{\frac{\ln(n)}{n}}$
- $\forall m \in \mathcal{M}, \|A_m\| \leq L_1$ and $\text{tr}(A_m^\top A_m) \leq \text{tr}(A_m) \leq n$
- *Gaussian noise:* $\varepsilon_j \sim \mathcal{N}(0, \sigma^2)$

Then, with probability at least $1 - 3 \text{Card}(\mathcal{M})n^{-\delta}$, if $n \geq n_0$, for every $\eta \in (0, 1)$,

$$\frac{1}{n} \left\| \widehat{F}_{\widehat{m}} - F \right\|^2 \leq (1 + \eta) \inf_{m \in \mathcal{M}} \left\{ \frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 \right\} + K(\delta, L_0, L_1) \frac{\ln(n)\sigma^2}{\eta n}$$

Note: ridge, $\lambda \in [0, +\infty] \Leftrightarrow \text{Card}(\mathcal{M}) \propto n$

Motivation (1): penalties known up to a constant factor

$$\text{Ex.:} \quad \text{pen}_{\text{Cp}}(m) = \frac{2\sigma^2 D_m}{n} \quad \text{pen}_{\text{CL}}(m) = \frac{2\sigma^2 \text{tr}(A_m)}{n}$$

Motivation (1): penalties known up to a constant factor

$$\text{Ex.:} \quad \text{pen}_{\text{Cp}}(m) = \frac{2\sigma^2 D_m}{n} \quad \text{pen}_{\text{CL}}(m) = \frac{2\sigma^2 \text{tr}(A_m)}{n}$$

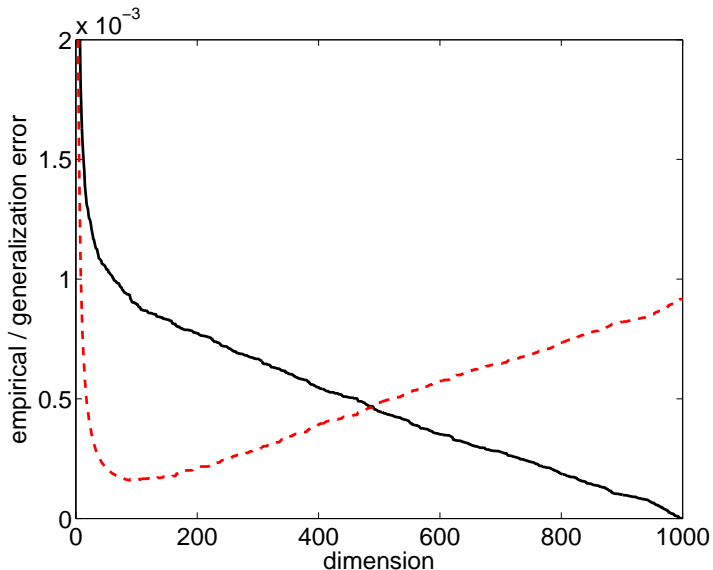
- Classical estimators of σ^2 :
 - $\hat{\sigma}_{m_0}^2 := \|Y - \hat{F}_{m_0}\|^2 / (n - D_{m_0})$ (OLS)
problem: choice of m_0 ?
 - $\hat{\sigma}_m^2 \Rightarrow$ **FPE** (Akaike, 1970) and **GCV** (Craven & Wahba, 1979)
problem: avoiding the largest models

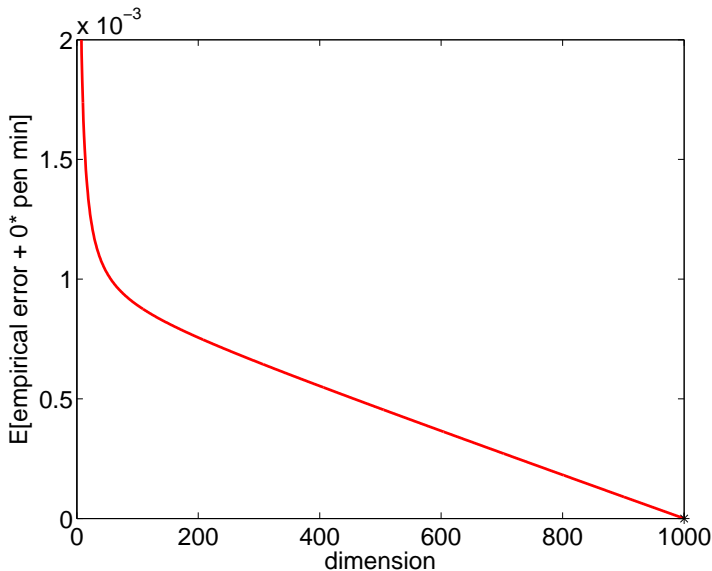
Motivation (1): penalties known up to a constant factor

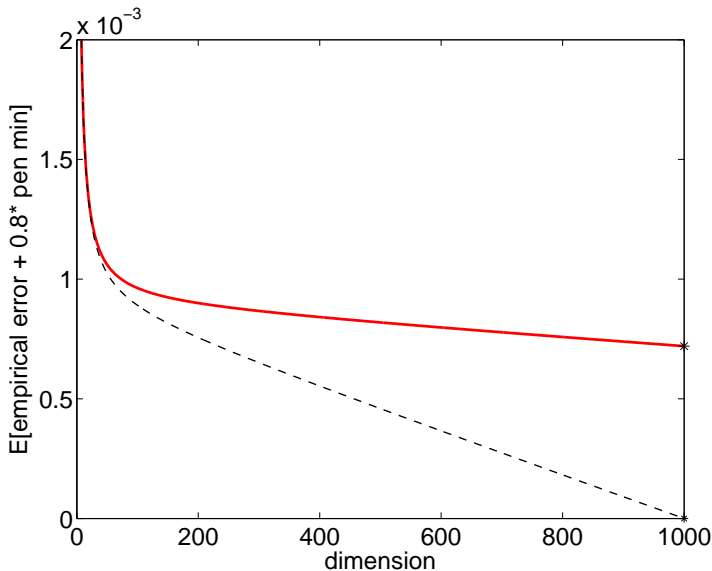
$$\text{Ex.:} \quad \text{pen}_{\text{Cp}}(m) = \frac{2\sigma^2 D_m}{n} \quad \text{pen}_{\text{CL}}(m) = \frac{2\sigma^2 \text{tr}(A_m)}{n}$$

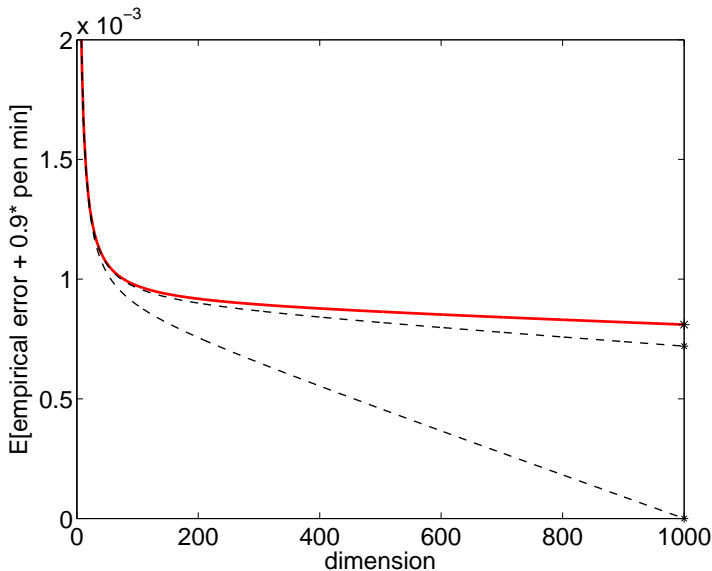
- Classical estimators of σ^2 :
 - $\hat{\sigma}_{m_0}^2 := \|Y - \hat{F}_{m_0}\|^2 / (n - D_{m_0})$ (OLS)
problem: choice of m_0 ?
 - $\hat{\sigma}_m^2 \Rightarrow$ **FPE** (Akaike, 1970) and **GCV** (Craven & Wahba, 1979)
problem: avoiding the largest models
- **Goals**: estimation of σ^2 for model selection, under minimal assumptions, without overfitting

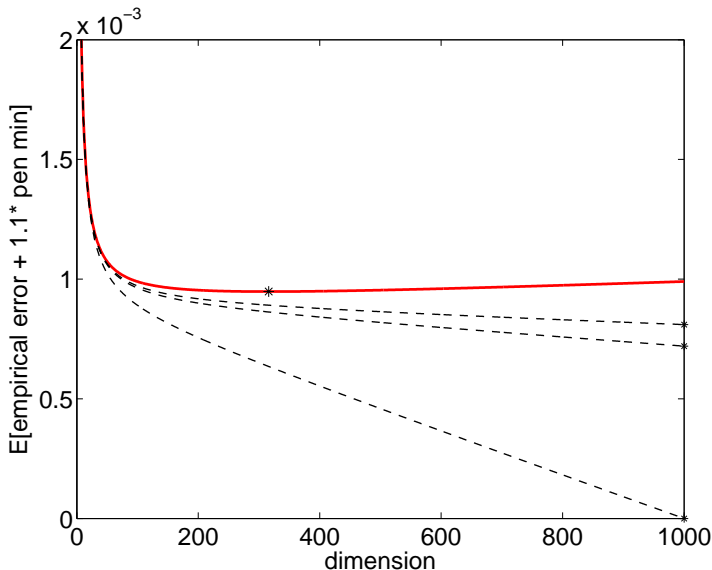
Motivation (2): “L-curve” and elbow heuristics?

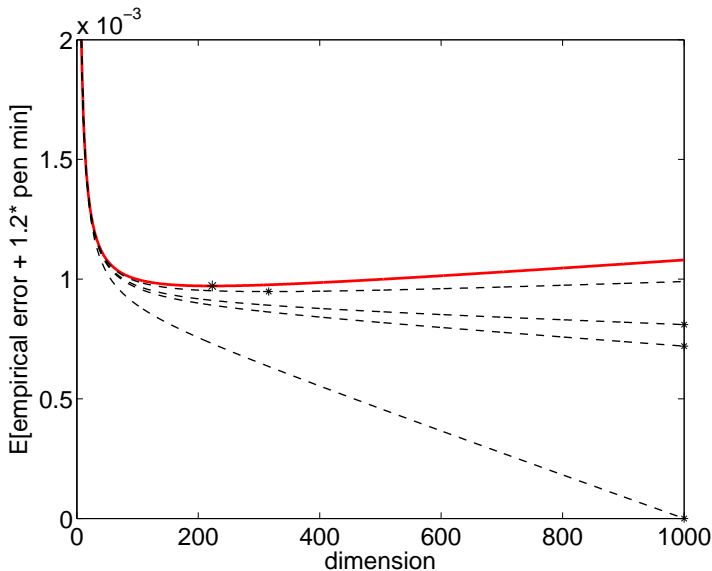


$\mathbb{E}[\text{Empirical risk}] + 0 \times \sigma^2 D_m n^{-1} \text{ (OLS)}$ 

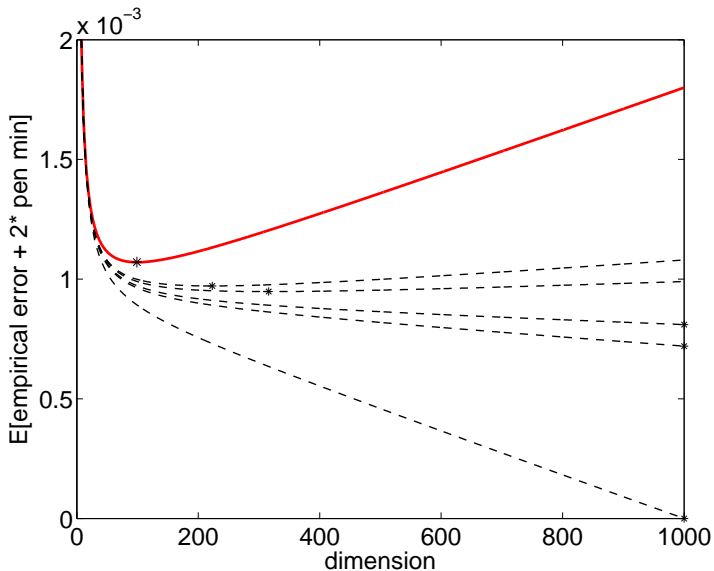
$\mathbb{E}[\text{Empirical risk}] + 0.8 \times \sigma^2 D_m n^{-1}$ (OLS)

$\mathbb{E}[\text{Empirical risk}] + 0.9 \times \sigma^2 D_m n^{-1} \text{ (OLS)}$ 

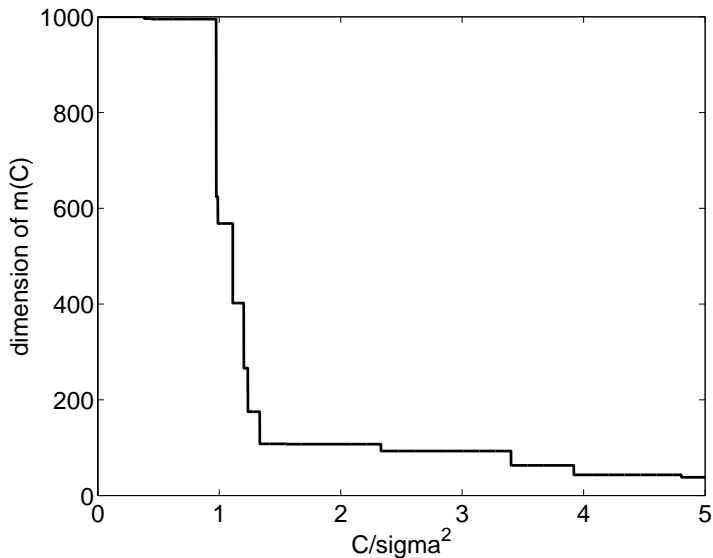
$\mathbb{E}[\text{Empirical risk}] + 1.1 \times \sigma^2 D_m n^{-1} \text{ (OLS)}$ 

$\mathbb{E}[\text{Empirical risk}] + 1.2 \times \sigma^2 D_m n^{-1} \text{ (OLS)}$ 

$$\mathbb{E}[\text{Empirical risk}] + 2 \times \sigma^2 D_m n^{-1} \text{ (OLS)}$$



OLS: Dimension jump



OLS: algorithm (Birgé & Massart 2007)

- 1 for every $C > 0$, compute

$$\hat{m}(C) \in \arg \min_{m \in \mathcal{M}_n} \left\{ \frac{1}{n} \left\| \hat{F}_m - Y \right\|^2 + C \frac{D_m}{n} \right\}$$

- 2 find \hat{C}_{\min} such that $D_{\hat{m}(C)}$ is “very large” when $C < \hat{C}_{\min}$ and “reasonably small” when $C > \hat{C}_{\min}$

- 3 select $\hat{m} = \hat{m} \left(2\hat{C}_{\min} \right)$

Practical use: CAPUSHE package (Baudry, Maugis & Michel, 2010)

<http://www.math.univ-toulouse.fr/~maugis/CAPUSHE.html>

Practical qualities of the algorithm

- **visual checking** of existence of a jump
- calibration **independent from the choice of some m_0**
- too strong **overfitting** almost impossible
- one remaining parameter: how to **localize the jump**

Theorem (1): Dimension jump / Minimal penalty

Theorem (Birgé & Massart 2007, A. & Bach 2009–2011)

Assumptions:

- $\exists m_0 \in \mathcal{M}$, $D_{m_0} \leq \sqrt{n}$ and $\frac{1}{n} \|F_{m_0} - F\|^2 \leq \sigma^2 \sqrt{\ln(n)/n}$.
- *Gaussian noise:* $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$

Then, with probability at least $1 - 3 \text{Card}(\mathcal{M})n^{-\delta}$, if $n \geq n_0(\delta)$,

$$\forall C < \left(1 - L_3 \delta \sqrt{\frac{\ln(n)}{n}}\right) \sigma^2, \quad D_{\hat{m}(C)} \geq \frac{n}{3}$$

$$\forall C > \left(1 + L_3 \delta \sqrt{\frac{\ln(n)}{n}}\right) \sigma^2, \quad D_{\hat{m}(C)} \leq \frac{n}{10}$$

and in the first case, $\frac{1}{n} \|\hat{F}_{\hat{m}(C)} - F\|^2 \gg \inf_{m \in \mathcal{M}_n} \left\{ \frac{1}{n} \|\hat{F}_m - F\|^2 \right\}$.

Theorem (2): Oracle inequality

Theorem (Birgé & Massart 2007, A. & Bach 2009–2011)

Assumptions:

- $\hat{m} \in \arg \min_{m \in \mathcal{M}} \left\{ \frac{1}{n} \|\hat{F}_m - Y\|^2 + 2\hat{C}_{\min} \frac{D_m}{n} \right\}$
- $\exists m_0 \in \mathcal{M}, D_{m_0} \leq \sqrt{n}$ and $\frac{1}{n} \|F_{m_0} - F\|^2 \leq \sigma^2 \sqrt{\ln(n)/n}$.
- *Gaussian noise:* $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$

Then, with probability at least $1 - 3 \text{Card}(\mathcal{M})n^{-\delta}$, if $n \geq n_0(\delta)$, for every $\eta > 0$,

$$\frac{1}{n} \|\hat{F}_{\hat{m}} - F\|^2 \leq (1 + \eta) \inf_{m \in \mathcal{M}} \left\{ \frac{1}{n} \|\hat{F}_m - F\|^2 \right\} + K(\delta) \frac{\ln(n)\sigma^2}{\eta n}$$

Generalization of minimal penalties to linear estimators?

OLS

$$\text{pen}_{\text{Cp}}(m) = \frac{2\sigma^2 D_m}{n}$$

$$\arg \min_{m \in \mathcal{M}} \left\{ \frac{1}{n} \|\widehat{F}_m - Y\|^2 + c \frac{D_m}{n} \right\}$$

$\Rightarrow D_{\widehat{m}(c)}$ “jumps” at $\widehat{C}_{\min} \approx \sigma^2$

\Rightarrow optimal choice with $\widehat{m}(2\widehat{C}_{\min})$

Generalization of minimal penalties to linear estimators?

OLS

$$\text{pen}_{\text{Cp}}(m) = \frac{2\sigma^2 D_m}{n}$$

$$\arg \min_{m \in \mathcal{M}} \left\{ \frac{1}{n} \|\widehat{F}_m - Y\|^2 + c \frac{D_m}{n} \right\}$$

$\Rightarrow D_{\widehat{m}(c)}$ “jumps” at $\widehat{C}_{\min} \approx \sigma^2$

\Rightarrow optimal choice with $\widehat{m}(2\widehat{C}_{\min})$

Linear estimators

$$\text{pen}_{\text{CL}}(m) = \frac{2\sigma^2 \text{tr}(A_m)}{n}$$

Generalization of minimal penalties to linear estimators?

OLS

$$\text{pen}_{\text{Cp}}(m) = \frac{2\sigma^2 D_m}{n}$$

$$\arg \min_{m \in \mathcal{M}} \left\{ \frac{1}{n} \|\widehat{F}_m - Y\|^2 + c \frac{D_m}{n} \right\}$$

$\Rightarrow D_{\widehat{m}(c)}$ “jumps” at $\widehat{C}_{\min} \approx \sigma^2$

\Rightarrow optimal choice with $\widehat{m}(2\widehat{C}_{\min})$

Linear estimators

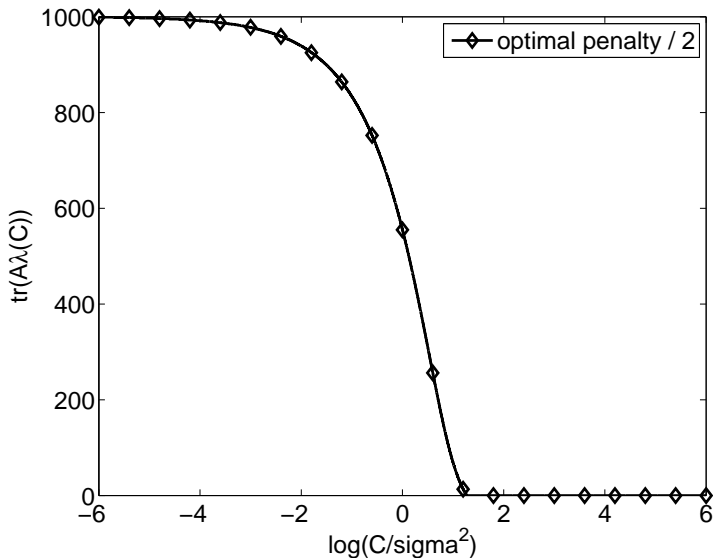
$$\text{pen}_{\text{CL}}(m) = \frac{2\sigma^2 \text{tr}(A_m)}{n}$$

$$\arg \min_{m \in \mathcal{M}} \left\{ \frac{1}{n} \|\widehat{F}_m - Y\|^2 + c \frac{\text{tr}(A_m)}{n} \right\}$$

Does $\text{tr}(A_{\widehat{m}(c)})$ jump at $\widehat{C}_{\min} \approx \sigma^2$?

optimal choice with $\widehat{m}(2\widehat{C}_{\min})$?

No dimension jump with a penalty $\propto \text{tr}(A_m)$



Minimal penalties for linear estimators

$$\mathbb{E} \left[\frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 \right] = \frac{1}{n} \left\| (I - A_m) F \right\|^2 + \frac{\text{tr}(A_m^\top A_m) \sigma^2}{n} = \text{bias} + \text{variance}$$

Minimal penalties for linear estimators

$$\mathbb{E} \left[\frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 \right] = \frac{1}{n} \left\| (I - A_m) F \right\|^2 + \frac{\text{tr}(A_m^\top A_m) \sigma^2}{n} = \text{bias} + \text{variance}$$

$$\mathbb{E} \left[\frac{1}{n} \left\| \widehat{F}_m - Y \right\|^2 \right] = \sigma^2 + \frac{1}{n} \left\| (I - A_m) F \right\|^2 - \frac{(2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m)) \sigma^2}{n}$$

Minimal penalties for linear estimators

$$\mathbb{E} \left[\frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 \right] = \frac{1}{n} \left\| (I - A_m) F \right\|^2 + \frac{\text{tr}(A_m^\top A_m) \sigma^2}{n} = \text{bias} + \text{variance}$$

$$\mathbb{E} \left[\frac{1}{n} \left\| \widehat{F}_m - Y \right\|^2 \right] = \sigma^2 + \frac{1}{n} \left\| (I - A_m) F \right\|^2 - \frac{(2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m)) \sigma^2}{n}$$

$$\Rightarrow \text{optimal penalty } \frac{(2 \text{tr}(A_m)) \sigma^2}{n}$$

Minimal penalties for linear estimators

$$\mathbb{E} \left[\frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 \right] = \frac{1}{n} \left\| (I - A_m) F \right\|^2 + \frac{\text{tr}(A_m^\top A_m) \sigma^2}{n} = \text{bias} + \text{variance}$$

$$\mathbb{E} \left[\frac{1}{n} \left\| \widehat{F}_m - Y \right\|^2 \right] = \sigma^2 + \frac{1}{n} \left\| (I - A_m) F \right\|^2 - \frac{(2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m)) \sigma^2}{n}$$

$$\Rightarrow \text{optimal penalty} \quad \frac{(2 \text{tr}(A_m)) \sigma^2}{n}$$

$$\Rightarrow \text{minimal penalty} \quad \frac{(2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m)) \sigma^2}{n}$$

Minimal penalties for linear estimators

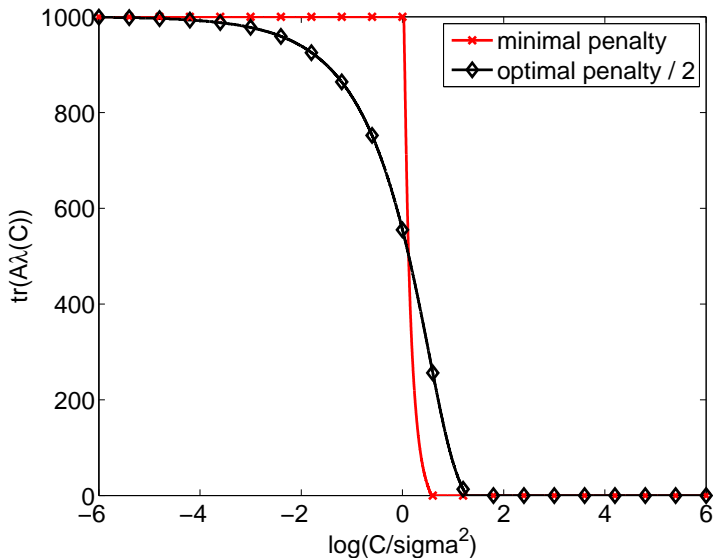
$$\mathbb{E} \left[\frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 \right] = \frac{1}{n} \left\| (I - A_m) F \right\|^2 + \frac{\text{tr}(A_m^\top A_m) \sigma^2}{n} = \text{bias} + \text{variance}$$

$$\mathbb{E} \left[\frac{1}{n} \left\| \widehat{F}_m - Y \right\|^2 \right] = \sigma^2 + \frac{1}{n} \left\| (I - A_m) F \right\|^2 - \frac{(2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m)) \sigma^2}{n}$$

$$\Rightarrow \text{optimal penalty } \frac{(2 \text{tr}(A_m)) \sigma^2}{n}$$

$$\widehat{m}(C) \in \arg \min_{\lambda \in \Lambda} \left\{ \frac{1}{n} \left\| \widehat{F}_m - Y \right\|^2 + C \times \frac{2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m)}{n} \right\}$$

“Dimension” jump (ridge regression)



Penalty calibration algorithm (A. & Bach 2009)

- 1 for every $C > 0$, compute

$$\hat{m}_{\min}(C) \in \arg \min_{m \in \mathcal{M}} \left\{ \frac{1}{n} \|Y - \hat{F}_m\|^2 + \frac{C (2 \operatorname{tr}(A_m) - \operatorname{tr}(A_m^\top A_m))}{n} \right\}$$

- 2 find \hat{C}_{\min} such that $\operatorname{tr}(A_{\hat{m}_{\min}(C)})$ is “too large” when $C < \hat{C}_{\min}$ and “reasonably small” when $C > \hat{C}_{\min}$,

- 3 select

$$\hat{m} \in \arg \min_{m \in \mathcal{M}} \left\{ \frac{1}{n} \|Y - \hat{F}_m\|^2 + \frac{2\hat{C}_{\min} \operatorname{tr}(A_m)}{n} \right\}$$

Penalty calibration algorithm (A. & Bach 2009)

- 1 for every $C > 0$, compute

$$\hat{m}_{\min}(C) \in \arg \min_{m \in \mathcal{M}} \left\{ \frac{1}{n} \|Y - \hat{F}_m\|^2 + \frac{C (2 \operatorname{tr}(A_m) - \operatorname{tr}(A_m^\top A_m))}{n} \right\}$$

- 2 find \hat{C}_{\min} such that $\operatorname{tr}(A_{\hat{m}_{\min}(C)})$ is “too large” when $C < \hat{C}_{\min}$ and “reasonably small” when $C > \hat{C}_{\min}$,

- 3 select

$$\hat{m} \in \arg \min_{m \in \mathcal{M}} \left\{ \frac{1}{n} \|Y - \hat{F}_m\|^2 + \frac{2\hat{C}_{\min} \operatorname{tr}(A_m)}{n} \right\}$$

$\Rightarrow \left| \hat{C}_{\min} \sigma^{-2} - 1 \right| \leq L_4 \sqrt{\ln(n)/n}$ and oracle inequality (same assumptions as before).

Comparison with least-squares

- Linear estimators:

$$\text{pen}_{\min}(m) = \frac{\sigma^2 (2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m))}{n}$$

$$\text{pen}_{\text{opt}}(m) = \frac{\sigma^2 (2 \text{tr}(A_m))}{n}$$

$$\frac{\text{pen}_{\text{opt}}(m)}{\text{pen}_{\min}(m)} = \frac{2 \text{tr}(A_m)}{2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m)} \in (1, 2]$$

Comparison with least-squares

- Linear estimators:

$$\text{pen}_{\min}(m) = \frac{\sigma^2 (2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m))}{n}$$

$$\text{pen}_{\text{opt}}(m) = \frac{\sigma^2 (2 \text{tr}(A_m))}{n}$$

$$\frac{\text{pen}_{\text{opt}}(m)}{\text{pen}_{\min}(m)} = \frac{2 \text{tr}(A_m)}{2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m)} \in (1, 2]$$

- Least-squares case:

$$A_m^\top A_m = A_m \quad \Rightarrow \quad \frac{\text{pen}_{\text{opt}}(m)}{\text{pen}_{\min}(m)} = 2 \quad \Rightarrow \quad \text{Slope heuristics}$$

The k -nearest neighbours case

$$\forall i, j \in \{1, \dots, n\}, \quad A_{i,j} \in \left\{ 0, \frac{1}{k} \right\}$$
$$\forall i \in \{1, \dots, n\}, \quad A_{i,i} = \frac{1}{k} \quad \text{and} \quad \sum_{j=1}^n A_{i,j} = 1$$

The k -nearest neighbours case

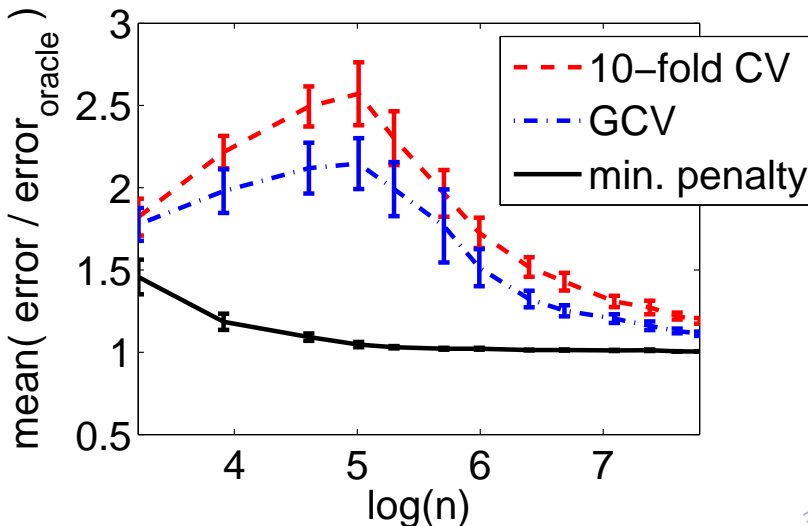
$$\forall i, j \in \{1, \dots, n\}, \quad A_{i,j} \in \left\{ 0, \frac{1}{k} \right\}$$

$$\forall i \in \{1, \dots, n\}, \quad A_{i,i} = \frac{1}{k} \quad \text{and} \quad \sum_{j=1}^n A_{i,j} = 1$$

$$\Rightarrow \quad \text{tr}(A) = \frac{n}{k} = \text{tr}(A^\top A)$$

$$\Rightarrow \quad \text{pen}_{\text{opt}} = 2 \text{pen}_{\text{min}}$$

Simulation study (ridge regression, choice of λ)



Multi-task regression

- We want to solve $p \geq 2$ regression problems simultaneously

Multi-task regression

- We want to solve $p \geq 2$ regression problems simultaneously
- Observations: $Y_1, \dots, Y_n \in \mathbb{R}^p$

$$Y_i^j = F_i^j + \varepsilon_i^j \quad j = 1, \dots, p \quad (\text{e.g., } F_i^j = F^j(x_i))$$

with $(\varepsilon_i)_{1 \leq i \leq n}$ i.i.d. $\mathcal{N}(0, \Sigma)$, $\Sigma \in \mathcal{S}_p^+(\mathbb{R})$

Multi-task regression

- We want to solve $p \geq 2$ regression problems simultaneously
- Observations: $Y_1, \dots, Y_n \in \mathbb{R}^p$

$$Y_i^j = F_i^j + \varepsilon_i^j \quad j = 1, \dots, p \quad (\text{e.g., } F_i^j = F^j(x_i))$$

with $(\varepsilon_i)_{1 \leq i \leq n}$ i.i.d. $\mathcal{N}(0, \Sigma)$, $\Sigma \in \mathcal{S}_p^+(\mathbb{R})$

- Implicit assumption: the p problems are “similar”

Multi-task regression

- We want to solve $p \geq 2$ regression problems simultaneously
- Observations: $Y_1, \dots, Y_n \in \mathbb{R}^p$

$$Y_i^j = F_i^j + \varepsilon_i^j \quad j = 1, \dots, p \quad (\text{e.g., } F_i^j = F^j(x_i))$$

with $(\varepsilon_i)_{1 \leq i \leq n}$ i.i.d. $\mathcal{N}(0, \Sigma)$, $\Sigma \in \mathcal{S}_p^+(\mathbb{R})$

- Implicit assumption: the p problems are “similar”
- **Least-squares loss** of a predictor $t \in \mathbb{R}^{np}$ (“ $t_i^j = t^j(x_i)$ ”):

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

Multi-task regression

- We want to solve $p \geq 2$ regression problems simultaneously
- Observations: $Y_1, \dots, Y_n \in \mathbb{R}^p$

$$Y_i^j = F_i^j + \varepsilon_i^j \quad j = 1, \dots, p \quad (\text{e.g., } F_i^j = F^j(x_i))$$

with $(\varepsilon_i)_{1 \leq i \leq n}$ i.i.d. $\mathcal{N}(0, \Sigma)$, $\Sigma \in \mathcal{S}_p^+(\mathbb{R})$

- Implicit assumption: the p problems are “similar”
- **Least-squares loss** of a predictor $t \in \mathbb{R}^{np}$ (“ $t_i^j = t^j(x_i)$ ”):

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

⇒ **Estimator** $\hat{F}(Y_1, \dots, Y_n) \in \mathbb{R}^{np}$?

Ridge multi-task regression

$\hat{F} = (\hat{F}_i^j)_{1 \leq i \leq n, 1 \leq j \leq p}$ with $\hat{F}_i^j = \hat{f}^j(x_i)$ and \hat{f} defined by:

- If we consider the tasks separately:

$$\arg \min_{f \in \mathcal{F}_K^p} \left\{ \frac{1}{np} \sum_{i=1}^n \sum_{j=1}^p \left(Y_i^j - f^j(x_i) \right)^2 + \sum_{j=1}^p \lambda^j \|f^j\|_{\mathcal{F}_K}^2 \right\}$$

Ridge multi-task regression

$\widehat{F} = (\widehat{F}_i^j)_{1 \leq i \leq n, 1 \leq j \leq p}$ with $\widehat{F}_i^j = \widehat{f}^j(x_i)$ and \widehat{f} defined by:

- If we consider the tasks separately:

$$\arg \min_{f \in \mathcal{F}_K^p} \left\{ \frac{1}{np} \sum_{i=1}^n \sum_{j=1}^p \left(Y_i^j - f^j(x_i) \right)^2 + \sum_{j=1}^p \lambda^j \|f^j\|_{\mathcal{F}_K}^2 \right\}$$

- A possible multi-task approach (Evgeniou *et al.*, 2005):

$$\arg \min_{f \in \mathcal{F}_K^p} \left\{ \frac{1}{np} \sum_{i=1}^n \sum_{j=1}^p \left(Y_i^j - f^j(x_i) \right)^2 + \lambda \sum_{j=1}^p \|f^j\|_{\mathcal{F}_K}^2 + \mu \sum_{j \neq \ell} \|f^j - f^\ell\|_{\mathcal{F}_K}^2 \right\}$$

Ridge multi-task regression

$\widehat{F} = (\widehat{F}_i^j)_{1 \leq i \leq n, 1 \leq j \leq p}$ with $\widehat{F}_i^j = \widehat{f}^j(x_i)$ and \widehat{f} defined by:

- If we consider the tasks separately:

$$\arg \min_{f \in \mathcal{F}_K^p} \left\{ \frac{1}{np} \sum_{i=1}^n \sum_{j=1}^p \left(Y_i^j - f^j(x_i) \right)^2 + \sum_{j=1}^p \lambda^j \|f^j\|_{\mathcal{F}_K}^2 \right\}$$

- A possible multi-task approach (Evgeniou *et al.*, 2005):

$$\arg \min_{f \in \mathcal{F}_K^p} \left\{ \frac{1}{np} \sum_{i=1}^n \sum_{j=1}^p \left(Y_i^j - f^j(x_i) \right)^2 + \lambda \sum_{j=1}^p \|f^j\|_{\mathcal{F}_K}^2 + \mu \sum_{j \neq \ell} \|f^j - f^\ell\|_{\mathcal{F}_K}^2 \right\}$$

- More generally: for $M \in \mathcal{S}_p^+(\mathbb{R})$,

$$\arg \min_{f \in \mathcal{F}_K^p} \left\{ \frac{1}{np} \sum_{i=1}^n \sum_{j=1}^p \left(Y_i^j - f^j(x_i) \right)^2 + \sum_{j,\ell} M_{j,\ell} \langle f^j, f^\ell \rangle_{\mathcal{F}_K} \right\}$$

Multi-task estimator selection

⇒ Estimators collection $(\hat{F}_M)_{M \in \mathcal{M}}$, $\mathcal{M} \subset \mathcal{S}_p^+(\mathbb{R})$,

with $\hat{F}_M = A_M Y$ and $A_M = (M^{-1} \otimes K) ((M^{-1} \otimes K) + npl_{np})^{-1}$

Multi-task estimator selection

⇒ Estimators collection $(\hat{F}_M)_{M \in \mathcal{M}}$, $\mathcal{M} \subset \mathcal{S}_p^+(\mathbb{R})$,

with $\hat{F}_M = A_M Y$ and $A_M = (M^{-1} \otimes K) ((M^{-1} \otimes K) + npl_{np})^{-1}$

- Goal: select $\hat{M} \in \mathcal{M}$ such that $\frac{1}{np} \|\hat{F}_{\hat{M}} - F\|^2$ is minimal

Multi-task estimator selection

⇒ Estimators collection $(\hat{F}_M)_{M \in \mathcal{M}}$, $\mathcal{M} \subset \mathcal{S}_p^+(\mathbb{R})$,

with $\hat{F}_M = A_M Y$ and $A_M = (M^{-1} \otimes K) ((M^{-1} \otimes K) + npI_{np})^{-1}$

- Goal: select $\hat{M} \in \mathcal{M}$ such that $\frac{1}{np} \|\hat{F}_{\hat{M}} - F\|^2$ is minimal
- Expectation of the ideal penalty:

$$\mathbb{E}[\text{pen}_{\text{id}}(M)] = \frac{2}{np} \text{tr}(A_M (\Sigma \otimes I_n))$$

Multi-task estimator selection

⇒ Estimators collection $(\hat{F}_M)_{M \in \mathcal{M}}$, $\mathcal{M} \subset \mathcal{S}_p^+(\mathbb{R})$,

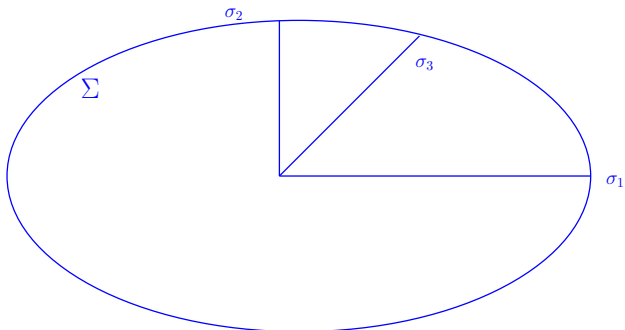
with $\hat{F}_M = A_M Y$ and $A_M = (M^{-1} \otimes K) ((M^{-1} \otimes K) + npI_{np})^{-1}$

- Goal: select $\hat{M} \in \mathcal{M}$ such that $\frac{1}{np} \|\hat{F}_{\hat{M}} - F\|^2$ is minimal
- Expectation of the ideal penalty:

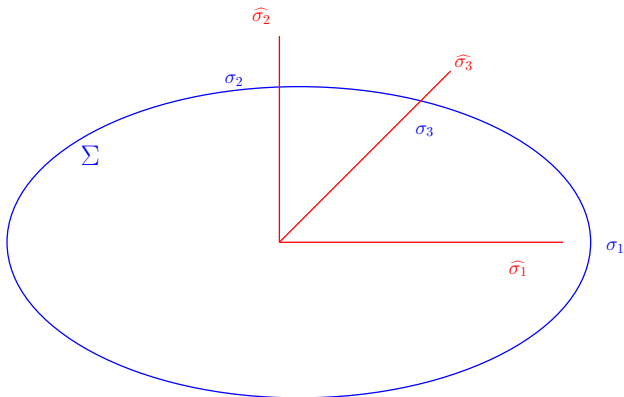
$$\mathbb{E}[\text{pen}_{\text{id}}(M)] = \frac{2}{np} \text{tr}(A_M (\Sigma \otimes I_n))$$

- Problem: How to estimate Σ ?

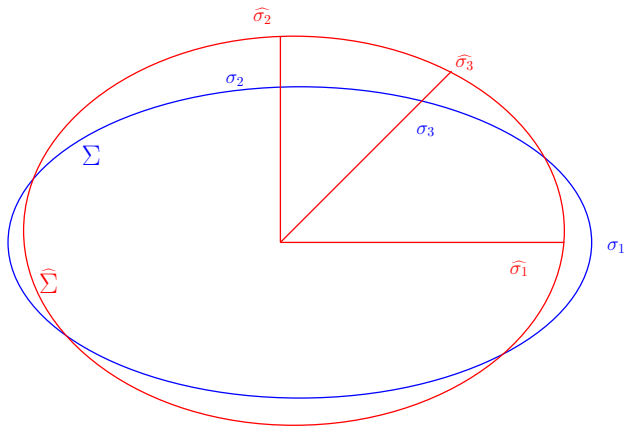
Estimating the covariance matrix: idea ($p = 2$)



Estimating the covariance matrix: idea ($p = 2$)



Estimating the covariance matrix: idea ($p = 2$)



Estimating the covariance matrix: algorithm

- for every $j \in \{1, \dots, p\}$, apply the “minimal penalties” algorithm to the data set $(Y_i^j)_{1 \leq i \leq n}$
⇒ estimator $a(e_j)$ of $\Sigma_{j,j}$
- for every $j \neq \ell \in \{1, \dots, p\}$, apply the “minimal penalties” algorithm to the data set $(Y_i^j + Y_i^\ell)_{1 \leq i \leq n}$
⇒ estimator $a(e_j + e_\ell)$ of $\Sigma_{j,j} + \Sigma_{\ell,\ell} + 2\Sigma_{j,\ell}$

Estimating the covariance matrix: algorithm

- for every $j \in \{1, \dots, p\}$, apply the “minimal penalties” algorithm to the data set $(Y_i^j)_{1 \leq i \leq n}$
 \Rightarrow estimator $a(e_j)$ of $\Sigma_{j,j}$
- for every $j \neq \ell \in \{1, \dots, p\}$, apply the “minimal penalties” algorithm to the data set $(Y_i^j + Y_i^\ell)_{1 \leq i \leq n}$
 \Rightarrow estimator $a(e_j + e_\ell)$ of $\Sigma_{j,j} + \Sigma_{\ell,\ell} + 2\Sigma_{j,\ell}$
- Recover an estimator $\hat{\Sigma}$ of Σ :

$$\hat{\Sigma} = J(a(e_1), \dots, a(e_p), a(e_1 + e_2), \dots, a(e_{p-1} + e_p))$$

where J is the unique linear application $R^{p(p+1)/2} \mapsto \mathcal{S}_p(\mathbb{R})$ such that

$$\Sigma = J(\Sigma_{1,1}, \dots, \Sigma_{p,p}, \Sigma_{1,1} + \Sigma_{2,2} + 2\Sigma_{1,2}, \dots, \Sigma_{p-1,p-1} + \Sigma_{p,p} + 2\Sigma_{p-1,p})$$

Theorem: Estimating the covariance matrix

Theorem (Solnon, A. & Bach, 2011)

If for every $j = 1, \dots, p$, some $\lambda_j > 0$ exists such that $\text{tr}(A_{\lambda_j}) \leq \sqrt{n}$ and

$$\frac{1}{n} \|(I_n - A_{\lambda_j})F^j\|^2 \leq \Sigma_{j,j} \sqrt{\frac{\ln(n)}{n}} \quad \text{where} \quad A_{\lambda_j} = K(K + n\lambda_j I_n)^{-1},$$

Then, with probability $1 - L_5 p^2 n^{-\delta}$, if $n \geq n_0(\delta)$,

$$(1 - \eta)\Sigma \preceq \hat{\Sigma} \preceq (1 + \eta)\Sigma \quad \text{with} \quad \eta := L(2 + \delta)c(\Sigma)^2 p \sqrt{\frac{\ln(n)}{n}}$$

where $c(\Sigma) = \max(\text{Sp}(\Sigma)) / \min(\text{Sp}(\Sigma))$.

⇒ sufficient condition for consistency

Theorem: Oracle inequality

Theorem (Solnon, A. & Bach, 2011)

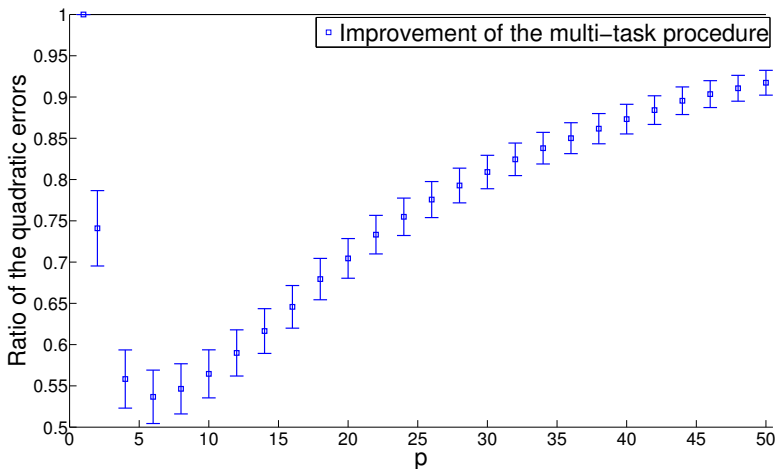
If moreover matrices $M \in \mathcal{M}$ can be diagonalized in the same orthogonal basis, and if

$$\hat{M} \in \arg \min_{M \in \mathcal{M}} \left\{ \frac{1}{np} \left\| \hat{F}_M - Y \right\|^2 + \frac{2}{np} \operatorname{tr} \left(A_M \left(\hat{\Sigma} \otimes I_n \right) \right) \right\},$$

Then, with probability $1 - L_5 p^2 n^{-\delta}$, if $n \geq n_0(\delta)$,

$$\begin{aligned} \frac{1}{np} \left\| \hat{F}_{\hat{M}} - F \right\|^2 &\leq \left(1 + \frac{1}{\ln(n)} \right)^2 \inf_{M \in \mathcal{M}} \left\{ \frac{1}{np} \left\| \hat{F}_M - F \right\|^2 \right\} \\ &\quad + L(2 + \delta)^2 c(\Sigma)^4 \frac{\operatorname{tr}(\Sigma)}{p} \frac{p^3 \ln(n)^3}{n} \end{aligned}$$

Simulations: $n = 100$, $2 \leq p \leq 50$, $1.1 \leq c(\Sigma) \leq 22.5$



Summary

- **Minimal penalties**: efficient for data-driven calibration of **multiplicative constants** in penalties

Summary

- **Minimal penalties**: efficient for data-driven calibration of **multiplicative constants** in penalties
- $\text{pen}_{\text{opt}} / \text{pen}_{\text{min}} \approx 2$ for least-squares estimators

Summary

- **Minimal penalties**: efficient for data-driven calibration of **multiplicative constants** in penalties
- $\text{pen}_{\text{opt}} / \text{pen}_{\text{min}} \approx 2$ for least-squares estimators
- $\text{pen}_{\text{opt}} / \text{pen}_{\text{min}} \in (1; 2]$ for linear estimators

Summary

- **Minimal penalties**: efficient for data-driven calibration of **multiplicative constants** in penalties
- $\text{pen}_{\text{opt}} / \text{pen}_{\text{min}} \approx 2$ for least-squares estimators
- $\text{pen}_{\text{opt}} / \text{pen}_{\text{min}} \in (1; 2]$ for linear estimators
- Can be applied with a data-driven shape of penalty (e.g., if data are heteroscedastic): **V-fold/resampling penalties** (OLS: A. 2008, 2009; Lerasle, 2009)

Minimal penalties: which frameworks?

- Theoretical results:
 - OLS, homoscedastic Gaussian regression (Birgé & Massart, 2007)
 - regressograms, heteroscedastic (A. & Massart, 2009)
 - Least-squares density estimation, i.i.d. (Lerasle, 2009) or mixing (Lerasle, 2010) data
 - Linear estimators, regression (A. & Bach, 2009–2011)
 - **Minimum contrast estimator, regular contrast** (Saumard, 2010)
 - Multitask regression (Solnon, A. & Bach, 2011)
 - ...
- Empirical results:
 - **Change-point detection** (Lebarbier, 2005)
 - Gaussian mixture models (Maugis & Michel, 2008–2010)
 - Unsupervised classification (Baudry, 2009)
 - Computational geometry (Caillerie & Michel, 2009)
 - **Lasso** (Connault, 2011)
 - ... (see Baudry, Maugis & Michel, 2011)