

Model selection and estimator selection for statistical learning

Sylvain Arlot

¹CNRS

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

Scuola Normale Superiore di Pisa, 14–23 February 2011

Outline of the 5 lectures

- 1 Statistical learning
- 2 Model selection for least-squares regression
- 3 **Linear estimator selection for least-squares regression**
- 4 Resampling and model selection
- 5 Cross-validation and model/estimator selection

Part III

Linear estimator selection for least-squares regression

Outline

- 1 Linear estimators for regression
- 2 Minimal penalties for linear estimators selection in regression
- 3 Minimal penalties and calibration in the general framework

Outline

- 1 Linear estimators for regression
- 2 Minimal penalties for linear estimators selection in regression
- 3 Minimal penalties and calibration in the general framework

Linear estimators for regression

- Definition:

$$\hat{F} = AY$$

with A **deterministic** (so, A can depend on X)

Linear estimators for regression

- Definition:

$$\hat{F} = AY$$

with A deterministic (so, A can depend on X)

- Example: **Least-squares estimator**: A orthogonal projection matrix
- Example: **"Regularized" estimator**: A diagonalizable with eigenvalues $1 \geq \lambda_1 \geq \dots \geq \lambda_n \geq 0$

Linear estimators for regression

- Definition:

$$\hat{F} = AY$$

with A deterministic (so, A can depend on X)

- Example: Least-squares estimator: A orthogonal projection matrix
- Example: “Regularized” estimator: A diagonalizable with eigenvalues $1 \geq \lambda_1 \geq \dots \geq \lambda_n \geq 0$
- Assumptions:

$$\|A\| \leq 1 \quad \text{or} \quad \|A\| \leq B$$

$$\text{tr}(A^T A) \leq \text{tr}(A)$$

Kernel ridge regression (1/2)

- Idea: look for an estimator having a small empirical risk **and** a small norm in some functional space \mathcal{F}
- $\mathcal{F} \subset \mathbb{S}$ the Reproducing Kernel Hilbert Space (RKHS) associated with a positive-definite kernel $k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$

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

- **Representer theorem** $\Rightarrow \exists \hat{\alpha} \in \mathbb{R}^n$ such that $\hat{f} = \sum_{i=1}^n \hat{\alpha}_i k(X_i, \cdot)$

Kernel ridge regression (2/2)

$$\hat{f} = \sum_{i=1}^n \hat{\alpha}_i k(X_i, \cdot) \in \arg \min_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2 + \lambda \|f\|_{\mathcal{F}}^2 \right\}$$

- Fixed design: $\hat{F} = (\hat{f}(x_i))_{1 \leq i \leq n} = K\hat{\alpha}$ with $K = (k(x_i, x_j))_{1 \leq i, j \leq n}$ and

$$\hat{\alpha} \in \arg \min_{\alpha \in \mathbb{R}^n} \left\{ \frac{1}{n} \|Y - K\alpha\|^2 + \lambda \alpha^\top K \alpha \right\}$$

$$\Rightarrow \hat{F} = K(K + n\lambda I_n)^{-1} Y$$

Kernel ridge regression (2/2)

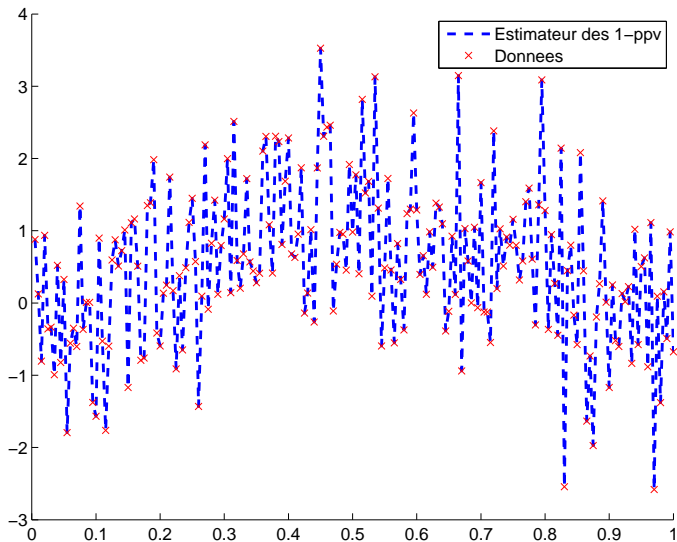
- Fixed design: $\widehat{F} = (\widehat{f}(x_i))_{1 \leq i \leq n} = K\widehat{\alpha}$ with $K = (k(x_i, x_j))_{1 \leq i, j \leq n}$ and

$$\widehat{\alpha} \in \arg \min_{\alpha \in \mathbb{R}^n} \left\{ \frac{1}{n} \|Y - K\alpha\|^2 + \lambda \alpha^\top K \alpha \right\}$$

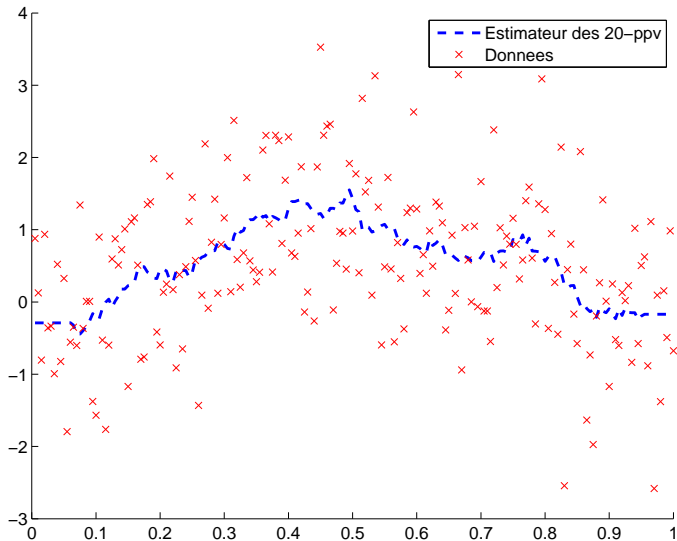
$$\Rightarrow \widehat{F} = K(K + n\lambda I_n)^{-1} Y$$

- If $(e_j)_{1 \leq j \leq n}$ orthonormal basis of eigenvectors of K (with eigenvalues $(\kappa_j)_{1 \leq j \leq n}$):

$$\widehat{F} e_j = \frac{\kappa_j}{\kappa_j + n\lambda} e_j$$

Nearest-neighbour estimator: $\hat{F} = Y$ 

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



k -nearest-neighbours estimator

For all $i \in \{1, \dots, n\}$,

$$\begin{aligned}\hat{F}_i &= \frac{1}{k} \sum_{x_j \in \{k\text{-NN of } x_i\}} Y_j \\ &= \sum_{1 \leq j \leq n} \left(\frac{1}{k} \mathbb{1}_{x_j \in \{k\text{-NN of } x_i\}} Y_j \right)\end{aligned}$$

k -nearest-neighbours estimator

For all $i \in \{1, \dots, n\}$,

$$\begin{aligned}\widehat{F}_i &= \frac{1}{k} \sum_{x_j \in \{k\text{-NN of } x_i\}} Y_j \\ &= \sum_{1 \leq j \leq n} \left(\frac{1}{k} \mathbb{1}_{x_j \in \{k\text{-NN of } x_i\}} Y_j \right)\end{aligned}$$

$$\Rightarrow \widehat{F} = A(k)Y$$

$$\text{with } A(k)_{i,j} = \frac{1}{k} \mathbb{1}_{x_j \in \{k\text{-NN of } x_i\}}$$

Nadaraya-Watson estimator

Let $K : \mathcal{X} \times \mathcal{X} \mapsto [0, +\infty)$ be some “kernel”.

$$\forall i, j \in \{1, \dots, n\}, \quad A(K)_{i,j} = \frac{K(x_i, x_j)}{\sum_{1 \leq \ell \leq n} K(x_i, x_\ell)}$$

$$\hat{F} = A(K)Y$$

Nadaraya-Watson estimator

Let $K : \mathcal{X} \times \mathcal{X} \mapsto [0, +\infty)$ be some “kernel”.

$$\forall i, j \in \{1, \dots, n\}, \quad A(K)_{i,j} = \frac{K(x_i, x_j)}{\sum_{1 \leq \ell \leq n} K(x_i, x_\ell)}$$

$$\hat{F} = A(K)Y$$

- Typically, $K(x, y) = g(d(x, y)/h)$ for some **distance** d over \mathcal{X} , a **bandwidth** $h > 0$ and a non-increasing function $g : [0, +\infty) \mapsto [0, +\infty)$.

Nadaraya-Watson estimator

Let $K : \mathcal{X} \times \mathcal{X} \mapsto [0, +\infty)$ be some “kernel”.

$$\forall i, j \in \{1, \dots, n\}, \quad A(K)_{i,j} = \frac{K(x_i, x_j)}{\sum_{1 \leq \ell \leq n} K(x_i, x_\ell)}$$

$$\hat{F} = A(K)Y$$

- Typically, $K(x, y) = g(d(x, y)/h)$ for some **distance** d over \mathcal{X} , a **bandwidth** $h > 0$ and a non-increasing function $g : [0, +\infty) \mapsto [0, +\infty)$.
- **Gaussian kernel** $g(t) = \exp(-t^2)$.
- **Window kernel** $g(t) = \mathbb{1}_{t \in [0,1]}$

Outline

- 1 Linear estimators for regression
- 2 Minimal penalties for linear estimators selection in regression
- 3 Minimal penalties and calibration in the general framework

Risk, empirical risk, ideal penalty

$$Y = F + \varepsilon \quad \text{with} \quad \mathbb{E}[\varepsilon_i^2] = \sigma^2$$

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

Risk, empirical risk, ideal penalty

$$Y = F + \varepsilon \quad \text{with} \quad \mathbb{E}[\varepsilon_i^2] = \sigma^2$$

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

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

Risk, empirical risk, ideal penalty

$$Y = F + \varepsilon \quad \text{with} \quad \mathbb{E}[\varepsilon_i^2] = \sigma^2$$

$$\widehat{F}_m = A_m Y$$

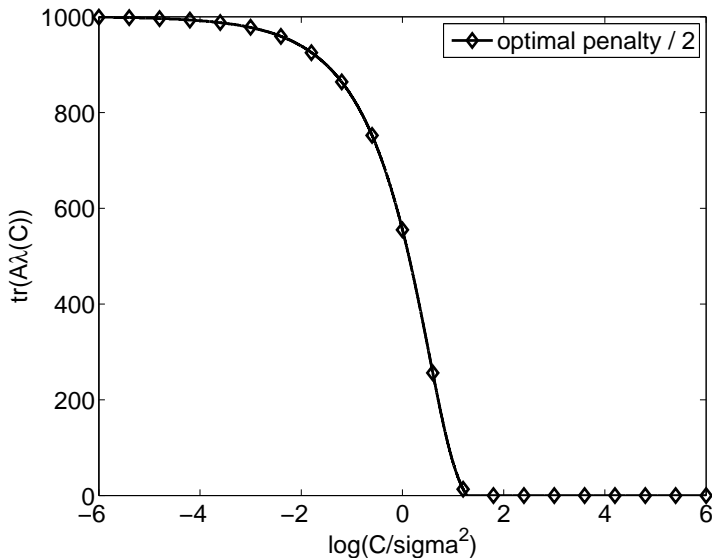
$$\mathbb{E} \left[\frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 \right] = \frac{1}{n} \|(A_m - I)F\|^2 + \frac{\sigma^2 \text{tr}(A_m^\top A_m)}{n}$$

$$\text{pen}_{\text{id}}(m) = \frac{2}{n} \langle A_m \varepsilon, \varepsilon \rangle + \frac{2}{n} \langle (A_m - I_n)F, \varepsilon \rangle$$

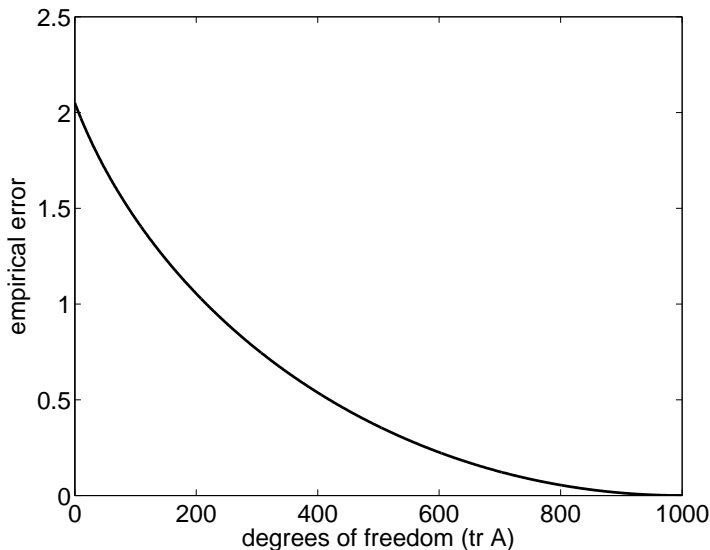
$$\mathbb{E}[\text{pen}_{\text{id}}(m)] = \frac{2\sigma^2 \text{tr}(A_m)}{n} \Rightarrow C_L \text{ (Mallows, 1973)}$$

generalized degrees of freedom: $\text{tr}(A_m)$

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



The minimal penalty is not proportional to $\text{tr}(A_m)$



What is the minimal penalty?

$$\begin{aligned}
 p_2(m) &= \frac{1}{n} \|Y - F_m\|^2 - \frac{1}{n} \|Y - \hat{F}_m\|^2 \\
 &= \frac{2}{n} \langle \varepsilon, A_m \varepsilon \rangle - \frac{1}{n} \|A_m \varepsilon\|^2 - \frac{2}{n} \langle \varepsilon, A_m^\top (I_n - A_m) F \rangle
 \end{aligned}$$

What is the minimal penalty?

$$\begin{aligned} p_2(m) &= \frac{1}{n} \|Y - F_m\|^2 - \frac{1}{n} \|Y - \hat{F}_m\|^2 \\ &= \frac{2}{n} \langle \varepsilon, A_m \varepsilon \rangle - \frac{1}{n} \|A_m \varepsilon\|^2 - \frac{2}{n} \langle \varepsilon, A_m^\top (I_n - A_m) F \rangle \end{aligned}$$

so that

$$\text{pen}_{\min}(m) = \mathbb{E}[p_2(m)] = \frac{\sigma^2 (2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m))}{n}$$

What is the minimal penalty?

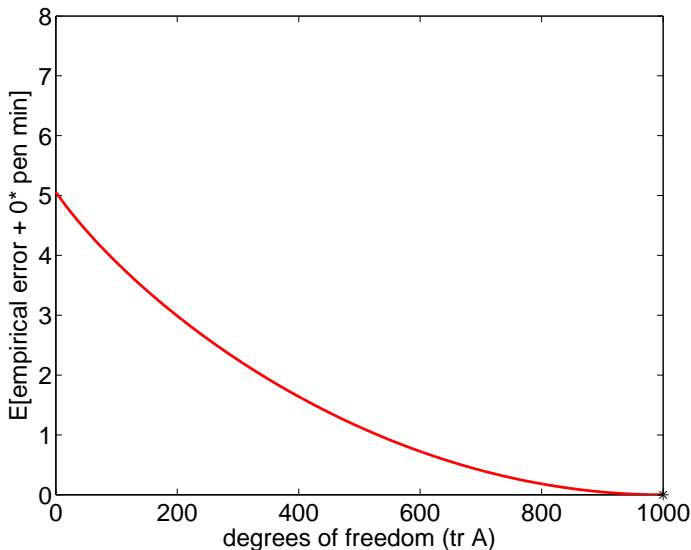
$$\begin{aligned}
 p_2(m) &= \frac{1}{n} \|Y - F_m\|^2 - \frac{1}{n} \left\| Y - \hat{F}_m \right\|^2 \\
 &= \frac{2}{n} \langle \varepsilon, A_m \varepsilon \rangle - \frac{1}{n} \|A_m \varepsilon\|^2 - \frac{2}{n} \left\langle \varepsilon, A_m^\top (I_n - A_m) F \right\rangle
 \end{aligned}$$

so that

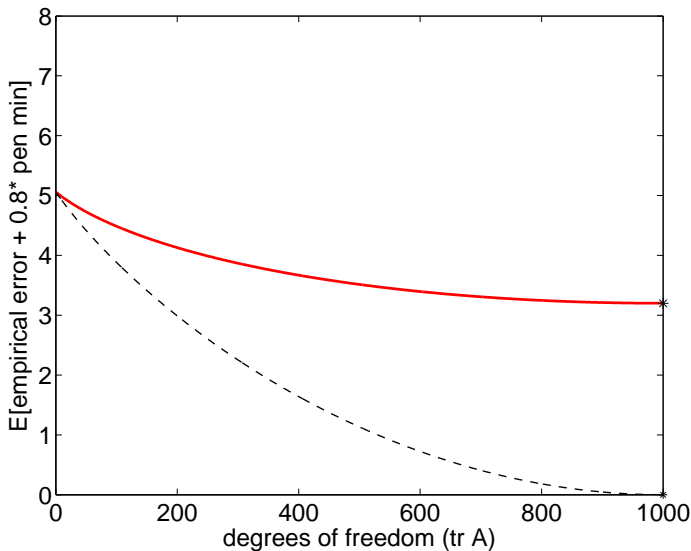
$$\text{pen}_{\min}(m) = \mathbb{E}[p_2(m)] = \frac{\sigma^2 (2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m))}{n}$$

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

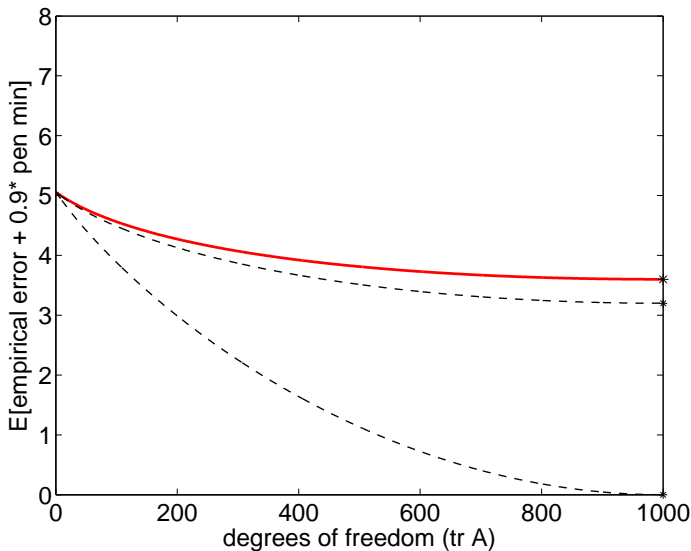
$$\mathbb{E}[n^{-1} \|Y - \hat{F}_m\|^2] + 0 \times \sigma^2 (2 \operatorname{tr}(A_m) - \operatorname{tr}(A_m^\top A_m)) n^{-1}$$



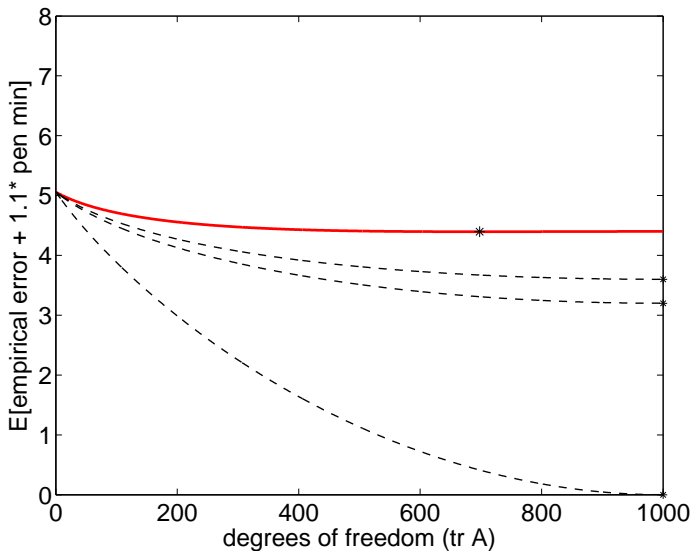
$$\mathbb{E}[n^{-1} \|Y - \widehat{F}_m\|^2] + 0.8 \times \sigma^2 (2 \operatorname{tr}(A_m) - \operatorname{tr}(A_m^\top A_m)) n^{-1}$$



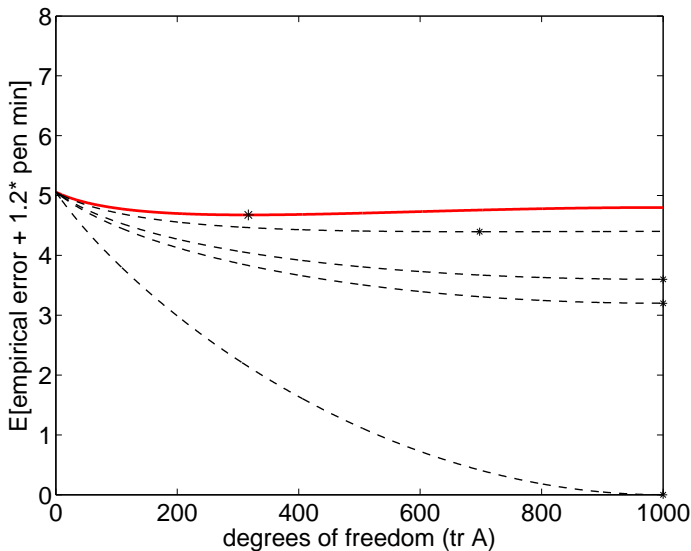
$$\mathbb{E}[n^{-1} \|Y - \widehat{F}_m\|^2] + 0.9 \times \sigma^2 (2 \operatorname{tr}(A_m) - \operatorname{tr}(A_m^\top A_m)) n^{-1}$$



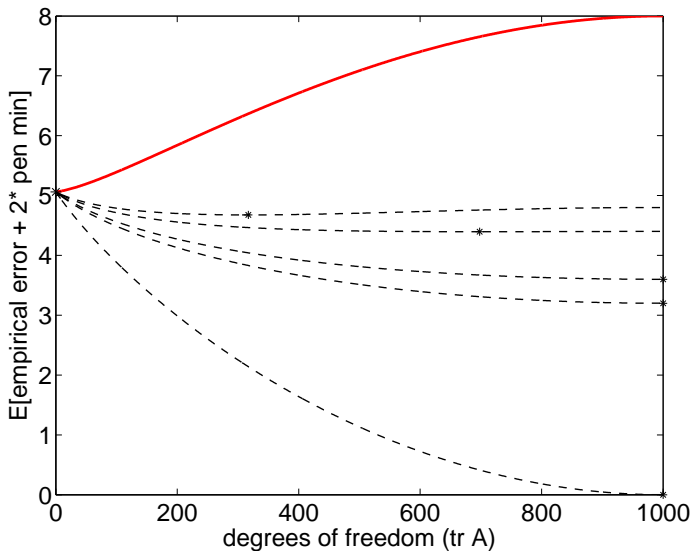
$$\mathbb{E}[n^{-1} \|Y - \widehat{F}_m\|^2] + 1.1 \times \sigma^2 (2 \operatorname{tr}(A_m) - \operatorname{tr}(A_m^\top A_m)) n^{-1}$$



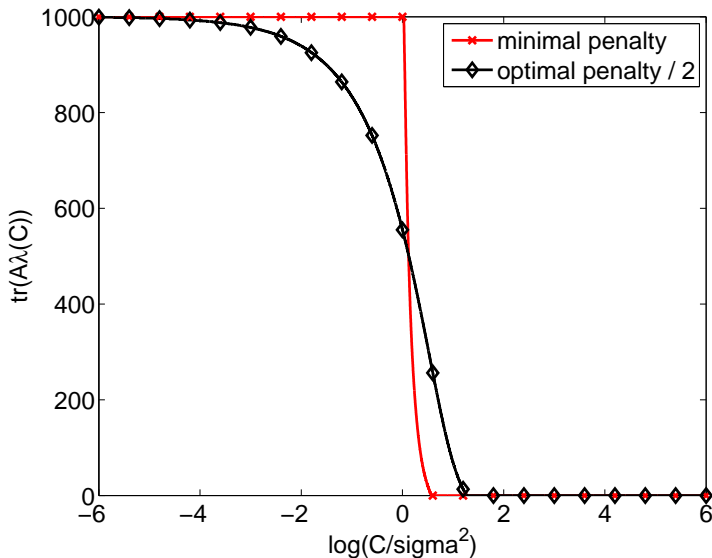
$$\mathbb{E}[n^{-1} \|Y - \widehat{F}_m\|^2] + 1.2 \times \sigma^2 (2 \operatorname{tr}(A_m) - \operatorname{tr}(A_m^\top A_m)) n^{-1}$$



$$\mathbb{E}[n^{-1} \|Y - \widehat{F}_m\|^2] + 2 \times \sigma^2 (2 \operatorname{tr}(A_m) - \operatorname{tr}(A_m^\top A_m)) n^{-1}$$



“Dimension” jump (ridge regression)



Penalty calibration algorithm (A. & Bach 2009)

- 1 for every $C > 0$, compute

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

- 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

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

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

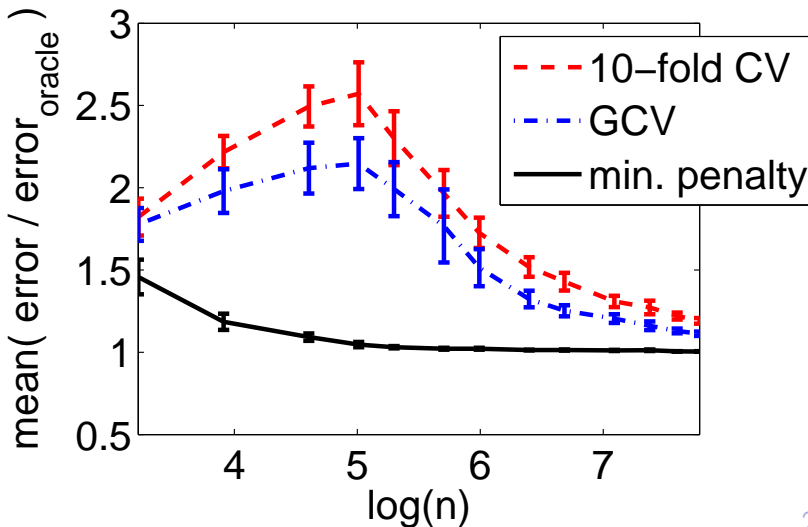
$$\forall i, j \in \{1, \dots, n\}, \quad A_{i,j} \in \left\{ 0, \frac{1}{k} \right\}$$

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

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

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

Simulation study (ridge regression, choice of λ)



Theorem (1): dimension jump (A. & Bach, 2009)

- polynomial complexity: $\text{Card}(\mathcal{M}_n) \leq C_{\mathcal{M}} n^{\alpha}$
- homoscedastic Gaussian noise, fixed design
- $\exists m_1, m_2 \in \mathcal{M}_n$ s.t. $\text{tr}(A_{m_1}) \geq n/2$, $\text{tr}(A_{m_2}) \leq \sqrt{n}$ and $\forall i \in \{1, 2\}$, $n^{-1} \|(I - A_{m_i})F\|^2 \leq \sigma^2 \sqrt{\ln(n)/n}$

Theorem (Minimal penalty; A. & Bach 2009)

With probability at least $1 - 8C_{\mathcal{M}}n^{-2}$, if $n \geq n_0(\alpha)$,

$$\forall C < \left(1 - L_{\alpha} \sqrt{\frac{\ln(n)}{n}}\right) \sigma^2, \quad \text{tr}(A_{\hat{m}_{\min}(C)}) \geq \frac{n}{3}$$

$$\forall C > \left(1 + L_{\alpha} \frac{\sqrt{\ln(n)}}{n^{1/4}}\right) \sigma^2, \quad \text{tr}(A_{\hat{m}_{\min}(C)}) \leq n^{3/4}$$

Theorem (2): oracle inequality (A. & Bach, 2009)

Additional assumption:

$$\exists \kappa \geq 1, \forall m \in \mathcal{M}_n, \quad \text{tr}(A_m)\sigma^2 \leq \kappa \mathbb{E} \left[\left\| F - \widehat{F}_m \right\|^2 \right]$$

Theorem (Oracle inequality; A. & Bach 2009)

Then, with probability at least $1 - 8C_{\mathcal{M}}n^{-2}$, si $n \geq n_0(\alpha)$,

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

$$\frac{1}{n} \left\| F - \widehat{F}_{\widehat{m}} \right\|^2 \leq \left(1 + \frac{40\kappa}{\ln(n)} \right) \inf_{m \in \mathcal{M}_n} \left\{ \frac{1}{n} \left\| F - \widehat{F}_m \right\|^2 \right\}$$

$$+ \frac{36(\kappa + \alpha + 2) \ln(n)\sigma^2}{n}$$

Outline

- 1 Linear estimators for regression
- 2 Minimal penalties for linear estimators selection in regression
- 3 Minimal penalties and calibration in the general framework

Minimal penalties: general case

$$\hat{m}_{\min}(C) \in \arg \min_{m \in \mathcal{M}_n} \{P_n \gamma(\hat{s}_m) + C \text{pen}_0(m)\}$$

- **risk burst** for $C < C_{\min}^*$
- **oracle inequality** for $C > C_{\min}^*$

Minimal penalties: general case

$$\hat{m}_{\min}(C) \in \arg \min_{m \in \mathcal{M}_n} \{P_n \gamma(\hat{s}_m) + C \text{pen}_0(m)\}$$

- risk burst for $C < C_{\min}^*$
- oracle inequality for $C > C_{\min}^*$
- **jump of some complexity** $\mathcal{C}_{\hat{m}(C)}$ around $C = C_{\min}^*$
- $\text{pen}_{\min} = C_{\min}^* \text{pen}_0$ with **pen₀ known (or can be estimated)**

Minimal penalties: general case

$$\hat{m}_{\min}(C) \in \arg \min_{m \in \mathcal{M}_n} \{ P_n \gamma(\hat{s}_m) + C \text{pen}_0(m) \}$$

- risk burst for $C < C_{\min}^*$
- oracle inequality for $C > C_{\min}^*$
- jump of some complexity $\mathcal{C}_{\hat{m}(C)}$ around $C = C_{\min}^*$
- $\text{pen}_{\min} = C_{\min}^* \text{pen}_0$ with pen_0 known (or can be estimated)
- **optimal oracle inequality** when $\text{pen} = \text{pen}_{\text{opt}}$
- **pen_{opt} known (that can be estimated) up to C^***

Minimal penalties: general case

$$\hat{m}_{\min}(C) \in \arg \min_{m \in \mathcal{M}_n} \{ P_n \gamma(\hat{s}_m) + C \text{pen}_0(m) \}$$

- risk burst for $C < C_{\min}^*$
- oracle inequality for $C > C_{\min}^*$
- jump of some complexity $\mathcal{C}_{\hat{m}(C)}$ around $C = C_{\min}^*$
- $\text{pen}_{\min} = C_{\min}^* \text{pen}_0$ with pen_0 known (or can be estimated)
- optimal oracle inequality when $\text{pen} = \text{pen}_{\text{opt}}$
- pen_{opt} known (that can be estimated) up to C^*
- **known relationship between C_{\min}^* and C^***

Calibration with minimal penalties

Inputs: $(\text{pen}_0(m))_{m \in \mathcal{M}_n}$ $(\text{pen}_1(m))_{m \in \mathcal{M}_n}$ $(C_m)_{m \in \mathcal{M}_n}$ f

Assumptions: $\text{pen}_0 = \frac{1}{C_{\min}^*} \text{pen}_{\min}$ $\text{pen}_1 = \frac{1}{f(C_{\min}^*)} \text{pen}_{\text{opt}}$

- 1 for every $C > 0$, compute

$$\hat{m}_{\min}(C) \in \arg \min_{m \in \mathcal{M}_n} \{ P_n \gamma(\hat{s}_m) + C \text{pen}_0(m) \}$$

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

- 3 select

$$\hat{m} \in \arg \min_{m \in \mathcal{M}_n} \left\{ P_n \gamma(\hat{s}_m) + f(\hat{C}_{\min}) \text{pen}_1(m) \right\}$$

Known mathematical results

- **Concentration of p_2** : minimum contrast estimators with bounded contrast (Boucheron & Massart, 2010)
- **Burst of the risk if $C < C_{\min}^*$** :
 - OLS, Gaussian noise, exponential complexity of \mathcal{M}_n (Birgé & Massart, 2007)
 - OLS, Gaussian noise, medium range complexity of \mathcal{M}_n and $s^* = 0$ (Birgé & Massart, 2007)
 - Density estimation, Dantzig estimator, special case (Bertin, Le Pennec & Rivoirard, 2009)
 - Estimation of intensity of a Poisson process by thresholding (Reynaud-Bouret & Rivoirard, 2009), with $C^*/C_{\min}^* \in [1, 12]$
- **Burst of the dimension if $C < C_{\min}^*$** : OLS, Gaussian noise, multiplicative penalties, polynomial or exponential complexity of \mathcal{M}_n and $s^* = 0$ (Baraud, Giraud & Huet, 2009)