

# Model selection and estimator selection for statistical learning

Sylvain Arlot

<sup>1</sup>CNRS

<sup>2</sup>École 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 II

# Model selection for least-squares regression

# Outline

- 1 An oracle inequality for model selection
- 2 The penalty calibration problem
- 3 Slope heuristics in homoscedastic regression
- 4 The slope heuristics
- 5 Practical issues
- 6 Conclusion

# Outline

- 1 An oracle inequality for model selection
- 2 The penalty calibration problem
- 3 Slope heuristics in homoscedastic regression
- 4 The slope heuristics
- 5 Practical issues
- 6 Conclusion

# A key lemma

## Lemma

Let  $\text{pen} : \mathcal{M}_n \mapsto \mathbb{R}$  some penalty (possibly data-dependent).  
On the event  $\Omega$  on which for every  $m, m' \in \mathcal{M}_n$ ,

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

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

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

# Oracle inequality for Gaussian regression (1)

Assumptions:

- Fixed design regression, least-squares contrast
- **Gaussian homoscedastic noise**:  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$
- Model collection of **polynomial complexity**:  $\text{Card}(\mathcal{M}_n) \leq Cn^\alpha$
- For all  $m \in \mathcal{M}_n$ ,  $\hat{F}_m = A_m Y = \Pi_{S_m} Y$  (least-squares estimator)
- Penalty

$$\text{pen}(m) = \frac{K\sigma^2 \dim(S_m)}{n} \quad \text{with } K > 1$$

# Oracle inequality for Gaussian regression (2)

$$\begin{aligned}
 & -B(m) \leq \text{pen}(m) - \text{pen}_{\text{id}}(m, D_n) \leq A(m) \\
 \Rightarrow & \frac{1}{n} \left\| \widehat{F}_{\widehat{m}} - F \right\|^2 - B(\widehat{m}) \leq \inf_{m \in \mathcal{M}_n} \left\{ \frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 + A(m) \right\}
 \end{aligned}$$

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

First term has expectation  $\frac{2\sigma^2 \dim(S_m)}{n}$ , the second term is centered.



# Oracle inequality for Gaussian regression (3)

Two **Gaussian concentration** results (see Massart 2007):

## Proposition

Let  $\xi$  be some standard Gaussian vector in  $\mathbb{R}^n$ ,  $\alpha \in \mathbb{R}^n$ ,  $M \in \mathcal{M}_n(\mathbb{R})$ . Then, for every  $x \geq 0$ ,

$$\mathbb{P} \left( |\langle \xi, \alpha \rangle| \leq \sqrt{2x} \|\alpha\|_2 \right) \geq 1 - 2e^{-x}$$

$$\mathbb{P} \left( |\langle \xi, M\xi \rangle - \text{tr}(M)| \leq 2\sqrt{x \text{tr}(M^T M)} + 2\|M\| x \right) \geq 1 - 2e^{-x}$$

# Oracle inequality for Gaussian regression (4)

Sketch of the proof:

- For all  $m \in \mathcal{M}_n$ ,  
**concentrate**  $\langle A_m \varepsilon, \varepsilon \rangle$  around  $\sigma^2 \dim(S_m)$   
and  $\langle (A_m - I_n)F, \varepsilon \rangle$  around 0
- Apply the **Lemma** on the intersection of these  $\text{Card}(\mathcal{M}_n)$  events
- Control the **remainder terms**

# Oracle inequality for Gaussian regression (5)

## Theorem (Birgé & Massart 2007)

For every  $x \geq 0$ , *with probability at least  $1 - 4 \text{Card}(\mathcal{M}_n)e^{-x}$* , for every

$$\hat{m} \in \arg \min_{m \in \mathcal{M}_n} \left\{ \frac{1}{n} \left\| Y - \hat{F}_m \right\|^2 + \frac{K\sigma^2 \dim(S_m)}{n} \right\},$$

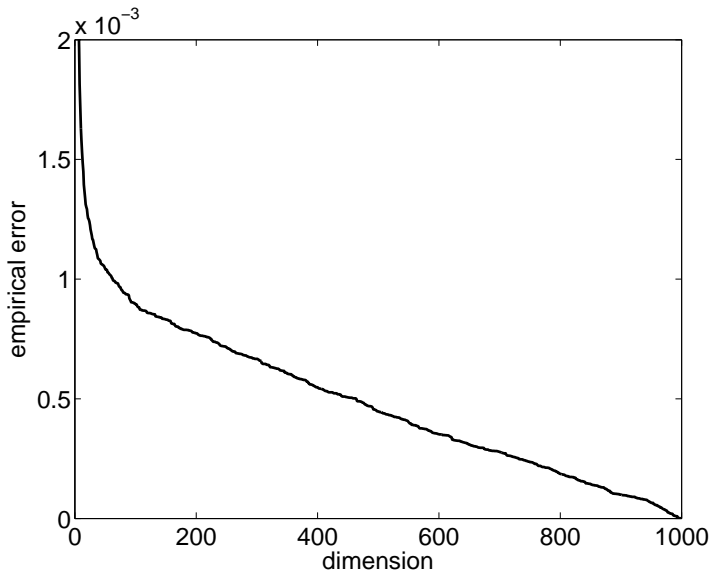
we get the oracle inequality  $\forall \delta > 0$ ,

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

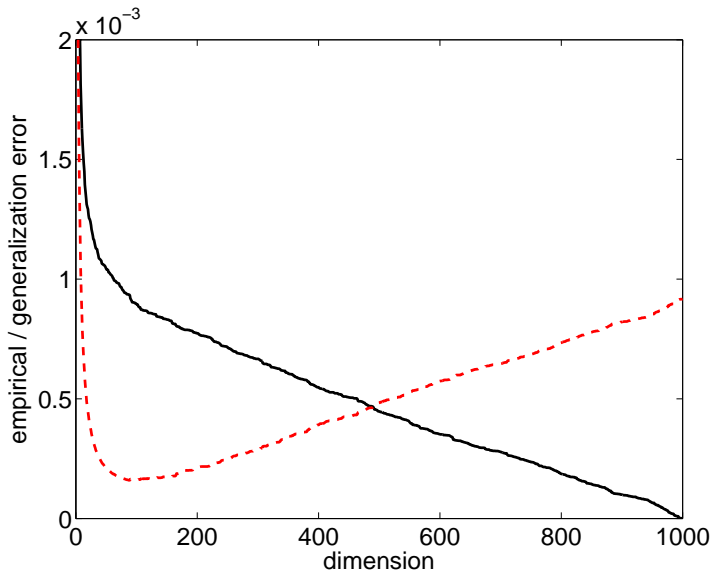
# Outline

- 1 An oracle inequality for model selection
- 2 The penalty calibration problem**
- 3 Slope heuristics in homoscedastic regression
- 4 The slope heuristics
- 5 Practical issues
- 6 Conclusion

# Motivation (1): L-curve and elbow heuristics?



# Motivation (1): L-curve and elbow heuristics?



# Motivation (2): what if $K \leq 1$ ?

## Theorem (Birgé & Massart 2007)

If  $K > 1$ , for every  $x \geq 0$ , with probability  $1 - 4 \text{Card}(\mathcal{M}_n)e^{-x}$ , for every

$$\hat{m} \in \arg \min_{m \in \mathcal{M}_n} \left\{ \frac{1}{n} \|Y - \hat{F}_m\|^2 + \frac{K\sigma^2 \dim(S_m)}{n} \right\},$$

we get the concentration inequality  $\forall \delta > 0$ ,

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

## Motivation (3): penalty calibration

- $C_p$  and  $C_L$  (Mallows, 1973):

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

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

- Penalties proportional to  $D_m$  with **the optimal multiplying factor unknown**: change-point detection (Birgé & Massart, 2001; Lebarbier, 2005), mixture models (Maugis & Michel, 2008), and so on
- Rademacher penalties

$$\text{pen}(m) = 2 \times \mathbb{E} \left[ \sup_{t \in S_m} \left\{ \frac{1}{n} \sum_{i=1}^n (\varepsilon_i \gamma(t; \xi_i)) \mid D_n \right\} \right]$$

- ...



# Naive estimator of $\sigma^2$

Example: homoscedastic regression on a fixed design  
Computation of the empirical risk

# Naive estimator of $\sigma^2$

Example: homoscedastic regression on a fixed design

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

Naive estimator of  $\sigma^2$ :

$$\widehat{\sigma}_m^2 := \frac{1}{n - D_m} \left\| Y - \widehat{F}_m \right\|^2$$

**Bias** of this estimator:

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

# Naive estimator of $\sigma^2$

Example: homoscedastic regression on a fixed design

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

Naive estimator of  $\sigma^2$ :

$$\widehat{\sigma}_m^2 := \frac{1}{n - D_m} \left\| Y - \widehat{F}_m \right\|^2$$

Bias of this estimator:

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

⇒ Using it inside the penalty  $2\sigma^2 D_m/n$ ?

# Naive estimator of $\sigma^2$

Naive estimator of  $\sigma^2$ :

$$\hat{\sigma}_m^2 := \frac{1}{n - D_m} \left\| Y - \hat{F}_m \right\|^2$$

Bias of this estimator:

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

⇒ Using it inside the penalty  $2\sigma^2 D_m/n$ ?

First idea:

$$\text{crit}(m) = \frac{1}{n} \left\| Y - \hat{F}_m \right\|^2 + \frac{2\hat{\sigma}_{m_0}^2 D_m}{n}$$

**Drawbacks:** we have to know/choose  $m_0$ , overpenalization by an unknown factor

# FPE (Akaike, 1970) and GCV (Craven & Wahba, 1979)

$$\text{crit}_{\text{FPE}}(m) = \frac{1}{n} \left\| Y - \hat{F}_m \right\|^2 + \frac{2\hat{\sigma}_m^2 D_m}{n} = \frac{1}{n} \left\| Y - \hat{F}_m \right\|^2 \left( 1 + \frac{2D_m}{n - D_m} \right)$$

(Akaike, 1970; see also Baraud, Giraud & Huet, 2009)

# FPE (Akaike, 1970) and GCV (Craven & Wahba, 1979)

$$\text{crit}_{\text{FPE}}(m) = \frac{1}{n} \left\| Y - \hat{F}_m \right\|^2 + \frac{2\hat{\sigma}_m^2 D_m}{n} = \frac{1}{n} \left\| Y - \hat{F}_m \right\|^2 \left( 1 + \frac{2D_m}{n - D_m} \right)$$

(Akaike, 1970; see also Baraud, Giraud & Huet, 2009)

**Generalized cross-validation** (GCV, Craven & Wahba, 1979)

$$\text{crit}_{\text{GCV}}(m) = \frac{1}{n} \frac{\left\| Y - \hat{F}_m \right\|^2}{\left( 1 - \frac{D_m}{n} \right)^2}$$

# FPE (Akaike, 1970) and GCV (Craven & Wahba, 1979)

$$\text{crit}_{\text{FPE}}(m) = \frac{1}{n} \left\| Y - \hat{F}_m \right\|^2 + \frac{2\hat{\sigma}_m^2 D_m}{n} = \frac{1}{n} \left\| Y - \hat{F}_m \right\|^2 \left( 1 + \frac{2D_m}{n - D_m} \right)$$

(Akaike, 1970; see also Baraud, Giraud & Huet, 2009)

Generalized cross-validation (GCV, Craven & Wahba, 1979)

$$\text{crit}_{\text{GCV}}(m) = \frac{1}{n} \frac{\left\| Y - \hat{F}_m \right\|^2}{\left( 1 - \frac{D_m}{n} \right)^2}$$

If  $D_m \ll n$ ,

$$\text{crit}_{\text{GCV}}(m) \approx \frac{1}{n} \left\| Y - \hat{F}_m \right\|^2 \frac{n + D_m}{n - D_m} = \frac{1}{n} \left\| Y - \hat{F}_m \right\|^2 \left( 1 + \frac{2D_m}{n - D_m} \right)$$

# FPE (Akaike, 1970) and GCV (Craven & Wahba, 1979)

$$\text{crit}_{\text{FPE}}(m) = \frac{1}{n} \left\| Y - \hat{F}_m \right\|^2 + \frac{2\hat{\sigma}_m^2 D_m}{n} = \frac{1}{n} \left\| Y - \hat{F}_m \right\|^2 \left( 1 + \frac{2D_m}{n - D_m} \right)$$

(Akaike, 1970; see also Baraud, Giraud & Huet, 2009)

Generalized cross-validation (GCV, Craven & Wahba, 1979)

$$\text{crit}_{\text{GCV}}(m) = \frac{1}{n} \frac{\left\| Y - \hat{F}_m \right\|^2}{\left( 1 - \frac{D_m}{n} \right)^2}$$

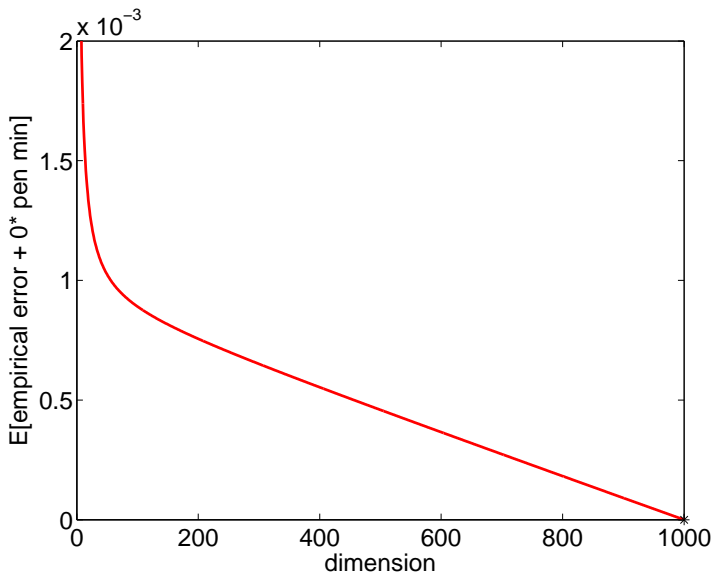
If  $D_m \ll n$ ,

$$\text{crit}_{\text{GCV}}(m) \approx \frac{1}{n} \left\| Y - \hat{F}_m \right\|^2 \frac{n + D_m}{n - D_m} = \frac{1}{n} \left\| Y - \hat{F}_m \right\|^2 \left( 1 + \frac{2D_m}{n - D_m} \right)$$

**Drawbacks:** for the largest models  $\left\| Y - \hat{F}_m \right\|^2 \approx 0$



$$\mathbb{E}[\|Y - \hat{F}_m\|^2] = n\sigma^2 + \|F - F_m\|^2 - \sigma^2 \times D_m$$



# Outline

- 1 An oracle inequality for model selection
- 2 The penalty calibration problem
- 3 Slope heuristics in homoscedastic regression**
- 4 The slope heuristics
- 5 Practical issues
- 6 Conclusion

# Minimal penalty: heuristics

For all  $C > 0$ ,

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

$\Rightarrow \exists C_{\min}$  s.t.:  
for  $C < C_{\min}$ ,  $\hat{F}_{\hat{m}(C)}$  overfits  
for  $C > C_{\min}$ , oracle-inequality for  $\hat{F}_{\hat{m}(C)}$ ?

# Minimal penalty: heuristics

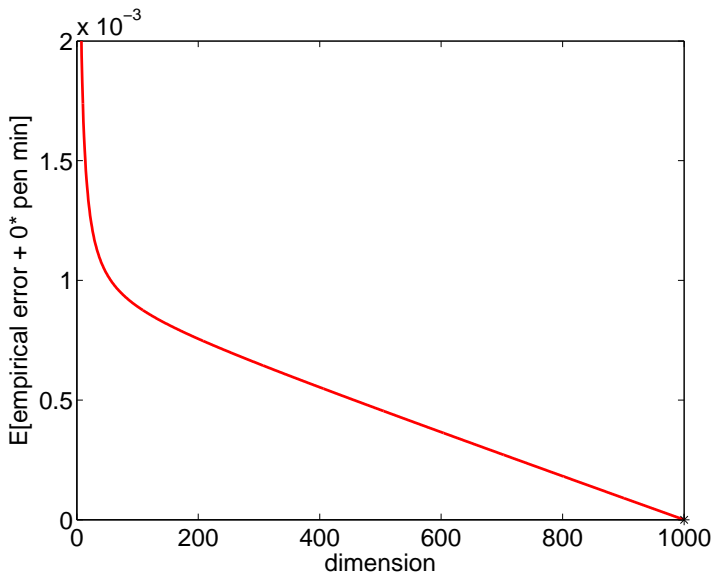
For all  $C > 0$ ,

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

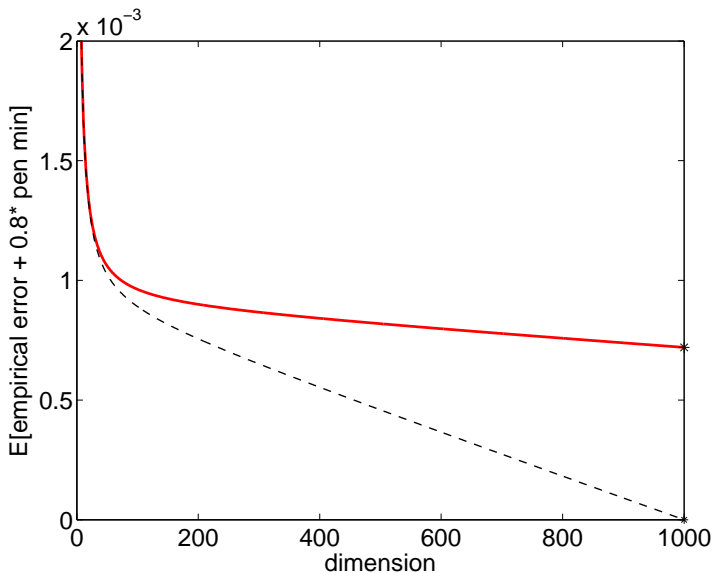
$\Rightarrow \exists C_{\min}$  s.t.:  
 for  $C < C_{\min}$ ,  $\hat{F}_{\hat{m}(C)}$  overfits  
 for  $C > C_{\min}$ , oracle-inequality for  $\hat{F}_{\hat{m}(C)}$ ?

$$\begin{aligned} m^*(C) &\in \arg \min_{m \in \mathcal{M}_n} \left\{ \mathbb{E} \left[ \frac{1}{n} \|Y - \hat{F}_m\|^2 + \frac{CD_m}{n} \right] \right\} \\ &= \arg \min_{m \in \mathcal{M}_n} \left\{ \frac{1}{n} \|F - F_m\|^2 + (C - \sigma^2) \frac{D_m}{n} \right\} \end{aligned}$$

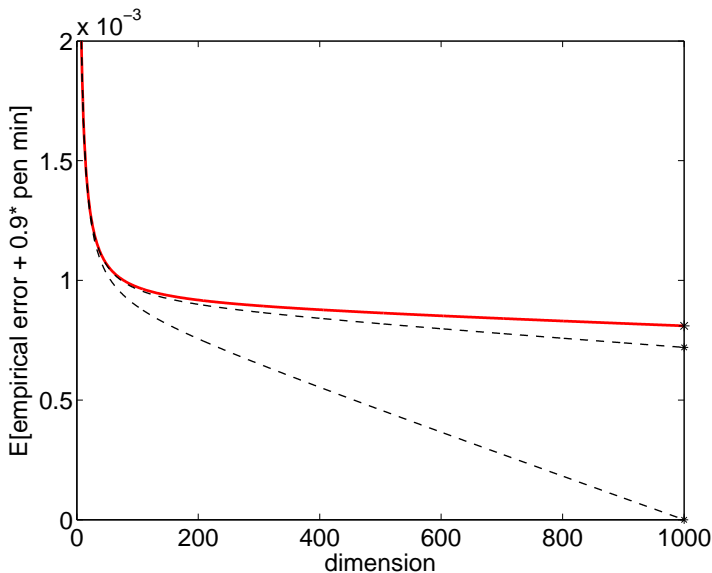
$$\mathbb{E}[n^{-1} \|Y - \widehat{F}_m\|^2] + 0 \times \sigma^2 D_m n^{-1}$$



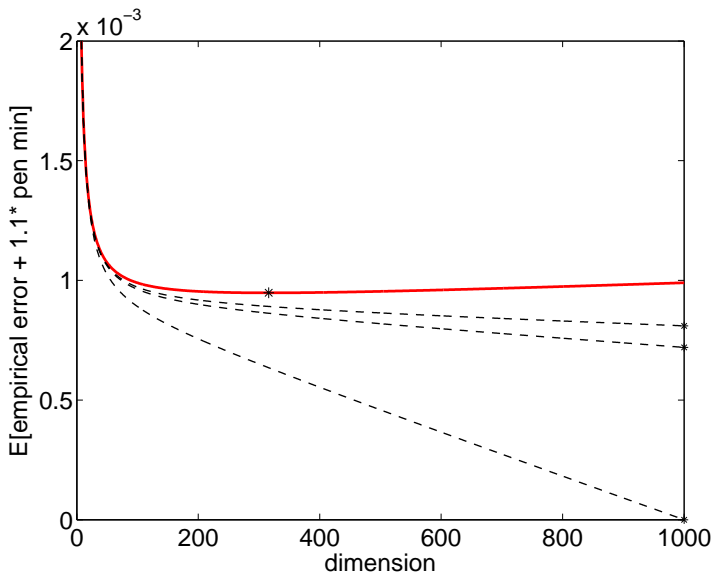
$$\mathbb{E}[n^{-1} \|Y - \hat{F}_m\|^2] + 0.8 \times \sigma^2 D_m n^{-1}$$



$$\mathbb{E}[n^{-1} \|Y - \hat{F}_m\|^2] + 0.9 \times \sigma^2 D_m n^{-1}$$

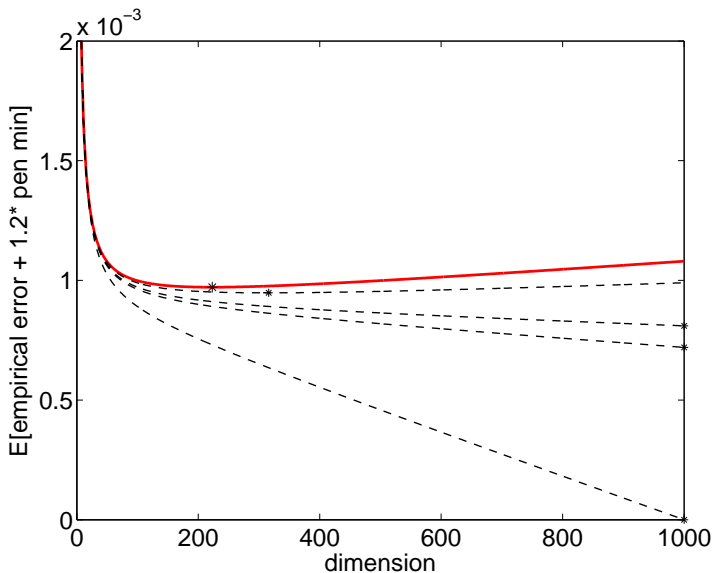


$$\mathbb{E}[n^{-1} \|Y - \widehat{F}_m\|^2] + 1.1 \times \sigma^2 D_m n^{-1}$$

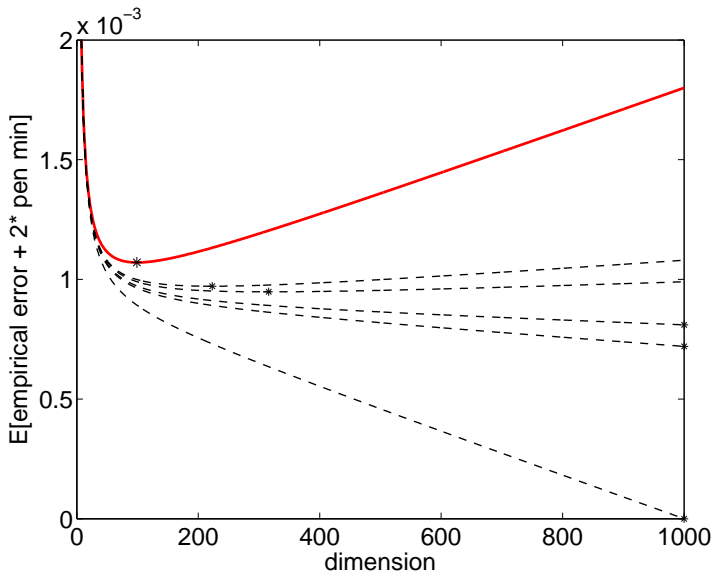




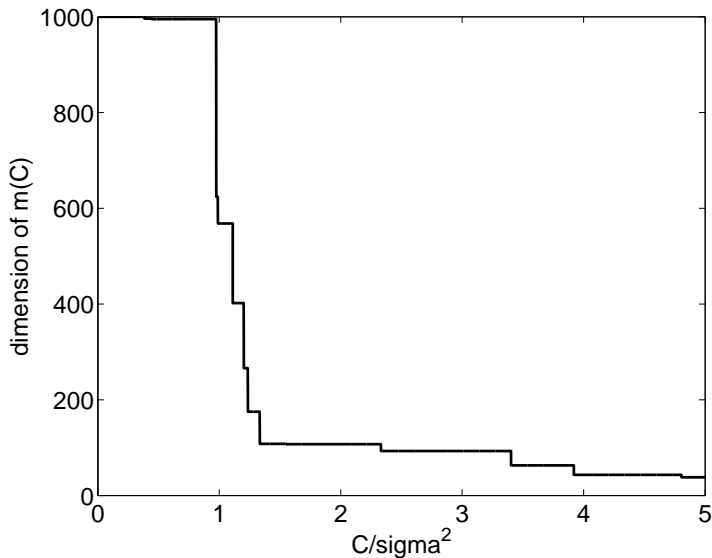
$$\mathbb{E}[n^{-1} \|Y - \widehat{F}_m\|^2] + 1.2 \times \sigma^2 D_m n^{-1}$$



$$\mathbb{E}[n^{-1} \|Y - \widehat{F}_m\|^2] + 2 \times \sigma^2 D_m n^{-1}$$



# Dimension jump



# Calibration of penalties (Birgé & Massart 2007)

- 1 for all  $C > 0$ , compute

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

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

# Proof: assumptions and concentration inequalities

## Assumptions:

- 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.  $D_{m_1} \geq n/2$ ,  $D_{m_2} \leq \sqrt{n}$  and  $\forall i \in \{1, 2\}$ ,  
 $n^{-1} \|F - F_{m_i}\|^2 \leq \sigma^2 \sqrt{\ln(n)/n}$

## Proposition

If  $\xi \sim \mathcal{N}(0, I_n)$ ,  $\alpha \in \mathbb{R}^n$ ,  $M \in \mathcal{M}_n(\mathbb{R})$ , for all  $x \geq 0$ ,

$$\mathbb{P} \left( |\langle \xi, \alpha \rangle| \leq \sqrt{2x} \|\alpha\|_2 \right) \geq 1 - 2e^{-x}$$

$$\mathbb{P} \left( |\langle \xi, M\xi \rangle - \text{tr}(M)| \leq 2\sqrt{x \text{tr}(M^T M)} + 2\|M\|x \right) \geq 1 - 2e^{-x}$$

# Theorem (1): Minimal penalty / Dimension jump

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

With probability at least  $1 - 4C_M n^{-2}$ , si  $n \geq n_0(\alpha)$ ,

$$\forall C < \left( 1 - 42 \sqrt{\frac{(\alpha + 2) \ln(n)}{n}} \right) \sigma^2, \quad D_{\hat{m}(C)} \geq \frac{n}{3}$$
$$\forall C > \left( 1 + 8 \frac{\sqrt{(\alpha + 2) \ln(n)}}{n^{1/4}} \right) \sigma^2, \quad D_{\hat{m}(C)} \leq n^{3/4}$$

and in the first case,

$$\|F - \hat{F}_{\hat{m}(C)}\|^2 \geq \ln(n) \inf_{m \in \mathcal{M}_n} \left\{ \|F - \hat{F}_m\|^2 \right\}$$

# Theorem (2): Oracle inequality

## Theorem (Birgé & Massart 2007)

For every  $x \geq 0$ , with probability  $1 - 4 \text{Card}(\mathcal{M}_n)e^{-x}$ , for every  $K > 1$ ,  $\delta > 0$ , and every

$$\hat{m}(K\sigma^2) \in \arg \min_{m \in \mathcal{M}_n} \left\{ \frac{1}{n} \|Y - \hat{F}_m\|^2 + \frac{K\sigma^2 \dim(S_m)}{n} \right\},$$

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

# Theorem (2): Oracle inequality

## Theorem (Birgé & Massart 2007)

If  $\mathbb{P}(2\hat{C} \in [(1 - \eta_-)2\sigma^2, (1 + \eta_+)2\sigma^2]) \geq 1 - 4C_{\mathcal{M}}n^{-2}$ .

For every  $x \geq 0$ , with probability  $1 - 4 \text{Card}(\mathcal{M}_n)e^{-x} - 4C_{\mathcal{M}}n^{-2}$ ,  
for every  $\delta > 0$ , and every

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

$$\frac{1}{n} \|\hat{F}_{\hat{m}(2\hat{C})} - F\|^2 \leq \left( \frac{1 + \eta_+}{1 - \eta_-} + \delta \right) \inf_{m \in \mathcal{M}_n} \left\{ \frac{1}{n} \|\hat{F}_m - F\|^2 \right\} + \frac{\max\{C(2 - \eta_-), C(2 + \eta_+)\} x \sigma^2}{\delta n}$$



# Theorem (2): Oracle inequality

## Theorem (Birgé & Massart 2007)

We take  $x = (\alpha + 2) \ln(n)$  and assume  $n \geq n_0(\alpha)$ .

With probability  $1 - 4C_{\mathcal{M}}n^{-2}$ , for every

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

$$\frac{1}{n} \left\| \hat{F}_{\hat{m}(2\hat{C})} - F \right\|^2 \leq \left( 1 + \frac{L_\alpha \sqrt{\ln(n)}}{n^{1/4}} + \delta \right) \inf_{m \in \mathcal{M}_n} \left\{ \frac{1}{n} \left\| \hat{F}_m - F \right\|^2 \right\} + \frac{L_\alpha \ln(n) \sigma^2}{\delta n}$$

# Outline

- 1 An oracle inequality for model selection
- 2 The penalty calibration problem
- 3 Slope heuristics in homoscedastic regression
- 4 The slope heuristics**
- 5 Practical issues
- 6 Conclusion

# The slope heuristics (Birgé & Massart, 2007)

- ① existence of a **minimal penalty**  $\text{pen}_{\min}(m)$ :

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

$$\frac{\ell(s^*, \hat{s}_{\hat{m}_{\min}(C)})}{\inf_{m \in \mathcal{M}_n} \{\ell(s^*, \hat{s}_m)\}} \quad \text{jumps at } C = 1$$

# The slope heuristics (Birgé & Massart, 2007)

- ① existence of a **minimal penalty**  $\text{pen}_{\min}(m)$ :

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

$$\frac{\ell(s^*, \hat{s}_{\hat{m}_{\min}(C)})}{\inf_{m \in \mathcal{M}_n} \{ \ell(s^*, \hat{s}_m) \}} \quad \text{jumps at } C = 1$$

- ② the minimal penalty **can be detected**:  
 $\mathcal{C}_{\hat{m}_{\min}(C)}$  “jumps” around  $C = 1$

# The slope heuristics (Birgé & Massart, 2007)

- ① existence of a **minimal penalty**  $\text{pen}_{\min}(m)$ :

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

$$\frac{\ell(s^*, \hat{S}_{\hat{m}_{\min}(C)})}{\inf_{m \in \mathcal{M}_n} \{ \ell(s^*, \hat{S}_m) \}} \quad \text{jumps at } C = 1$$

- ② the minimal penalty **can be detected**:  
 $\mathcal{C}_{\hat{m}_{\min}(C)}$  “jumps” around  $C = 1$
- ③ link between minimal and optimal penalty:

$$\text{pen}_{\text{opt}}(m) \approx 2 \text{pen}_{\min}(m)$$

# Data-driven penalties with the slope heuristics

**Inputs:**  $(\text{pen}_0(m))_{m \in \mathcal{M}_n}$   $(\mathcal{C}_m)_{m \in \mathcal{M}_n}$

**Assumption:**  $\text{pen}_0(m) \propto \text{pen}_{\min}(m)$

- 1 for every  $C > 0$ , compute

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

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

# Slope heuristics recipe

$$p_2(m) = P_n(\gamma(s_m^*) - \gamma(\hat{s}_m))$$

$$\text{pen}_{\min}(m) = \mathbb{E}[p_2(m)]$$

# Slope heuristics recipe

$$p_2(m) = P_n(\gamma(s_m^*) - \gamma(\hat{s}_m)) \quad \text{pen}_{\min}(m) = \mathbb{E}[p_2(m)]$$

$$p_1(m) = P(\gamma(\hat{s}_m) - \gamma(s_m^*)) \quad \delta(m) = (P - P_n)\gamma(s_m^*)$$

$$\text{pen}_{\text{id}}(m) = p_1(m) + p_2(m) - \delta(m)$$

$$\mathbb{E}[\text{pen}_{\text{id}}(m)] = \text{pen}_{\text{opt}}(m) = \mathbb{E}[p_1(m)] + \mathbb{E}[p_2(m)]$$



# Slope heuristics recipe

$$p_2(m) = P_n(\gamma(s_m^*) - \gamma(\widehat{s}_m)) \quad \text{pen}_{\min}(m) = \mathbb{E}[p_2(m)]$$

$$p_1(m) = P(\gamma(\widehat{s}_m) - \gamma(s_m^*)) \quad \delta(m) = (P - P_n)\gamma(s_m^*)$$

$$\text{pen}_{\text{id}}(m) = p_1(m) + p_2(m) - \delta(m)$$

$$\mathbb{E}[\text{pen}_{\text{id}}(m)] = \text{pen}_{\text{opt}}(m) = \mathbb{E}[p_1(m)] + \mathbb{E}[p_2(m)]$$

Heuristics:  $p_1(m) \approx p_2(m)$

- concentration of  $p_1$ ,  $p_2$ ,  $\delta$
- $\mathbb{E}[p_1(m)] \approx \mathbb{E}[p_2(m)]$
- increase of the expectation for compensating the bias

# Known results

- Least-squares, regression, **homoscedastic Gaussian noise** (Birgé & Massart, 2007)
- **Heteroscedastic regressograms** (A. & Massart, 2009)
- **Least-squares density estimation**, i.i.d. (Lerasle, 2009) or **mixing** data (Lerasle, 2010)
- **Minimum contrast estimators, regular contrast** (Saumard, 2010)

# Minimizer of a regular contrast (Saumard, 2010)

- **Regular contrast** on some convex model  $S_m$ :
    - $s_m^* \in \arg \min_{t \in S_m} P\gamma(t)$  exists
    - $t \in S_m \mapsto P\gamma(t)$  strictly convex
    - $\exists c > 0, t \in B_\infty(s_m^*, c) \mapsto \gamma(t; \cdot) \in L_\infty(P)$  is  $\mathcal{C}^3$
  - **Concentration of  $p_1(m)$  and  $p_2(m)$**  around the same deterministic quantity  $D_m \mathcal{K}_m^2 / (4n)$  (unobservable in general)
- + control of  $\|\widehat{s}_m - s_m^*\|_\infty$
- ⇒ validates the slope heuristics for:
- heteroscedastic regression (histograms, **piecewise polynomials**)
  - least-squares density estimation
  - **log-likelihood** density estimation on histograms

# Experimental results

- Change-point detection (Lebarbier, 2005)
- Gaussian mixture models (Maugis & Michel, 2008)
- Unsupervised classification (choice of the number of clusters) (Baudry, 2009)
- Computational geometry (Caillerie & Michel, 2009)
- Lasso (Connault, 2011)
- ...

for a complete list, see Baudry, Maugis & Michel, 2010

# Outline

- 1 An oracle inequality for model selection
- 2 The penalty calibration problem
- 3 Slope heuristics in homoscedastic regression
- 4 The slope heuristics
- 5 Practical issues**
- 6 Conclusion

# 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**

## How to localize the jump in practice?

- **Complexity jump**: largest jump? largest relative jump? complexity threshold?

# How to localize the jump in practice?

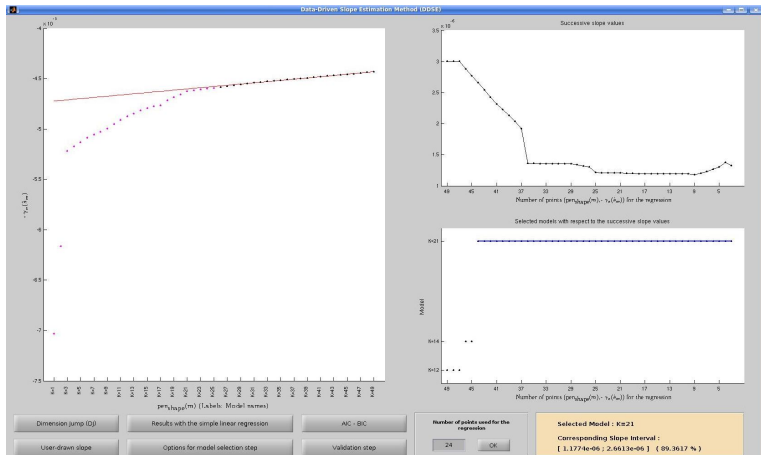
- Complexity jump: largest jump? largest relative jump? complexity threshold?
- Estimation of the slope of the empirical risk as a function of the complexity:  
computed with which models? robust regression?



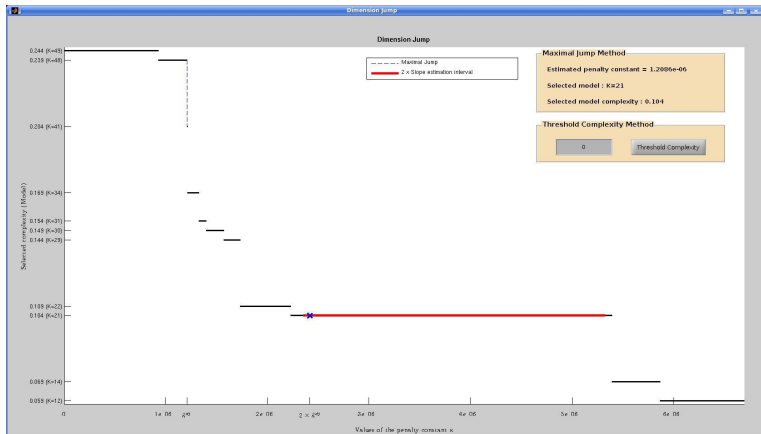
# How to localize the jump in practice?

- Complexity jump: largest jump? largest relative jump? complexity threshold?
- Estimation of the slope of the empirical risk as a function of the complexity:  
computed with which models? robust regression?
- **Jump vs. slope? Take both!**  
⇒ package CAPUSHE (Baudry, Maugis & Michel, 2010)  
<http://www.math.univ-toulouse.fr/~maugis/CAPUSHE.html>

# CAPUSHE (Baudry, Maugis & Michel, 2010): slope



# CAPUSHE (Baudry, Maugis & Michel, 2010): jump



# Outline

- 1 An oracle inequality for model selection
- 2 The penalty calibration problem
- 3 Slope heuristics in homoscedastic regression
- 4 The slope heuristics
- 5 Practical issues
- 6 Conclusion**

# Overpenalization

