

# V-fold cross-validation improved: V-fold penalization

Sylvain Arlot

<sup>1</sup>University Paris-Sud XI, Orsay

<sup>2</sup>Inria Saclay, Projet Select

Journées Statistiques du Sud, INSA Toulouse  
June 18, 2008

# Statistical framework: regression on a random design

$(X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \mathcal{Y}$  i.i.d.       $(X_i, Y_i) \sim P$  unknown

$Y = s(X) + \sigma(X)\epsilon$        $X \in \mathcal{X} \subset \mathbb{R}^d$ ,       $Y \in \mathcal{Y} = [0; 1]$  or  $\mathbb{R}$

noise  $\epsilon$  :       $\mathbb{E}[\epsilon|X] = 0$       noise level       $\sigma(X)$

predictor       $t : \mathcal{X} \mapsto \mathcal{Y}$       ?

# Statistical framework: regression on a random design

$(X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \mathcal{Y}$  i.i.d.       $(X_i, Y_i) \sim P$  unknown

$Y = s(X) + \sigma(X)\epsilon$        $X \in \mathcal{X} \subset \mathbb{R}^d$ ,       $Y \in \mathcal{Y} = [0; 1]$  or  $\mathbb{R}$

noise  $\epsilon$  :       $\mathbb{E}[\epsilon|X] = 0$       noise level       $\sigma(X)$

predictor       $t : \mathcal{X} \mapsto \mathcal{Y}$       ?

# Statistical framework: regression on a random design

$(X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \mathcal{Y}$  i.i.d.       $(X_i, Y_i) \sim P$  unknown

$Y = s(X) + \sigma(X)\epsilon$        $X \in \mathcal{X} \subset \mathbb{R}^d$ ,       $Y \in \mathcal{Y} = [0; 1]$  or  $\mathbb{R}$

noise  $\epsilon$  :       $\mathbb{E}[\epsilon|X] = 0$       noise level       $\sigma(X)$

predictor       $t : \mathcal{X} \mapsto \mathcal{Y}$       ?

# Statistical framework: regression on a random design

$(X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \mathcal{Y}$  i.i.d.       $(X_i, Y_i) \sim P$  unknown

$Y = s(X) + \sigma(X)\epsilon$        $X \in \mathcal{X} \subset \mathbb{R}^d$ ,       $Y \in \mathcal{Y} = [0; 1]$  or  $\mathbb{R}$

noise  $\epsilon$  :       $\mathbb{E}[\epsilon|X] = 0$       noise level       $\sigma(X)$

predictor       $t : \mathcal{X} \mapsto \mathcal{Y}$       ?

# Loss function, least-square estimator

- Least-square risk:

$$\mathbb{E}\gamma(t, (X, Y)) = P\gamma(t, \cdot)$$

$$\text{with } \gamma(t, (x, y)) = (t(x) - y)^2$$

- Empirical risk minimizer on  $S_m$  (= model):

$$\hat{s}_m \in \arg \min_{t \in S_m} P_n \gamma(t, \cdot) = \arg \min_{t \in S_m} \frac{1}{n} \sum_{i=1}^n (t(X_i) - Y_i)^2.$$

- e.g. histograms on a partition  $(I_\lambda)_{\lambda \in \Lambda_m}$  of  $\mathcal{X}$ .

$$\hat{s}_m = \sum_{\lambda \in \Lambda_m} \hat{\beta}_\lambda \mathbf{1}_{I_\lambda} \quad \hat{\beta}_\lambda = \frac{1}{\text{Card}\{X_i \in I_\lambda\}} \sum_{X_i \in I_\lambda} Y_i.$$

# Loss function, least-square estimator

- **Loss function:**

$$\ell(s, t) = P\gamma(t, \cdot) - P\gamma(s, \cdot) = \mathbb{E} [(t(X) - s(X))^2]$$

$$\text{with } \gamma(t, (x, y)) = (t(x) - y)^2$$

- Empirical risk minimizer on  $S_m$  (= model):

$$\hat{s}_m \in \arg \min_{t \in S_m} P_n \gamma(t, \cdot) = \arg \min_{t \in S_m} \frac{1}{n} \sum_{i=1}^n (t(X_i) - Y_i)^2.$$

- e.g. histograms on a partition  $(I_\lambda)_{\lambda \in \Lambda_m}$  of  $\mathcal{X}$ .

$$\hat{s}_m = \sum_{\lambda \in \Lambda_m} \hat{\beta}_\lambda \mathbf{1}_{I_\lambda} \quad \hat{\beta}_\lambda = \frac{1}{\text{Card}\{X_i \in I_\lambda\}} \sum_{X_i \in I_\lambda} Y_i.$$

# Loss function, least-square estimator

- Loss function:

$$\ell(s, t) = P\gamma(t, \cdot) - P\gamma(s, \cdot) = \mathbb{E} [(t(X) - s(X))^2]$$

with  $\gamma(t, (x, y)) = (t(x) - y)^2$

- Empirical risk minimizer on  $S_m$  (= model):

$$\hat{s}_m \in \arg \min_{t \in S_m} P_n \gamma(t, \cdot) = \arg \min_{t \in S_m} \frac{1}{n} \sum_{i=1}^n (t(X_i) - Y_i)^2 .$$

- e.g. histograms on a partition  $(I_\lambda)_{\lambda \in \Lambda_m}$  of  $\mathcal{X}$ .

$$\hat{s}_m = \sum_{\lambda \in \Lambda_m} \hat{\beta}_\lambda \mathbb{1}_{I_\lambda} \quad \hat{\beta}_\lambda = \frac{1}{\text{Card}\{X_i \in I_\lambda\}} \sum_{X_i \in I_\lambda} Y_i .$$



# Loss function, least-square estimator

- Loss function:

$$\ell(s, t) = P\gamma(t, \cdot) - P\gamma(s, \cdot) = \mathbb{E} [(t(X) - s(X))^2]$$

$$\text{with } \gamma(t, (x, y)) = (t(x) - y)^2$$

- Empirical risk minimizer on  $S_m$  (= model):

$$\hat{s}_m \in \arg \min_{t \in S_m} P_n \gamma(t, \cdot) = \arg \min_{t \in S_m} \frac{1}{n} \sum_{i=1}^n (t(X_i) - Y_i)^2 .$$

- e.g. histograms on a partition  $(I_\lambda)_{\lambda \in \Lambda_m}$  of  $\mathcal{X}$ .

$$\hat{s}_m = \sum_{\lambda \in \Lambda_m} \hat{\beta}_\lambda \mathbb{1}_{I_\lambda} \quad \hat{\beta}_\lambda = \frac{1}{\text{Card}\{X_i \in I_\lambda\}} \sum_{X_i \in I_\lambda} Y_i .$$

# Model selection

$$(\mathcal{S}_m)_{m \in \mathcal{M}} \longrightarrow (\widehat{\mathcal{S}}_m)_{m \in \mathcal{M}} \longrightarrow \widehat{\mathcal{S}}_{\widehat{m}} \quad ???$$

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

$$\ell(s, \widehat{\mathcal{S}}_{\widehat{m}}) \leq C \inf_{m \in \mathcal{M}} \{\ell(s, \widehat{\mathcal{S}}_m) + R(m, n)\}$$

- Adaptivity (e.g.,  $\alpha$  if  $s$  is  $\alpha$ -hölder,  $\sigma(X)$  in the heteroscedastic framework)

# Model selection

$$(S_m)_{m \in \mathcal{M}} \longrightarrow (\widehat{S}_m)_{m \in \mathcal{M}} \longrightarrow \widehat{S}_{\widehat{m}} \quad ???$$

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

$$\ell(s, \widehat{S}_{\widehat{m}}) \leq C \inf_{m \in \mathcal{M}} \{\ell(s, \widehat{S}_m) + R(m, n)\}$$

- Adaptivity (e.g.,  $\alpha$  if  $s$  is  $\alpha$ -holder,  $\sigma(X)$  in the heteroscedastic framework)

# Model selection

$$(S_m)_{m \in \mathcal{M}} \longrightarrow (\widehat{S}_m)_{m \in \mathcal{M}} \longrightarrow \widehat{s}_{\widehat{m}} \quad ???$$

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

$$\ell(s, \widehat{s}_{\widehat{m}}) \leq C \inf_{m \in \mathcal{M}} \{\ell(s, \widehat{s}_m) + R(m, n)\}$$

- Adaptivity (e.g.,  $\alpha$  if  $s$  is  $\alpha$ -hölder,  $\sigma(X)$  in the heteroscedastic framework)

# Cross-validation

$$\underbrace{(X_1, Y_1), \dots, (X_q, Y_q)}_{\text{Training}}, \underbrace{(X_{q+1}, Y_{q+1}), \dots, (X_n, Y_n)}_{\text{Validation}}$$

$$\hat{s}_m^{(t)} \in \arg \min_{t \in \mathcal{S}_m} \left\{ \frac{1}{q} \sum_{i=1}^q \gamma(t, (X_i, Y_i)) \right\}$$

$$P_n^{(v)} = \frac{1}{n-q} \sum_{i=q+1}^n \delta_{(X_i, Y_i)} \quad \Rightarrow \quad P_n^{(v)} \gamma \left( \hat{s}_m^{(t)} \right)$$

V-fold cross-validation:  $(B_j)_{1 \leq j \leq V}$  partition of  $\{1, \dots, n\}$

$$\Rightarrow \hat{m} \in \arg \min_{m \in \mathcal{M}} \left\{ \frac{1}{V} \sum_{j=1}^V P_n^{(j)} \gamma \left( \hat{s}_m^{(-j)} \right) \right\} \quad \tilde{s} = \hat{s}_{\hat{m}}$$

# Cross-validation

$$\underbrace{(X_1, Y_1), \dots, (X_q, Y_q)}_{\text{Training}}, \underbrace{(X_{q+1}, Y_{q+1}), \dots, (X_n, Y_n)}_{\text{Validation}}$$

$$\hat{s}_m^{(t)} \in \arg \min_{t \in \mathcal{S}_m} \left\{ \frac{1}{q} \sum_{i=1}^q \gamma(t, (X_i, Y_i)) \right\}$$

$$P_n^{(v)} = \frac{1}{n-q} \sum_{i=q+1}^n \delta_{(X_i, Y_i)} \quad \Rightarrow \quad P_n^{(v)} \gamma \left( \hat{s}_m^{(t)} \right)$$

V-fold cross-validation:  $(B_j)_{1 \leq j \leq V}$  partition of  $\{1, \dots, n\}$

$$\Rightarrow \hat{m} \in \arg \min_{m \in \mathcal{M}} \left\{ \frac{1}{V} \sum_{j=1}^V P_n^{(j)} \gamma \left( \hat{s}_m^{(-j)} \right) \right\} \quad \tilde{s} = \hat{s}_{\hat{m}}$$

# Cross-validation

$$\underbrace{(X_1, Y_1), \dots, (X_q, Y_q)}_{\text{Training}}, \underbrace{(X_{q+1}, Y_{q+1}), \dots, (X_n, Y_n)}_{\text{Validation}}$$

$$\hat{s}_m^{(t)} \in \arg \min_{t \in \mathcal{S}_m} \left\{ \frac{1}{q} \sum_{i=1}^q \gamma(t, (X_i, Y_i)) \right\}$$

$$P_n^{(v)} = \frac{1}{n-q} \sum_{i=q+1}^n \delta_{(X_i, Y_i)} \quad \Rightarrow \quad P_n^{(v)} \gamma \left( \hat{s}_m^{(t)} \right)$$

**V-fold cross-validation:**  $(B_j)_{1 \leq j \leq V}$  partition of  $\{1, \dots, n\}$

$$\Rightarrow \hat{m} \in \arg \min_{m \in \mathcal{M}} \left\{ \frac{1}{V} \sum_{j=1}^V P_n^{(j)} \gamma \left( \hat{s}_m^{(-j)} \right) \right\} \quad \tilde{s} = \hat{s}_{\hat{m}}$$

# Bias of cross-validation

**Ideal criterion:**  $P\gamma(\hat{s}_m)$

Regression on an histogram model of dimension  $D_m$ , when  $\sigma(X) \equiv \sigma$ :

$$\mathbb{E} [P\gamma(\hat{s}_m)] \approx P\gamma(s_m) + \frac{D_m\sigma^2}{n}$$

$$\mathbb{E} \left[ P_n^{(j)} \gamma \left( \hat{s}_m^{(-j)} \right) \right] = \mathbb{E} \left[ P\gamma \left( \hat{s}_m^{(-j)} \right) \right] \approx P\gamma(s_m) + \frac{V}{V-1} \frac{D_m\sigma^2}{n}$$

$\Rightarrow$  **bias** if  $V$  is fixed



# Bias of cross-validation

Ideal criterion:  $P\gamma(\hat{s}_m)$

Regression on an histogram model of dimension  $D_m$ , when  $\sigma(X) \equiv \sigma$ :

$$\mathbb{E} [P\gamma(\hat{s}_m)] \approx P\gamma(s_m) + \frac{D_m\sigma^2}{n}$$

$$\mathbb{E} \left[ P_n^{(j)} \gamma \left( \hat{s}_m^{(-j)} \right) \right] = \mathbb{E} \left[ P\gamma \left( \hat{s}_m^{(-j)} \right) \right] \approx P\gamma(s_m) + \frac{V}{V-1} \frac{D_m\sigma^2}{n}$$

⇒ **bias** if  $V$  is fixed

# Suboptimality of $V$ -fold cross-validation

- $Y = X + \sigma\epsilon$  with  $\epsilon$  bounded and  $\sigma > 0$
- $\mathcal{M}_n$ : family of regular histograms on  $\mathcal{X} = [0, 1]$
- $V$  fixed

## Theorem

*With probability at least  $1 - \diamond n^{-2}$ ,*

$$l(s, \widehat{s}_m) \geq (1 + \kappa(V)) \inf_{m \in \mathcal{M}} \{l(s, \widehat{s}_m)\}$$

*with  $\kappa(V) > 0$ .*

# Choice of $V$

- Bias: decreases with  $V$  (can be corrected: Burman 1989)
- Variability: large if  $V$  is small ( $V = 2$ ), or sometimes when  $V$  is very large ( $V = n$ , unstable algorithms)
- Computation time: complexity proportional to  $V$

⇒ trade-off

⇒ classical conclusion: “ $V = 10$  is fine”

# Choice of $V$

- Bias: decreases with  $V$  (can be corrected: Burman 1989)
- Variability: large if  $V$  is small ( $V = 2$ ), or sometimes when  $V$  is very large ( $V = n$ , unstable algorithms)
- Computation time: complexity proportional to  $V$

⇒ trade-off

⇒ classical conclusion: “ $V = 10$  is fine”

# Choice of $V$

- Bias: decreases with  $V$  (can be corrected: Burman 1989)
- Variability: large if  $V$  is small ( $V = 2$ ), or sometimes when  $V$  is very large ( $V = n$ , unstable algorithms)
- Computation time: complexity proportional to  $V$

⇒ trade-off

⇒ classical conclusion: “ $V = 10$  is fine”

# Choice of $V$

- Bias: decreases with  $V$  (can be corrected: Burman 1989)
- Variability: large if  $V$  is small ( $V = 2$ ), or sometimes when  $V$  is very large ( $V = n$ , unstable algorithms)
- Computation time: complexity proportional to  $V$

⇒ trade-off

⇒ classical conclusion: “ $V = 10$  is fine”

# Simulation framework

$$Y_i = s(X_i) + \sigma(X_i)\epsilon_i \quad X_i \sim^{\text{i.i.d.}} \mathcal{U}([0; 1]) \quad \epsilon_i \sim^{\text{i.i.d.}} \mathcal{N}(0, 1)$$

$$\mathcal{M}_n = \left\{ \text{regular histograms with } D \text{ pieces, } 1 \leq D \leq \frac{n}{\log(n)} \right. \\ \left. \text{and s.t. } \min_{\lambda \in \Lambda_m} \text{Card}\{X_i \in I_\lambda\} \geq 2 \right\}$$

⇒ Benchmark:

$$C_{\text{classical}} = \frac{\mathbb{E}[\ell(s, \hat{s}_{\hat{m}})]}{\mathbb{E}[\inf_{m \in \mathcal{M}} \ell(s, \hat{s}_m)]} \quad \text{computed with } N = 1000 \text{ samples}$$

# Simulation framework

$$Y_i = s(X_i) + \sigma(X_i)\epsilon_i \quad X_i \sim^{\text{i.i.d.}} \mathcal{U}([0; 1]) \quad \epsilon_i \sim^{\text{i.i.d.}} \mathcal{N}(0, 1)$$

$$\mathcal{M}_n = \left\{ \text{regular histograms with } D \text{ pieces, } 1 \leq D \leq \frac{n}{\log(n)} \right. \\ \left. \text{and s.t. } \min_{\lambda \in \Lambda_m} \text{Card}\{X_i \in I_\lambda\} \geq 2 \right\}$$

⇒ Benchmark:

$$C_{\text{classical}} = \frac{\mathbb{E}[\ell(s, \hat{s}_{\hat{m}})]}{\mathbb{E}[\inf_{m \in \mathcal{M}} \ell(s, \hat{s}_m)]} \quad \text{computed with } N = 1000 \text{ samples}$$



# Simulation framework

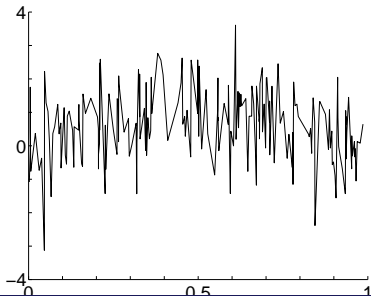
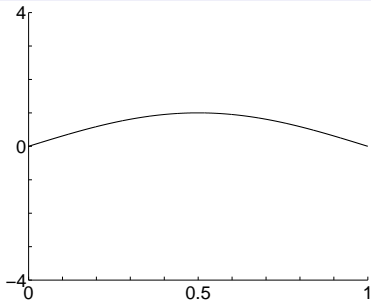
$$Y_i = s(X_i) + \sigma(X_i)\epsilon_i \quad X_i \sim^{\text{i.i.d.}} \mathcal{U}([0; 1]) \quad \epsilon_i \sim^{\text{i.i.d.}} \mathcal{N}(0, 1)$$

$$\mathcal{M}_n = \left\{ \text{regular histograms with } D \text{ pieces, } 1 \leq D \leq \frac{n}{\log(n)} \right. \\ \left. \text{and s.t. } \min_{\lambda \in \Lambda_m} \text{Card}\{X_i \in I_\lambda\} \geq 2 \right\}$$

⇒ Benchmark:

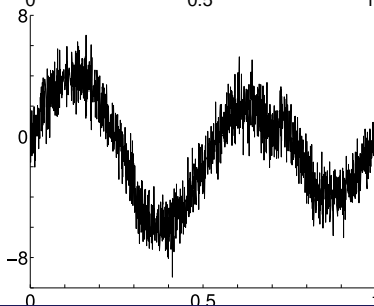
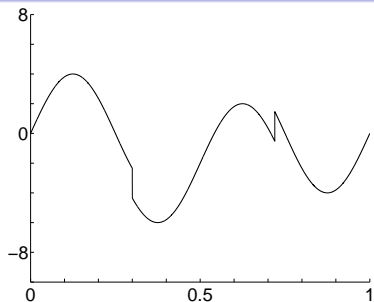
$$C_{\text{classical}} = \frac{\mathbb{E}[\ell(s, \hat{s}_{\hat{m}})]}{\mathbb{E}[\inf_{m \in \mathcal{M}} \ell(s, \hat{s}_m)]} \quad \text{computed with } N = 1000 \text{ samples}$$

# Simulations: $s(x) = \sin(\pi x)$ , $n = 200$ , $\sigma \equiv 1$



2-fold	$2.08 \pm 0.04$
5-fold	$2.14 \pm 0.04$
10-fold	$2.10 \pm 0.05$
20-fold	$2.09 \pm 0.04$
leave-one-out	$2.08 \pm 0.04$

# Simulations: HeaviSine, $n = 2048$ , $\sigma \equiv 1$



2-fold	$1.002 \pm 0.003$
5-fold	$1.014 \pm 0.003$
10-fold	$1.021 \pm 0.003$
20-fold	$1.029 \pm 0.004$
leave-one-out	$1.034 \pm 0.004$

# The penalization viewpoint

- penalization:  $\hat{m} \in \arg \min_{m \in \mathcal{M}} \{P_n \gamma(\hat{S}_m) + \text{pen}(m)\}$
- ideal penalty:  $\text{pen}_{\text{id}}(m) = P \gamma(\hat{S}_m) - P_n \gamma(\hat{S}_m)$
- V-fold cross-validation is **overpenalizing**:

$$\frac{\mathbb{E} \left[ \frac{1}{V} \sum_{j=1}^V P_n^{(j)} \gamma(\hat{S}_m^{(-j)}) - P_n \gamma(\hat{S}_m) \right]}{\mathbb{E} [\text{pen}_{\text{id}}(m)]} \approx 1 + \frac{1}{2(V-1)}$$

- **non-asymptotic** phenomenon:  
better to **overpenalize** when the signal-to-noise ratio  $n/\sigma^2$  is small.

# The penalization viewpoint

- penalization:  $\hat{m} \in \arg \min_{m \in \mathcal{M}} \{P_n \gamma(\hat{S}_m) + \text{pen}(m)\}$
- ideal penalty:  $\text{pen}_{\text{id}}(m) = P \gamma(\hat{S}_m) - P_n \gamma(\hat{S}_m)$
- V-fold cross-validation is **overpenalizing**:

$$\frac{\mathbb{E} \left[ \frac{1}{V} \sum_{j=1}^V P_n^{(j)} \gamma(\hat{S}_m^{(-j)}) - P_n \gamma(\hat{S}_m) \right]}{\mathbb{E} [\text{pen}_{\text{id}}(m)]} \approx 1 + \frac{1}{2(V-1)}$$

- **non-asymptotic** phenomenon:  
better to **overpenalize** when the signal-to-noise ratio  $n/\sigma^2$  is small.

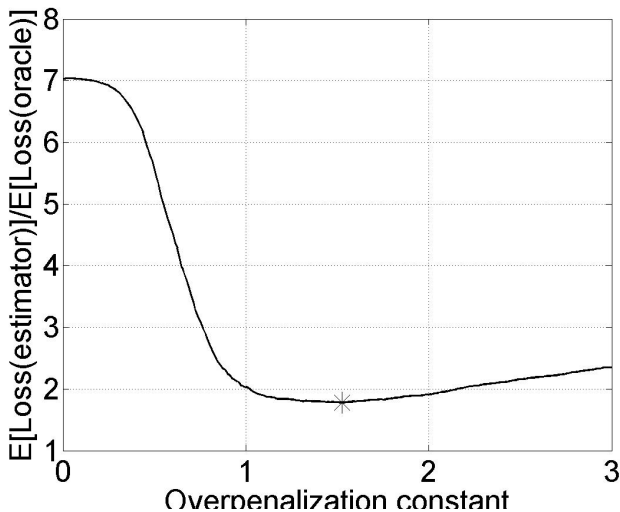
# The penalization viewpoint

- penalization:  $\hat{m} \in \arg \min_{m \in \mathcal{M}} \{P_n \gamma(\hat{S}_m) + \text{pen}(m)\}$
- ideal penalty:  $\text{pen}_{\text{id}}(m) = P \gamma(\hat{S}_m) - P_n \gamma(\hat{S}_m)$
- V-fold cross-validation is **overpenalizing**:

$$\frac{\mathbb{E} \left[ \frac{1}{V} \sum_{j=1}^V P_n^{(j)} \gamma(\hat{S}_m^{(-j)}) - P_n \gamma(\hat{S}_m) \right]}{\mathbb{E} [\text{pen}_{\text{id}}(m)]} \approx 1 + \frac{1}{2(V-1)}$$

- **non-asymptotic** phenomenon:  
better to **overpenalize** when the signal-to-noise ratio  $n/\sigma^2$  is small.

# Overpenalization ( $s = \sin$ , $\sigma \equiv 1$ , $n = 200$ , Mallows' $C_p$ )



# Conclusions on $V$ -fold cross-validation

- asymptotically suboptimal if  $V$  fixed
- optimal  $V^*$ : trade-off **variability–overpenalization**
- $V^* = 2$  can happen for prediction
  
- **difficult** to find  $V^*$  from the data (+ complexity issue)
- low signal-to-noise ratio  $\Rightarrow V^*$  **unsatisfactory** (highly **variable**)
- large signal-to-noise ratio  $\Rightarrow V^*$  too large (**computation time**)



# Conclusions on $V$ -fold cross-validation

- asymptotically suboptimal if  $V$  fixed
- optimal  $V^*$ : trade-off **variability–overpenalization**
- $V^* = 2$  can happen for prediction
  
- **difficult** to find  $V^*$  from the data (+ complexity issue)
- low signal-to-noise ratio  $\Rightarrow V^*$  **unsatisfactory** (highly **variable**)
- large signal-to-noise ratio  $\Rightarrow V^*$  too large (**computation time**)

# Penalization

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

Ideal penalty:  $\text{pen}_{\text{id}}(m) = (P - P_n)(\gamma(\hat{S}_m, \cdot))$

$$\text{pen}(m) = \frac{2\sigma^2 D_m}{n} \quad (\text{Mallows 1973}) \quad \text{pen}(m) = \frac{2\hat{\sigma}^2 D_m}{n}$$

# Penalization

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

Ideal penalty:  $\text{pen}_{\text{id}}(m) = (P - P_n)(\gamma(\hat{S}_m, \cdot))$

$$\text{pen}(m) = \frac{2\sigma^2 D_m}{n} \quad (\text{Mallows 1973}) \quad \text{pen}(m) = \frac{2\hat{\sigma}^2 D_m}{n}$$

# Penalization

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

$$\text{Ideal penalty: } \text{pen}_{\text{id}}(m) = (P - P_n)(\gamma(\hat{s}_m, \cdot))$$

## Theorem (Suboptimality of linear penalties, A., 2008)

$\mathcal{X} = [0, 1]$ ,  $Y = X + \sigma(X)\epsilon$ ,  $\sigma(x) = \mathbb{1}_{x \leq 1/2} + 3\mathbb{1}_{x > 1/2}$

$\mathcal{M}_n$ : Regular histograms on  $[0; 1/2]$  and  $[1/2; 1]$

With a probability at least  $1 - \diamond n^{-2}$ , for every  $K \geq 0$  and

$$\hat{m}(K) \in \arg \min_{m \in \mathcal{M}_n} \{P_n \gamma(\hat{s}_m) + KD_m\} \quad ,$$

$$\ell(s, \hat{s}_{\hat{m}(K)}) \geq (1 + \kappa) \inf_{m \in \mathcal{M}_n} \{\ell(s, \hat{s}_m)\} \quad \text{with } \kappa > 0 \quad .$$

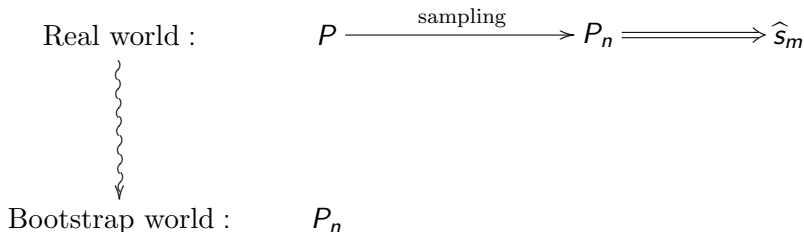
## Resampling heuristics (bootstrap, Efron 1979)

Real world :  $P \xrightarrow{\text{sampling}} P_n \Longrightarrow \hat{S}_m$

$$\text{pen}_{\text{id}}(m) = (P - P_n)\gamma(\hat{S}_m) = F(P, P_n)$$

V-fold:  $P_n^W = \frac{1}{n - \text{Card}(B_J)} \sum_{i \notin B_J} \delta_{(X_i, Y_i)}$  with  $J \sim \mathcal{U}(1, \dots, V)$

# Resampling heuristics (bootstrap, Efron 1979)



$$\text{pen}_{\text{id}}(m) = (P - P_n)\gamma(\hat{S}_m) = F(P, P_n)$$

V-fold:  $P_n^W = \frac{1}{n - \text{Card}(B_J)} \sum_{i \notin B_J} \delta_{(X_i, Y_i)} \quad \text{with } J \sim \mathcal{U}(1, \dots, V)$

## Resampling heuristics (bootstrap, Efron 1979)

Real world :

$$P \xrightarrow{\text{sampling}} P_n \xRightarrow{\quad\quad\quad} \widehat{S}_m$$



Bootstrap world :

$$P_n \xrightarrow{\text{resampling}} P_n^W \xRightarrow{\quad\quad\quad} \widehat{S}_m^W$$

$$(P - P_n)\gamma(\widehat{S}_m) = F(P, P_n) \rightsquigarrow F(P_n, P_n^W) = (P_n - P_n^W)\gamma(\widehat{S}_m^W)$$

$$\text{V-fold: } P_n^W = \frac{1}{n - \text{Card}(B_J)} \sum_{i \notin B_J} \delta_{(X_i, Y_i)} \quad \text{with } J \sim \mathcal{U}(1, \dots, V)$$

## Resampling heuristics (bootstrap, Efron 1979)

Real world :

$$P \xrightarrow{\text{sampling}} P_n \xRightarrow{\quad\quad\quad} \widehat{S}_m$$



Bootstrap world :

$$P_n \xrightarrow{\text{subsampling}} P_n^W \xRightarrow{\quad\quad\quad} \widehat{S}_m^W$$

$$(P - P_n)\gamma(\widehat{S}_m) = F(P, P_n) \rightsquigarrow F(P_n, P_n^W) = (P_n - P_n^W)\gamma(\widehat{S}_m^W)$$

$$\text{V-fold: } P_n^W = \frac{1}{n - \text{Card}(B_J)} \sum_{i \notin B_J} \delta_{(X_i, Y_i)} \quad \text{with } J \sim \mathcal{U}(1, \dots, V)$$



# V-fold penalization

- Ideal penalty:

$$(P - P_n)(\gamma(\hat{s}_m))$$

- V-fold penalty:

$$\text{pen}(m) = \frac{C}{V} \sum_{j=1}^V \left[ (P_n - P_n^{(-j)})(\gamma(\hat{s}_m^{(-j)})) \right]$$

$$\hat{s}_m^{(-j)} \in \arg \min_{t \in \mathcal{S}_m} P_n^{(-j)} \gamma(t)$$

with  $C \geq V - 1$  to be chosen

( $C = V - 1 \Rightarrow$  we recover Burman's corrected V-fold, 1989)

- The final estimator is  $\hat{s}_{\hat{m}}$  with

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

# V-fold penalization

- Ideal penalty:

$$(P - P_n)(\gamma(\hat{s}_m))$$

- V-fold penalty:

$$\text{pen}(m) = \frac{C}{V} \sum_{j=1}^V \left[ (P_n - P_n^{(-j)})(\gamma(\hat{s}_m^{(-j)})) \right]$$

$$\hat{s}_m^{(-j)} \in \arg \min_{t \in S_m} P_n^{(-j)} \gamma(t)$$

with  $C \geq V - 1$  to be chosen

( $C = V - 1 \Rightarrow$  we recover Burman's corrected V-fold, 1989)

- The final estimator is  $\hat{s}_{\hat{m}}$  with

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

# V-fold penalization

- Ideal penalty:

$$(P - P_n)(\gamma(\hat{s}_m))$$

- V-fold penalty:

$$\text{pen}(m) = \frac{C}{V} \sum_{j=1}^V \left[ (P_n - P_n^{(-j)})(\gamma(\hat{s}_m^{(-j)})) \right]$$

$$\hat{s}_m^{(-j)} \in \arg \min_{t \in S_m} P_n^{(-j)} \gamma(t)$$

with  $C \geq V - 1$  to be chosen

( $C = V - 1 \Rightarrow$  we recover Burman's corrected V-fold, 1989)

- The final estimator is  $\hat{s}_{\hat{m}}$  with

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

# Some references on model selection and resampling

- **Hold-out, Cross-validation, Leave-one-out, V-fold cross-validation:**

$I \subset \{1, \dots, n\}$  random sub-sample of size  $q$  (VFCV:  
 $q = \frac{n(V-1)}{V}$ ).

- **Efron's bootstrap penalties** (Efron 1983, Shibata 1997):

$$\text{pen}(m) = \mathbb{E} \left[ (P_n - P_n^W)(\gamma(\hat{s}_m^W)) \middle| (X_i, Y_i)_{1 \leq i \leq n} \right]$$

- **Rademacher complexities** (Koltchinskii 2001 ; Bartlett, Boucheron, Lugosi 2002): subsampling

$$\text{pen}_{\text{id}}(m) \leq \text{pen}_{\text{id}}^{\text{glo}}(m) = \sup_{t \in S_m} (P - P_n)\gamma(t, \cdot)$$

- idem with general exchangeable weights (Fromont 2004)
- **Local Rademacher complexities** (Bartlett, Bousquet, Mendelson 2004 ; Koltchinskii 2004)

# Non-asymptotic pathwise oracle inequality

- $C \approx V - 1$
- Histogram regression on a random design
- Small number of models (at most polynomial in  $n$ )
- Model pre-selection: remove  $m$  when

$$\min_{\lambda \in \Lambda_m} \{\text{Card} \{X_i \in I_\lambda\}\} \leq 1$$

- Fixed  $V$  or  $V = n$

## Theorem

*Under a “reasonable” set of assumptions on  $P$ , with probability at least  $1 - \diamond n^{-2}$ ,*

$$l(s, \widehat{s}_{\widehat{m}}) \leq \left(1 + \ln(n)^{-1/5}\right) \inf_{m \in \mathcal{M}} \{l(s, \widehat{s}_m)\}$$

# Non-asymptotic pathwise oracle inequality

- $C \approx V - 1$
- Histogram regression on a random design
- Small number of models (at most polynomial in  $n$ )
- Model pre-selection: remove  $m$  when

$$\min_{\lambda \in \Lambda_m} \{\text{Card} \{X_i \in I_\lambda\}\} \leq 1$$

- Fixed  $V$  or  $V = n$

## Theorem

*Under a “reasonable” set of assumptions on  $P$ , with probability at least  $1 - \diamond n^{-2}$ ,*

$$\ell(s, \widehat{s}_{\widehat{m}}) \leq \left(1 + \ln(n)^{-1/5}\right) \inf_{m \in \mathcal{M}} \{\ell(s, \widehat{s}_m)\}$$

# Sufficient assumptions

Reminder: *the procedure does not use any of these assumptions.*

- Bounded data:  $\|Y\|_\infty \leq A < \infty$
- Minimal noise-level:

$$0 < \sigma_{\min} \leq \sigma(X)$$

- Smoothness of the regression function  $s$ : non-constant, belongs to some hölderian ball  $\mathcal{H}_\alpha(R)$
- Regularity of the partition:  $\min_\lambda \mathbb{P}(X \in I_\lambda) \geq \diamond D_m^{-1}$

and they can be relaxed...

# Sufficient assumptions

Reminder: *the procedure does not use any of these assumptions.*

- Bounded data:  $\|Y\|_\infty \leq A < \infty$
- Minimal noise-level:

$$0 < \sigma_{\min} \leq \sigma(X)$$

- Smoothness of the regression function  $s$ : non-constant, belongs to some hölderian ball  $\mathcal{H}_\alpha(R)$
- Regularity of the partition:  $\min_\lambda \mathbb{P}(X \in I_\lambda) \geq \diamond D_m^{-1}$

*and they can be relaxed...*



## Corollaries

- Classical oracle inequality:

$$\mathbb{E} [\ell(s, \widehat{s}_{\widehat{m}})] \leq \left(1 + \ln(n)^{-1/5}\right) \mathbb{E} \left[ \inf_{m \in \mathcal{M}} \{\ell(s, \widehat{s}_m)\} \right] + \diamond n^{-2}$$

- **Asymptotic optimality** if  $C \sim_{n \rightarrow +\infty} V - 1$  :

$$\frac{\ell(s, \widehat{s}_{\widehat{m}})}{\inf_{m \in \mathcal{M}} \{\ell(s, \widehat{s}_m)\}} \xrightarrow[n \rightarrow +\infty]{a.s.} 1$$

- **Adaptation** to hölderian regularity in an heteroscedastic framework (regular histograms):  
 $s \in \mathcal{H}(\alpha, R)$ ,  $\alpha \in (0, 1]$ ,  $\mathcal{X} \subset \mathbb{R}^k$ ,  $(\dots) \mathcal{Y}$  bounded

$$\Rightarrow \text{rate } \|\sigma\|_{L^2(\text{Leb})}^{\frac{4\alpha}{2\alpha+k}} R^{\frac{2k}{2\alpha+k}} n^{\frac{-2\alpha}{2\alpha+k}} .$$

# Simulation framework

$$Y_i = s(X_i) + \sigma(X_i)\epsilon_i \quad X_i \sim^{\text{i.i.d.}} \mathcal{U}([0; 1]) \quad \epsilon_i \sim^{\text{i.i.d.}} \mathcal{N}(0, 1)$$

$$\mathcal{M}_n = \left\{ \text{regular histograms with } D \text{ pieces, } 1 \leq D \leq \frac{n}{\log(n)} \right. \\ \left. \text{and s.t. } \min_{\lambda \in \Lambda_m} \text{Card}\{X_i \in I_\lambda\} \geq 2 \right\}$$

⇒ Benchmark:

$$C_{\text{classical}} = \frac{\mathbb{E}[\ell(s, \hat{s}_m)]}{\mathbb{E}[\inf_{m \in \mathcal{M}} \ell(s, \hat{s}_m)]} \quad \text{computed with } N = 1000 \text{ samples}$$

# Model selection methods

- Mallows:

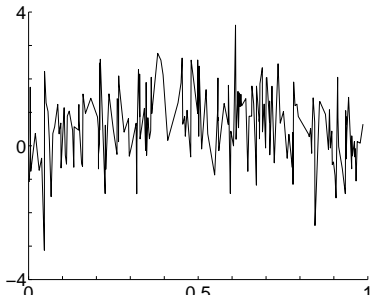
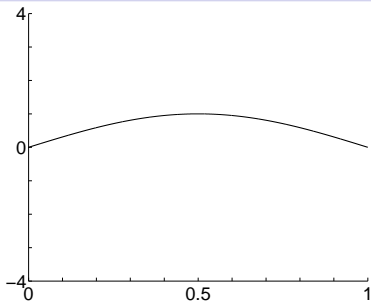
$$\text{pen}(m) = 2\hat{\sigma}^2 D_m n^{-1}$$

- “Classical”  $V$ -fold cross-validation ( $V \in \{2, 5, 10, 20, n\}$ ):

$$\hat{m} \in \arg \min_{m \in \mathcal{M}} \left\{ \frac{1}{V} \sum_{j=1}^V P_n^j \gamma \left( \hat{s}_m^{(-j)}, \cdot \right) \right\} \quad \tilde{s} = \hat{s}_{\hat{m}}$$

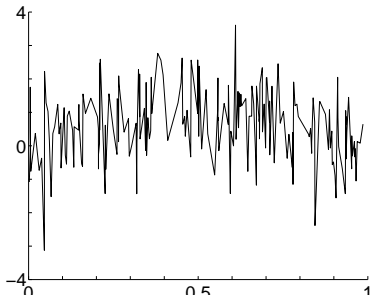
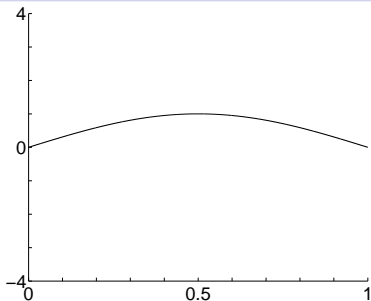
- $V$ -fold penalties ( $V \in \{2, 5, 10, n\}$ ),  $C = V - 1$

# Simulations: $s(x) = \sin(\pi x)$ , $n = 200$ , $\sigma \equiv 1$



<b>Mallows</b>	<b><math>1.93 \pm 0.04</math></b>
2-fold	$2.08 \pm 0.04$
5-fold	$2.14 \pm 0.04$
10-fold	$2.10 \pm 0.05$
20-fold	$2.09 \pm 0.04$
leave-one-out	$2.08 \pm 0.04$
pen 2-f	$2.58 \pm 0.06$
pen 5-f	$2.22 \pm 0.05$
pen 10-f	$2.12 \pm 0.05$
pen Loo	$2.08 \pm 0.05$

# Simulations: $s(x) = \sin(\pi x)$ , $n = 200$ , $\sigma \equiv 1$

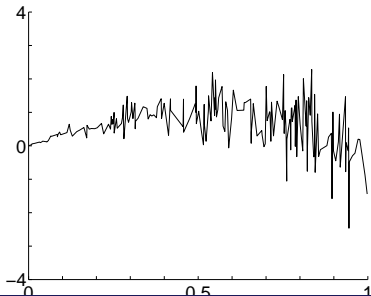
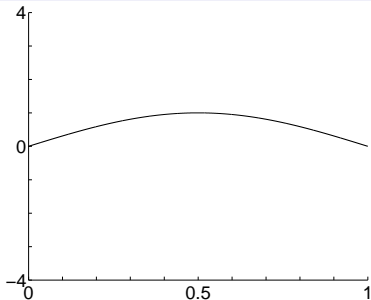


Mallows	$1.93 \pm 0.04$
2-fold	$2.08 \pm 0.04$
5-fold	$2.14 \pm 0.04$
10-fold	$2.10 \pm 0.05$
20-fold	$2.09 \pm 0.04$
leave-one-out	$2.08 \pm 0.04$

pen 2-f	$2.58 \pm 0.06$
pen 5-f	$2.22 \pm 0.05$
pen 10-f	$2.12 \pm 0.05$
pen Loo	$2.08 \pm 0.05$

<b>Mallows <math>\times 1.25</math></b>	<b><math>1.80 \pm 0.03</math></b>
pen 2-f $\times 1.25$	$2.17 \pm 0.05$
pen 5-f $\times 1.25$	$1.91 \pm 0.05$
pen 10-f $\times 1.25$	$1.87 \pm 0.03$
<b>pen Loo <math>\times 1.25</math></b>	<b><math>1.84 \pm 0.03</math></b>

# Simulations: $\sin$ , $n = 200$ , $\sigma(x) = x$ , 2 bin sizes



Mallows	$3.69 \pm 0.07$
2-fold	$2.54 \pm 0.05$
5-fold	$2.58 \pm 0.06$
10-fold	$2.60 \pm 0.06$
20-fold	$2.58 \pm 0.06$
leave-one-out	$2.59 \pm 0.06$

pen 2-f	$3.06 \pm 0.07$
pen 5-f	$2.75 \pm 0.06$
pen 10-f	$2.65 \pm 0.06$
pen Loo	$2.59 \pm 0.06$

Mallows $\times 1.25$	$3.17 \pm 0.07$
pen 2-f $\times 1.25$	$2.75 \pm 0.06$
pen 5-f $\times 1.25$	$2.38 \pm 0.06$
pen 10-f $\times 1.25$	$2.28 \pm 0.05$
pen Loo $\times 1.25$	$2.21 \pm 0.05$

25/26

# Conclusions on $V$ -fold penalization

- **asymptotically optimal**, even if  $V$  fixed
- **optimal  $V^*$** : the largest possible one  
⇒ easier to balance with the computational cost
- low signal-to-noise ratio ⇒ easy to **overpenalize and decrease variability** (keep  $V$  large)
- large signal-to-noise ratio ⇒ possible to stay **unbiased with a small  $V$**  (for computational reasons)
- **flexibility** improves  $V$ -fold cross-validation (according to both **theoretical** results and **simulations**)
- theory can be extended to **exchangeable weighted bootstrap penalties** (e.g. bootstrap, i.i.d. Rademacher, leave-one-out, leave- $p$ -out with  $p = \alpha n$ ).
- Some open problems: consistency when  $C \gg V - 1$ , prediction in a general framework, automatic choice of the **overpenalization constant**.

# Conclusions on $V$ -fold penalization

- **asymptotically optimal**, even if  $V$  fixed
- **optimal  $V^*$** : the largest possible one  
⇒ easier to balance with the computational cost
- low signal-to-noise ratio ⇒ easy to **overpenalize and decrease variability** (keep  $V$  large)
- large signal-to-noise ratio ⇒ possible to stay **unbiased with a small  $V$**  (for computational reasons)
  
- **flexibility** improves  $V$ -fold cross-validation (according to both **theoretical** results and **simulations**)
- theory can be extended to **exchangeable weighted bootstrap penalties** (e.g. bootstrap, i.i.d. Rademacher, leave-one-out, leave- $p$ -out with  $p = \alpha n$ ).
- Some open problems: consistency when  $C \gg V - 1$ , prediction in a general framework, automatic choice of the **overpenalization constant**.



# Conclusions on $V$ -fold penalization

- **asymptotically optimal**, even if  $V$  fixed
- **optimal  $V^*$** : the largest possible one  
⇒ easier to balance with the computational cost
- low signal-to-noise ratio ⇒ easy to **overpenalize and decrease variability** (keep  $V$  large)
- large signal-to-noise ratio ⇒ possible to stay **unbiased with a small  $V$**  (for computational reasons)
  
- **flexibility** improves  $V$ -fold cross-validation (according to both **theoretical** results and **simulations**)
- theory can be extended to **exchangeable weighted bootstrap penalties** (e.g. bootstrap, i.i.d. Rademacher, leave-one-out, leave- $p$ -out with  $p = \alpha n$ ).
- Some open problems: consistency when  $C \gg V - 1$ , prediction in a general framework, automatic choice of the **overpenalization constant**.

# Conclusions on $V$ -fold penalization

- **asymptotically optimal**, even if  $V$  fixed
- **optimal  $V^*$** : the largest possible one  
⇒ easier to balance with the computational cost
- low signal-to-noise ratio ⇒ easy to **overpenalize and decrease variability** (keep  $V$  large)
- large signal-to-noise ratio ⇒ possible to stay **unbiased with a small  $V$**  (for computational reasons)
  
- **flexibility** improves  $V$ -fold cross-validation (according to both **theoretical** results and **simulations**)
- theory can be extended to **exchangeable weighted bootstrap penalties** (e.g. bootstrap, i.i.d. Rademacher, leave-one-out, leave- $p$ -out with  $p = \alpha n$ ).
- Some open problems: consistency when  $C \gg V - 1$ , prediction in a general framework, automatic choice of the overpenalization constant.

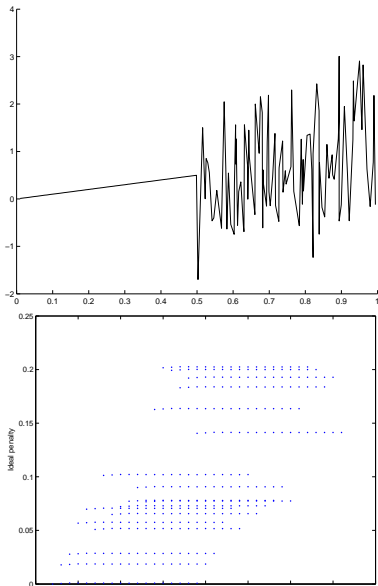
Thank you for your attention !

Preprint: [arXiv:0802.0566](https://arxiv.org/abs/0802.0566)

## Part I

# Appendix

# Limitations of a linear penalty

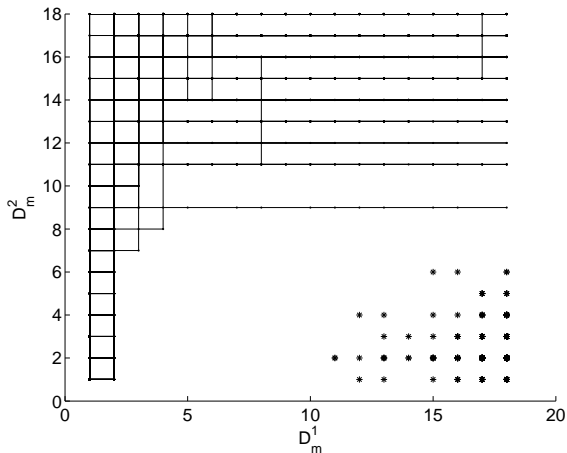


$$Y = X + \sigma(X)\epsilon$$

$$\sigma(X) = \mathbb{1}_{X \geq \frac{1}{2}} \quad \epsilon \sim \mathcal{N}(0, 1)$$

Regular histograms on  $[0; \frac{1}{2}]$  ( $D_{m,1}$  pieces), then regular histograms on  $[\frac{1}{2}; 1]$  ( $D_{m,2}$  pieces).

$\Rightarrow \text{pen}_{\text{id}}(m)$  is not a linear function of  $D_m$ .

Limitations of a linear penalty:  $\hat{m}(K) \neq m^*$ 

# Sketch of the proof

For each  $m \in \mathcal{M}_n$ ,

$$\text{pen}_{\text{id}}(m) \approx \mathbb{E}[\text{pen}_{\text{id}}(m)] \propto \mathbb{E}[\text{pen}(m)] \approx \text{pen}(m)$$

with remainders  $\ll \ell(s, \hat{s}_m)$  when  $D_m \rightarrow +\infty$ :

- Explicit computation of  $\text{pen}_{\text{id}}$  and  $\text{pen}$
- Comparison of expectations:  $\mathbb{E}(\text{pen}_{\text{id}}) \propto \mathbb{E}(\text{pen})$  (if  $\min_{\lambda \in \Lambda_m} \{n\mathbb{P}(X \in I_\lambda)\} \rightarrow +\infty$ )
- Moment inequalities (Boucheron, Bousquet, Lugosi, Massart 2003)  
 $\Rightarrow$  concentration inequalities (for  $\text{pen}_{\text{id}}$  and  $\text{pen}$ )
- Assumptions  $\Rightarrow$  control of the remainders in terms of  $\ell(s, \hat{s}_m)$ .

# Sketch of the proof

For each  $m \in \mathcal{M}_n$ ,

$$\text{pen}_{\text{id}}(m) \approx \mathbb{E}[\text{pen}_{\text{id}}(m)] \propto \mathbb{E}[\text{pen}(m)] \approx \text{pen}(m)$$

with remainders  $\ll \ell(s, \hat{s}_m)$  when  $D_m \rightarrow +\infty$ :

- **Explicit computation of  $\text{pen}_{\text{id}}$  and  $\text{pen}$**
- Comparison of expectations:  $\mathbb{E}(\text{pen}_{\text{id}}) \propto \mathbb{E}(\text{pen})$  (if  $\min_{\lambda \in \Lambda_m} \{n\mathbb{P}(X \in I_\lambda)\} \rightarrow +\infty$ )
- Moment inequalities (Boucheron, Bousquet, Lugosi, Massart 2003)  
 $\Rightarrow$  concentration inequalities (for  $\text{pen}_{\text{id}}$  and  $\text{pen}$ )
- Assumptions  $\Rightarrow$  control of the remainders in terms of  $\ell(s, \hat{s}_m)$ .



# Sketch of the proof

For each  $m \in \mathcal{M}_n$ ,

$$\text{pen}_{\text{id}}(m) \approx \mathbb{E}[\text{pen}_{\text{id}}(m)] \propto \mathbb{E}[\text{pen}(m)] \approx \text{pen}(m)$$

with remainders  $\ll \ell(s, \hat{s}_m)$  when  $D_m \rightarrow +\infty$ :

- Explicit computation of  $\text{pen}_{\text{id}}$  and  $\text{pen}$
- Comparison of expectations:  $\mathbb{E}(\text{pen}_{\text{id}}) \propto \mathbb{E}(\text{pen})$  (if  $\min_{\lambda \in \Lambda_m} \{n\mathbb{P}(X \in I_\lambda)\} \rightarrow +\infty$ )
- Moment inequalities (Boucheron, Bousquet, Lugosi, Massart 2003)  
 $\Rightarrow$  concentration inequalities (for  $\text{pen}_{\text{id}}$  and  $\text{pen}$ )
- Assumptions  $\Rightarrow$  control of the remainders in terms of  $\ell(s, \hat{s}_m)$ .

# Sketch of the proof

For each  $m \in \mathcal{M}_n$ ,

$$\text{pen}_{\text{id}}(m) \approx \mathbb{E}[\text{pen}_{\text{id}}(m)] \propto \mathbb{E}[\text{pen}(m)] \approx \text{pen}(m)$$

with remainders  $\ll \ell(s, \hat{s}_m)$  when  $D_m \rightarrow +\infty$ :

- Explicit computation of  $\text{pen}_{\text{id}}$  and  $\text{pen}$
- Comparison of expectations:  $\mathbb{E}(\text{pen}_{\text{id}}) \propto \mathbb{E}(\text{pen})$  (if  $\min_{\lambda \in \Lambda_m} \{n\mathbb{P}(X \in I_\lambda)\} \rightarrow +\infty$ )
- **Moment inequalities (Boucheron, Bousquet, Lugosi, Massart 2003)**  
 $\Rightarrow$  **concentration inequalities (for  $\text{pen}_{\text{id}}$  and  $\text{pen}$ )**
- Assumptions  $\Rightarrow$  control of the remainders in terms of  $\ell(s, \hat{s}_m)$ .

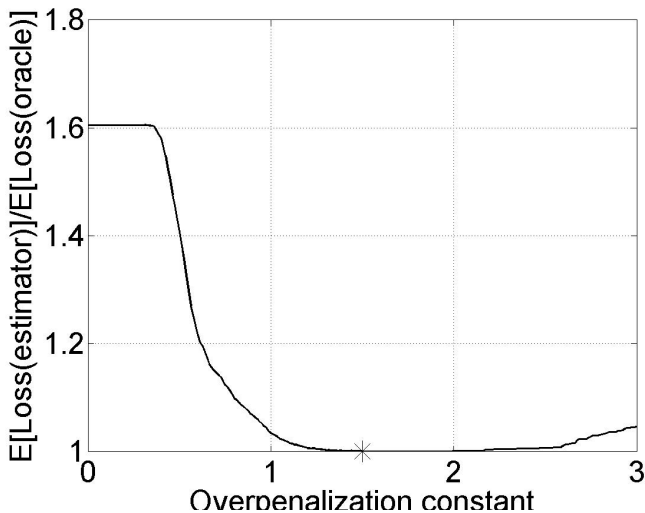
# Sketch of the proof

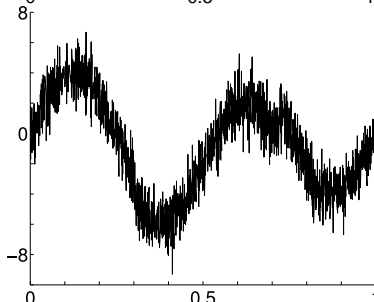
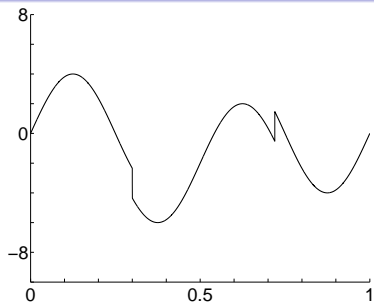
For each  $m \in \mathcal{M}_n$ ,

$$\text{pen}_{\text{id}}(m) \approx \mathbb{E}[\text{pen}_{\text{id}}(m)] \propto \mathbb{E}[\text{pen}(m)] \approx \text{pen}(m)$$

with remainders  $\ll \ell(s, \widehat{s}_m)$  when  $D_m \rightarrow +\infty$ :

- Explicit computation of  $\text{pen}_{\text{id}}$  and  $\text{pen}$
- Comparison of expectations:  $\mathbb{E}(\text{pen}_{\text{id}}) \propto \mathbb{E}(\text{pen})$  (if  $\min_{\lambda \in \Lambda_m} \{n\mathbb{P}(X \in I_\lambda)\} \rightarrow +\infty$ )
- Moment inequalities (Boucheron, Bousquet, Lugosi, Massart 2003)  
 $\Rightarrow$  concentration inequalities (for  $\text{pen}_{\text{id}}$  and  $\text{pen}$ )
- Assumptions  $\Rightarrow$  control of the remainders in terms of  $\ell(s, \widehat{s}_m)$ .

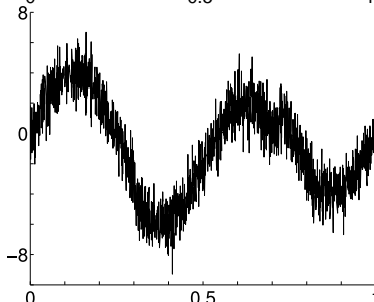
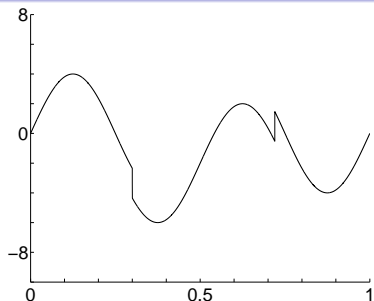
Overpenalization (HeaviSine,  $n = 2048$ ,  $\sigma \equiv 1$ )

Simulations: HeaviSine,  $n = 2048$ ,  $\sigma \equiv 1$ 

Mallows	$1.015 \pm 0.003$
<b>2-fold</b>	<b><math>1.002 \pm 0.003</math></b>
5-fold	$1.014 \pm 0.003$
10-fold	$1.021 \pm 0.003$
20-fold	$1.029 \pm 0.004$
leave-one-out	$1.034 \pm 0.004$

pen 2-f	$1.038 \pm 0.004$
pen 5-f	$1.037 \pm 0.004$
pen 10-f	$1.034 \pm 0.004$
pen Loo	$1.034 \pm 0.004$

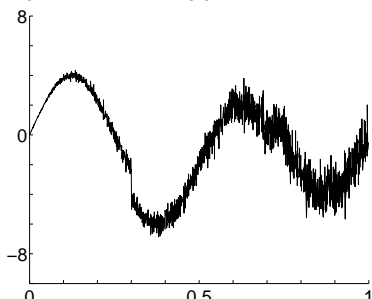
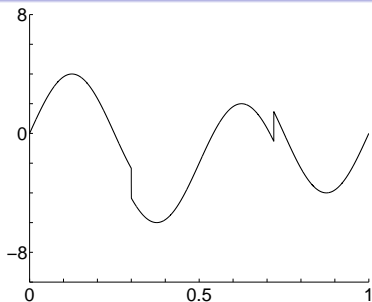
Mallows $\times 1.25$	$1.002 \pm 0.003$
pen 2-f $\times 1.25$	$1.011 \pm 0.003$
pen 5-f $\times 1.25$	$1.006 \pm 0.003$
pen 10-f $\times 1.25$	$1.005 \pm 0.003$
pen Loo $\times 1.25$	$1.004 \pm 0.003$

Simulations: HeaviSine,  $n = 2048$ ,  $\sigma \equiv 1$ 

Mallows	$1.015 \pm 0.003$
2-fold	$1.002 \pm 0.003$
5-fold	$1.014 \pm 0.003$
10-fold	$1.021 \pm 0.003$
20-fold	$1.029 \pm 0.004$
leave-one-out	$1.034 \pm 0.004$

pen 2-f	$1.038 \pm 0.004$
pen 5-f	$1.037 \pm 0.004$
pen 10-f	$1.034 \pm 0.004$
pen Loo	$1.034 \pm 0.004$

Mallows $\times 1.25$	$1.002 \pm 0.003$
pen 2-f $\times 1.25$	$1.011 \pm 0.003$
pen 5-f $\times 1.25$	$1.006 \pm 0.003$
pen 10-f $\times 1.25$	$1.005 \pm 0.003$
pen Loo $\times 1.25$	$1.004 \pm 0.003$

Simulations: HeaviSine,  $n = 2048$ ,  $\sigma(x) = x$ , 2 bin sizes

Mallows	$1.373 \pm 0.010$
2-fold	$1.184 \pm 0.004$
5-fold	$1.115 \pm 0.005$
10-fold	$1.109 \pm 0.004$
20-fold	$1.105 \pm 0.004$
leave-one-out	$1.105 \pm 0.004$
pen 2-f	$1.103 \pm 0.005$
pen 5-f	$1.104 \pm 0.004$
pen 10-f	$1.104 \pm 0.004$
pen Loo	$1.105 \pm 0.004$
Mallows $\times 1.25$	$1.411 \pm 0.008$
pen 2-f $\times 1.25$	$1.106 \pm 0.004$
pen 5-f $\times 1.25$	$1.102 \pm 0.004$
pen 10-f $\times 1.25$	$1.098 \pm 0.004$
pen Loo $\times 1.25$	$1.096 \pm 0.004$

Simulations: sin, variable  $n$  and  $\sigma$ , regular histograms

$n$	200	1000	200
$\sigma$	1	1	0.1
Mallows ( $K = 2$ )	1.93 ± 0.04	1.67 ± 0.04	1.40 ± 0.02
2-fold	2.08 ± 0.04	1.67 ± 0.04	1.39 ± 0.02
10-fold	2.10 ± 0.05	1.75 ± 0.04	1.38 ± 0.02
pen 10-fold	2.12 ± 0.05	1.78 ± 0.05	1.37 ± 0.02



Simulations: sin, variable  $n$  and  $\sigma$ , regular histograms

$n$	200	1000	200
$\sigma$	1	1	0.1
Mallows ( $K = 2$ )	1.93 $\pm$ 0.04	1.67 $\pm$ 0.04	1.40 $\pm$ 0.02
2-fold	2.08 $\pm$ 0.04	1.67 $\pm$ 0.04	1.39 $\pm$ 0.02
10-fold	2.10 $\pm$ 0.05	1.75 $\pm$ 0.04	1.38 $\pm$ 0.02
pen 10-fold	2.12 $\pm$ 0.05	1.78 $\pm$ 0.05	1.37 $\pm$ 0.02
Mallows ( $K = 2.5$ )	1.80 $\pm$ 0.03	1.62 $\pm$ 0.03	1.43 $\pm$ 0.02
pen 10-fold $\times 1.25$	1.87 $\pm$ 0.03	1.63 $\pm$ 0.04	1.38 $\pm$ 0.02