

## V-fold selection of kernel estimators

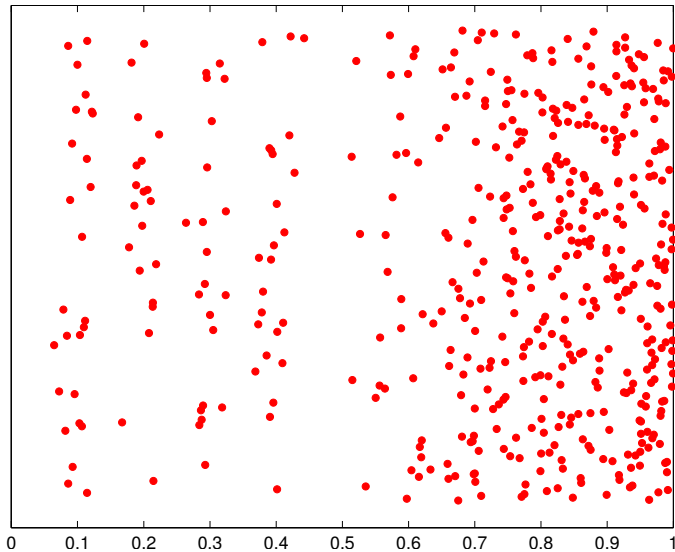
Sylvain Arlot (joint work with Matthieu Lerasle & Nelo Magalhães)

<sup>1</sup>CNRS

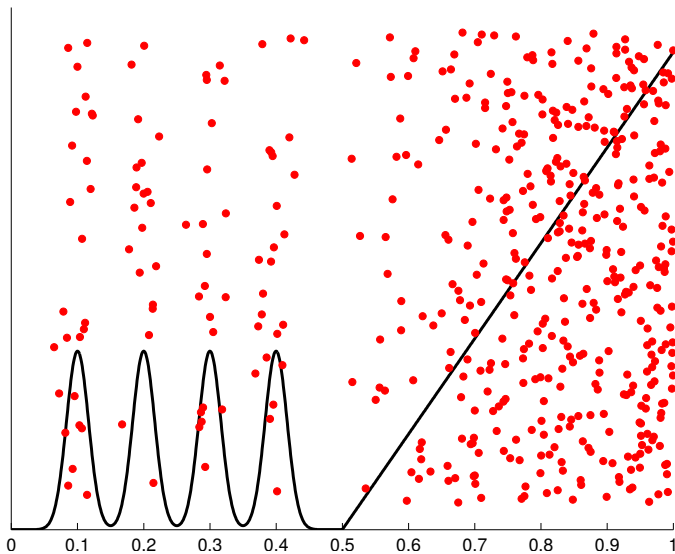
<sup>2</sup>École Normale Supérieure (Paris), DI/ENS, Équipe SIERRA

EMS 2015, Amsterdam, July 6th, 2015

# Density estimation: data $X_1, \dots, X_n$



Goal: estimate the common density  $f^*$  of  $X_i$



# Problem: density estimation

- **Data**  $D_n$ :  $X_1, \dots, X_n \in \mathbb{X}$  (i.i.d.  $\sim P$ , density  $f^*$  w.r.t.  $\mu$ )  
Assumption:  $f^* \in L^\infty$
- **Least-squares loss**  $\gamma(t, x) = \|t\|_{L^2(\mu)}^2 - 2t(x)$
- **Goal**: learn  $t \in \mathbb{S} = \{\text{measurable functions } \mathbb{X} \rightarrow \mathbb{R}\}$  s.t.  
 $\mathbb{E}_{X \sim P}[\gamma(t; X)] =: P\gamma(t)$  is minimal.

# Problem: density estimation

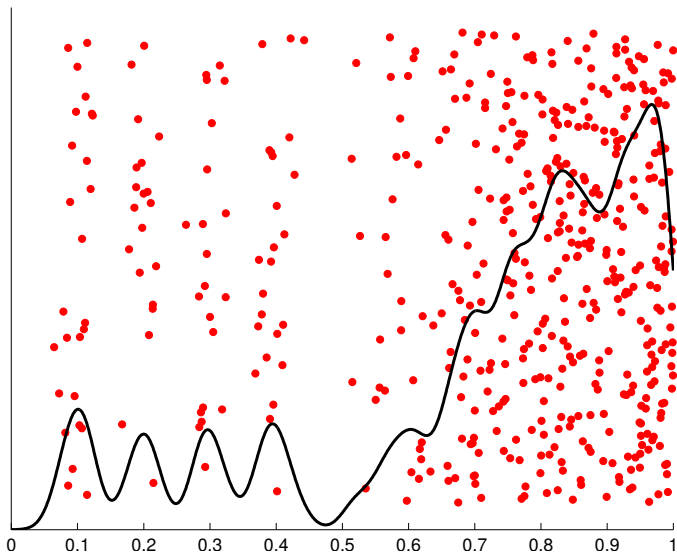
- **Data**  $D_n$ :  $X_1, \dots, X_n \in \mathbb{X}$  (i.i.d.  $\sim P$ , density  $f^*$  w.r.t.  $\mu$ )  
Assumption:  $f^* \in L^\infty$
- **Least-squares loss**  $\gamma(t, x) = \|t\|_{L^2(\mu)}^2 - 2t(x)$
- **Goal**: learn  $t \in \mathbb{S} = \{\text{measurable functions } \mathbb{X} \rightarrow \mathbb{R}\}$  s.t.  $\mathbb{E}_{X \sim P}[\gamma(t; X)] =: P\gamma(t)$  is minimal.

$$P\gamma(t) = \int t^2 d\mu - 2 \int tf^* d\mu = \int (t - f^*)^2 d\mu - \|f^*\|_{L^2(\mu)}^2$$

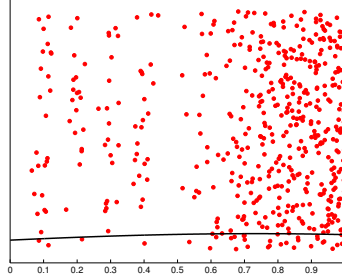
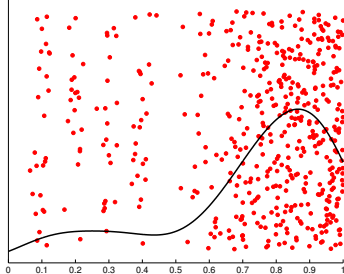
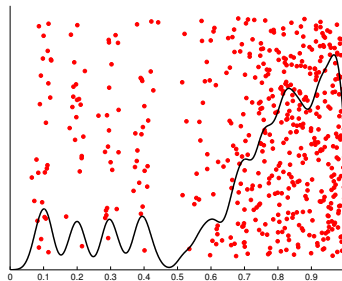
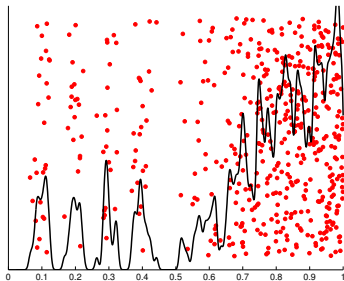
$\Rightarrow$  the true density  $f^* \in \operatorname{argmin}_{t \in \mathbb{S}} P\gamma(t)$  and the **excess risk** is

$$P\gamma(t) - P\gamma(f^*) = \|t - f^*\|_{L^2(\mu)}^2$$

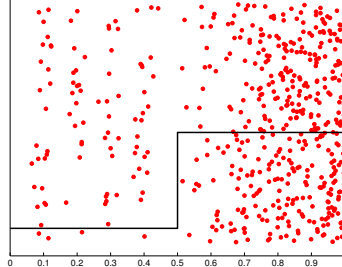
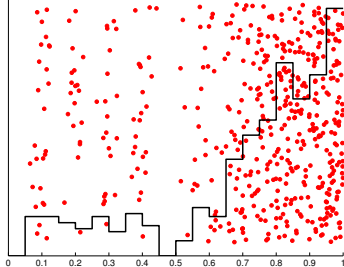
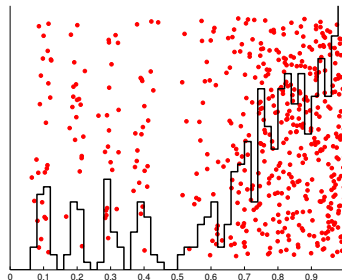
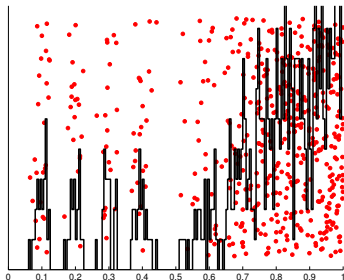
# Estimators: example: Parzen, Gaussian kernel



# Estimator selection: Parzen, Gaussian kernel



# Estimator selection: regular histograms





# Selection of linear estimators

- **Linear estimator** (a.k.a. additive / delta-sequence estimator):  
 $P_n \mapsto \hat{f}_m$  linear, i.e.,

$$\hat{f}_m(D_n) : x \mapsto \frac{1}{n} \sum_{i=1}^n K_m(x, X_i) .$$

$K_m : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$  symmetric (the **kernel**).

Assumptions in this talk:  $K_m$  bounded and  $(K_m(x, \cdot))_{x \in \mathbb{X}}$  bounded in  $L^2(\mu)$ .

# Selection of linear estimators

- Linear estimator (a.k.a. additive / delta-sequence estimator):  
 $P_n \mapsto \hat{f}_m$  linear, i.e.,

$$\hat{f}_m(D_n) : x \mapsto \frac{1}{n} \sum_{i=1}^n K_m(x, X_i) .$$

$K_m : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$  symmetric (the kernel).

Assumptions in this talk:  $K_m$  bounded and  $(K_m(x, \cdot))_{x \in \mathbb{X}}$  bounded in  $L^2(\mu)$ .

- Estimator collection  $(\hat{f}_m)_{m \in \mathcal{M}} \Rightarrow$  choose  $\hat{m} = \hat{m}(D_n)$ ?  
i.e., family of kernels  $(K_m)_{m \in \mathcal{M}} \Rightarrow$  family of linear estimators

# Selection of linear estimators

- Linear estimator (a.k.a. additive / delta-sequence estimator):  
 $P_n \mapsto \hat{f}_m$  linear, i.e.,

$$\hat{f}_m(D_n) : x \mapsto \frac{1}{n} \sum_{i=1}^n K_m(x, X_i) .$$

$K_m : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$  symmetric (the kernel).

Assumptions in this talk:  $K_m$  bounded and  $(K_m(x, \cdot))_{x \in \mathbb{X}}$  bounded in  $L^2(\mu)$ .

- Estimator collection  $(\hat{f}_m)_{m \in \mathcal{M}} \Rightarrow$  choose  $\hat{m} = \hat{m}(D_n)$ ?  
i.e., family of kernels  $(K_m)_{m \in \mathcal{M}} \Rightarrow$  family of linear estimators
- Goal: minimize the risk, i.e.,  
**Oracle inequality** (in expectation or with a large probability):

$$\|\hat{f}_{\hat{m}} - f^*\|^2 \leq C \inf_{m \in \mathcal{M}} \left\{ \|\hat{f}_m - f^*\|^2 \right\} + R_n$$

# Linear estimators: examples

- Projection on  $S_m = \text{span}(\psi_\lambda)_{\lambda \in \Lambda_m}$  (orthonormal in  $L^2(\mu)$ )

$$\hat{f}_m \in \underset{t \in S_m}{\text{argmin}} \{P_n \gamma(t)\} \quad \text{where} \quad P_n \gamma(t) := \frac{1}{n} \sum_{x \in D_n} \gamma(t; x)$$

$$K_m(x, y) = \sum_{\lambda \in \Lambda_m} \psi_\lambda(x) \psi_\lambda(y) .$$

# Linear estimators: examples

- Projection on  $S_m = \text{span}(\psi_\lambda)_{\lambda \in \Lambda_m}$  (orthonormal in  $L^2(\mu)$ )

$$\hat{f}_m \in \underset{t \in S_m}{\text{argmin}} \{P_n \gamma(t)\} \quad \text{where} \quad P_n \gamma(t) := \frac{1}{n} \sum_{x \in D_n} \gamma(t; x)$$

$$K_m(x, y) = \sum_{\lambda \in \Lambda_m} \psi_\lambda(x) \psi_\lambda(y) .$$

- Weighted projection estimator (e.g., Pinsker):

$$K_m(x, y) = \sum_{\lambda \in \Lambda_m} m_\lambda \psi_\lambda(x) \psi_\lambda(y) \quad \text{with} \quad m_\lambda \in [0, 1] .$$

# Linear estimators: examples

- Projection on  $S_m = \text{span}(\psi_\lambda)_{\lambda \in \Lambda_m}$  (orthonormal in  $L^2(\mu)$ )

$$\hat{f}_m \in \operatorname{argmin}_{t \in S_m} \{P_n \gamma(t)\} \quad \text{where} \quad P_n \gamma(t) := \frac{1}{n} \sum_{x \in D_n} \gamma(t; x)$$

$$K_m(x, y) = \sum_{\lambda \in \Lambda_m} \psi_\lambda(x) \psi_\lambda(y) .$$

- Weighted projection estimator (e.g., Pinsker):

$$K_m(x, y) = \sum_{\lambda \in \Lambda_m} m_\lambda \psi_\lambda(x) \psi_\lambda(y) \quad \text{with} \quad m_\lambda \in [0, 1] .$$

- Kernel estimator on  $\mathbb{X} = \mathbb{R}$  (or  $\mathbb{R}^d$ ), Parzen-Rosenblatt:

$$K_{(k,h)}(x, y) = \frac{1}{h} k\left(\frac{x-y}{h}\right) \quad \text{with } k \in L^2(\mu) \text{ symmetric and } h > 0 .$$

Example: Gaussian kernel

$$k(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$$

# Expectation of the loss

$$\mathbb{E} \left[ \left\| \hat{f}_m - f^* \right\|^2 \right] = \underbrace{\left\| f_m^* - f^* \right\|^2}_{\text{approximation error}} + \underbrace{\frac{1}{n} \mathbb{E} [A_m(X, X) - A_m(X, Y)]}_{\text{estimation error}}$$

$$f_m^*(x) = \mathbb{E} [K_m(x, X)]$$

$$A_m(x, y) = \int_{\mathbb{X}} K_m(x, z) K_m(y, z) d\mu(z)$$

# Expectation of the loss

$$\mathbb{E} \left[ \left\| \hat{f}_m - f^* \right\|^2 \right] = \underbrace{\left\| f_m^* - f^* \right\|^2}_{\text{approximation error}} + \underbrace{\frac{1}{n} \mathbb{E} [A_m(X, X) - A_m(X, Y)]}_{\text{estimation error}}$$

$$f_m^*(x) = \mathbb{E} [K_m(x, X)]$$

$$A_m(x, y) = \int_{\mathbb{X}} K_m(x, z) K_m(y, z) d\mu(z)$$

- Projection estimators:

$$f_m^* \in \operatorname{argmin}_{t \in S_m} \|t - f^*\|^2 \quad \text{and} \quad A_m = K_m .$$

Regular histograms on  $\mathbb{R}$  with bin size  $d_m^{-1}$ : estim. error  $\approx \frac{d_m}{n}$



# Expectation of the loss

$$\mathbb{E} \left[ \left\| \hat{f}_m - f^* \right\|^2 \right] = \underbrace{\left\| f_m^* - f^* \right\|^2}_{\text{approximation error}} + \underbrace{\frac{1}{n} \mathbb{E} [A_m(X, X) - A_m(X, Y)]}_{\text{estimation error}}$$

$$f_m^*(x) = \mathbb{E} [K_m(x, X)]$$

$$A_m(x, y) = \int_{\mathbb{X}} K_m(x, z) K_m(y, z) d\mu(z)$$

- Projection estimators:

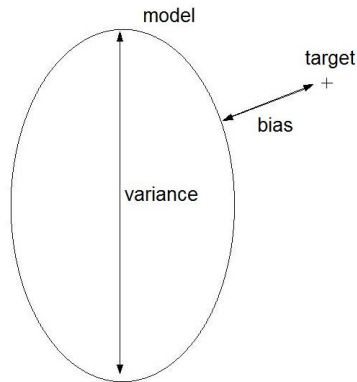
$$f_m^* \in \operatorname{argmin}_{t \in \mathcal{S}_m} \|t - f^*\|^2 \quad \text{and} \quad A_m = K_m .$$

Regular histograms on  $\mathbb{R}$  with bin size  $d_m^{-1}$ : estim. error  $\approx \frac{d_m}{n}$

- Parzen:  $A_{(k,h)}(x, y) = \frac{1}{h} (k * k) \left( \frac{x-y}{h} \right)$ ; estim. error  $\approx \frac{\|k\|_{L^2(\mu)}^2}{nh}$

# Bias-variance trade-off

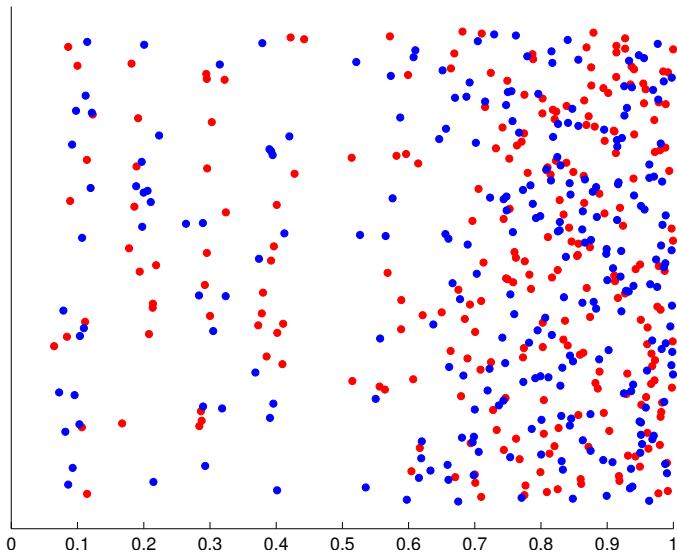
$$\begin{aligned} \mathbb{E} \left[ \left\| \hat{f}_m - f^* \right\|^2 \right] &= \text{Approx. error} \\ &\quad + \text{Estim. error} \\ &= \text{Bias} + \text{Variance} \end{aligned}$$



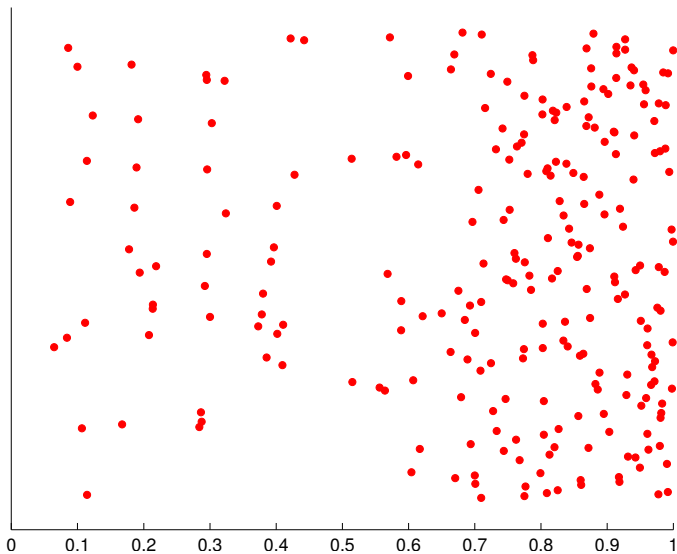
**Bias-variance trade-off**

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

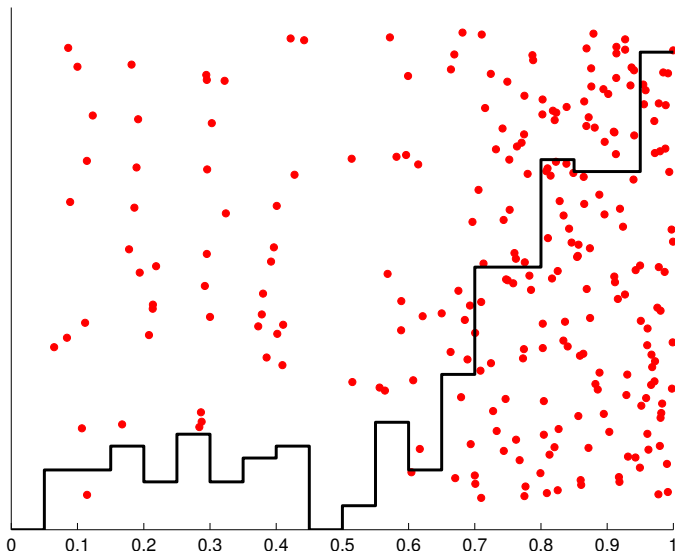
# Validation principle



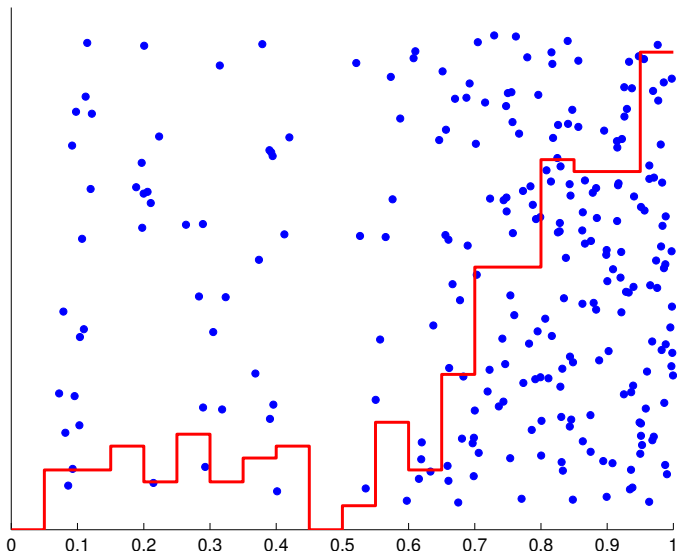
# Validation principle: training sample



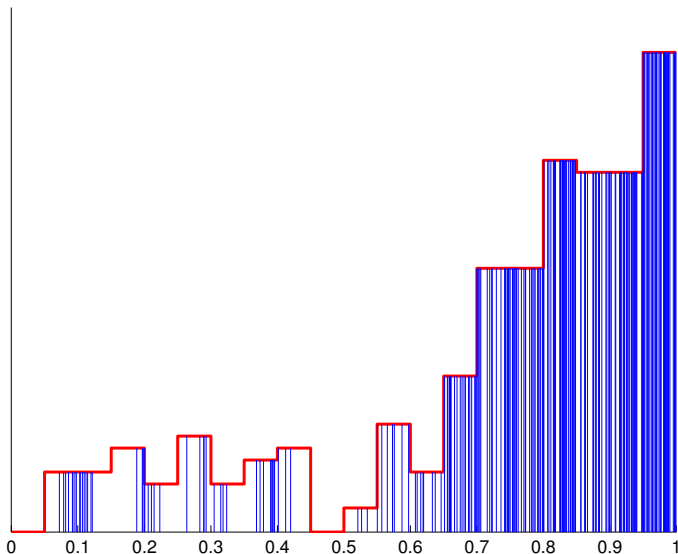
# Validation principle: training sample



# Validation principle: validation sample



# Validation principle: validation sample



# V-fold cross-validation

$$\underbrace{X_1, \dots, X_q}_{\text{Training}}, \underbrace{X_{q+1}, \dots, X_n}_{\text{Validation}}$$

$$\hat{f}_m^{(t)} = \hat{f}_m((X_i)_{1 \leq i \leq q}) = \frac{1}{q} \sum_{i=1}^q K_m(\cdot, X_i)$$

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

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

$$\Rightarrow \hat{\mathcal{R}}^{\text{vf}}(\hat{f}_m; D_n; \mathcal{B}) = \frac{1}{v} \sum_{j=1}^v P_n^j \gamma(\hat{f}_m^{(-j)}) \quad \hat{m} \in \arg \min_{m \in \mathcal{M}} \left\{ \hat{\mathcal{R}}^{\text{vf}}(\hat{f}_m) \right\}$$



# Expectation of cross-validation criteria

- In this talk, we always assume  $\text{Card}(B_j) = n/V$  for all  $j$ .

# Expectation of cross-validation criteria

- In this talk, we always assume  $\text{Card}(B_j) = n/V$  for all  $j$ .
- **Ideal criterion:**  $P_\gamma(\hat{f}_m)$
- General analysis for the bias:

$$\begin{aligned} \mathbb{E} \left[ P_\gamma \left( \hat{f}_m(D_n) \right) \right] &= \alpha(m) + \frac{\beta(m)}{n} \\ \Rightarrow \mathbb{E} \left[ \hat{\mathcal{R}}^{\text{vf}} \left( \hat{f}_m; D_n; \mathcal{B} \right) \right] &= \mathbb{E} \left[ P_n^{(j)} \gamma \left( \hat{f}_m^{(-j)} \right) \right] = \mathbb{E} \left[ P_\gamma \left( \hat{f}_m^{(-j)} \right) \right] \\ &= \alpha(m) + \frac{V}{V-1} \frac{\beta(m)}{n} \end{aligned}$$

# Expectation of cross-validation criteria

- In this talk, we always assume  $\text{Card}(B_j) = n/V$  for all  $j$ .
- **Ideal criterion:**  $P_\gamma(\hat{f}_m)$
- General analysis for the bias:

$$\begin{aligned}\mathbb{E} \left[ P_\gamma \left( \hat{f}_m(D_n) \right) \right] &= \alpha(m) + \frac{\beta(m)}{n} \\ \Rightarrow \mathbb{E} \left[ \hat{\mathcal{R}}^{\text{vf}} \left( \hat{f}_m; D_n; \mathcal{B} \right) \right] &= \mathbb{E} \left[ P_n^{(j)} \gamma \left( \hat{f}_m^{(-j)} \right) \right] = \mathbb{E} \left[ P_\gamma \left( \hat{f}_m^{(-j)} \right) \right] \\ &= \alpha(m) + \frac{V}{V-1} \frac{\beta(m)}{n}\end{aligned}$$

$\Rightarrow$  **bias** decreases with  $V$ , vanishes as  $V \rightarrow \infty$

# Expectation of cross-validation criteria

- In this talk, we always assume  $\text{Card}(B_j) = n/V$  for all  $j$ .
- **Ideal criterion:**  $P_\gamma(\hat{f}_m)$
- General analysis for the bias:

$$\begin{aligned} \mathbb{E} \left[ P_\gamma \left( \hat{f}_m(D_n) \right) \right] &= \alpha(m) + \frac{\beta(m)}{n} \\ \Rightarrow \mathbb{E} \left[ \hat{\mathcal{R}}^{\text{vf}} \left( \hat{f}_m; D_n; \mathcal{B} \right) \right] &= \mathbb{E} \left[ P_n^{(j)} \gamma \left( \hat{f}_m^{(-j)} \right) \right] = \mathbb{E} \left[ P_\gamma \left( \hat{f}_m^{(-j)} \right) \right] \\ &= \alpha(m) + \frac{V}{V-1} \frac{\beta(m)}{n} \end{aligned}$$

$\Rightarrow$  **bias** decreases with  $V$ , vanishes as  $V \rightarrow \infty$

- Same result for the leave- $p$ -out with  $V$  replaced by  $n/p$ .

# Bias and model selection

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \hat{\mathcal{R}}^{\text{vf}}(\hat{f}_m) \right\} \quad \text{vs.} \quad m^* \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ P\gamma(\hat{f}_m(D_n)) \right\}$$

- Perfect ranking among  $(\hat{f}_m)_{m \in \mathcal{M}} \Leftrightarrow \forall m, m' \in \mathcal{M},$

$$\operatorname{sign} \left( \hat{\mathcal{R}}^{\text{vf}}(\hat{f}_m) - \hat{\mathcal{R}}^{\text{vf}}(\hat{f}_{m'}) \right) = \operatorname{sign} \left( P\gamma(\hat{f}_m) - P\gamma(\hat{f}_{m'}) \right)$$

# Bias and model selection

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \hat{\mathcal{R}}^{\text{vf}}(\hat{f}_m) \right\} \quad \text{vs.} \quad m^* \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ P\gamma(\hat{f}_m(D_n)) \right\}$$

- Perfect ranking among  $(\hat{f}_m)_{m \in \mathcal{M}} \Leftrightarrow \forall m, m' \in \mathcal{M}$ ,

$$\operatorname{sign} \left( \hat{\mathcal{R}}^{\text{vf}}(\hat{f}_m) - \hat{\mathcal{R}}^{\text{vf}}(\hat{f}_{m'}) \right) = \operatorname{sign} \left( P\gamma(\hat{f}_m) - P\gamma(\hat{f}_{m'}) \right)$$

- **Key quantities:**

$$\mathbb{E} \left[ P\gamma(\hat{f}_m) - P\gamma(\hat{f}_{m'}) \right] = \alpha(m) - \alpha(m') + \frac{\beta(m) - \beta(m')}{n}$$

$$\mathbb{E} \left[ \hat{\mathcal{R}}^{\text{vf}}(\hat{f}_m) - \hat{\mathcal{R}}^{\text{vf}}(\hat{f}_{m'}) \right] = \alpha(m) - \alpha(m') + \frac{V}{V-1} \frac{\beta(m) - \beta(m')}{n}$$

- **V-fold CV favours  $m$  with smaller complexity  $\beta(m)$**

# Bias-corrected VFCV / V-fold penalization

- **Bias-corrected V-fold CV** (Burman, 1989):

$$\widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{f}_m; D_n; \mathcal{B}) := \widehat{\mathcal{R}}^{\text{vf}}(\widehat{f}_m; D_n; \mathcal{B}) + P_n \gamma(\widehat{f}_m) - \frac{1}{V} \sum_{j=1}^V P_n \gamma(\widehat{f}_m^{(-j)})$$

- Resampling heuristics (Efron, 1983), V-fold subsampling and penalization principle  $\Rightarrow$  **V-fold penalty** (A. 2008)

$$\text{pen}_{\text{VF}}(\widehat{f}_m; D_n; \mathcal{B}) := \frac{V-1}{V} \sum_{j=1}^V (P_n - P_n^{(-j)}) \gamma(\widehat{f}_m^{(-j)})$$

# Bias-corrected VFCV / V-fold penalization

- **Bias-corrected V-fold CV** (Burman, 1989):

$$\widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{f}_m; D_n; \mathcal{B}) := \widehat{\mathcal{R}}^{\text{vf}}(\widehat{f}_m; D_n; \mathcal{B}) + P_n \gamma(\widehat{f}_m) - \frac{1}{V} \sum_{j=1}^V P_n \gamma(\widehat{f}_m^{(-j)})$$

- Resampling heuristics (Efron, 1983), V-fold subsampling and penalization principle  $\Rightarrow$  **V-fold penalty** (A. 2008)

$$\text{pen}_{\text{VF}}(\widehat{f}_m; D_n; \mathcal{B}) := \frac{V-1}{V} \sum_{j=1}^V (P_n - P_n^{(-j)}) \gamma(\widehat{f}_m^{(-j)})$$

- $\widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{f}_m; D_n; \mathcal{B}) = P_n \gamma(\widehat{f}_m(D_n)) + \text{pen}_{\text{VF}}(\widehat{f}_m; D_n; \mathcal{B})$
- Projection estimators:

$$\widehat{\mathcal{R}}^{\text{vf}}(\widehat{f}_m; D_n; \mathcal{B}) = P_n \gamma(\widehat{f}_m(D_n)) + \frac{V-1/2}{V-1} \text{pen}_{\text{VF}}(\widehat{f}_m; D_n; \mathcal{B})$$



# Expectations

Bias-corrected V-fold CV:

$$\begin{aligned} & \mathbb{E} \left[ \widehat{\mathcal{R}}^{\text{vf,corr}} \left( \widehat{f}_m; D_n; \mathcal{B} \right) \right] \\ &= \mathbb{E} \left[ P_\gamma \left( \widehat{f}_m(D_n) \right) \right] \\ &= \|f_m^* - f^*\|^2 + \frac{1}{n} \mathbb{E} [A_m(X, X) - A_m(X, Y)] - \|f^*\|^2 \end{aligned}$$

# Expectations

Bias-corrected V-fold CV:

$$\begin{aligned}
 & \mathbb{E} \left[ \widehat{\mathcal{R}}^{\text{vf,corr}} \left( \widehat{f}_m; D_n; \mathcal{B} \right) \right] \\
 &= \mathbb{E} \left[ P_\gamma \left( \widehat{f}_m(D_n) \right) \right] \\
 &= \|f_m^* - f^*\|^2 + \frac{1}{n} \mathbb{E} [A_m(X, X) - A_m(X, Y)] - \|f^*\|^2
 \end{aligned}$$

V-fold penalization:

$$\begin{aligned}
 & \mathbb{E} \left[ P_n \gamma \left( \widehat{f}_m(D_n) \right) + C \text{pen}_{\text{VF}} \left( \widehat{f}_m; D_n; \mathcal{B} \right) \right] \\
 &= \mathbb{E} \left[ P_\gamma \left( \widehat{f}_m(D_n) \right) \right] + \frac{C-1}{n} \mathbb{E} [K_m(X, X) - K_m(X, Y)]
 \end{aligned}$$

# Oracle inequality for bias-corrected V-fold

Theorem (A., Lerasle & Magalhães 2015)

Under some “reasonable” assumptions, *with probability at least  $1 - e^{-x}$* , for all  $\varepsilon \in (0, 1)$ , for any

$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \{P_n \gamma(\hat{f}_m(D_n)) + \operatorname{pen}_{\text{VF}}(\hat{f}_m; D_n; \mathcal{B})\}$ ,

$$\begin{aligned} \|\hat{f}_{\hat{m}} - f^*\|^2 &\leq \frac{1 + \varepsilon}{1 - \varepsilon} \inf_{m \in \mathcal{M}} \left\{ \|\hat{f}_m - f^*\|^2 \right\} \\ &\quad + \frac{L (\log \operatorname{Card}(\mathcal{M}) + x)^2}{(1 - \varepsilon) \varepsilon^3 n} . \end{aligned}$$

Related result: A. & Lerasle (2012) for projection estimators.

# Oracle inequality for CV procedures

## Theorem (A., Lerasle & Magalhães 2015)

Under some “reasonable” assumptions, with probability at least  $1 - e^{-x}$ , for all  $\varepsilon > 0$ , for any  $\hat{m}$  selected by *V-fold penalization* ( $C \times \text{pen}_{\text{VF}}$ ), *V-fold CV* or *Leave-p-out*, for any  $m \in \mathcal{M}$ ,

$$\frac{1 - \delta_{\hat{m}}^-}{1 + \delta_{\hat{m}}^+} \left\| \hat{f}_{\hat{m}} - f^* \right\|^2 \leq \inf_{m \in \mathcal{M}} \left\{ \left\| \hat{f}_m - f^* \right\|^2 \right\} + R_{n,x,\varepsilon}$$

where  $\gamma_m := \frac{\mathbb{E}[K_m(X,X)]}{A_m(X,X)}$ ,  $R_{n,x,\varepsilon} = \frac{L(C^2 \vee 1)(\log \text{Card}(\mathcal{M}) + x)^2}{\varepsilon^3 n}$  and

$\delta_m^+ =$	$C \times \text{pen}_{\text{VF}}$	$\text{VFCV}$	$\text{LPO}$
$\delta_m^- =$	$2(C-1)_+ \gamma_m + \varepsilon$	$2/(V-1) + \varepsilon$	$2p/(n-p) + \varepsilon$
	$2(C-1)_- \gamma_m + \varepsilon$	$\varepsilon$	$\varepsilon$

Related results: van der Laan, Dudoit & Keles (2004); Celisse (2014) for the leave-p-out; A. & Lerasle (2012).

# Assumptions

General assumption set for linear estimators, holds true for instance if, for some numerical constant  $\kappa > 0$ ,

- **projection estimators** ( $\gamma_m = 1$ ):

$$\kappa L \geq 1 \vee \frac{\|\sum_{\lambda \in \Lambda_m} \psi_\lambda^2\|_\infty}{n} \vee \|f_m^*\|_\infty .$$

# Assumptions

General assumption set for linear estimators, holds true for instance if, for some numerical constant  $\kappa > 0$ ,

- projection estimators ( $\gamma_m = 1$ ):

$$\kappa L \geq 1 \vee \frac{\|\sum_{\lambda \in \Lambda_m} \psi_\lambda^2\|_\infty}{n} \vee \|f_m^*\|_\infty .$$

- kernel (Parzen) estimators ( $\gamma_m = k(0) / \|k\|^2$ ):

$$h \geq \frac{\|k\|_{L^2(\mu)}^2 \vee \|k\|_\infty}{n} \quad \text{and} \quad \kappa L \geq 1 \vee \frac{k(0)}{\|k\|_{L^2(\mu)}^2} \vee \|f^*\|_\infty \|k\|_{L^1(\mu)}^2 .$$

# Assumptions

General assumption set for linear estimators, holds true for instance if, for some numerical constant  $\kappa > 0$ ,

- projection estimators ( $\gamma_m = 1$ ):

$$\kappa L \geq 1 \vee \frac{\|\sum_{\lambda \in \Lambda_m} \psi_\lambda^2\|_\infty}{n} \vee \|f_m^*\|_\infty .$$

- kernel (Parzen) estimators ( $\gamma_m = k(0) / \|k\|^2$ ):

$$h \geq \frac{\|k\|_{L^2(\mu)}^2 \vee \|k\|_\infty}{n} \quad \text{and} \quad \kappa L \geq 1 \vee \frac{k(0)}{\|k\|_{L^2(\mu)}^2} \vee \|f^*\|_\infty \|k\|_{L^1(\mu)}^2 .$$

Also possible to get **sharp minimax Sobolev adaptivity** with V-fold calibrated Pinsker estimators.

# Variance of the (corrected)-VFCV criterion

- Exact non-asymptotic formula for

$$\text{var}(\mathcal{C}(m)) \quad \text{and} \quad \text{var}(\mathcal{C}(m) - \mathcal{C}(m'))$$

with

$$\mathcal{C}(m) \in \left\{ \widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{f}_m; D_n; \mathcal{B}), \widehat{\mathcal{R}}^{\text{vf}}(\widehat{f}_m; D_n; \mathcal{B}), \widehat{\mathcal{R}}^{\text{lp0}}(\widehat{f}_m; D_n) \right\} .$$



# Variance of the (corrected)-VFCV criterion

- Exact non-asymptotic formula for

$$\text{var}(\mathcal{C}(m)) \quad \text{and} \quad \text{var}(\mathcal{C}(m) - \mathcal{C}(m'))$$

with

$$\mathcal{C}(m) \in \left\{ \widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{f}_m; D_n; \mathcal{B}), \widehat{\mathcal{R}}^{\text{vf}}(\widehat{f}_m; D_n; \mathcal{B}), \widehat{\mathcal{R}}^{\ell_{\text{po}}}(\widehat{f}_m; D_n) \right\} .$$

- For bias-corrected  $V$ -fold (and  $V$ -fold penalization), the variance always **decreases with  $V$** .

# Variance of the (corrected)-VFCV criterion

- Exact non-asymptotic formula for

$$\text{var}(\mathcal{C}(m)) \quad \text{and} \quad \text{var}(\mathcal{C}(m) - \mathcal{C}(m'))$$

with

$$\mathcal{C}(m) \in \left\{ \widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{f}_m; D_n; \mathcal{B}), \widehat{\mathcal{R}}^{\text{vf}}(\widehat{f}_m; D_n; \mathcal{B}), \widehat{\mathcal{R}}^{\text{lp0}}(\widehat{f}_m; D_n) \right\} .$$

- For bias-corrected  $V$ -fold (and  $V$ -fold penalization), the variance always **decreases with  $V$** .
- For projection estimators (A. & Lerasle 2012), the variance **decreases by a constant factor** (at most).

# Conclusion

- Non-asymptotic oracle inequality for  $V$ -fold cross-validation and the leave- $p$ -out, for any  $V, p$ .  
Optimal up to a constant factor (bias).

# Conclusion

- Non-asymptotic oracle inequality for  $V$ -fold cross-validation and the leave- $p$ -out, for **any  $V, p$** .  
Optimal up to a constant factor (bias).
- **Bias-corrected  $V$ -fold (or  $V$ -fold penalization)  $\Rightarrow$  first-order optimal non-asymptotic oracle inequality, for any  $V$ .**

# Conclusion

- Non-asymptotic oracle inequality for  $V$ -fold cross-validation and the leave- $p$ -out, for **any  $V, p$** .  
Optimal up to a constant factor (bias).
- Bias-corrected  $V$ -fold (or  $V$ -fold penalization)  $\Rightarrow$  **first-order optimal** non-asymptotic oracle inequality, for **any  $V$** .
- **Result holds for various settings, such as: model selection, bandwidth and kernel choice (Parzen), Pinsker estimators, mix of various kinds of these.**

# Conclusion

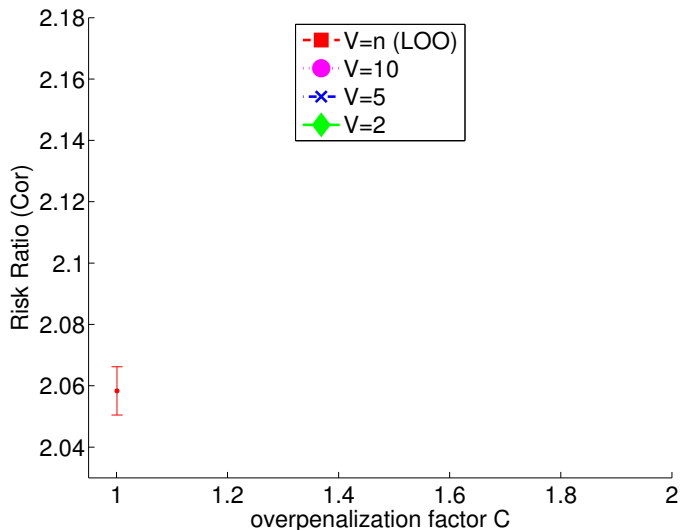
- Non-asymptotic oracle inequality for  $V$ -fold cross-validation and the leave- $p$ -out, for **any  $V, p$** .  
Optimal up to a constant factor (bias).
  - Bias-corrected  $V$ -fold (or  $V$ -fold penalization)  $\Rightarrow$  **first-order optimal** non-asymptotic oracle inequality, for **any  $V$** .
  - Result holds for **various settings**, such as: model selection, bandwidth and kernel choice (Parzen), Pinsker estimators, mix of various kinds of these.
  - **Choice of  $V$  for  $V$ -fold?**
- $\Rightarrow$  **must take into account second-order terms:**
- **variance: decreases with  $V$**
  - **small bias can be benefic (open problem).**

<https://tel.archives-ouvertes.fr/tel-01164581> Chap. 3–4

# Part I

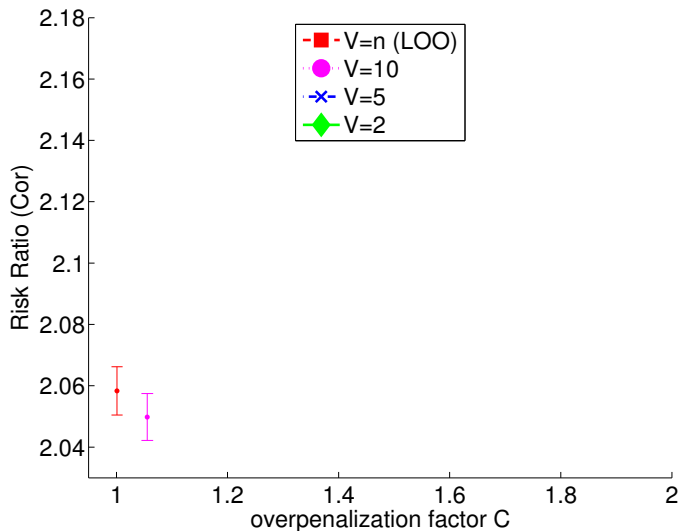
## Appendix

# Experiment: $V$ -fold CV

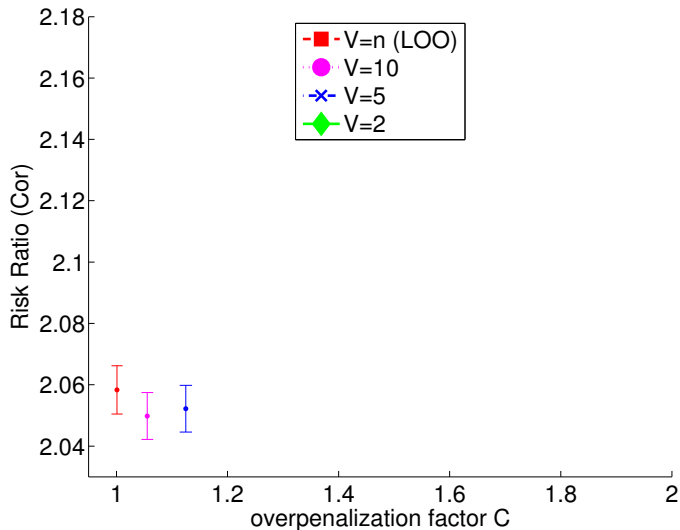




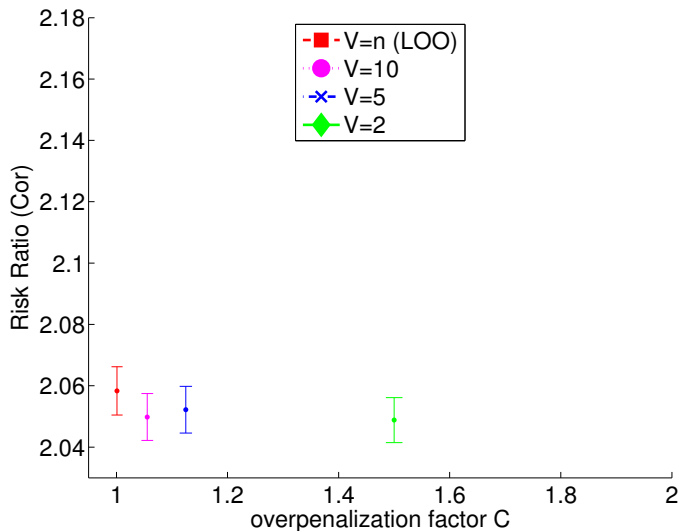
# Experiment: $V$ -fold CV



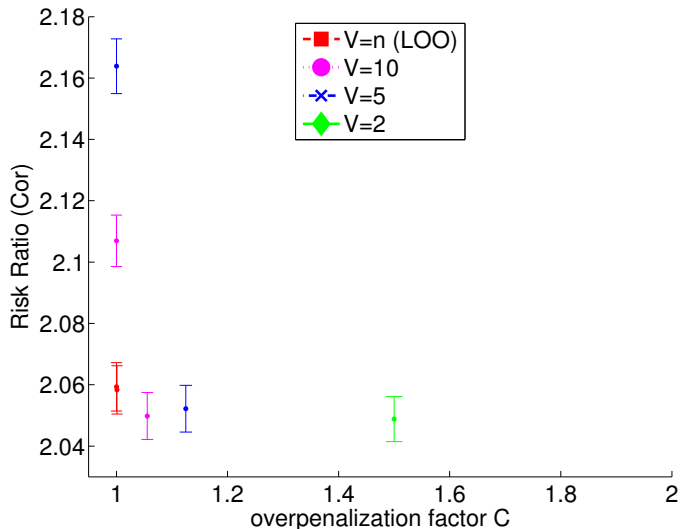
# Experiment: $V$ -fold CV



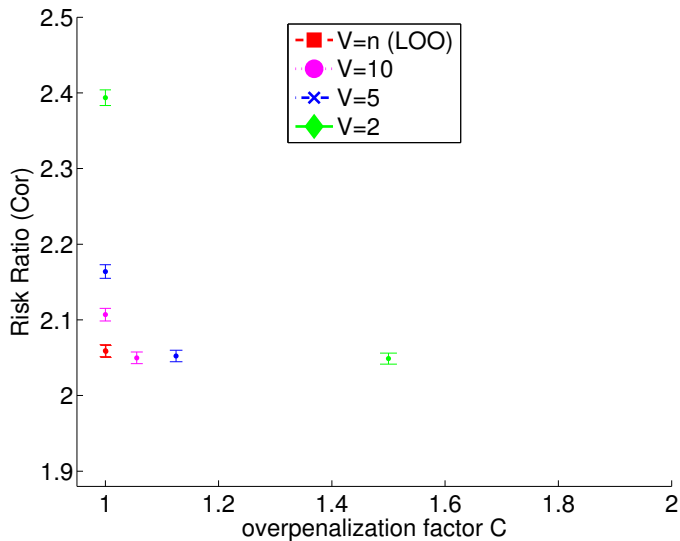
# Experiment: $V$ -fold CV



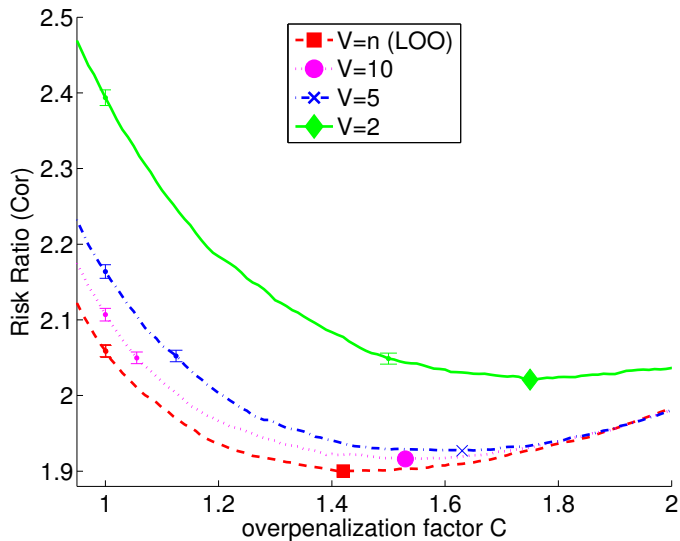
# Experiment: $V$ -fold penalization



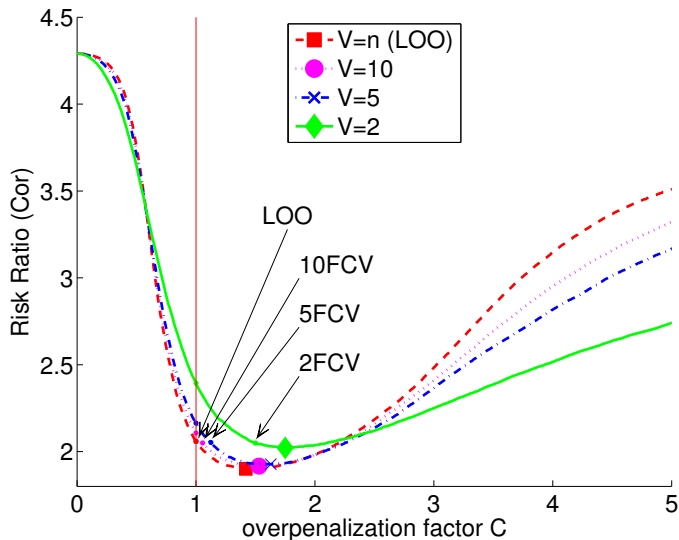
# Experiment: $V$ -fold penalization



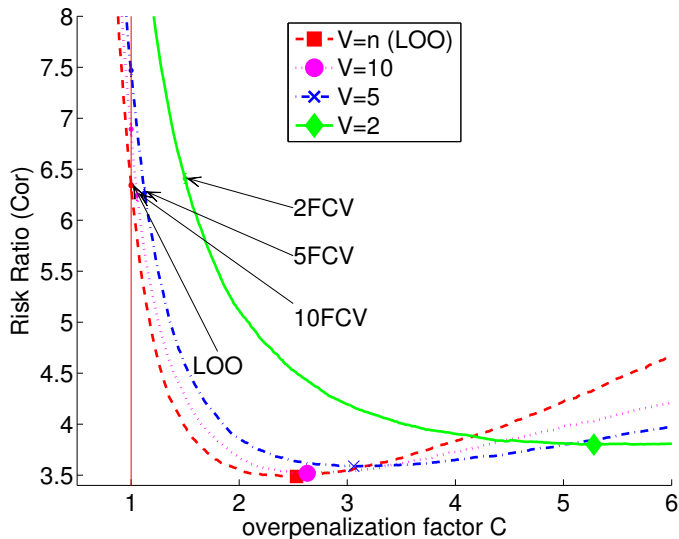
# Experiment: overpenalization



# Experiment: conclusion (setting S)

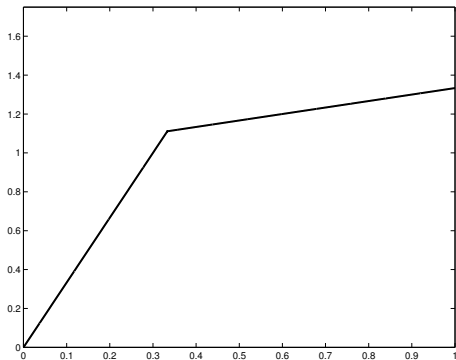


# Experiment: “practically parametric” setting (L)

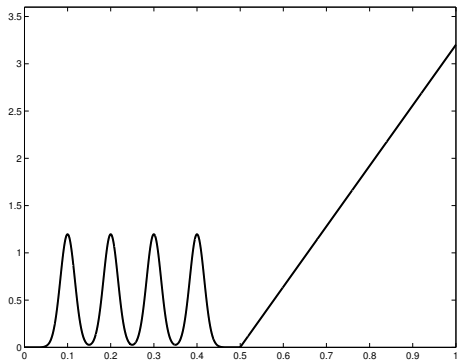




## Simulation setting: densities

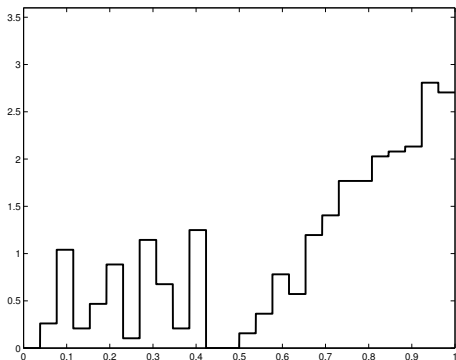


L

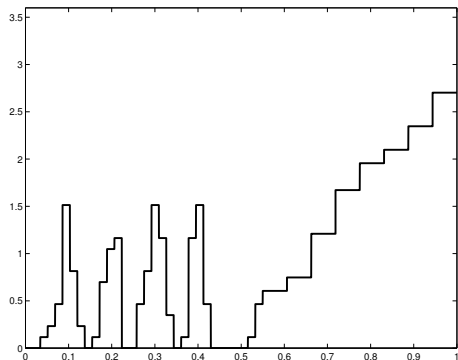


S

# Simulation setting: model families



Regu



Dya2