

Cross-validation for estimator selection

Sylvain Arlot (joint works with Alain Celisse, Matthieu Lerasle,
Nelo Magalhães)

¹CNRS

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

IHP, Paris

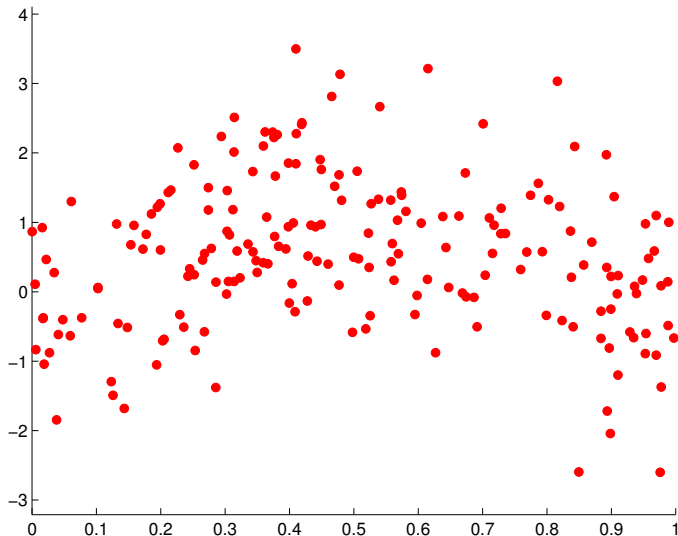
April, 9th, 2015

Main reference (survey): [arXiv:0907.4728](https://arxiv.org/abs/0907.4728)

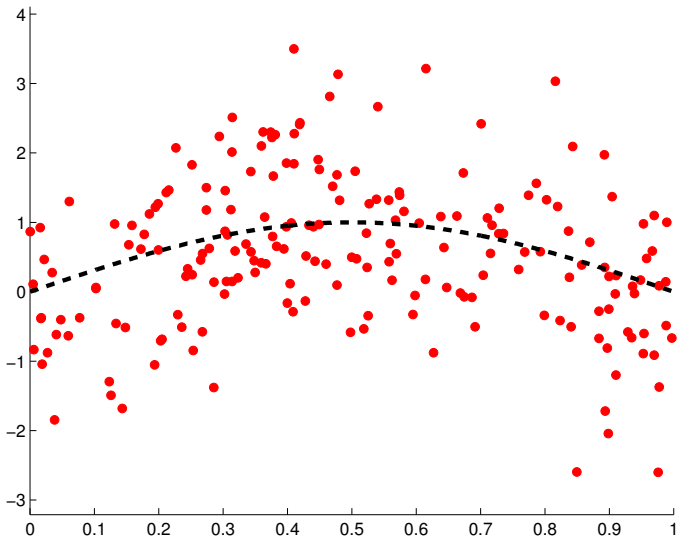
Outline

- 1 Estimator selection
- 2 Cross-validation
- 3 Cross-validation for risk estimation
- 4 Cross-validation for estimator selection
- 5 Large \mathcal{M}
- 6 Conclusion

Regression: data $(X_1, Y_1), \dots, (X_n, Y_n)$



Goal: predict Y given X , i.e., denoising



Prediction problem / regression

- **Data** $D_n: (X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \mathcal{Y}$ (i.i.d. $\sim P$)
- **Contrast** $\gamma(t; (x, y))$ measures how well $t(x)$ “predicts” y
- **Goal:** learn $t \in \mathbb{S} = \{ \text{measurable functions } \mathcal{X} \rightarrow \mathcal{Y} \}$ s.t.
 $\mathbb{E}_{(X, Y) \sim P} [\gamma(t; (X, Y))] =: P\gamma(t)$ is minimal.

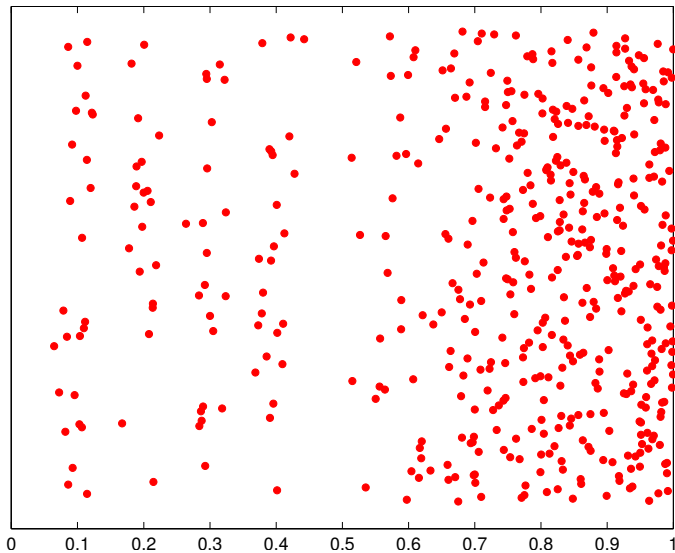
Prediction problem / regression

- **Data** D_n : $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \mathcal{Y}$ (i.i.d. $\sim P$)
- **Contrast** $\gamma(t; (x, y))$ measures how well $t(x)$ “predicts” y
- **Goal**: learn $t \in \mathbb{S} = \{ \text{measurable functions } \mathcal{X} \rightarrow \mathcal{Y} \}$ s.t.
 $\mathbb{E}_{(X, Y) \sim P} [\gamma(t; (X, Y))] =: P\gamma(t)$ is minimal.
- **Example: regression** $\mathcal{Y} = \mathbb{R}$,
least-squares contrast $\gamma(t; (x, y)) = (t(x) - y)^2$
 $s^* \in \operatorname{argmin}_{t \in \mathbb{S}} P\gamma(t)$ is the regression function:
 $s^*(X) = \mathbb{E}[Y | X]$

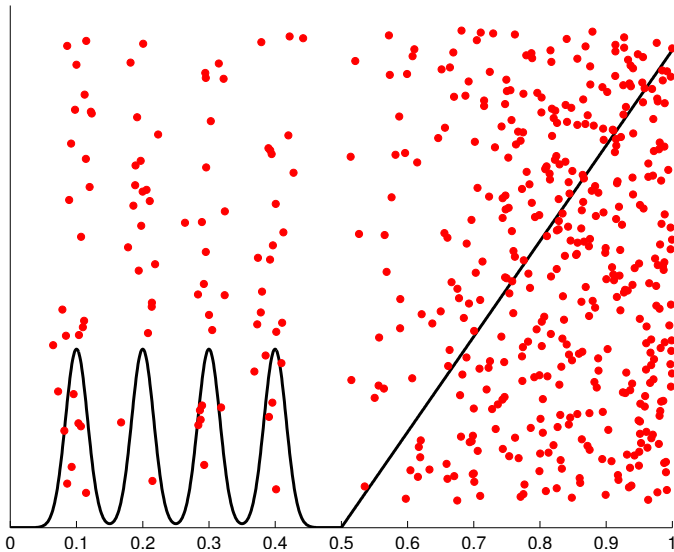
⇒ excess loss

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

Density estimation: data ξ_1, \dots, ξ_n



Goal: estimate the common density s^* of ξ_i



Problem: density estimation

- **Data** D_n : $\xi_1, \dots, \xi_n \in \Xi$ (i.i.d. $\sim P$, density s^* w.r.t. μ)
- **Least-squares contrast** $\gamma(t, \xi) = \|t\|_{L^2(\mu)}^2 - 2t(\xi)$
- **Goal**: learn $t \in \mathbb{S} = \{\text{measurable functions } \Xi \rightarrow \mathbb{R}\}$ s.t.
 $\mathbb{E}_{\xi \sim P} [\gamma(t; \xi)] =: P\gamma(t)$ is minimal.

Problem: density estimation

- **Data** D_n : $\xi_1, \dots, \xi_n \in \Xi$ (i.i.d. $\sim P$, density s^* w.r.t. μ)
- **Least-squares contrast** $\gamma(t, \xi) = \|t\|_{L^2(\mu)}^2 - 2t(\xi)$
- **Goal**: learn $t \in \mathbb{S} = \{\text{measurable functions } \Xi \rightarrow \mathbb{R}\}$ s.t. $\mathbb{E}_{\xi \sim P} [\gamma(t; \xi)] =: P\gamma(t)$ is minimal.

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

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

$$\ell(s^*, t) := P\gamma(t) - P\gamma(s^*) = \|t - s^*\|_{L^2(\mu)}^2$$

General setting

- Data $\xi_1, \dots, \xi_n \in \Xi$ i.i.d. with distribution P
prediction: $\xi_i = (X_i, Y_i) \in \mathcal{X} \times \mathcal{Y}$
- Goal: Estimate some feature $s^* \in \mathbb{S}$ of P
density, regression function, Bayes predictor...
- Contrast function $\gamma : \mathbb{S} \times \Xi \rightarrow \mathbb{R}$ such that

$$s^* \in \underset{t \in \mathbb{S}}{\operatorname{argmin}} \{P\gamma(t)\} \quad \text{with} \quad P\gamma(t) := \mathbb{E}_{\xi \sim P} [\gamma(t; \xi)]$$

- Excess loss

$$\ell(s^*, t) := P\gamma(t) - P\gamma(s^*) \geq 0$$

Examples

- **Prediction:** $\xi_i = (X_i, Y_i)$
 $X_{n+1} \rightsquigarrow$ “predict” Y_{n+1} with $t(X_{n+1})$?
 $\gamma(t; (x, y))$ quantifies the “distance” between $t(x)$ and y

- **Regression** ($\mathcal{Y} = \mathbb{R}$), least squares:

$$\gamma(t; (x, y)) = (t(x) - y)^2 \quad s^*(X) = \mathbb{E}[Y|X]$$

- **Binary classification** ($\mathcal{Y} = \{0, 1\}$), 0–1 contrast:

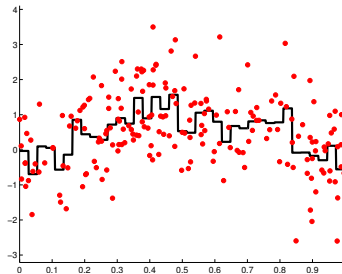
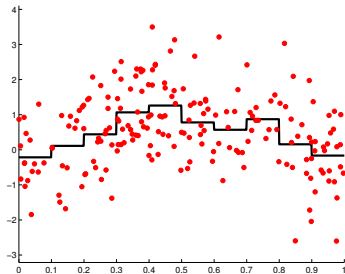
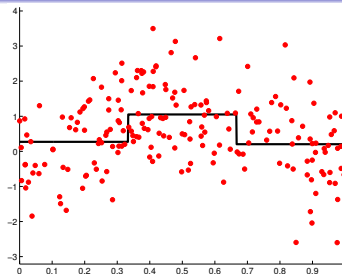
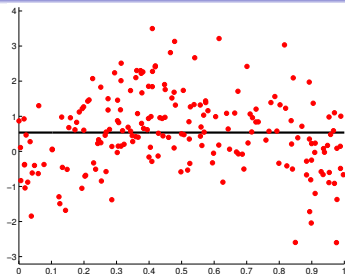
$$\gamma(t; (x, y)) = \mathbf{1}_{t(x) \neq y}$$

- **Density estimation** (reference measure μ):

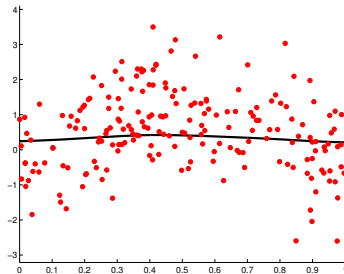
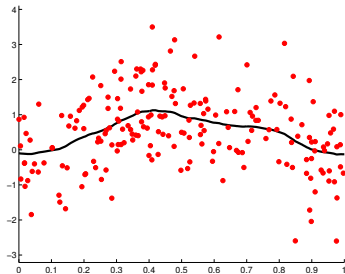
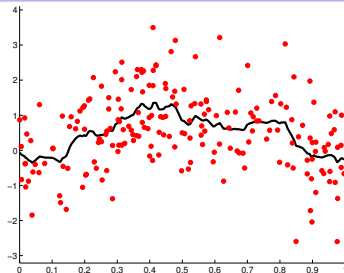
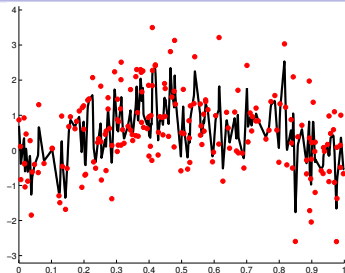
least squares: $\gamma(t; \xi) = \|t\|_{L^2(\mu)}^2 - 2t(\xi)$

log-likelihood: $\gamma(t; \xi) = -\log(t(\xi))$

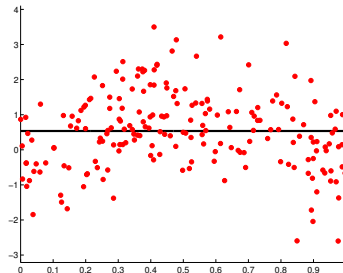
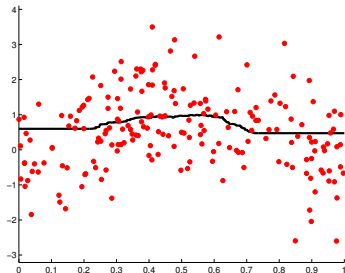
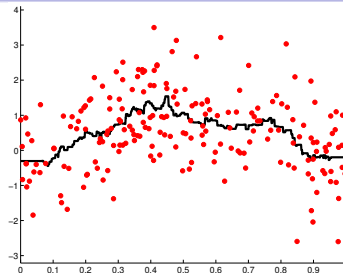
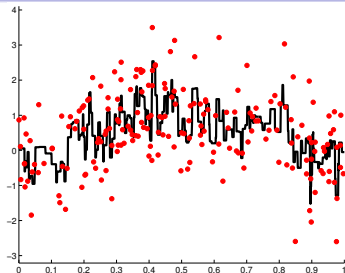
Estimator selection (regression): regular regressograms



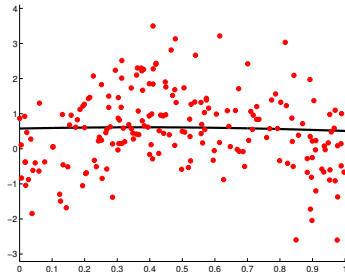
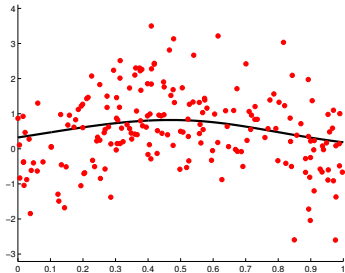
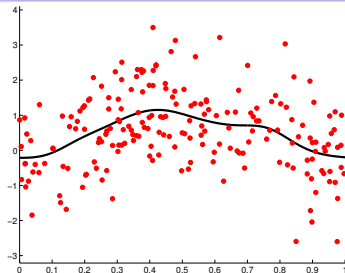
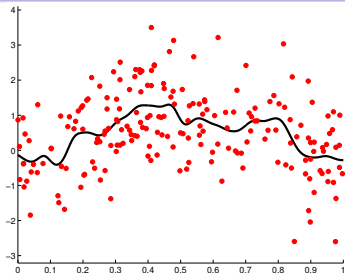
Estimator selection (regression): kernel ridge



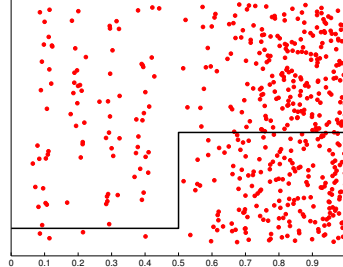
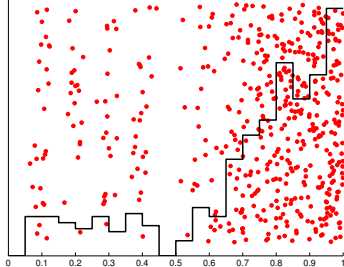
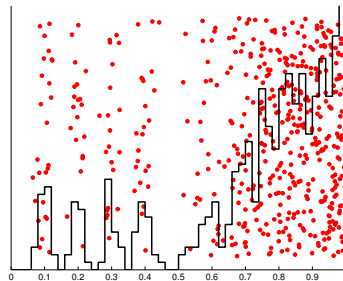
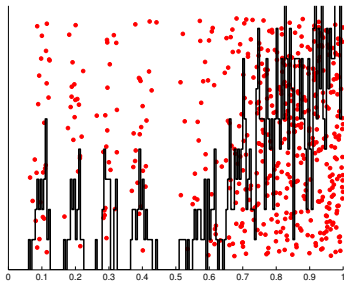
Estimator selection (regression): k nearest neighbours



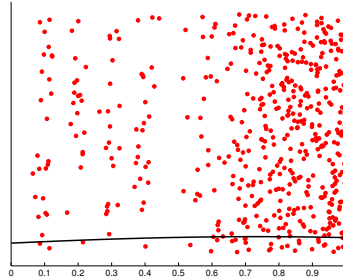
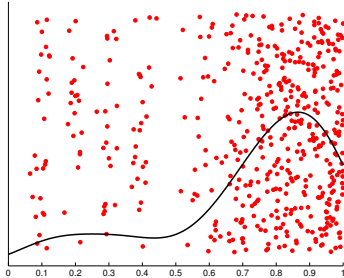
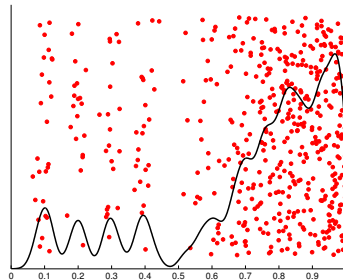
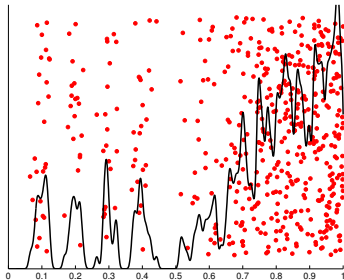
Estimator selection (regression): Nadaraya-Watson



Estimator selection (density): regular histograms



Estimator selection (density): Parzen, Gaussian kernel



Estimator selection

- **Estimator/Learning algorithm:** $\hat{s} : D_n \mapsto \hat{s}(D_n) \in \mathbb{S}$
- Example: **least-squares estimator** on some **model** $S_m \subset \mathbb{S}$

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

Examples of models: histograms, $\operatorname{span}\{\varphi_1, \dots, \varphi_D\}$

Estimator selection

- Estimator/Learning algorithm: $\hat{s} : D_n \mapsto \hat{s}(D_n) \in \mathbb{S}$
- Example: least-squares estimator on some model $S_m \subset \mathbb{S}$

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

Examples of models: histograms, $\operatorname{span}\{\varphi_1, \dots, \varphi_D\}$

- Estimator collection $(\hat{s}_m)_{m \in \mathcal{M}} \Rightarrow$ choose $\hat{m} = \hat{m}(D_n)$?

Estimator selection

- Estimator/Learning algorithm: $\hat{s} : D_n \mapsto \hat{s}(D_n) \in \mathbb{S}$
- Example: least-squares estimator on some model $S_m \subset \mathbb{S}$

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

Examples of models: histograms, $\operatorname{span}\{\varphi_1, \dots, \varphi_D\}$

- Estimator collection $(\hat{s}_m)_{m \in \mathcal{M}} \Rightarrow$ choose $\hat{m} = \hat{m}(D_n)?$
- Examples:
 - **model selection**
 - **calibration of tuning parameters** (choosing k or the distance for k -NN, choice of a regularization parameter, choice of a kernel, etc.)
 - choice between **different methods**
ex.: k -NN vs. smoothing splines?

Estimator selection: two possible goals

- **Estimation goal:** minimize the risk of the final estimator, i.e., **Oracle inequality** (in expectation or with a large probability):

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

Estimator selection: two possible goals

- **Estimation goal:** minimize the risk of the final estimator, i.e., Oracle inequality (in expectation or with a large probability):

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

- **Identification goal:** select the (asymptotically) best model/estimator, assuming it is well-defined, i.e., Selection consistency:

$$\mathbb{P}(\widehat{m}(D_n) = m^*) \xrightarrow[n \rightarrow \infty]{} 1.$$

Equivalent to estimation in the **parametric** setting.

Estimator selection: two possible goals

- **Estimation goal:** minimize the risk of the final estimator, i.e., Oracle inequality (in expectation or with a large probability):

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

- **Identification goal:** select the (asymptotically) best model/estimator, assuming it is well-defined, i.e., Selection consistency:

$$\mathbb{P}(\widehat{m}(D_n) = m^*) \xrightarrow{n \rightarrow \infty} 1.$$

Equivalent to estimation in the **parametric** setting.

- Both goals with the same procedure (AIC-BIC dilemma)?
No in general (Yang, 2005). Sometimes possible.

Estimation goal: Bias-variance trade-off

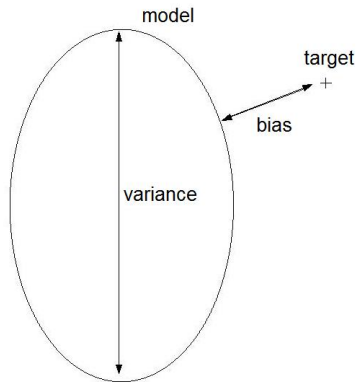
$$\mathbb{E}[\ell(s^*, \hat{s}_m)] = \text{Bias} + \text{Variance}$$

Bias or Approximation error

$$\ell(s^*, s_m^*) = \inf_{t \in S_m} \ell(s^*, t)$$

Variance or Estimation error

$$\text{OLS in regression: } \frac{\sigma^2 \dim(S_m)}{n}$$



Estimation goal: Bias-variance trade-off

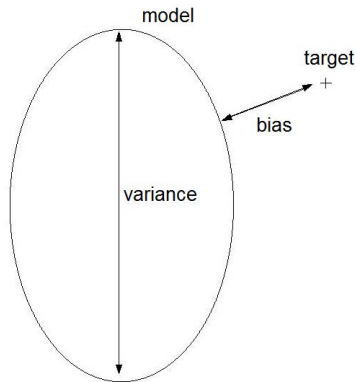
$$\mathbb{E}[\ell(s^*, \hat{s}_m)] = \text{Bias} + \text{Variance}$$

Bias or Approximation error

$$\ell(s^*, s_m^*) = \inf_{t \in S_m} \ell(s^*, t)$$

Variance or Estimation error

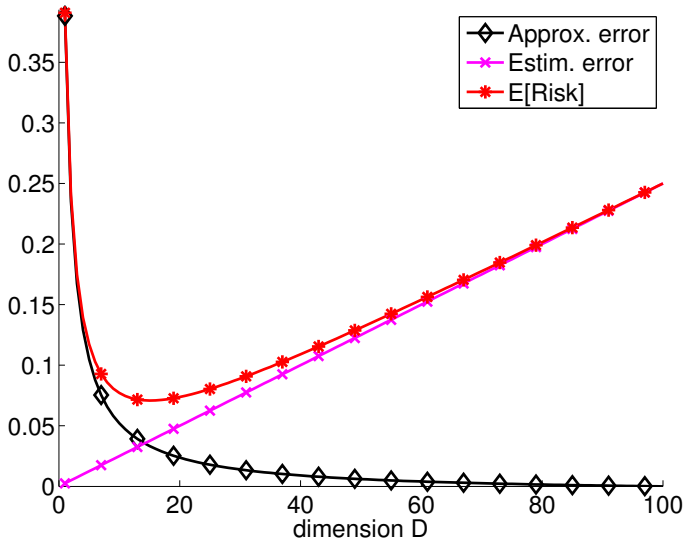
$$\text{OLS in regression: } \frac{\sigma^2 \dim(S_m)}{n}$$



Bias-variance trade-off

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

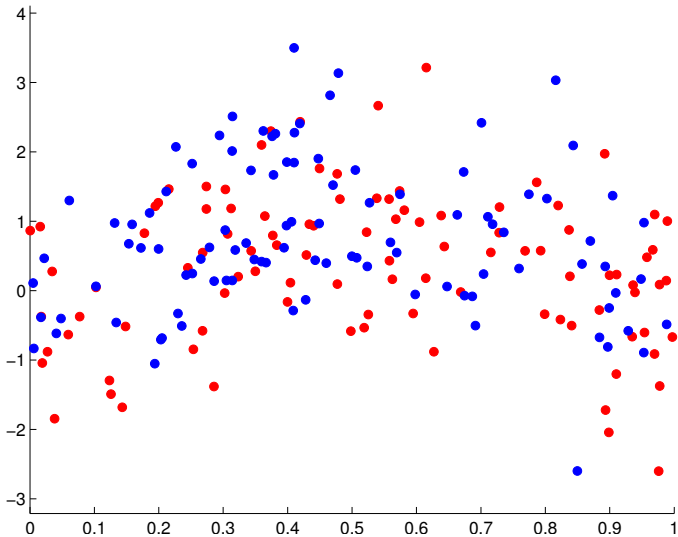
Estimation goal: Bias-variance trade-off



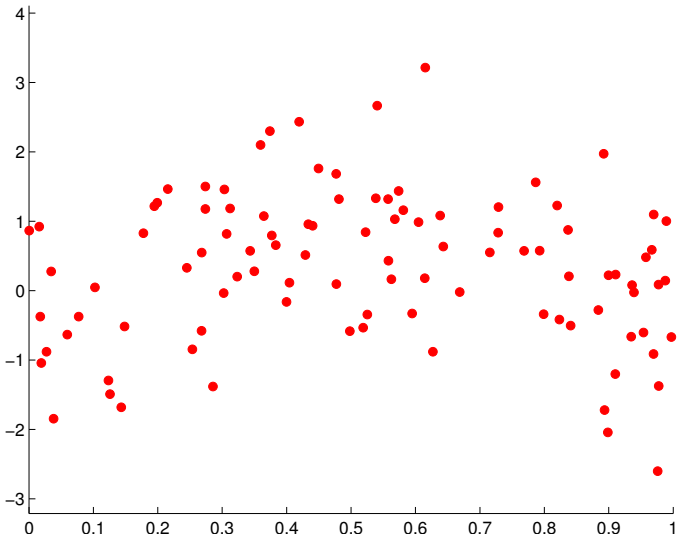
Outline

- 1 Estimator selection
- 2 Cross-validation
- 3 Cross-validation for risk estimation
- 4 Cross-validation for estimator selection
- 5 Large \mathcal{M}
- 6 Conclusion

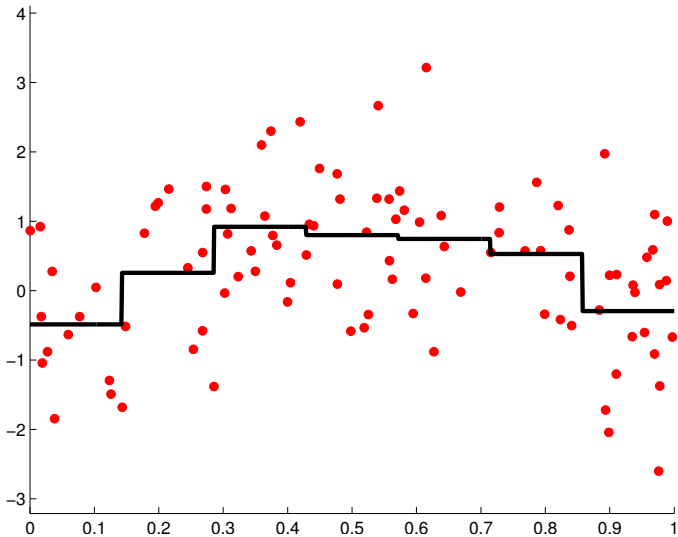
Validation principle



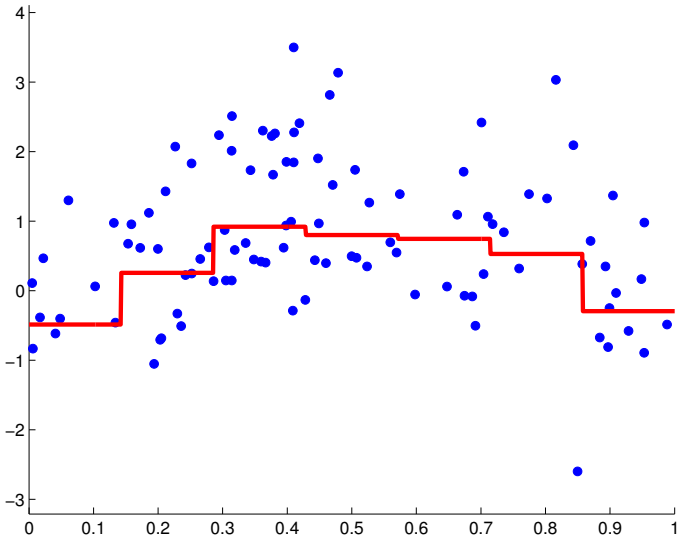
Validation principle: learning sample



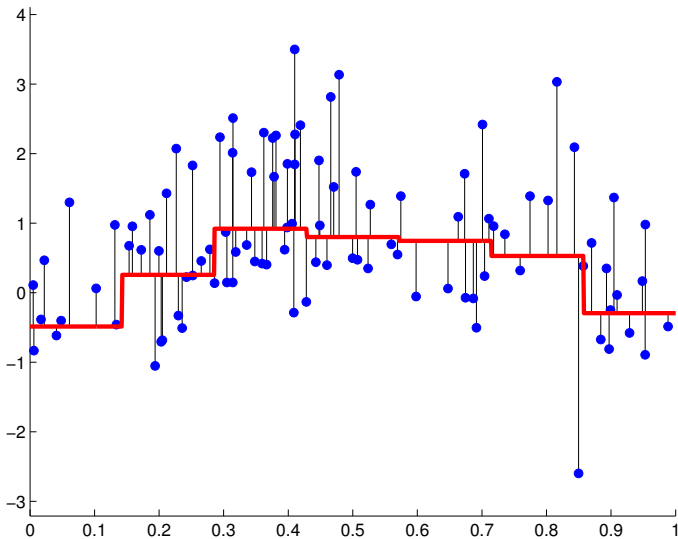
Validation principle: learning sample



Validation principle: validation sample



Validation principle: validation sample



Cross-validation

$(X_1, Y_1), \dots, (X_{n_t}, Y_{n_t})$

Training set $D_n^{(t)} \Rightarrow \hat{s}_m^{(t)} = \hat{s}_m(D_n^{(t)})$

$(X_{n_t+1}, Y_{n_t+1}), \dots, (X_n, Y_n)$

Validation set $D_n^{(v)} \Rightarrow$ evaluate risk

Cross-validation

$(X_1, Y_1), \dots, (X_{n_t}, Y_{n_t})$

Training set $D_n^{(t)} \Rightarrow \hat{s}_m^{(t)} = \hat{s}_m(D_n^{(t)})$

$(X_{n_t+1}, Y_{n_t+1}), \dots, (X_n, Y_n)$

Validation set $D_n^{(v)} \Rightarrow$ evaluate risk

- hold-out estimator of the risk:

$$P_n^{(v)} \gamma(\hat{s}_m^{(t)}) = \frac{1}{n_v} \sum_{\xi \in D_n^{(v)}} \gamma(\hat{s}_m^{(t)}; \xi)$$

$$n_v = |D_n^{(v)}| = n - n_t$$

Cross-validation

$$\underbrace{(X_1, Y_1), \dots, (X_{n_t}, Y_{n_t})}_{\text{Training set}}$$

Training set $D_n^{(t)} \Rightarrow \hat{s}_m^{(t)} = \hat{s}_m(D_n^{(t)})$

$$\underbrace{(X_{n_t+1}, Y_{n_t+1}), \dots, (X_n, Y_n)}_{\text{Validation set}}$$

Validation set $D_n^{(v)} \Rightarrow$ evaluate risk

- **hold-out** estimator of the risk:

$$P_n^{(v)} \gamma(\hat{s}_m^{(t)}) = \frac{1}{n_v} \sum_{\xi \in D_n^{(v)}} \gamma(\hat{s}_m^{(t)}; \xi) \quad n_v = |D_n^{(v)}| = n - n_t$$

- **cross-validation**: average several hold-out estimators

$$\hat{\mathcal{R}}^{\text{cv}}(\hat{s}_m; D_n; (I_j^{(t)})_{1 \leq j \leq B}) = \frac{1}{B} \sum_{j=1}^B P_n^{(v,j)} \gamma(\hat{s}_m^{(t,j)}) \quad D_n^{(t,j)} = (\xi_i)_{i \in I_j^{(t)}}$$

Cross-validation

$(X_1, Y_1), \dots, (X_{n_t}, Y_{n_t})$

$(X_{n_t+1}, Y_{n_t+1}), \dots, (X_n, Y_n)$

Training set $D_n^{(t)} \Rightarrow \hat{s}_m^{(t)} = \hat{s}_m(D_n^{(t)})$

Validation set $D_n^{(v)} \Rightarrow$ evaluate risk

- **hold-out** estimator of the risk:

$$P_n^{(v)} \gamma(\hat{s}_m^{(t)}) = \frac{1}{n_v} \sum_{\xi \in D_n^{(v)}} \gamma(\hat{s}_m^{(t)}; \xi) \quad n_v = |D_n^{(v)}| = n - n_t$$

- **cross-validation**: average several hold-out estimators

$$\hat{\mathcal{R}}^{cv}(\hat{s}_m; D_n; (I_j^{(t)})_{1 \leq j \leq B}) = \frac{1}{B} \sum_{j=1}^B P_n^{(v,j)} \gamma(\hat{s}_m^{(t,j)}) \quad D_n^{(t,j)} = (\xi_i)_{i \in I_j^{(t)}}$$

- **estimator selection**:

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \hat{\mathcal{R}}^{cv}(\hat{s}_m; D_n) \right\}$$

Cross-validation: examples

- Exhaustive data splitting: all possible subsets of size n_t
 \Rightarrow leave-one-out ($n_t = n - 1$)

$$\widehat{\mathcal{R}}^{\text{loo}}(\widehat{s}_m; D_n) = \frac{1}{n} \sum_{j=1}^n \gamma(\widehat{s}_m^{(-j)}; \xi_j)$$

\Rightarrow leave- p -out ($n_t = n - p$)

Cross-validation: examples

- Exhaustive data splitting: all possible subsets of size n_t
 \Rightarrow leave-one-out ($n_t = n - 1$)

$$\widehat{\mathcal{R}}^{\text{loo}}(\widehat{s}_m; D_n) = \frac{1}{n} \sum_{j=1}^n \gamma(\widehat{s}_m^{(-j)}; \xi_j)$$

\Rightarrow leave- p -out ($n_t = n - p$)

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

$$\Rightarrow \widehat{\mathcal{R}}^{\text{vf}}(\widehat{s}_m; D_n; \mathcal{B}) = \frac{1}{V} \sum_{j=1}^V P_n^j \gamma(\widehat{s}_m^{(-j)})$$

Cross-validation: examples

- Exhaustive data splitting: all possible subsets of size n_t
 \Rightarrow leave-one-out ($n_t = n - 1$)

$$\widehat{\mathcal{R}}^{\text{loo}}(\widehat{s}_m; D_n) = \frac{1}{n} \sum_{j=1}^n \gamma(\widehat{s}_m^{(-j)}; \xi_j)$$

\Rightarrow leave- p -out ($n_t = n - p$)

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

$$\Rightarrow \widehat{\mathcal{R}}^{\text{vf}}(\widehat{s}_m; D_n; \mathcal{B}) = \frac{1}{V} \sum_{j=1}^V P_n^j \gamma(\widehat{s}_m^{(-j)})$$

- Monte-Carlo CV / Repeated learning testing:

$$I_1^{(t)}, \dots, I_B^{(t)} \text{ i.i.d. uniform}$$

Outline

- 1 Estimator selection
- 2 Cross-validation
- 3 Cross-validation for risk estimation
- 4 Cross-validation for estimator selection
- 5 Large \mathcal{M}
- 6 Conclusion

Bias of cross-validation

- In this talk, we always assume: $\forall j, \text{Card}(D_n^{(t,j)}) = n_t$
For V -fold CV: $\text{Card}(B_j) = n/V$.
- Ideal criterion: $P\gamma(\hat{s}_m(D_n))$

Bias of cross-validation

- In this talk, we always assume: $\forall j, \text{Card}(D_n^{(t,j)}) = n_t$
For V -fold CV: $\text{Card}(B_j) = n/V$.
- Ideal criterion: $P_\gamma(\hat{s}_m(D_n))$
- General analysis for the bias:

$$\mathbb{E} \left[\hat{\mathcal{R}}^{\text{cv}} \left(\hat{s}_m; D_n; \left(I_j^{(t)} \right)_{1 \leq j \leq B} \right) \right] = \mathbb{E} \left[P_\gamma(\hat{s}_m(D_{n_t})) \right]$$

Bias of cross-validation

- In this talk, we always assume: $\forall j, \text{Card}(D_n^{(t,j)}) = n_t$
For V -fold CV: $\text{Card}(B_j) = n/V$.
- Ideal criterion: $P_\gamma(\hat{s}_m(D_n))$
- General analysis for the bias:

$$\mathbb{E} \left[\hat{\mathcal{R}}^{\text{cv}} \left(\hat{s}_m; D_n; \left(I_j^{(t)} \right)_{1 \leq j \leq B} \right) \right] = \mathbb{E} \left[P_\gamma(\hat{s}_m(D_{n_t})) \right]$$

\Rightarrow everything depends on $n \rightarrow \mathbb{E} \left[P_\gamma(\hat{s}_m(D_n)) \right]$

Bias of cross-validation

- In this talk, we always assume: $\forall j, \text{Card}(D_n^{(t,j)}) = n_t$
For V -fold CV: $\text{Card}(B_j) = n/V$.
- Ideal criterion: $P\gamma(\hat{s}_m(D_n))$
- General analysis for the bias:

$$\mathbb{E} \left[\hat{\mathcal{R}}^{\text{cv}} \left(\hat{s}_m; D_n; \left(I_j^{(t)} \right)_{1 \leq j \leq B} \right) \right] = \mathbb{E} \left[P\gamma(\hat{s}_m(D_{n_t})) \right]$$

⇒ everything depends on $n \rightarrow \mathbb{E} \left[P\gamma(\hat{s}_m(D_n)) \right]$

- Note: **bias can be corrected** in some settings (Burman, 1989).

Bias of cross-validation

- In this talk, we always assume: $\forall j, \text{Card}(D_n^{(t,j)}) = n_t$
For V -fold CV: $\text{Card}(B_j) = n/V$.
- Ideal criterion: $P\gamma(\hat{s}_m(D_n))$
- General analysis for the bias:

$$\mathbb{E} \left[\hat{\mathcal{R}}^{\text{cv}} \left(\hat{s}_m; D_n; \left(I_j^{(t)} \right)_{1 \leq j \leq B} \right) \right] = \mathbb{E} \left[P\gamma(\hat{s}_m(D_{n_t})) \right]$$

\Rightarrow everything depends on $n \rightarrow \mathbb{E} \left[P\gamma(\hat{s}_m(D_n)) \right]$

- Note: **bias can be corrected** in some settings (Burman, 1989).
- Note: $D_n \rightarrow \hat{s}_m(D_n)$ must be fixed **before seeing any data**; otherwise, stronger bias.

Bias of cross-validation: generic example

Assume

$$\mathbb{E} \left[P\gamma(\hat{s}_m(D_n)) \right] = \alpha(m) + \frac{\beta(m)}{n}$$

(e.g., LS/ridge/ k -NN regression, LS/kernel density estimation).

Bias of cross-validation: generic example

Assume

$$\mathbb{E} \left[P\gamma(\hat{s}_m(D_n)) \right] = \alpha(m) + \frac{\beta(m)}{n}$$

(e.g., LS/ridge/ k -NN regression, LS/kernel density estimation).

$$\Rightarrow \mathbb{E} \left[\hat{\mathcal{R}}^{\text{cv}} \left(\hat{s}_m; D_n; \left(I_j^{(t)} \right)_{1 \leq j \leq B} \right) \right] = \alpha(m) + \frac{n}{n_t} \frac{\beta(m)}{n}$$

Bias of cross-validation: generic example

Assume

$$\mathbb{E} \left[P\gamma(\hat{s}_m(D_n)) \right] = \alpha(m) + \frac{\beta(m)}{n}$$

(e.g., LS/ridge/ k -NN regression, LS/kernel density estimation).

$$\Rightarrow \mathbb{E} \left[\hat{\mathcal{R}}^{\text{cv}} \left(\hat{s}_m; D_n; \left(I_j^{(t)} \right)_{1 \leq j \leq B} \right) \right] = \alpha(m) + \frac{n}{n_t} \frac{\beta(m)}{n}$$

\Rightarrow Bias:

- decreases as a function of n_t ,
- minimal for $n_t = n - 1$,
- negligible if $n_t \sim n$.

Bias of cross-validation: generic example

Assume

$$\mathbb{E} \left[P\gamma(\hat{s}_m(D_n)) \right] = \alpha(m) + \frac{\beta(m)}{n}$$

(e.g., LS/ridge/ k -NN regression, LS/kernel density estimation).

$$\Rightarrow \mathbb{E} \left[\hat{\mathcal{R}}^{\text{cv}} \left(\hat{s}_m; D_n; \left(I_j^{(t)} \right)_{1 \leq j \leq B} \right) \right] = \alpha(m) + \frac{n}{n_t} \frac{\beta(m)}{n}$$

\Rightarrow Bias:

- decreases as a function of n_t ,
- minimal for $n_t = n - 1$,
- negligible if $n_t \sim n$.

\Rightarrow V -fold: bias decreases when V increases, vanishes as $V \rightarrow +\infty$.

Variance of cross-validation

- **Hold-out** (Nadeau & Bengio, 2003):

$$\begin{aligned} \text{var} \left(P_n^{(v)} \gamma \left(\hat{s}_m^{(t)} \right) \right) &= \frac{1}{n_v} \mathbb{E} \left[\text{var} \left(\gamma(u; \xi) \mid u = \hat{s}_m^{(t)} \right) \right] \\ &\quad + \text{var} \left(P \gamma \left(\hat{s}_m(D_{n_t}) \right) \right) \end{aligned}$$

Variance of cross-validation

- Hold-out (Nadeau & Bengio, 2003):

$$\begin{aligned} \text{var} \left(P_n^{(v)} \gamma \left(\hat{s}_m^{(t)} \right) \right) &= \frac{1}{n_v} \mathbb{E} \left[\text{var} \left(\gamma(u; \xi) \mid u = \hat{s}_m^{(t)} \right) \right] \\ &\quad + \text{var} \left(P \gamma \left(\hat{s}_m(D_{n_t}) \right) \right) \end{aligned}$$

- Monte-Carlo CV and number of splits: ($p = n - n_t$)

$$\begin{aligned} \text{var} \left(\hat{\mathcal{R}}^{\text{cv}} \left(\hat{s}_m; D_n; \left(I_j^{(t)} \right)_{1 \leq j \leq B} \right) \right) &= \text{var} \left(\hat{\mathcal{R}}^{\text{lp0}} \left(\hat{s}_m; D_n \right) \right) \\ &\quad + \underbrace{\frac{1}{B} \mathbb{E} \left[\text{var}_{I^{(t)}} \left(P_n^{(v)} \gamma \left(\hat{s}_m^{(t)} \right) \mid D_n \right) \right]}_{\text{permutation variance}} \end{aligned}$$

Variance of cross-validation

- Hold-out (Nadeau & Bengio, 2003):

$$\begin{aligned} \text{var} \left(P_n^{(v)} \gamma \left(\hat{s}_m^{(t)} \right) \right) &= \frac{1}{n_v} \mathbb{E} \left[\text{var} \left(\gamma(u; \xi) \mid u = \hat{s}_m^{(t)} \right) \right] \\ &\quad + \text{var} \left(P \gamma \left(\hat{s}_m(D_{n_t}) \right) \right) \end{aligned}$$

- Monte-Carlo CV and number of splits: ($p = n - n_t$)

$$\begin{aligned} \text{var} \left(\hat{\mathcal{R}}^{\text{cv}} \left(\hat{s}_m; D_n; \left(I_j^{(t)} \right)_{1 \leq j \leq B} \right) \right) &= \text{var} \left(\hat{\mathcal{R}}^{\text{lp0}} \left(\hat{s}_m; D_n \right) \right) \\ &\quad + \underbrace{\frac{1}{B} \mathbb{E} \left[\text{var}_{I^{(t)}} \left(P_n^{(v)} \gamma \left(\hat{s}_m^{(t)} \right) \mid D_n \right) \right]}_{\text{permutation variance}} \end{aligned}$$

- **V-fold CV**: B, n_t, n_v related
leave-one-out: related to stability? (empirical results)

Variance of the V -fold CV criterion

- **Least-squares density estimation** (A. & Lerasle 2012), exact computation (non-asymptotic):

$$\begin{aligned} \text{var} \left(\widehat{\mathcal{R}}^{\text{vf}}(\widehat{s}_m; D_n; \mathcal{B}) \right) &= \frac{1 + \mathcal{O}(1)}{n} \text{var}_P(s_m^*) \\ &+ \frac{2}{n^2} \left[1 + \frac{4}{V-1} + \mathcal{O}\left(\frac{1}{V} + \frac{1}{n}\right) \right] A(m) \end{aligned}$$

(simplified formula, histogram model with bin size d_m^{-1} , $A(m) \approx d_m$)

- Linear regression, specific setting, asymptotic formula (Burman, 1989):

$$\text{var} \left(\widehat{\mathcal{R}}^{\text{vf}}(\widehat{s}_m; D_n; \mathcal{B}) \right) = \frac{2\sigma^2}{n} + \frac{4\sigma^4}{n^2} \left[4 + \frac{4}{V-1} + \frac{2}{(V-1)^2} + \frac{1}{(V-1)^3} \right] + o(n^{-2})$$

⇒ decreasing with V , dependence only in second order terms.

Outline

- 1 Estimator selection
- 2 Cross-validation
- 3 Cross-validation for risk estimation
- 4 Cross-validation for estimator selection
- 5 Large \mathcal{M}
- 6 Conclusion

Risk estimation and estimator selection are different goals

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \widehat{\mathcal{R}}^{\text{cv}}(\hat{s}_m) \right\} \quad \text{vs.} \quad m^* \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ P_\gamma(\hat{s}_m(D_n)) \right\}$$

- For any Z (deterministic or random),

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \widehat{\mathcal{R}}^{\text{cv}}(\hat{s}_m) + Z \right\}$$

⇒ bias and variance meaningless.

Risk estimation and estimator selection are different goals

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

- For any Z (deterministic or random),

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \widehat{\mathcal{R}}^{\text{cv}}(\hat{s}_m) + Z \right\}$$

⇒ bias and variance meaningless.

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

$$\operatorname{sign}(\widehat{\mathcal{R}}^{\text{cv}}(\hat{s}_m) - \widehat{\mathcal{R}}^{\text{cv}}(\hat{s}_{m'})) = \operatorname{sign}(P\gamma(\hat{s}_m) - P\gamma(\hat{s}_{m'}))$$

Risk estimation and estimator selection are different goals

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

- For any Z (deterministic or random),

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \widehat{\mathcal{R}}^{\text{cv}}(\hat{s}_m) + Z \right\}$$

⇒ bias and variance meaningless.

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

$$\operatorname{sign}(\widehat{\mathcal{R}}^{\text{cv}}(\hat{s}_m) - \widehat{\mathcal{R}}^{\text{cv}}(\hat{s}_{m'})) = \operatorname{sign}(P\gamma(\hat{s}_m) - P\gamma(\hat{s}_{m'}))$$

⇒ $\mathbb{E} \left[\widehat{\mathcal{R}}^{\text{cv}}(\hat{s}_m) - \widehat{\mathcal{R}}^{\text{cv}}(\hat{s}_{m'}) \right]$ should be of the good sign (unbiased risk estimation heuristic: AIC, C_p , leave-one-out...)

Risk estimation and estimator selection are different goals

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

- For any Z (deterministic or random),

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \widehat{\mathcal{R}}^{\text{cv}}(\hat{s}_m) + Z \right\}$$

\Rightarrow bias and variance meaningless.

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

$$\operatorname{sign}(\widehat{\mathcal{R}}^{\text{cv}}(\hat{s}_m) - \widehat{\mathcal{R}}^{\text{cv}}(\hat{s}_{m'})) = \operatorname{sign}(P\gamma(\hat{s}_m) - P\gamma(\hat{s}_{m'}))$$

$\Rightarrow \mathbb{E} \left[\widehat{\mathcal{R}}^{\text{cv}}(\hat{s}_m) - \widehat{\mathcal{R}}^{\text{cv}}(\hat{s}_{m'}) \right]$ should be of the good sign (unbiased risk estimation heuristic: AIC, C_p , leave-one-out...)

$\Rightarrow \operatorname{var} \left(\widehat{\mathcal{R}}^{\text{vf}}(\hat{s}_m) - \widehat{\mathcal{R}}^{\text{vf}}(\hat{s}_{m'}) \right)$ should be minimal (detailed heuristic: A. & Lerasle 2012)

Bias and estimator selection: generic example

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

- Assume

$$\mathbb{E} \left[P\gamma(\widehat{s}_m(D_n)) \right] = \alpha(m) + \frac{\beta(m)}{n}$$

(e.g., LS/ridge/ k NN regression, LS/kernel density estimation).

Bias and estimator selection: generic example

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

- Assume

$$\mathbb{E} \left[P\gamma(\widehat{s}_m(D_n)) \right] = \alpha(m) + \frac{\beta(m)}{n}$$

(e.g., LS/ridge/ k NN regression, LS/kernel density estimation).

- Key quantities:**

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

$$\mathbb{E} \left[\widehat{\mathcal{R}}^{\text{cv}}(\widehat{s}_m) - \widehat{\mathcal{R}}^{\text{cv}}(\widehat{s}_{m'}) \right] = \alpha(m) - \alpha(m') + \frac{n}{n_t} \frac{\beta(m) - \beta(m')}{n}$$

\Rightarrow CV favours m with smaller complexity $\beta(m)$, more and more as n_t decreases.

CV with an estimation goal: the big picture (\mathcal{M} “small”)

- At first order, the **bias drives the performance** of:
 - leave- p -out, V -fold CV,
 - Monte-Carlo CV if $B \gg n^2$
or if n_v large enough (including hold-out)
- CV performs similarly to

$$\operatorname{argmin}_{m \in \mathcal{M}} \left\{ \mathbb{E} \left[P_\gamma(\hat{s}_m(D_{n_t})) \right] \right\}$$

CV with an estimation goal: the big picture (\mathcal{M} "small")

- At first order, the bias drives the performance of:
 - leave- p -out, V -fold CV,
 - Monte-Carlo CV if $B \gg n^2$
 - or if n_v large enough (including hold-out)
- CV performs similarly to

$$\operatorname{argmin}_{m \in \mathcal{M}} \left\{ \mathbb{E} \left[P_\gamma(\widehat{s}_m(D_{n_t})) \right] \right\}$$

\Rightarrow first-order optimality if $n_t \sim n$

\Rightarrow suboptimal otherwise

e.g., V -fold CV with V fixed.

- Theoretical results for least-squares regression and density estimation at least.

Bias-corrected VFCV / V-fold penalization

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

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

Bias-corrected VFCV / V-fold penalization

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

$$\begin{aligned}\widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{s}_m; D_n; \mathcal{B}) &:= \widehat{\mathcal{R}}^{\text{vf}}(\widehat{s}_m; D_n; \mathcal{B}) + P_n\gamma(\widehat{s}_m) - \frac{1}{V} \sum_{j=1}^V P_n\gamma(\widehat{s}_m^{(-j)}) \\ &= P_n\gamma(\widehat{s}_m) + \underbrace{\text{pen}_{\text{VF}}(\widehat{s}_m; D_n; \mathcal{B})}_{\text{V-fold penalty (A. 2008)}}\end{aligned}$$

Bias-corrected VFCV / V-fold penalization

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

$$\begin{aligned}\widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{s}_m; D_n; \mathcal{B}) &:= \widehat{\mathcal{R}}^{\text{vf}}(\widehat{s}_m; D_n; \mathcal{B}) + P_n\gamma(\widehat{s}_m) - \frac{1}{V} \sum_{j=1}^V P_n\gamma(\widehat{s}_m^{(-j)}) \\ &= P_n\gamma(\widehat{s}_m) + \underbrace{\text{pen}_{\text{VF}}(\widehat{s}_m; D_n; \mathcal{B})}_{\text{V-fold penalty (A. 2008)}}\end{aligned}$$

- In least-squares density estimation (A. & Lerasle, 2012):

$$\widehat{\mathcal{R}}^{\text{vf}}(\widehat{s}_m; D_n; \mathcal{B}) = P_n\gamma(\widehat{s}_m(D_n)) + \underbrace{\left(1 + \frac{1}{2(V-1)}\right)}_{\text{overpenalization factor}} \text{pen}_{\text{VF}}(\widehat{s}_m; D_n; \mathcal{B})$$

$$\widehat{\mathcal{R}}^{\ell\text{po}}(\widehat{s}_m; D_n; \mathcal{B}) = P_n\gamma(\widehat{s}_m(D_n)) + \underbrace{\left(1 + \frac{1}{2\left(\frac{n}{p} - 1\right)}\right)}_{\text{overpenalization factor}} \text{pen}_{\text{VF}}(\widehat{s}_m; D_n; \mathcal{B}_{\text{loo}})$$

Variance and estimator selection

$$\Delta(m, m', V) = \widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{s}_m) - \widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{s}_{m'})$$

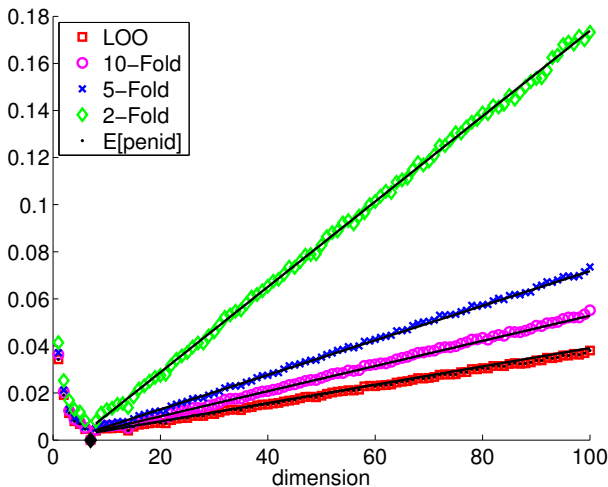
Theorem (A. & Lerasle 2012, least-squares density estimation)

$$\begin{aligned} \text{var}(\Delta(m, m', V)) &= 4 \left(1 + \frac{2}{n} + \frac{1}{n^2} \right) \frac{\text{var}_\rho(s_m^* - s_{m'}^*)}{n} \\ &\quad + 2 \left(1 + \frac{4}{V-1} - \frac{1}{n} \right) \underbrace{\frac{B(m, m')}{n^2}}_{\geq 0} \end{aligned}$$

If $S_m \subset S_{m'}$ are two histogram models with constant bin sizes $d_m^{-1}, d_{m'}^{-1}$, then, $B(m, m') \propto \|s_m^* - s_{m'}^*\| d_m$.

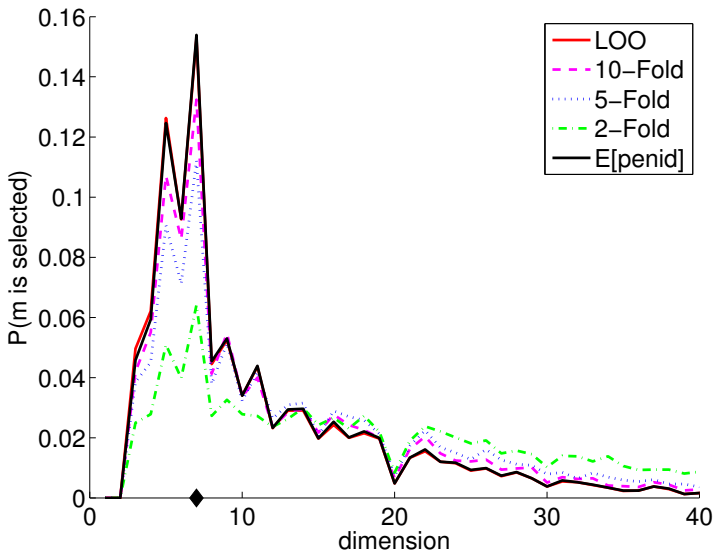
The two terms are of the same order if $\|s_m^* - s_{m'}^*\| \approx d_m/n$.

Variance of $\widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{s}_m) - \widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{s}_{m^*})$ vs. (d_m, V)

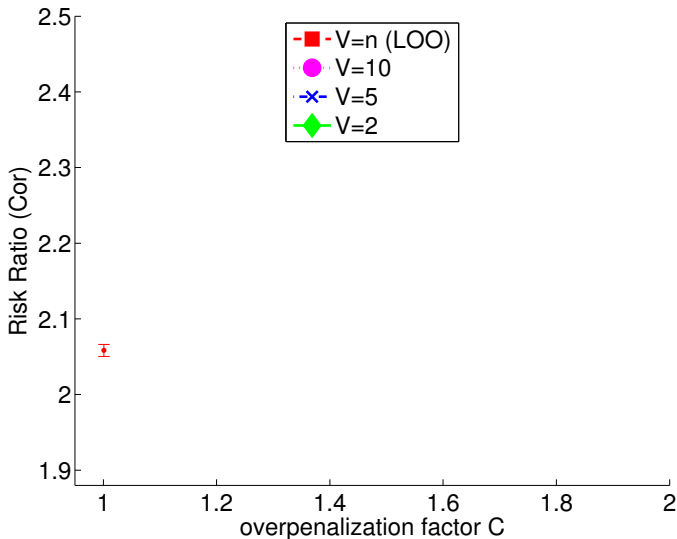


$$\text{var}(\Delta(m, m', V)) \approx n^{-2} \left[29 \left(1 + \frac{0.8}{V-1} \right) + 3.7 \left(1 + \frac{3.8}{V-1} \right) (d_m - d_{m^*}) \right]$$

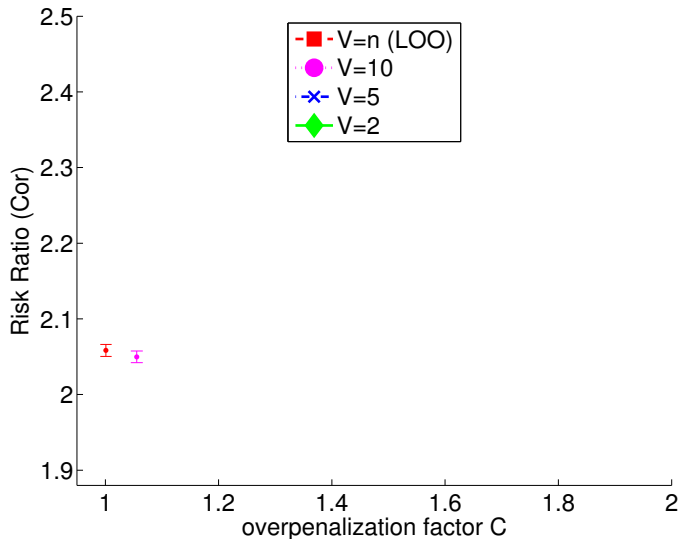
Probability of selection of every m



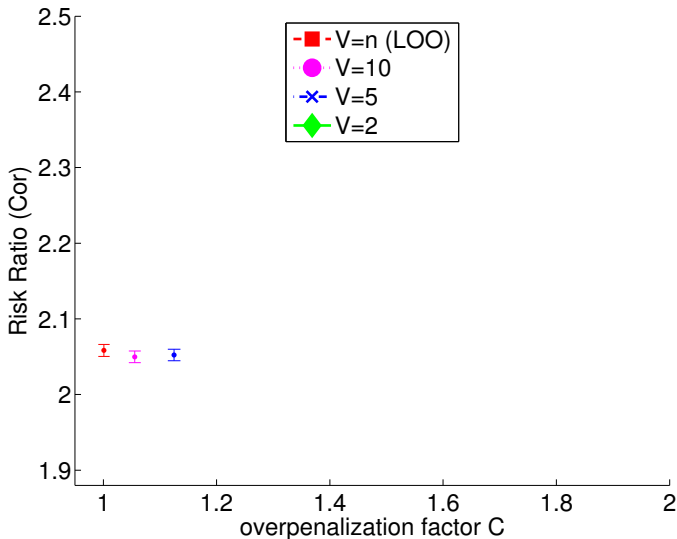
Experiment (LS density estimation): V -fold CV



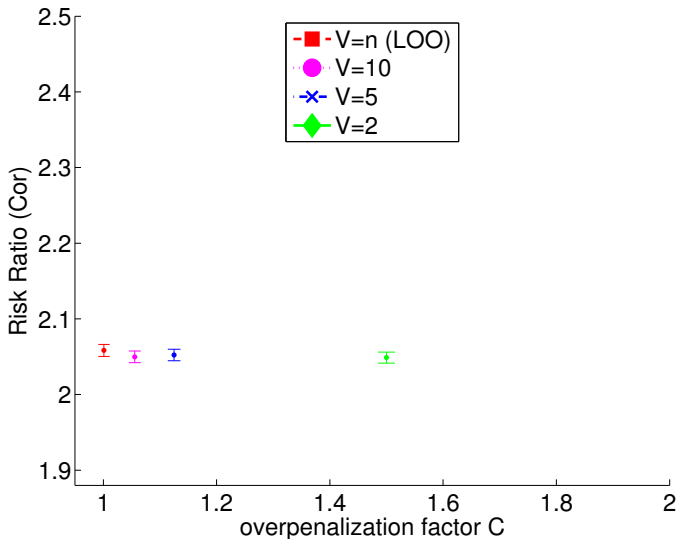
Experiment (LS density estimation): V -fold CV



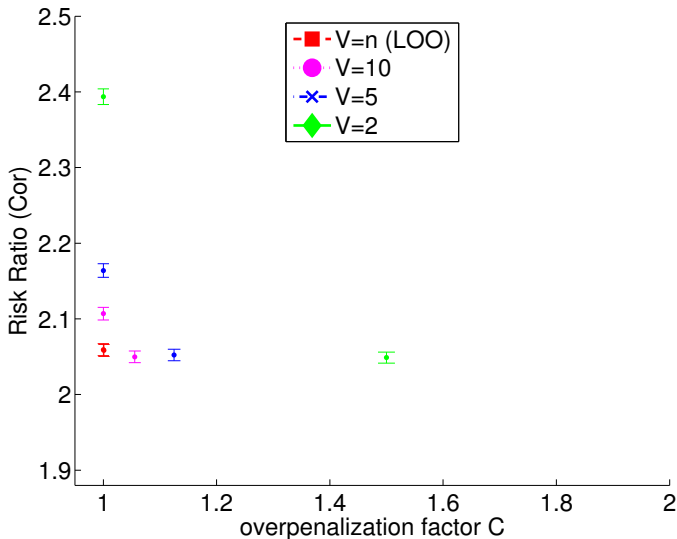
Experiment (LS density estimation): V -fold CV



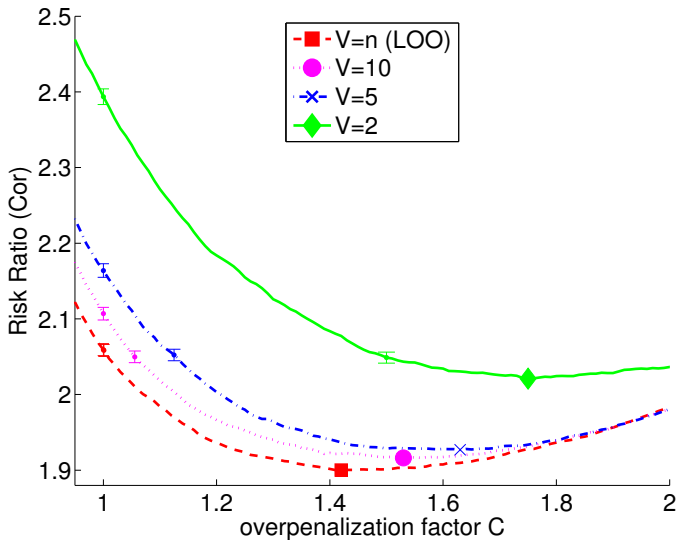
Experiment (LS density estimation): V -fold CV



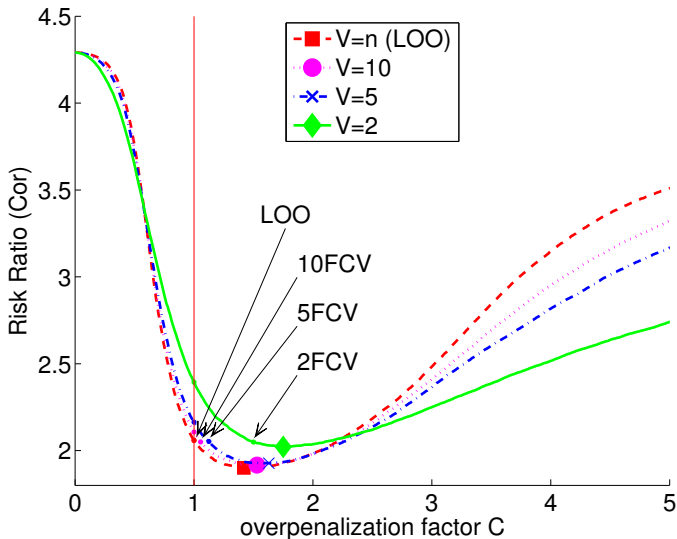
Experiment (LS density estimation): V -fold penalization



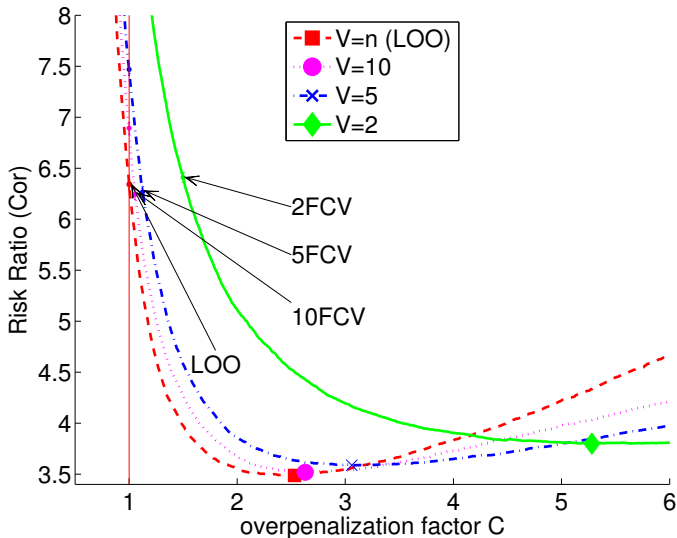
Experiment (LS density estimation): overpenalization



Experiment (LS density estimation): conclusion



Experiment (LS density estimation): other setting



Estimator selection with V -fold: conclusion

- **Computational complexity:** $\mathcal{O}(V)$ in general

Estimator selection with V -fold: conclusion

- **Computational complexity:** $\mathcal{O}(V)$ in general
 - **V -fold cross-validation:**
 - **Bias:** decreases with V / can be removed
 - **Variance:** decreases with V / almost minimal with $V \in [5, 10]$
- ⇒ best performance for the largest V and **almost optimal with $V = 10$...**

Estimator selection with V -fold: conclusion

- **Computational complexity:** $\mathcal{O}(V)$ in general
 - **V -fold cross-validation:**
 - Bias: decreases with V / can be removed
 - Variance: decreases with V / almost minimal with $V \in [5, 10]$
- ⇒ best performance for the largest V and **almost optimal with $V = 10$...**
- ... **if optimal overpenalization factor $C^* \approx 1$ (various behaviours possible).**

Estimator selection with V -fold: conclusion

- **Computational complexity:** $\mathcal{O}(V)$ in general
- **V -fold cross-validation:**
 - Bias: decreases with V / can be removed
 - Variance: decreases with V / almost minimal with $V \in [5, 10]$

⇒ best performance for the largest V and almost optimal with $V = 10...$

... if optimal overpenalization factor $C^* \approx 1$ (various behaviours possible).
- **V -fold penalization:**
 - **Decoupling** of bias and variance ⇒ **easier to understand.**
 - Bias: **chosen directly** through C , **without any constraint.**
 - Variance: decreases with V / **almost minimal with $V \in [5, 10]$.**

Outline

- 1 Estimator selection
- 2 Cross-validation
- 3 Cross-validation for risk estimation
- 4 Cross-validation for estimator selection
- 5 Large \mathcal{M}
- 6 Conclusion

Large collection of estimators/models

- Estimator/model selection with an “exponential” collection (implicitly excluded in all results above).
⇒ Expectations do not drive the first order!

Large collection of estimators/models

- Estimator/model selection with an “exponential” collection (implicitly excluded in all results above).
 ⇒ Expectations do not drive the first order!
- Examples: variable selection with $p \geq n$ variables, change-point detection.

Large collection of estimators/models

- Estimator/model selection with an “exponential” collection (implicitly excluded in all results above).
⇒ Expectations do not drive the first order!
- Examples: variable selection with $p \geq n$ variables, change-point detection.
- Solution: group the models \Rightarrow one estimator per dimension (e.g., empirical risk minimizer)
works for change-point detection (A. & Celisse, 2010).

Change-point detection and model selection

$$Y_i = \eta(t_i) + \sigma(t_i)\varepsilon_i \quad \text{with} \quad \mathbb{E}[\varepsilon_i] = 0 \quad \mathbb{E}[\varepsilon_i^2] = 1$$

- Goal: detect the **change-points of the mean η** of the signal Y
- ⇒ Model selection, collection of regressograms with $\mathcal{M}_n = \mathfrak{P}_{\text{interv}}(\{t_1, \dots, t_n\})$ (partitions of \mathcal{X} into intervals)
- No assumption on the variance $\sigma(t_i)^2$

Classical approach (Lebarbier, 2005; ...)

- “Birgé-Massart” penalty (assumes $\sigma(t_i) \equiv \sigma$):

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}_n} \left\{ P_n \gamma(\hat{s}_m) + \frac{C\sigma^2 D_m}{n} \left(5 + 2 \log \left(\frac{n}{D_m} \right) \right) \right\}$$

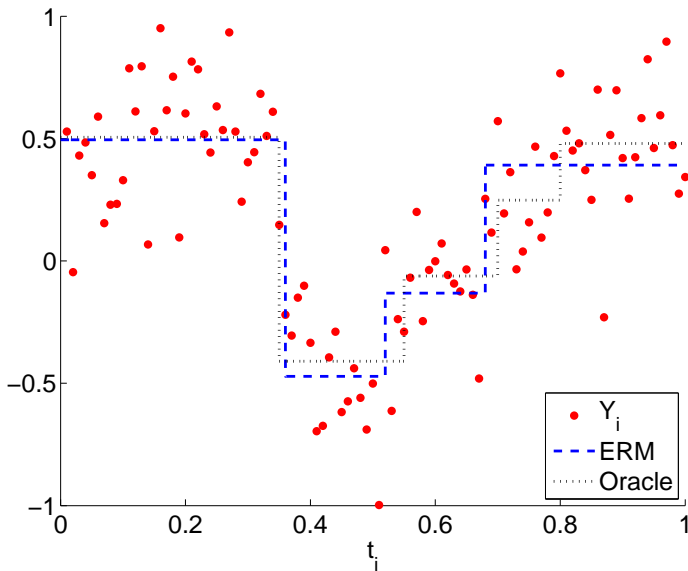
- Equivalent to aggregating models of the same dimension:

$$\tilde{S}_D := \bigcup_{m \in \mathcal{M}_n, D_m = D} S_m$$

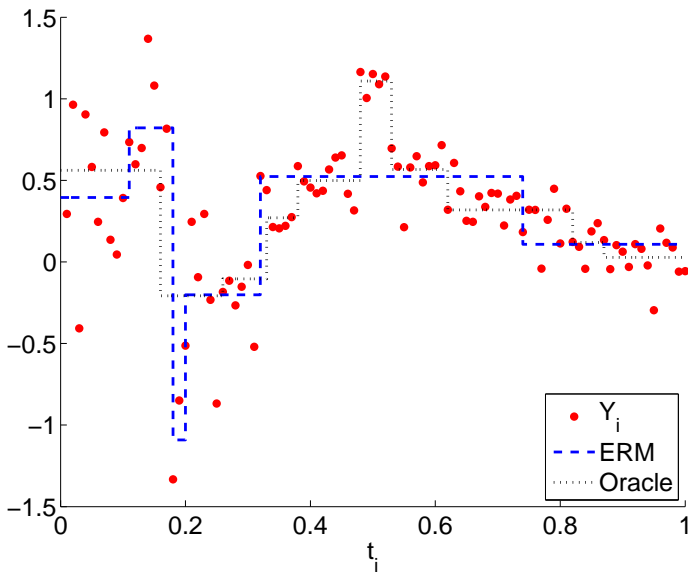
$$\hat{s}_D \in \operatorname{argmin}_{t \in \tilde{S}_D} \{ P_n \gamma(t) \} \quad \text{dynamic programming}$$

$$\hat{D} \in \operatorname{argmin}_{1 \leq D \leq n} \left\{ P_n \gamma(\hat{s}_D) + \frac{C\sigma^2 D}{n} \left(5 + 2 \log \left(\frac{n}{D} \right) \right) \right\}$$

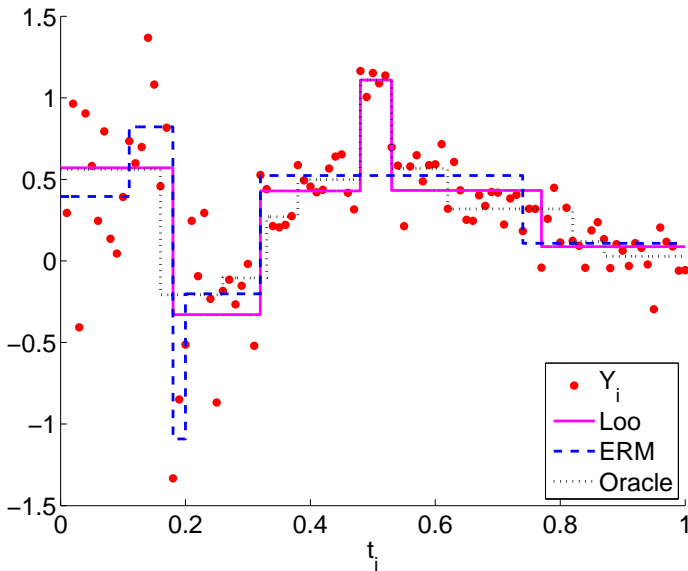
$D = 4$, homoscedastic; $n = 100$, $\sigma = 0.25$



$D = 6$, heteroscedastic; $n = 100$, $\|\sigma\| = 0.30$



$D = 6$, heteroscedastic; $n = 100$, $\|\sigma\| = 0.30$



Change-point detection algorithms (A. & Celisse, 2010)

- ① $\forall D \in \{1, \dots, D_{\max}\}$, **select**

$$\hat{m}(D) \in \operatorname{argmin}_{m \in \mathcal{M}_n, D_m = D} \left\{ \text{crit}_1(m; (t_i, Y_i)_i) \right\}$$

Examples for crit_1 : **empirical risk**, or **leave- p -out** or **V -fold estimators** of the risk (**dynamic programming**)

Change-point detection algorithms (A. & Celisse, 2010)

- ① $\forall D \in \{1, \dots, D_{\max}\}$, **select**

$$\hat{m}(D) \in \operatorname{argmin}_{m \in \mathcal{M}_n, D_m = D} \left\{ \operatorname{crit}_1(m; (t_i, Y_i)_i) \right\}$$

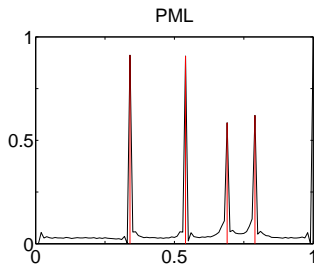
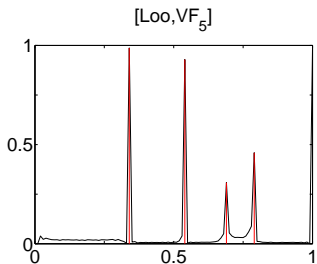
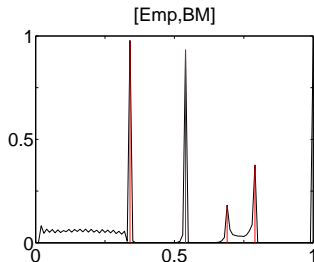
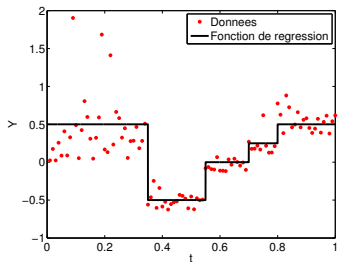
Examples for crit_1 : empirical risk, or leave- p -out or V -fold estimators of the risk (**dynamic programming**)

- ② **Select**

$$\hat{D} \in \operatorname{argmin}_{D \in \{1, \dots, D_{\max}\}} \left\{ \operatorname{crit}_2(D; (t_i, Y_i)_i; \operatorname{crit}_1(\cdot)) \right\}$$

Examples for crit_2 : **penalized empirical criterion, V -fold estimator of the risk**

Simulations: position of the change-points



Outline

- 1 Estimator selection
- 2 Cross-validation
- 3 Cross-validation for risk estimation
- 4 Cross-validation for estimator selection
- 5 Large \mathcal{M}
- 6 **Conclusion**

Generality of the results

- At least valid for least-square regression / density estimation, kernel density estimation.

Generality of the results

- At least valid for least-square regression / density estimation, kernel density estimation.
- Bias-correction / V -fold penalization: valid if

$$\mathbb{E} \left[(P - P_n) \gamma(\hat{S}_m) \right] \approx \frac{\gamma(m)}{n} .$$

Otherwise: use repeated V -fold or Monte-Carlo CV with a well-chosen n_t .

Generality of the results

- At least valid for least-square regression / density estimation, kernel density estimation.
- Bias-correction / V -fold penalization: valid if

$$\mathbb{E} \left[(P - P_n) \gamma(\hat{s}_m) \right] \approx \frac{\gamma(m)}{n} .$$

Otherwise: use repeated V -fold or Monte-Carlo CV with a well-chosen n_t .

- **Variance: different behaviours can occur in other settings (experiments).**

Generality of the results

- At least valid for least-square regression / density estimation, kernel density estimation.
- Bias-correction / V -fold penalization: valid if

$$\mathbb{E} \left[(P - P_n) \gamma(\hat{s}_m) \right] \approx \frac{\gamma(m)}{n} .$$

Otherwise: use repeated V -fold or Monte-Carlo CV with a well-chosen n_t .

- Variance: different behaviours can occur in other settings (experiments).
- Everything can be checked on synthetic data: plot

$$n \rightarrow \mathbb{E} \left[P \gamma(\hat{s}_m(D_n)) \right] \quad \text{and} \quad m \rightarrow \text{var} \left(\hat{\mathcal{R}}^{\text{cv}}(\hat{s}_m) - \hat{\mathcal{R}}^{\text{cv}}(\hat{s}_{m^*}) \right) .$$

Cross-validation with an identification goal

- **Main change:** value of the optimal overpenalization factor C^* , often $C^* \rightarrow +\infty$ when $n \rightarrow +\infty$.

Cross-validation with an identification goal

- Main change: value of the optimal overpenalization factor C^* , often $C^* \rightarrow +\infty$ when $n \rightarrow +\infty$.
- ⇔ **Cross-validation paradox** (Yang, 2006, 2007): $n_t \ll n$ can be necessary!
- Why? Smaller $n_t \Rightarrow$ easier to distinguish the two best procedures...

Cross-validation with an identification goal

- Main change: value of the optimal overpenalization factor C^* , often $C^* \rightarrow +\infty$ when $n \rightarrow +\infty$.
- ⇔ **Cross-validation paradox** (Yang, 2006, 2007): $n_t \ll n$ can be necessary!
- Why? Smaller $n_t \Rightarrow$ easier to distinguish the two best procedures... **if** n_t large enough (asymptotic regime).

Cross-validation with an identification goal

- Main change: value of the optimal overpenalization factor C^* , often $C^* \rightarrow +\infty$ when $n \rightarrow +\infty$.
- ⇔ Cross-validation paradox (Yang, 2006, 2007): $n_t \ll n$ can be necessary!
 - Why? Smaller $n_t \Rightarrow$ easier to distinguish the two best procedures... if n_t large enough (asymptotic regime).
 - Remark: **estimation goal, parametric setting** \Rightarrow similar behaviour.

Dependent data

- $D_n^{(t)}, D_n^{(v)}$ dependent \Rightarrow CV heuristic fails!

\Rightarrow possible troubles for risk estimation (Hart & Wehrly, 1986; Opsomer et al., 2001).

Dependent data

- $D_n^{(t)}, D_n^{(v)}$ dependent \Rightarrow CV heuristic fails!

\Rightarrow possible troubles for risk estimation (Hart & Wehrly, 1986; Opsomer et al., 2001).

- **Solution for short-term dependence:**
remove some data at each split \Rightarrow gap between training and validation samples.

Estimator selection
oooooooooooooooooooo

Cross-validation
oo

CV for risk estimation
oooo

CV for estimator selection
oooooooooooooooo

Large \mathcal{M}
oooooooo

Conclusion
ooo●

Questions?

Part I

Appendix

Outline

1 Change-point detection

Competitors

- **[Emp, BM]**: assume $\sigma(\cdot) \equiv \sigma$

$$\operatorname{argmin}_{m \in \mathcal{M}_n} \left\{ P_n \gamma(\hat{s}_m) + \frac{C\hat{\sigma}^2 D_m}{n} \left(5 + 2 \log \left(\frac{n}{D_m} \right) \right) \right\}$$

- **BGH** (Baraud, Giraud & Huet 2009): multiplicative penalty, $\sigma(\cdot) \equiv \sigma$

$$\operatorname{argmin}_{m \in \mathcal{M}_n} \left\{ P_n \gamma(\hat{s}_m) \left[1 + \frac{\text{pen}_{\text{BGH}}(m)}{n - D_m} \right] \right\}$$

- **ZS** (Zhang & Siegmund, 2007): modified BIC, $\sigma(\cdot) \equiv \sigma$
- **PML** (Picard *et al.*, 2005): penalized maximum likelihood, looks for **change-points of (η, σ)** , assuming a Gaussian model

$$\operatorname{argmin}_{m \in \mathcal{M}_n} \left\{ \sum_{\lambda \in m} n \hat{p}_\lambda \log \left(\frac{1}{n \hat{p}_\lambda} \sum_{t_i \in \lambda} (Y_i - \hat{s}_m(t_i))^2 \right) + \hat{C}'' D_m \right\}$$

Simulations: comparison to the oracle (quadratic risk)

$$\frac{\mathbb{E}[\ell(s^*, \widehat{s}_m)]}{\mathbb{E}[\inf_{m \in \mathcal{M}_n} \{\ell(s^*, \widehat{s}_m)\}]}$$

 $N = 10\,000$ sample

$\mathcal{L}(\varepsilon)$	Gaussian	Gaussian	Gaussian
$\sigma(\cdot)$	homosc.	heterosc.	heterosc.
η	s_2	s_2	s_3
[Loo, VF ₅]	4.02 ± 0.02	4.95 ± 0.05	5.59 ± 0.02
[Emp, VF ₅]	3.99 ± 0.02	5.62 ± 0.05	6.13 ± 0.02
[Emp, BM]	3.58 ± 0.02	9.25 ± 0.06	6.24 ± 0.02
BGH	3.52 ± 0.02	10.13 ± 0.07	6.31 ± 0.02
ZS	3.62 ± 0.02	6.50 ± 0.05	6.61 ± 0.02
PML	4.34 ± 0.02	2.73 ± 0.03	4.99 ± 0.03

Simulations: comparison to the oracle (quadratic risk)

$$\frac{\mathbb{E}[\ell(s^*, \widehat{s}_m)]}{\mathbb{E}[\inf_{m \in \mathcal{M}_n} \{\ell(s^*, \widehat{s}_m)\}]}$$

 $N = 10\,000$ sample

$\mathcal{L}(\varepsilon)$ $\sigma(\cdot)$ η	Gaussian homosc. s_2	Exponential heterosc. s_2	Exponential heterosc. s_3
[Loo, VF ₅]	4.02 ± 0.02	4.47 ± 0.05	5.11 ± 0.03
[Emp, VF ₅]	3.99 ± 0.02	5.98 ± 0.07	6.22 ± 0.04
[Emp, BM]	3.58 ± 0.02	10.81 ± 0.09	6.45 ± 0.04
BGH	3.52 ± 0.02	11.67 ± 0.09	6.42 ± 0.04
ZS	3.62 ± 0.02	9.34 ± 0.09	6.83 ± 0.04
PML	4.34 ± 0.02	5.04 ± 0.06	5.40 ± 0.03