

Model selection and estimator selection for statistical learning

Sylvain Arlot

¹CNRS

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

Scuola Normale Superiore di Pisa, 14–23 February 2011

Outline of the 5 lectures

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

Part IV

Resampling and model selection

Outline

- 1 Regressograms in heteroscedastic regression
- 2 The shape of the penalty must be estimated
- 3 Resampling
- 4 Theoretical guarantees for regressograms
- 5 Least-squares density estimation
- 6 Conclusion

Outline

- 1 Regressograms in heteroscedastic regression
- 2 The shape of the penalty must be estimated
- 3 Resampling
- 4 Theoretical guarantees for regressograms
- 5 Least-squares density estimation
- 6 Conclusion

Heteroscedastic regression framework

- **Random** design: $(X_1, Y_1), \dots, (X_n, Y_n)$ i.i.d.

$$Y_i = \eta(X_i) + \varepsilon_i$$

$$\mathbb{E}[\varepsilon_i | X_i] = 0 \quad \text{and} \quad \mathbb{E}[\varepsilon_i^2 | X_i] = \sigma^2(X_i)$$

Heteroscedastic regression framework

- Random design: $(X_1, Y_1), \dots, (X_n, Y_n)$ i.i.d.

$$Y_i = \eta(X_i) + \varepsilon_i$$

$$\mathbb{E}[\varepsilon_i | X_i] = 0 \quad \text{and} \quad \mathbb{E}[\varepsilon_i^2 | X_i] = \sigma^2(X_i)$$

- Quadratic loss:

$$P\gamma(t) = \mathbb{E}_{(X,Y) \sim P} [\gamma(t; (X, Y))] = \mathbb{E}_{(X,Y) \sim P} \left[(t(X) - Y)^2 \right]$$

- Excess loss: $\eta = s^*$ and

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

Regressograms

For any finite partition m of \mathcal{X}

$$S_m := \left\{ \sum_{\lambda \in m} \alpha_\lambda \mathbb{1}_\lambda \text{ s.t. } \alpha \in \mathbb{R}^m \right\}$$

\Rightarrow least-squares estimator over S_m (**regressogram**):

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

Regressograms

For any finite partition m of \mathcal{X}

$$S_m := \left\{ \sum_{\lambda \in m} \alpha_\lambda \mathbb{1}_\lambda \text{ s.t. } \alpha \in \mathbb{R}^m \right\}$$

\Rightarrow least-squares estimator over S_m (regressogram):

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

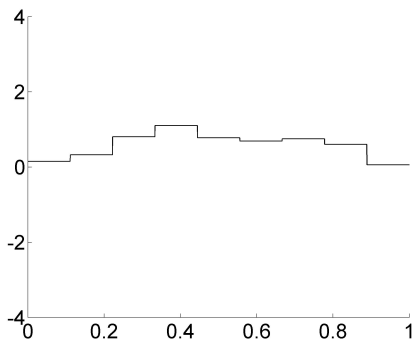
If for every $\lambda \in m$

$$\hat{p}_\lambda = \hat{p}_\lambda(D_n) = \frac{1}{n} \text{Card} \{ i \text{ s.t. } X_i \in \lambda \} > 0$$

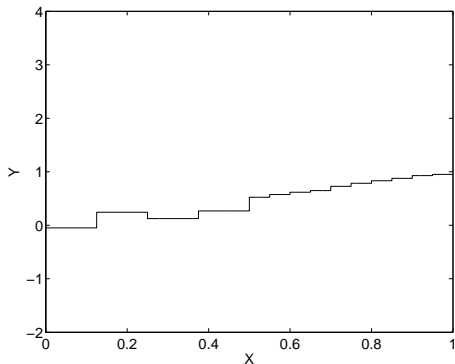
$$\hat{s}_m = \sum_{\lambda \in m} \hat{\beta}_\lambda \mathbb{1}_\lambda \quad \hat{\beta}_\lambda := \frac{1}{n \hat{p}_\lambda} \sum_{i \text{ s.t. } X_i \in \lambda} Y_i$$

Regressograms: examples ($\mathcal{X} = [0, 1]$)

$\mathcal{M}_n^{(\text{reg})}$ (regular partitions)



$\mathcal{M}_n^{(\text{reg}, 1/2)}$ (regular partitions on $[0, 1/2]$ and on $[1/2, 1]$)



Regressograms: bias, ideal penalty

Regressograms: bias, ideal penalty

$$s_m^* = \sum_{\lambda \in m} \beta_\lambda \mathbb{1}_\lambda \quad \beta_\lambda := \mathbb{E}_{(X,Y) \sim P} [Y | X \in \lambda]$$

$$\ell(s^*, s_m^*) = \sum_{\lambda \in m} p_\lambda \left(\sigma_\lambda^{(d)} \right)^2 \quad \left(\sigma_\lambda^{(d)} \right)^2 := \mathbb{E} \left[(\beta_\lambda - s^*(X))^2 \mid X \in \lambda \right]$$

Regressograms: bias, ideal penalty

$$s_m^* = \sum_{\lambda \in m} \beta_\lambda \mathbb{1}_\lambda \quad \beta_\lambda := \mathbb{E}_{(X,Y) \sim P} [Y \mid X \in \lambda]$$

$$\ell(s^*, s_m^*) = \sum_{\lambda \in m} p_\lambda \left(\sigma_\lambda^{(d)} \right)^2 \quad \left(\sigma_\lambda^{(d)} \right)^2 := \mathbb{E} \left[(\beta_\lambda - s^*(X))^2 \mid X \in \lambda \right]$$

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

$$p_1(m) = P(\gamma(\widehat{s}_m) - \gamma(s_m^*)) = \sum_{\lambda \in m} p_\lambda \left(\widehat{\beta}_\lambda - \beta_\lambda \right)^2$$

$$p_2(m) = P_n(\gamma(s_m^*) - \gamma(\widehat{s}_m)) = \sum_{\lambda \in m} \widehat{p}_\lambda \left(\widehat{\beta}_\lambda - \beta_\lambda \right)^2$$

$$\delta(m) = (P_n - P)\gamma(s_m^*)$$

Regressograms: conditional expectations

$$\mathcal{P}_m := (\mathbb{1}_{X_i \in \lambda})_{1 \leq i \leq n, \lambda \in m}$$

$$\mathbb{E}[p_1(m) \mid \mathcal{P}_m] = \frac{1}{n} \sum_{\lambda \in m} \frac{p_\lambda}{\hat{p}_\lambda} \sigma_\lambda^2$$

$$\mathbb{E}[p_2(m) \mid \mathcal{P}_m] = \frac{1}{n} \sum_{\lambda \in m} \sigma_\lambda^2$$

$$\sigma_\lambda^2 := \mathbb{E}_{(X,Y) \sim P} \left[(Y - \beta_\lambda)^2 \mid X \in \lambda \right] = \left(\sigma_\lambda^{(d)} \right)^2 + \left(\sigma_\lambda^{(r)} \right)^2$$

$$\left(\sigma_\lambda^{(r)} \right)^2 := \mathbb{E}_{(X,Y) \sim P} \left[(\sigma(X))^2 \mid X \in \lambda \right]$$

Regressograms: expectations

$$\mathbb{E} [p_1(m)] = \frac{1}{n} \sum_{\lambda \in m} \sigma_\lambda^2 (1 + \delta_{n,p_\lambda})$$

$$\mathbb{E} [p_2(m)] = \frac{1}{n} \sum_{\lambda \in m} \sigma_\lambda^2$$

$$\delta_{n,p_\lambda} := \mathbb{E} \left[\frac{p_\lambda}{\hat{p}_\lambda} \mid \hat{p}_\lambda > 0 \right] - 1$$

$$-\exp(-np) \leq \delta_{n,p} \leq \min \left\{ 1 + \frac{\kappa_1}{(np)^{1/4}}, \kappa_2 \right\}$$

Regressograms: risk, expectation of the ideal penalty

$$\mathbb{E}[\ell(s^*, \hat{s}_m)] = \sum_{\lambda \in m} p_\lambda \left(\sigma_\lambda^{(d)} \right)^2 + \frac{1}{n} \sum_{\lambda \in m} (1 + \delta_{n, p_\lambda}) \sigma_\lambda^2$$

$$\mathbb{E}[\text{pen}_{\text{id}}(m)] = \frac{1}{n} \sum_{\lambda \in m} (2 + \delta_{n, p_\lambda}) \sigma_\lambda^2$$

Outline

- 1 Regressograms in heteroscedastic regression
- 2 The shape of the penalty must be estimated
- 3 Resampling
- 4 Theoretical guarantees for regressograms
- 5 Least-squares density estimation
- 6 Conclusion

Drawbacks of pen = pen(D_m)

$$Y = s^*(X) + \varepsilon \quad \text{with} \quad X \sim \mathcal{U}([0, 1])$$

$$\mathbb{E} [\varepsilon^2 \mid X] = \sigma(X) \quad \text{and} \quad \int_0^{1/2} (\sigma(x))^2 dx \neq \int_{1/2}^1 (\sigma(x))^2 dx$$

$m \in \mathcal{M}_n^{(\text{reg}, 1/2)}$: $D_{m,1}$ pieces on $[0, \frac{1}{2}]$
 $D_{m,2}$ pieces on $[\frac{1}{2}, 1]$

Drawbacks of $\text{pen} = \text{pen}(D_m)$

$$Y = s^*(X) + \varepsilon \quad \text{with} \quad X \sim \mathcal{U}([0, 1])$$

$$\mathbb{E}[\varepsilon^2 | X] = \sigma(X) \quad \text{and} \quad \int_0^{1/2} (\sigma(x))^2 dx \neq \int_{1/2}^1 (\sigma(x))^2 dx$$

$$m \in \mathcal{M}_n^{(\text{reg}, 1/2)}: \quad D_{m,1} \text{ pieces on } [0, \frac{1}{2}]$$

$$D_{m,2} \text{ pieces on } [\frac{1}{2}, 1]$$

$$\mathbb{E}[\text{pen}_{\text{id}}(m)] \approx \frac{4}{n} \left[D_{m,1} \int_0^{1/2} (\sigma(x))^2 dx + D_{m,2} \int_{1/2}^1 (\sigma(x))^2 dx \right]$$

Drawbacks of $\text{pen} = \text{pen}(D_m)$: an example

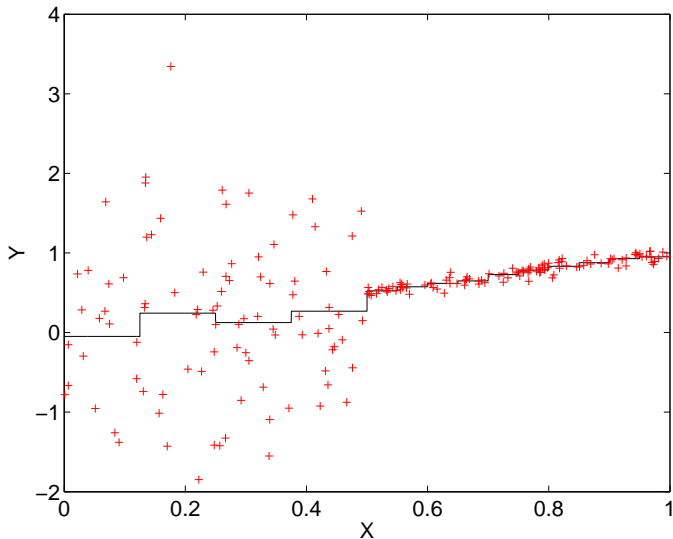
$$Y = s^*(X) + \varepsilon \quad \text{with} \quad X \sim \mathcal{U}([0, 1])$$

$$\mathcal{L}(\varepsilon | X) = \mathcal{N}(0, \sigma(X)^2)$$

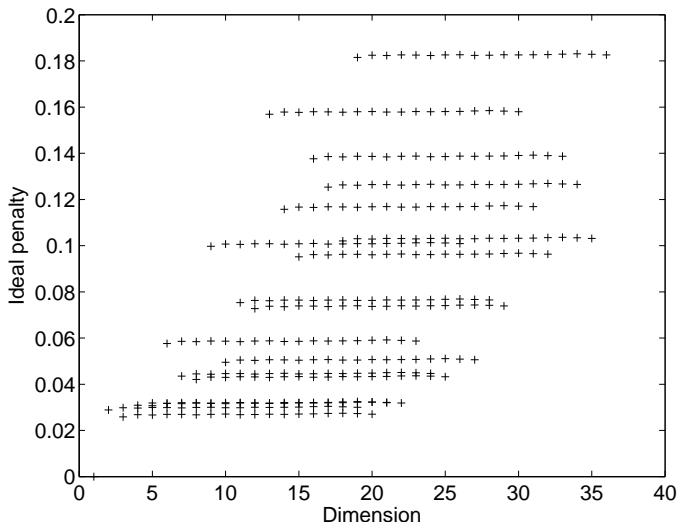
$$s^*(X) = X \quad \sigma(X) = \mathbb{1}_{X \leq \frac{1}{2}} + \frac{1}{20} \mathbb{1}_{X > 1/2}$$

$$\mathbb{E}[\text{pen}_{\text{id}}(m)] \approx \frac{2}{n} \left[D_{m,1} + \frac{D_{m,2}}{400} \right]$$

Example: data and oracle ($n = 200$)



Example: $\text{pen}_{\text{id}}(m)$ as a function of D_m



Penalties function of the dimension

Lemma

For any $D \in \mathcal{D}_n = \{D_m \text{ s.t. } m \in \mathcal{M}_n\}$

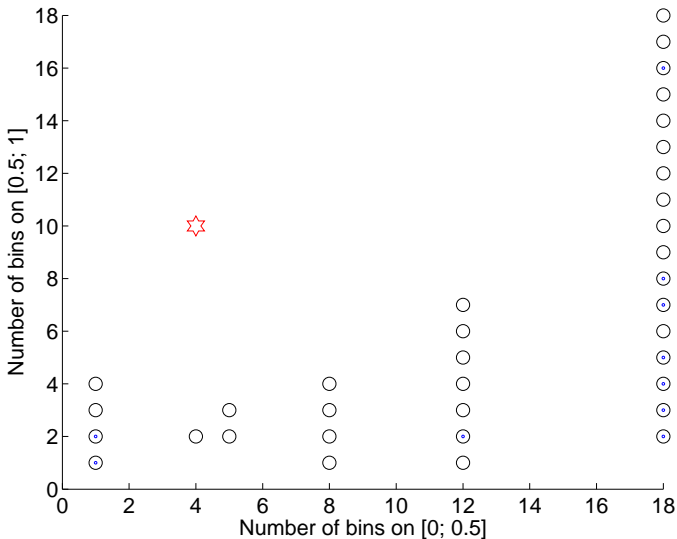
$$\mathcal{M}_{\dim}(D) := \operatorname{argmin}_{m \in \mathcal{M}_n \text{ s.t. } D_m = D} \{P_n \gamma(\widehat{s}_m)\}$$

$$\mathcal{M}_{\dim} := \bigcup_{D \in \mathcal{D}_n} \mathcal{M}_{\dim}(D)$$

Then, $\forall F : \mathcal{M}_n \mapsto \mathbb{R} \forall (X_i, Y_i)_{1 \leq i \leq n}$

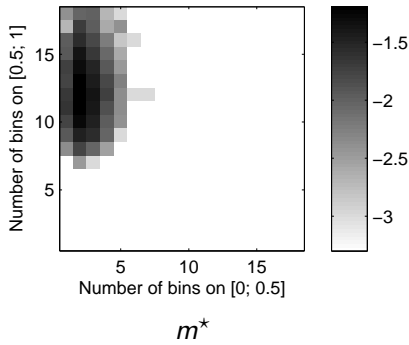
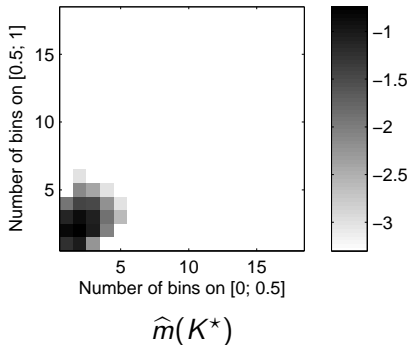
$$\operatorname{argmin}_{m \in \mathcal{M}_n} \{P_n \gamma(\widehat{s}_m) + F(D_m)\} \subset \mathcal{M}_{\dim}$$

Models that can be selected with $\text{pen}(D_m)$



Drawbacks of $\text{pen} = \text{pen}(D_m): \hat{m}(D^*) \neq m^*$

Densities of $(D_{\hat{m}(D^*),1}, D_{\hat{m}(D^*),2})$ and $(D_{m^*,1}, D_{m^*,2})$ over $N = 1000$ samples



Towards a proof: concentration of $p_{i;d}$

Assumption: $\|Y\|_\infty \leq A < \infty$ and $\sigma(\cdot) \geq \sigma_{\min} > 0$

- Concentration of p_1 and p_2 :
if $\min_{\lambda \in m} \{np_\lambda\} \geq \diamond \ln(n)$, with probability at least $1 - Ln^{-\gamma}$,
for $i = 1, 2$

$$|p_i(m) - \mathbb{E}[p_i(m)]| \leq \frac{L_{A, \sigma_{\min}, \gamma} (\ln(n))^2}{\sqrt{D_m}} \mathbb{E}[p_2(m)]$$

- Bernstein's inequality: with probability at least $1 - 2e^{-x}$,

$$\forall \theta \in (0, 1] \quad , \quad |(P_n - P)(\gamma(s_m^*) - \gamma(s^*))| \leq \theta \ell(s^*, s_m^*) + \frac{6A^2 x}{\theta n}$$

Heuristical proof: expectations

$$\mathbb{E}[p_1(m)] \approx \mathbb{E}[p_2(m)] \approx \frac{\beta_1 D_{m,1}}{n} + \frac{\beta_2 D_{m,1}}{n}$$

$$\beta_1 = 2 \int_0^{1/2} \sigma^2 \quad \beta_2 = 2 \int_{1/2}^1 \sigma^2$$

$$\ell(s^*, s_m^*) \approx \frac{\alpha_1}{D_{m,1}^2} + \frac{\alpha_2}{D_{m,2}^2}$$

$$\alpha_1 = \frac{1}{48} \int_0^{1/2} (s^{*'})^2 \quad \alpha_2 = \frac{1}{48} \int_{1/2}^1 (s^{*'})^2$$

Heuristical proof: expectations

$$\beta_1 = 2 \int_0^{1/2} \sigma^2 \quad \beta_2 = 2 \int_{1/2}^1 \sigma^2$$

$$\alpha_1 = \frac{1}{48} \int_0^{1/2} (s^{*t})^2 \quad \alpha_2 = \frac{1}{48} \int_{1/2}^1 (s^{*t})^2$$

$$P_n \gamma(\hat{s}_m) - P \gamma(s^*) \approx \frac{\alpha_1}{D_{m,1}^2} + \frac{\alpha_2}{D_{m,2}^2} - \frac{\beta_1 D_{m,1}}{n} - \frac{\beta_2 D_{m,1}}{n}$$

$$\ell(s^*, \hat{s}_m) \approx \frac{\alpha_1}{D_{m,1}^2} + \frac{\alpha_2}{D_{m,2}^2} + \frac{\beta_1 D_{m,1}}{n} + \frac{\beta_2 D_{m,1}}{n}$$

Heuristical proof: expectations

$$\beta_1 = 2 \int_0^{1/2} \sigma^2 > \beta_2 = 2 \int_{1/2}^1 \sigma^2$$

$$\alpha_1 = \frac{1}{48} \int_0^{1/2} (s^{*'})^2 \quad \alpha_2 = \frac{1}{48} \int_{1/2}^1 (s^{*'})^2$$

$$P_n \gamma(\widehat{s}_m) - P \gamma(s^*) \approx \frac{\alpha_1}{D_{m,1}^2} + \frac{\alpha_2}{D_{m,2}^2} - \frac{\beta_1 D_{m,1}}{n} - \frac{\beta_2 D_{m,1}}{n}$$

$$\ell(s^*, \widehat{s}_m) \approx \frac{\alpha_1}{D_{m,1}^2} + \frac{\alpha_2}{D_{m,2}^2} + \frac{\beta_1 D_{m,1}}{n} + \frac{\beta_2 D_{m,1}}{n}$$

$$m^* \approx \left(\left(\frac{2\alpha_1 n}{\beta_1} \right)^{1/3}, \left(\frac{2\alpha_2 n}{\beta_2} \right)^{1/3} \right)$$

Drawbacks of $\text{pen} = \text{pen}(D_m)$: theory

$$Y = s^*(X) + \varepsilon \quad \text{with} \quad X \sim \mathcal{U}([0, 1]) \quad , \quad \mathbb{E}[\varepsilon^2 \mid X] = \sigma(X)$$

$$\text{and} \quad \sigma_a^2 = \int_0^{1/2} (\sigma(x))^2 dx \neq \int_{1/2}^1 (\sigma(x))^2 dx = \sigma_b^2$$

Theorem (A. 2008)

If $\mathcal{M} = \mathcal{M}_n^{(\text{reg}, 1/2)}$, under “reasonable” assumptions on $(s^*, \varepsilon, \sigma)$,
 $\exists \eta(\sigma_a^2/\sigma_b^2) > 0$ such that with probability at least
 $1 - C(\|\varepsilon\|_\infty, \sigma_a^2, \sigma_b^2, \|s^*\|_\infty, \|s^{**}\|_\infty) n^{-2}$

$$\forall F \quad , \quad \forall \widehat{m}_F \in \arg \min_{m \in \mathcal{M}_n} \{P_n \gamma(\widehat{s}_m) + F(D_m)\} \quad ,$$

$$\ell(s^*, \widehat{s}_{\widehat{m}_F}) \geq \left(1 + \eta\left(\frac{\sigma_a^2}{\sigma_b^2}\right)\right) \inf_{m \in \mathcal{M}_n} \{\ell(s^*, \widehat{s}_m)\}$$

Why should we estimate the shape of the penalty?

- $\text{pen}(D) = F(D) \Rightarrow$ loss of a factor $(1 + \eta) > 1$

Why should we estimate the shape of the penalty?

- $\text{pen}(D) = F(D) \Rightarrow$ loss of a factor $(1 + \eta) > 1$
- $\text{pen}(m) = 2\mathbb{E} [\sigma(X)^2] D_m/n \Rightarrow$ possible burst of the risk

Why should we estimate the shape of the penalty?

- $\text{pen}(D) = F(D) \Rightarrow$ loss of a factor $(1 + \eta) > 1$
- $\text{pen}(m) = 2\mathbb{E} [\sigma(X)^2] D_m/n \Rightarrow$ possible burst of the risk
- $\text{pen}(m) = 2 \|\sigma\|_\infty^2 D_m/n \Rightarrow$ oracle-inequality with constant $\mathcal{O}(\max \sigma^2 / \min \sigma^2)$

Why should we estimate the shape of the penalty?

- $\text{pen}(D) = F(D) \Rightarrow$ loss of a **factor** $(1 + \eta) > 1$
 - $\text{pen}(m) = 2\mathbb{E} [\sigma(X)^2] D_m/n \Rightarrow$ possible **burst of the risk**
 - $\text{pen}(m) = 2 \|\sigma\|_{\infty}^2 D_m/n \Rightarrow$ oracle-inequality with **constant**
 $\mathcal{O}(\max \sigma^2 / \min \sigma^2)$
- \Rightarrow must estimate $\mathbb{E}[\text{pen}_{\text{id}}(m)]$ for an oracle inequality with **constant** $(1 + o(1))$ and for avoiding overfitting

Outline

- 1 Regressograms in heteroscedastic regression
- 2 The shape of the penalty must be estimated
- 3 Resampling**
- 4 Theoretical guarantees for regressograms
- 5 Least-squares density estimation
- 6 Conclusion

Resampling heuristics (bootstrap, Efron 1979)

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

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

Resampling heuristics (bootstrap, Efron 1979)

Real world :

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



Bootstrap world :

$$P_n$$

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

Resampling heuristics (bootstrap, Efron 1979)

Real world :

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



Bootstrap world :

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

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

Exchangeable weighted resampling

$$P_n^W := \frac{1}{n} \sum_{i=1}^n W_i \delta_{\xi_i}$$

- **Bootstrap:**

$$W \sim \mathcal{M} \left(n; \frac{1}{n}, \dots, \frac{1}{n} \right)$$

Exchangeable weighted resampling

$$P_n^W := \frac{1}{n} \sum_{i=1}^n W_i \delta_{\xi_i}$$

- Efron(m) or m out of n bootstrap:

$$\frac{m}{n} W \sim \mathcal{M} \left(m; \frac{1}{n}, \dots, \frac{1}{n} \right)$$

Exchangeable weighted resampling

$$P_n^W := \frac{1}{n} \sum_{i=1}^n W_i \delta_{\xi_i} \quad \text{or} \quad \frac{1}{\sum_k W_k} \sum_{i=1}^n W_i \delta_{\xi_i} = \frac{1}{n} \sum_{i=1}^n \frac{W_i}{\overline{W}} \delta_{\xi_i}$$

- **Efron(m)** or m out of n bootstrap:

$$\frac{m}{n} W \sim \mathcal{M} \left(m; \frac{1}{n}, \dots, \frac{1}{n} \right)$$

- **Subsampling:**

- **Random-hold out(q)**, $q \in \{1, \dots, n-1\}$:

$$W_i = \frac{n}{q} \mathbb{1}_{i \in I} \quad \text{with} \quad I \sim \mathcal{U}(\mathfrak{P}_q(\{1, \dots, n\}))$$

- **Rademacher(p)** or Bernoulli:

$$pW_1, \dots, pW_n \text{ i.i.d. } \sim \mathcal{B}(p)$$

Theoretical justification: asymptotics

Theorem (van der Vaart & Wellner, 1996)

Let $(W_{n,1}, \dots, W_{n,n}) \in \mathbb{R}^n$ be a non-negative random vector, exchangeable, independent from $\xi_{1\dots n}$, **bounded** and such that

$$n^{-1} \sum_{i=1}^n (W_{n,i} - \overline{W}_n)^2 \xrightarrow{(p)} c^2 > 0 .$$

Then, as n goes to infinity,

$$\sup_{h \in BL_1} \left| \mathbb{E}_W \left[h \left(\sqrt{n} \left(P_n^W - \overline{W}_n P_n \right) \right) \right] - \mathbb{E} \left[h(c\mathbb{G}) \right] \right| \xrightarrow{(p)} 0$$

where \mathbb{G} is a Gaussian process, **limit of $\sqrt{n}(P_n - P)$** , with zero mean and covariance function $\text{cov}(f, g) = P(fg) - P(f)P(g)$.

Classical uses of resampling

- estimating a **variance**, a **quadratic risk**
- estimation and/or **bias correction**
- **confidence intervals**, p -values
- estimation of **prediction error**, model selection
- **stabilization** (bagging, random forests)
- ...

A resampling-based estimator of variance

Framework:

$$\xi_1, \dots, \xi_n \text{ i.i.d. } \sim P \quad \mathbb{E}[\xi_i] = \mu \quad \mathbb{E}\left[(\xi_i - \mu)^2\right] = \sigma^2$$

$$\begin{aligned} \sigma^2 &= n\mathbb{E}\left[\left(\frac{1}{n}\sum_{i=1}^n \xi_i - \mu\right)^2\right] = n\mathbb{E}\left[\left(\mathbb{E}_{\xi \sim P_n} \xi - \mathbb{E}_{\xi \sim P} \xi\right)^2\right] \\ &= n\mathbb{E}\left[F(P, P_n)\right] \end{aligned}$$

⇒ resampling-based estimator

$$\widehat{\sigma}_W^2 = n\mathbb{E}_W\left[F(P_n, P_n^W)\right]$$

A resampling-based estimator of variance

$$\widehat{\sigma}_W^2 = n \mathbb{E}_W \left[F(P_n, P_n^W) \right]$$

$$\widehat{\sigma}_W^2 = \frac{R_V^{(W)}}{n} \left[\sum_{i=1}^n (\xi_i - \mu)^2 - \frac{1}{n-1} \sum_{i \neq j} (\xi_i - \mu) (\xi_j - \mu) \right]$$

$$R_V^{(W)} := \mathbb{E}_W \left[\left(\frac{nW_i}{\sum_{k=1}^n W_k} - 1 \right)^2 \right]$$

Comparison with the classical variance estimator

Classical unbiased estimator of variance:

$$\widehat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n \left(\xi_i - \frac{1}{n} \sum_{k=1}^n \xi_k \right)^2$$

Comparison with the classical variance estimator

Classical unbiased estimator of variance:

$$\widehat{\sigma^2} = \frac{1}{n-1} \sum_{i=1}^n \left(\xi_i - \frac{1}{n} \sum_{k=1}^n \xi_k \right)^2$$

$$\begin{aligned} \widehat{\sigma_W^2} &= R_V^{(W)} \widehat{\sigma^2} \\ \Rightarrow \mathbb{E} \left[\widehat{\sigma_W^2} \right] &= R_V^{(W)} \sigma^2 \end{aligned}$$

Resampling and structure

- Properties of $F(P, P_n) = (\mathbb{E}_{\xi \sim P_n} \xi - \mathbb{E}_{\xi \sim P} \xi)^2$:
 - exchangeable
 - translation-invariance
 - homogeneity
 - polynomial function of ξ_i and $\mathbb{E}_{\xi \sim P} \xi$

$\Rightarrow \mathbb{E}_W[F(P_n, P_n^W)]$ has similar properties

Resampling and structure

- Properties of $F(P, P_n) = (\mathbb{E}_{\xi \sim P_n} \xi - \mathbb{E}_{\xi \sim P} \xi)^2$:
 - exchangeable
 - translation-invariance
 - homogeneity
 - polynomial function of ξ_i and $\mathbb{E}_{\xi \sim P} \xi$

$\Rightarrow \mathbb{E}_W[F(P_n, P_n^W)]$ has similar properties

$$\Rightarrow \mathbb{E}_W \left[F(P_n, P_n^W) \right] \propto \widehat{\sigma}^2$$

Resampling and concentration

Over-concentration phenomenon for the resampling-based estimator:

$$\text{var} \left(n \left(\frac{1}{n} \sum_{i=1}^n \xi_i - \mu \right)^2 \right) = 2\sigma^4 + \frac{\mathbb{E} \left[(\xi_1 - \mu)^4 \right] - 3\sigma^4}{n}$$

$$\text{var} \left(\frac{1}{R_V^{(W)}} \widehat{\sigma_W^2} \right) = \frac{1}{n} \left(\mathbb{E} \left[(\xi_1 - \mu)^4 \right] - \sigma^4 \right) + \frac{2}{n(n-1)} \sigma^4$$

Computation of the multiplicative factor

$$R_V^{(W)} := \mathbb{E}_W \left[\left(\frac{nW_i}{\sum_{k=1}^n W_k} - 1 \right)^2 \right]$$

Computation of the multiplicative factor

$$R_V^{(W)} := \mathbb{E}_W \left[\left(\frac{nW_i}{\sum_{k=1}^n W_k} - 1 \right)^2 \right]$$

Efron(m): $R_V^{(W)} = \frac{n-1}{m}$

Rademacher(p): $R_V^{(W)} = \frac{1 + \delta_{n,p}}{p} - 1 \approx \frac{1}{p} - 1$

Random hold-out(q): $R_V^{(W)} = \frac{n}{q} - 1$

Leave-one-out = Rho($n-1$): $R_V^{(W)} = \frac{1}{n-1}$

Resampling-based estimator of $\text{pen}_{\text{id}}(m)$

- Ideal penalty:

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

- Resampling-based estimator of $\mathbb{E}[F(P, P_n)]$:

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

- bootstrap (Efron, 1983; Shibata, 1997), m out of n bootstrap for identification (Shao, 1996), general exchangeable weights (A. 2009)
- **Multiplicative factor C_W** : why? how can we estimate it?

Rademacher penalties

- Global penalties:

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

Rademacher penalties

- Global penalties:

$$\text{pen}_{\text{id}}(m) \leq \text{pen}_{\text{id}}^{\text{glo}}(m) = \sup_{t \in \mathcal{S}_m} (P - P_n)\gamma(t)$$

- **Global Rademacher penalties** in classification (Koltchinskii & Panchenko, 2001; Bartlett, Boucheron & Lugosi, 2002), exchangeable weights (Fromont, 2004)

$$\mathbb{E} \left[\sup_{t \in \mathcal{S}_m} \left\{ \frac{1}{n} \sum_{i=1}^n \varepsilon_i \gamma(t; \xi_i) \right\} \middle| P_n \right]$$

with $\varepsilon_1, \dots, \varepsilon_n$ i.i.d. $\sim \mathcal{U}(\{-1, +1\})$

- **Local Rademacher complexities** (Bartlett, Bousquet & Mendelson, 2004; Koltchinskii, 2006)

Outline

- 1 Regressograms in heteroscedastic regression
- 2 The shape of the penalty must be estimated
- 3 Resampling
- 4 Theoretical guarantees for regressograms**
- 5 Least-squares density estimation
- 6 Conclusion

Reminder

$$\hat{s}_m = \sum_{\lambda \in m} \hat{\beta}_\lambda \mathbf{1}_\lambda \quad \text{with} \quad \hat{\beta}_\lambda := \frac{1}{n\hat{p}_\lambda} \sum_{i \text{ s.t. } X_i \in \lambda} Y_i$$

$$\hat{p}_\lambda = \hat{p}_\lambda(D_n) = \frac{1}{n} \text{Card} \{ i \text{ s.t. } X_i \in \lambda \}$$

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

$$p_1(m) = P(\gamma(\hat{s}_m) - \gamma(s_m^*)) = \sum_{\lambda \in m} \left[p_\lambda (\hat{\beta}_\lambda - \beta_\lambda)^2 \right]$$

$$p_2(m) = P_n(\gamma(s_m^*) - \gamma(\hat{s}_m)) = \sum_{\lambda \in m} \left[\hat{p}_\lambda (\hat{\beta}_\lambda - \beta_\lambda)^2 \right]$$

$$\delta(m) = (P_n - P)\gamma(s_m^*)$$

Resampling-based penalty

$$\text{pen}_W(m) = \frac{C_W}{n} \sum_{\lambda \in m} \frac{R_{1,W} + R_{2,W}}{n\hat{p}_\lambda - 1} \left(S_{\lambda,2} - \frac{1}{n\hat{p}_\lambda} S_{\lambda,1}^2 \right) \mathbb{1}_{n\hat{p}_\lambda \geq 2}$$

$$\text{with } S_{\lambda,1} := \sum_{X_i \in \lambda} (Y_i - \beta_\lambda) \quad S_{\lambda,2} := \sum_{X_i \in \lambda} (Y_i - \beta_\lambda)^2$$

$$R_{1,W}(n, \hat{p}_\lambda) := \mathbb{E} \left[\frac{(W_1 - \widehat{W}_\lambda)^2}{\widehat{W}_\lambda^2} \middle| X_1 \in \lambda, \widehat{W}_\lambda > 0 \right]$$

$$\text{and } R_{2,W}(n, \hat{p}_\lambda) := \mathbb{E} \left[\frac{(W_1 - \widehat{W}_\lambda)^2}{\widehat{W}_\lambda} \middle| X_1 \in \lambda \right]$$

Value of R_1 and R_2 : examples

$$R_{1,W}(n, \hat{p}_\lambda) \sim R_{2,W}(n, \hat{p}_\lambda) \text{ as } n\hat{p}_\lambda \rightarrow \infty$$

$$C_{W,\infty}(n) := \lim_{n\hat{p}_\lambda \rightarrow \infty} \frac{1}{R_{2,W}(n, \hat{p}_\lambda)}$$

Efron(m):	$R_{2,W}(n, \hat{p}_\lambda) = \frac{n}{m} \left(1 - \frac{1}{n\hat{p}_\lambda} \right)$	$C_{W,\infty} = \frac{m}{n}$
Rademacher(p):	$R_{2,W}(n, \hat{p}_\lambda) = \frac{1}{p} - 1$	$C_{W,\infty} = \frac{p}{1-p}$
Random hold-out(q):	$R_{2,W}(n, \hat{p}_\lambda) = \frac{n}{q} - 1$	$C_{W,\infty} = \frac{q}{n-q}$
Leave-one-out:	$R_{2,W}(n, \hat{p}_\lambda) = \frac{1}{n-1}$	$C_{W,\infty} = n-1$

Expectations

$$\mathbb{E}[Y_i - \beta_\lambda \mid X_i \in \lambda] = 0 \quad \text{and} \quad \mathbb{E}\left[(Y_i - \beta_\lambda)^2 \mid X_i \in \lambda\right] = \sigma_\lambda^2$$

$$\mathbb{E}[\text{pen}_W(m) \mid \mathcal{P}_m] = \frac{C_W}{n} \sum_{\lambda \in m} (R_{1,W} + R_{2,W}) \sigma_\lambda^2 \mathbb{1}_{n\hat{p}_\lambda \geq 2}$$

$$\mathbb{E}[\text{pen}_W(m)] = \frac{C_W}{C_{W,\infty}} \frac{1}{n} \sum_{\lambda \in m} \left(2 + \bar{\delta}_{n,p_\lambda}^{(\text{pen}W)}\right) \sigma_\lambda^2$$

with $\bar{\delta}_{n,p_\lambda}^{(\text{pen}W)} \rightarrow 0$ quand $np_\lambda \rightarrow +\infty$

\Rightarrow adaptation to heteroscedasticity

Concentration

Proposition (A. 2009)

- *Bounded data*: $\|Y_i\|_\infty \leq A < \infty$
- *Lower-bounded noise*: $\sigma(X_i) \geq \sigma_{\min} > 0$
- $\mathcal{L}(W)$ among $Efr(n)$, $Rad(1/2)$, $Rho(n/2)$, Loo

For every $A_n \geq 2$, with probability at least $1 - L_1 n^{-\gamma}$,

$$\begin{aligned} & |\text{pen}_W(m) - \mathbb{E}[\text{pen}_W(m) \mid \mathcal{P}_m]| \mathbf{1}_{\min_{\lambda \in m} \{n\hat{p}_\lambda\} \geq A_n} \\ & \leq \frac{C_W}{C_{W,\infty}} \frac{L_2(A/\sigma_{\min}, \gamma) \ln(n)}{\sqrt{A_n D_m}} \mathbb{E}[p_2(m)] \end{aligned}$$

Pathwise non-asymptotic oracle inequality

- $\mathcal{L}(W)$ among Efr(n), Rad(1/2), Rho($n/2$), Loo
- $C_W \approx C_{W,\infty}$

Pathwise non-asymptotic oracle inequality

- $\mathcal{L}(W)$ among Efr(n), Rad(1/2), Rho($n/2$), Loo
- $C_W \approx C_{W,\infty}$
- $\text{Card}(\mathcal{M}_n) \leq C_{\mathcal{M}} n^{\alpha_{\mathcal{M}}}$
- **Bounded data:** $\|Y_i\|_{\infty} \leq A < \infty$
- **Lower-bounded noise:** $\sigma(X_i) \geq \sigma_{\min} > 0$

Pathwise non-asymptotic oracle inequality

- $\mathcal{L}(W)$ among Efr(n), Rad($1/2$), Rho($n/2$), Loo
- $C_W \approx C_{W,\infty}$
- $\text{Card}(\mathcal{M}_n) \leq C_{\mathcal{M}} n^{\alpha_{\mathcal{M}}}$
- **Bounded data:** $\|Y_i\|_{\infty} \leq A < \infty$
- **Lower-bounded noise:** $\sigma(X_i) \geq \sigma_{\min} > 0$
- $s^* \in \mathcal{H}(\alpha, R)$ non-constant
- **Pre-selected models:** $\forall m \in \mathcal{M}, \min_{\lambda \in m} n \hat{p}_{\lambda} \geq 3$

Pathwise non-asymptotic oracle inequality

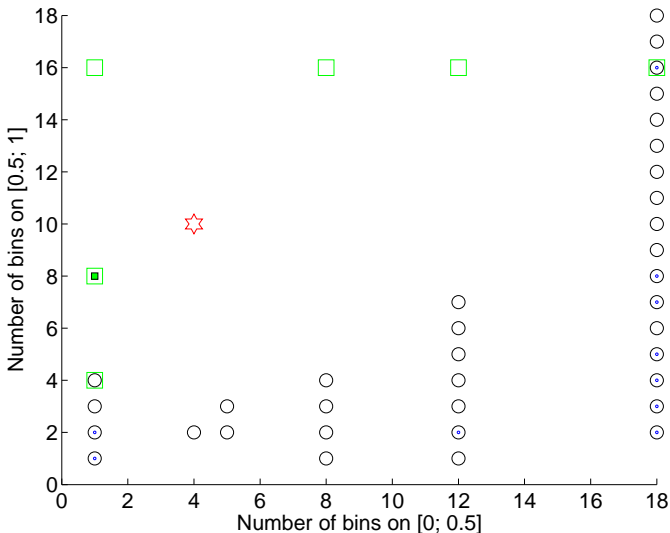
- $\mathcal{L}(W)$ among Efr(n), Rad(1/2), Rho($n/2$), Loo
- $C_W \approx C_{W,\infty}$
- $\text{Card}(\mathcal{M}_n) \leq C_{\mathcal{M}} n^{\alpha_{\mathcal{M}}}$
- **Bounded data**: $\|Y_i\|_{\infty} \leq A < \infty$
- **Lower-bounded noise**: $\sigma(X_i) \geq \sigma_{\min} > 0$
- $s^* \in \mathcal{H}(\alpha, R)$ non-constant
- **Pre-selected** models: $\forall m \in \mathcal{M}$, $\min_{\lambda \in m} n \hat{p}_{\lambda} \geq 3$

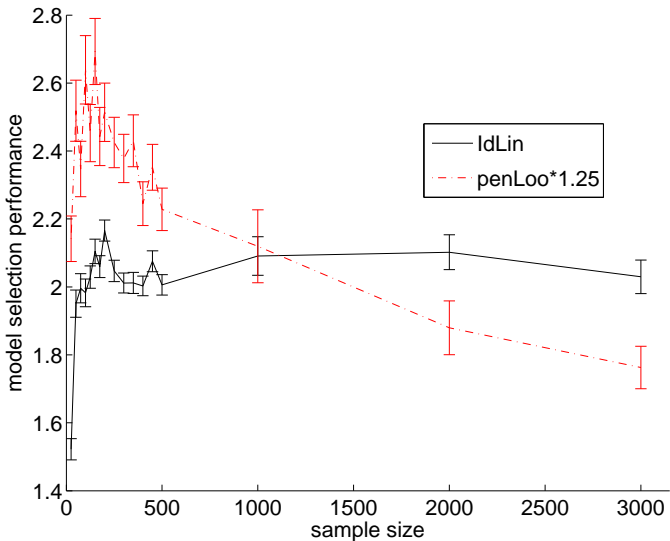
Theorem (A. 2009)

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

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

Models that can be selected: penLoo better than pen(D_m)



Simulations: $1.25 \times \text{pen}_{L_{00}}(m)$ vs. K^*D_m 

Adaptation

$$\tilde{s} := \hat{s}_{\hat{m}} \quad \text{with} \quad \hat{m} \in \underset{\substack{m \in \mathcal{M}_n^{(\text{reg})} \\ \min_{\lambda \in m} \{n\hat{p}_\lambda\} \geq 3}}{\text{argmin}} \{P_n \gamma(\hat{s}_m) + \text{pen}_W(m)\}$$

Assumptions:

- **Bounded data:** $\|Y_i\|_\infty \leq A < \infty$
- **Lower-bounded noise:** $\sigma(X_i) \geq \sigma_{\min} > 0$
- **Lower-bounded density of X :** $\forall I \subset \mathcal{X}$,
 $\mathbb{P}(X \in I) \geq c_X^{\min} \text{Leb}(I)$
- $s^* = \eta \in \mathcal{H}(\alpha, R)$ with $\alpha \in (0, 1]$:

$$\forall x_1, x_2 \in \mathcal{X}, \quad |s^*(x_1) - s^*(x_2)| \leq R \|x_1 - x_2\|_\infty^\alpha$$

Adaptation

$$\tilde{s} := \hat{s}_{\hat{m}} \quad \text{with} \quad \hat{m} \in \underset{m \in \mathcal{M}_n^{(\text{reg})}}{\text{argmin}} \{ P_n \gamma(\hat{s}_m) + \text{pen}_W(m) \}$$

$$\min_{\lambda \in m} \{ n \hat{p}_\lambda \} \geq 3$$

$$\mathbb{E}[\ell(s^*, \tilde{s})] \leq K_2 R^{\frac{2d}{2\alpha+d}} n^{\frac{-2\alpha}{2\alpha+d}} \|\sigma\|_\infty^{\frac{4\alpha}{2\alpha+d}} + \frac{K_3 A^2}{n^2}$$

Adaptation

$$\tilde{s} := \hat{s}_{\hat{m}} \quad \text{with} \quad \hat{m} \in \underset{m \in \mathcal{M}_n^{(\text{reg})}}{\text{argmin}} \{ P_n \gamma(\hat{s}_m) + \text{pen}_W(m) \}$$

$$\min_{\lambda \in m} \{ n \hat{p}_\lambda \} \geq 3$$

$$\mathbb{E}[\ell(s^*, \tilde{s})] \leq K_2 R^{\frac{2d}{2\alpha+d}} n^{\frac{-2\alpha}{2\alpha+d}} \|\sigma\|_\infty^{\frac{4\alpha}{2\alpha+d}} + \frac{K_3 A^2}{n^2}$$

and if $\sigma(\cdot)$ is K_σ -Lipschitz with at most J_σ jumps:

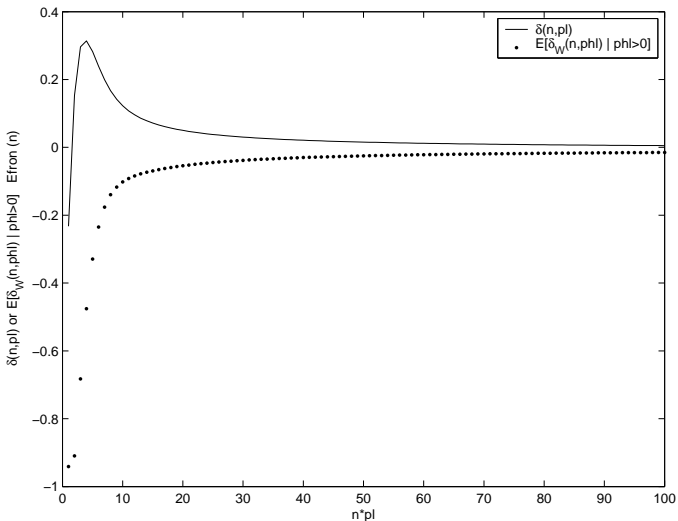
$$\mathbb{E}[\ell(s^*, \tilde{s})] \leq K_2 R^{\frac{2d}{2\alpha+d}} n^{\frac{-2\alpha}{2\alpha+d}} \|\sigma\|_{L^2(\text{Leb})}^{\frac{4\alpha}{2\alpha+d}} + \frac{K_4 A^2}{n^2}$$

Theoretical comparison of weights: reminder

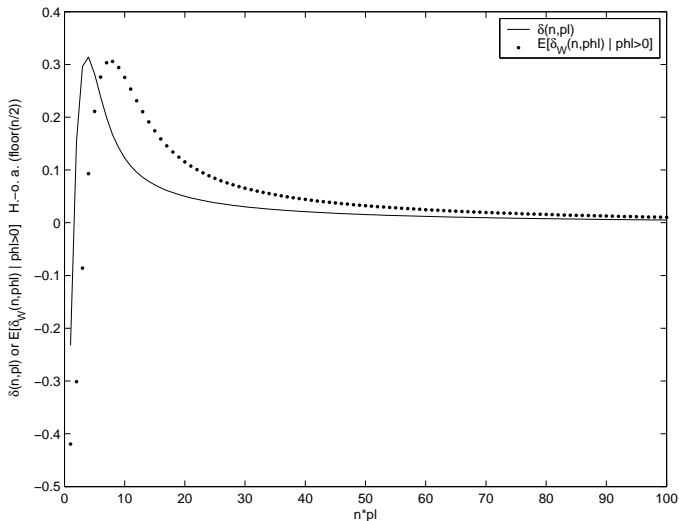
$$\mathbb{E}[\text{pen}_{\text{id}}(m)] = \frac{1}{n} \sum_{\lambda \in m} (2 + \delta_{n, p_\lambda}) \sigma_\lambda^2$$

and $\mathbb{E}[\text{pen}_W(m)] = \frac{1}{n} \sum_{\lambda \in m} \left(2 + \bar{\delta}_{n, \hat{p}_\lambda}^{(\text{pen}W)} \right) \sigma_\lambda^2$

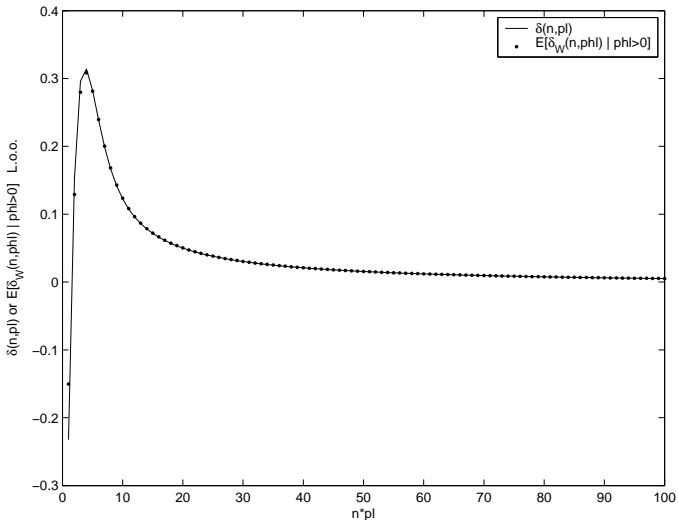
$\bar{\delta}_{n, \hat{p}_\lambda}(\text{pen}^W)$ vs. δ_{n, p_λ} : Efron(n) \approx Poisson(1)



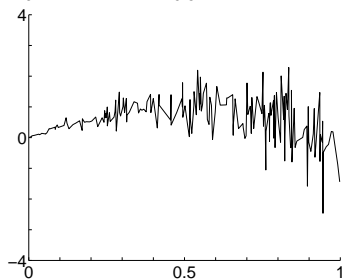
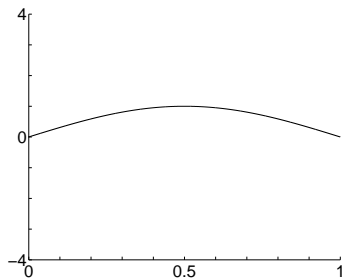
$\overline{\delta}_{n, \hat{p}_\lambda}^{(\text{pen}W)}$ vs. δ_{n, p_λ} : $\text{Rho}(n/2) \approx \text{Rad}(1/2)$



$\overline{\delta}_{n, \hat{p}_\lambda}(\text{pen}W)$ vs. δ_{n, p_λ} : Leave-one-out

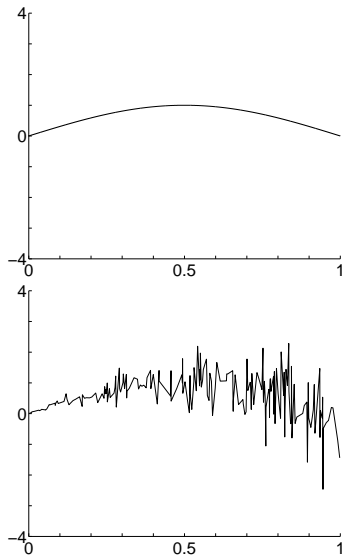


$$s^*(x) = \sin(2\pi x) \quad n = 200 \quad \sigma(x) = x \quad \mathcal{M} = \mathcal{M}^{(\text{reg}, 1/2)}$$



Mallows	3.69 ± 0.07
$\mathbb{E}[\text{pen}_{\text{id}}]$	2.30 ± 0.05
pen Efr	3.15 ± 0.07
pen Loo	2.59 ± 0.06
pen Rho	2.50 ± 0.06
pen Rad	2.49 ± 0.06

$$s^*(x) = \sin(2\pi x) \quad n = 200 \quad \sigma(x) = x \quad \mathcal{M} = \mathcal{M}^{(\text{reg}, 1/2)}$$



Mallows	3.69 ± 0.07
$\mathbb{E}[\text{pen}_{\text{id}}]$	2.30 ± 0.05
pen Efr	3.15 ± 0.07
pen Loo	2.59 ± 0.06
pen Rho	2.50 ± 0.06
pen Rad	2.49 ± 0.06
Mallows $\times 1.25$	3.17 ± 0.07
$\mathbb{E}[\text{pen}_{\text{id}}] \times 1.25$	2.03 ± 0.04
pen Efr $\times 1.25$	2.60 ± 0.06
pen Loo $\times 1.25$	2.22 ± 0.05
pen Rho $\times 1.25$	2.14 ± 0.05
pen Rad $\times 1.25$	2.14 ± 0.05

Outline

- 1 Regressograms in heteroscedastic regression
- 2 The shape of the penalty must be estimated
- 3 Resampling
- 4 Theoretical guarantees for regressograms
- 5 Least-squares density estimation**
- 6 Conclusion

Least-squares density estimation

- μ reference measure on Ξ
- $f = dP/d\mu \in \mathbb{S} = L^2(\mu)$

Least-squares density estimation

- μ reference measure on Ξ
- $f = dP/d\mu \in \mathbb{S} = L^2(\mu)$
- $\gamma(t; \xi) = \|t\|_{L^2(\mu)}^2 - 2t(\xi)$
 - $\Rightarrow P\gamma(t) = \|t\|_{L^2(\mu)}^2 - 2\langle t, f \rangle_{L^2(\mu)}$
 - $\Rightarrow s^* = f$ and $\ell(s^*, t) = \|t - s^*\|_{L^2(\mu)}^2$

Least-squares density estimation

- μ reference measure on Ξ
- $f = dP/d\mu \in \mathbb{S} = L^2(\mu)$
- $\gamma(t; \xi) = \|t\|_{L^2(\mu)}^2 - 2t(\xi)$
 - $\Rightarrow P\gamma(t) = \|t\|_{L^2(\mu)}^2 - 2\langle t, f \rangle_{L^2(\mu)}$
 - $\Rightarrow s^* = f$ and $\ell(s^*, t) = \|t - s^*\|_{L^2(\mu)}^2$
- $(\psi_\lambda)_{\lambda \in m}$ orthonormal basis of S_m
 - $\Rightarrow s_m^* = \sum_{\lambda \in m} (P\psi_\lambda)\psi_\lambda$ and $\hat{s}_m = \sum_{\lambda \in m} (P_n\psi_\lambda)\psi_\lambda$

Least-squares density estimation

- μ reference measure on Ξ
- $f = dP/d\mu \in \mathbb{S} = L^2(\mu)$
- $\gamma(t; \xi) = \|t\|_{L^2(\mu)}^2 - 2t(\xi)$
 $\Rightarrow P\gamma(t) = \|t\|_{L^2(\mu)}^2 - 2\langle t, f \rangle_{L^2(\mu)}$
 $\Rightarrow s^* = f$ and $\ell(s^*, t) = \|t - s^*\|_{L^2(\mu)}^2$
- $(\psi_\lambda)_{\lambda \in m}$ orthonormal basis of S_m
 $\Rightarrow s_m^* = \sum_{\lambda \in m} (P\psi_\lambda)\psi_\lambda$ and $\hat{s}_m = \sum_{\lambda \in m} (P_n\psi_\lambda)\psi_\lambda$

$$\begin{aligned} \text{pen}_{\text{id}}(m) &= (P - P_n)\gamma(\hat{s}_m) = 2(P_n - P)(\hat{s}_m) \\ &= 2\|s_m^* - \hat{s}_m\|_{L^2(\mu)}^2 + 2(P_n - P)(s_m^*) \end{aligned}$$

I.i.d. framework (Lerasle 2009)

$$\text{pen}_{\text{id}}(m) = 2(P_n - P)(\hat{s}_m)$$

$$\text{pen}_W(m) = C_W \mathbb{E}_W \left[2(P_n^W - \overline{W}P_n)(\hat{s}_m^W) \right]$$

I.i.d. framework (Lerasle 2009)

$$\text{pen}_W(m) = C_W \mathbb{E}_W \left[2(P_n^W - \overline{W}P_n)(\hat{s}_m^W) \right]$$

⇒ $\text{pen}_W(m)$ only depends on W through a **multiplicative factor**

I.i.d. framework (Lerasle 2009)

$$\text{pen}_W(m) = C_W \mathbb{E}_W \left[2(P_n^W - \overline{W}P_n)(\hat{s}_m^W) \right]$$

- ⇒ $\text{pen}_W(m)$ only depends on W through a multiplicative factor
- ⇒ $\mathbb{E}[\text{pen}_W(m)] = C_W \text{var}(W_1 - \overline{W}) \mathbb{E}[\text{pen}_{\text{id}}(m)]$
- + **concentration** of $\text{pen}_W(m)$ around its expectation (faster than $\text{pen}_{\text{id}}(m)$)

I.i.d. framework (Lerasle 2009)

$$\text{pen}_W(m) = C_W \mathbb{E}_W \left[2(P_n^W - \overline{W}P_n)(\hat{s}_m^W) \right]$$

- ⇒ $\text{pen}_W(m)$ only depends on W through a multiplicative factor
- ⇒ $\mathbb{E}[\text{pen}_W(m)] = C_W \text{var}(W_1 - \overline{W}) \mathbb{E}[\text{pen}_{\text{id}}(m)]$
 - + concentration of $\text{pen}_W(m)$ around its expectation (faster than $\text{pen}_{\text{id}}(m)$)
- ⇒ **oracle inequality with constant $1 + o(1)$** under well-chosen assumptions on P and \mathcal{M}_n

Dependent case (Lerasle 2010)

- **Mixing** (β or τ)
 - Split the data into several **blocks** \Rightarrow keep one every two blocks
 - **Resample the blocks** (which are almost independent)
- \Rightarrow Oracle inequality (with an oracle only based on part of the original sample)

Outline

- 1 Regressograms in heteroscedastic regression
- 2 The shape of the penalty must be estimated
- 3 Resampling
- 4 Theoretical guarantees for regressograms
- 5 Least-squares density estimation
- 6 Conclusion

Limits of resampling

- Computational complexity

⇒ alternative: non-exchangeable weights (e.g., V -fold)

- Non-asymptotic results: can we have some **without closed-form expressions?**