

Estimator selection

Christophe Giraud

Université Paris-Sud et Paris-Saclay

M2 MSV et MDA

What shall I do with these data ?

Classical steps

- 1 Elucidate the question(s) you want to answer to, and check your data
This requires some
 - ▶ deep discussions with specialists (biologists, physicians, etc),
 - ▶ low level analyses (PCA, LDA, etc) to detect key features, outliers, etc
 - ▶ and ... experience !
- 2 Choose and apply an estimation procedure
- 3 Check your results (residues, possible bias, stability, etc)

Setting

Gaussian regression with unknown variance:

- $Y_i = f_i^* + \varepsilon_i$ with $\varepsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$
- $f^* = (f_1^*, \dots, f_n^*)^T$ and σ^2 are unknown
- we want to estimate f^*

Ex1: sparse linear regression

- $f^* = \mathbf{X}\beta^*$ with β^* "sparse" in some sense and $\mathbf{X} \in \mathbb{R}^{n \times p}$ with possibly $p > n$

A plethora of estimators

Sparse linear regression

- **Coordinate sparsity:** Lasso, Dantzig, Elastic-Net, Exponential-Weighting, Projection on subspaces $\{V_\lambda : \lambda \in \Lambda\}$ given by PCA, Random Forest, PLS, etc.
- **Structured sparsity:** Group-lasso, Fused-Lasso, Bayesian estimators, etc

Important practical issues

Which estimator shall I use?

Lasso? Group-Lasso? Random-Forest? Exponential-Weighting?
Forward-Backward?

With which tuning parameter?

- which penalty level λ for the lasso?
- which beta for expo-weighting?
- etc

Difficulties

- No procedure is universally better than the others
- A sensible choice of the tuning parameters depends on
 - ▶ some unknown characteristics of f (sparsity, smoothness, etc)
 - ▶ the unknown variance σ^2 .

Even if you are a pure Lasso-enthusiast, you miss some key informations in order to apply properly the lasso procedure !

The objective

Formalization

We have a collection of estimation schemes (lasso, group-lasso, etc) and for each scheme we have a grid of different values for the tuning parameters.

At the end, putting all the estimators together we have a collection $\{\hat{f}_\lambda, \lambda \in \Lambda\}$ of estimators.

Ideal objective

Select the "best" estimator among the collection $\{\hat{f}_\lambda, \lambda \in \Lambda\}$.

Cross-Validation

The most popular technique for choosing tuning parameters

Principle

split the data into a **training set** and a **validation set**: the estimators are built on the *training* set and the *validation* set is used for estimating their prediction risk.

Most popular cross-validation scheme

- **Hold-out** : a single split of the data for *training* and *validation*.
- **V-fold CV** : the data is split into V subsamples. Each subsample is successively removed for *validation*, the remaining data being used for *training*.
- **Leave-one-out** : corresponds to n -fold CV.
- **Leave- q -out** : every possible subset of cardinality q of the data is removed for *validation*, the remaining data being used for *training*.

Classical choice of V : between 5 and 10 (remains tractable).

V-fold CV

train	train	train	train	test
train	train	train	test	train
train	train	test	train	train
train	test	train	train	train
test	train	train	train	train

Recursive data splitting for 5-fold Cross-Validation

Pros and Cons

- **Universality:** Cross-Validation can be implemented in most statistical frameworks and for most estimation procedures.
- Usually (but not always!) give good results in practice.
- But **limited theoretical garanties** in large dimensional settings.

Complexity selection (LinSelect)

Principle

To adapt the ideas of model selection to estimator selection.

Pros and Cons

- Strong theoretical guaranties,
- Computationally feasible,
- Good performances in the Gaussian setting,
- But relies on the Gaussian assumption

Reminder on BM model selection

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \|Y - \hat{f}_m\|^2 + \sigma^2 \operatorname{pen}_{BM}(m) \right\} \quad \text{with } \operatorname{pen}_{BM}(m) \approx 2 \log(1/\pi_m)$$

3 difficulties

- 1 We cannot explore a huge collection of models : we restrict to a subcollection $\{S_m, m \in \widehat{\mathcal{M}}\}$
- 2 A good model S_m for \hat{f}_λ must achieve a good balance between the approximation error $\|\hat{f}_\lambda - \operatorname{Proj}_{S_m} \hat{f}_\lambda\|^2$ and the complexity $\log(1/\pi_m)$
- 3 the criterion must not depend on the unknown variance σ^2 : We replace σ^2 in front of the penalty term by the estimator

$$\hat{\sigma}_m^2 = \frac{\|Y - \operatorname{Proj}_{S_m} Y\|^2}{n - \dim(S_m)}. \quad (1)$$

LinSelect procedure

Selection procedure

We select $\hat{f}_{\hat{\lambda}}$, with $\hat{\lambda} = \operatorname{argmin}_{\lambda} \operatorname{crit}(\hat{f}_{\lambda})$ where

$$\operatorname{crit}(\hat{f}_{\lambda}) = \inf_{m \in \hat{\mathcal{M}}} \left[\|Y - \operatorname{Proj}_{S_m} \hat{f}_{\lambda}\|^2 + \frac{1}{2} \|\hat{f}_{\lambda} - \operatorname{Proj}_{S_m} \hat{f}_{\lambda}\|^2 + \operatorname{pen}_{\pi}(m) \hat{\sigma}_m^2 \right]$$

where $\hat{\sigma}_m^2$ is given by (1) and $\operatorname{pen}_{\pi}(m) \approx \operatorname{pen}_{BM}(m)$.

Example : tuning the lasso

Collection of estimators: lasso estimators $\{\hat{f}_\lambda = \mathbf{X}\hat{\beta}_\lambda : \lambda > 0\}$

Collection of models $\{S_m, m \in \mathcal{M}\}$ and probability π : those for coordinate sparse regression

Subcollection: $\hat{\mathcal{M}} = \{\hat{m}(\lambda) : \lambda > 0 \text{ and } 1 \leq |\hat{m}(\lambda)| \leq n/(3 \log p)\}$ with $\hat{m}(\lambda) = \text{supp}(\hat{\beta}_\lambda)$

Theoretical garanty: under some suitable assumptions

$$\|\mathbf{X}(\hat{\beta}_\lambda - \beta^*)\|^2 \leq C \inf_{\beta \neq 0} \left\{ \|\mathbf{X}(\beta^* - \beta)\|^2 + \frac{|\beta|_0 \log(p)}{\kappa^2(\beta)} \sigma^2 \right\}$$

with probability at least $1 - C_1 p^{-C_2}$.

Scaled-Lasso

Automatic tuning of the Lasso

Scale invariance

The estimator $\hat{\beta}(Y, \mathbf{X})$ of β^* is scale-invariant if $\hat{\beta}(sY, \mathbf{X}) = s\hat{\beta}(Y, \mathbf{X})$ for any $s > 0$.

Example: the estimator

$$\hat{\beta}(Y, \mathbf{X}) \in \underset{\beta}{\operatorname{argmin}} \|Y - \mathbf{X}\beta\|^2 + \lambda\Omega(\beta),$$

where Ω is homogeneous with degree 1 is not scale-invariant unless λ is proportional to σ .

In particular the Lasso estimator is not scale-invariant when λ is not proportional to σ .

Rescaling

Idea:

- estimate σ with $\hat{\sigma} = \|Y - \mathbf{X}\beta\|/\sqrt{n}$.
- set $\lambda = \mu\hat{\sigma}$
- divide the criterion by $\hat{\sigma}$ to get a convex problem

Scale-invariant criterion

$$\hat{\beta}(Y, \mathbf{X}) \in \operatorname{argmin}_{\beta} \sqrt{n} \|Y - \mathbf{X}\beta\| + \mu\Omega(\beta).$$

Example: scaled-Lasso

$$\hat{\beta} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \{ \sqrt{n} \|Y - \mathbf{X}\beta\| + \mu|\beta|_1 \}.$$

Pros and Cons

- Universal choice $\mu = 5\sqrt{\log(p)}$
- strong theoretical guaranties (Corollary 5.5)
- computationally feasible
- but poor performances in practice

Numerical experiments (1/2)

Tuning the Lasso

- 165 examples extracted from the literature
- each example e is evaluated on the basis of 400 runs

Comparison to the oracle $\hat{\beta}_{\lambda^*}$

procedure	quantiles			
	0%	50%	75%	90%
Lasso 10-fold CV	1.03	1.11	1.15	1.19
Lasso LinSelect	0.97	1.03	1.06	1.19
Square-Root Lasso	1.32	2.61	3.37	11.2

For each procedure ℓ , quantiles of $\mathcal{R} [\hat{\beta}_{\lambda_\ell}; \beta_0] / \mathcal{R} [\hat{\beta}_{\lambda^*}; \beta_0]$, for $e = 1, \dots, 165$.

Numerical experiments (2/2)

Computation time

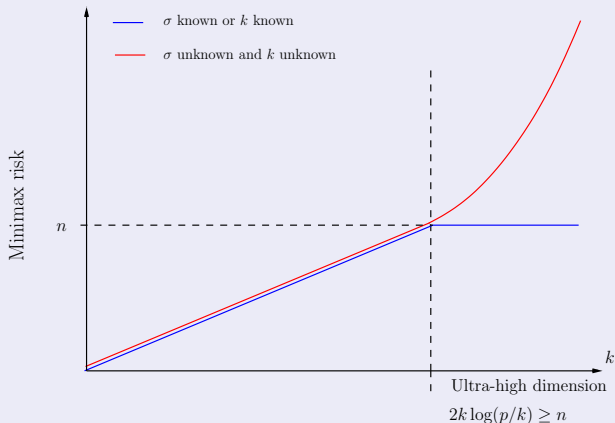
n	p	10-fold CV	LinSelect	Square-Root
100	100	4 s	0.21 s	0.18 s
100	500	4.8 s	0.43 s	0.4 s
500	500	300 s	11 s	6.3 s

Packages:

- `enet` for 10-fold CV and LinSelect
- `lars` for Square-Root Lasso (procedure of Sun & Zhang)

Impact of the unknown variance?

Case of coordinate-sparse linear regression



Minimax prediction risk over k -sparse signal as a function of k