Adaptive Validation, a possible alternative to Cross-Validation

Joseph Salmon

http://josephsalmon.eu Télécom Paristech, Institut Mines-Télécom

Joint work with: Didier Chételat (Cornell University) Johannes Lederer (Cornell University)





Sparsity of signals is all around

Signals can often be represented through a combination of a few elements / atoms :

Fourier decomposition for sounds





Sparsity of signals is all around

Signals can often be represented through a combination of a few elements / atoms :

- Fourier decomposition for sounds
- ▶ Wavelet for images (1990's)





Sparsity of signals is all around

Signals can often be represented through a combination of a few elements / atoms :

- Fourier decomposition for sounds
- ▶ Wavelet for images (1990's)
- Dictionary learning for images (late 2000's)





Sparse linear model

Let $y \in \mathbb{R}^n$ be a signal, *e.g.*, an image

Let $X = [\mathbf{x}_1, \dots, \mathbf{x}_p] \in \mathbb{R}^{n \times p}$ be a collection of (normalized) atoms: corresponds to a **dictionary**

 $\begin{array}{l} X \text{ well suited if one can} \\ \text{approximate the signal } y \approx X\beta \\ \text{with a sparse vector } \beta \in \mathbb{R}^p \end{array}$





The Lasso and its other names

Possible way to get a sparse vector when the dictionary is known:



<u>Rem</u>: Convex optimization problem, can be solved efficiently <u>Rem</u>: Did I mention you have to tune/choose λ ?

Vocabulary:

- Statistics: Lasso Tibshirani (1996)
- Signal processing: Basis Pursuit Chen, Donoho and Saunders (1998)

Dictionary learning: last motivation word

One typically observes T training signals $y_1, \ldots y_T$. start with an initial dictionary X_0 , then alternates over $t \in [T]$:

$$\begin{cases} \hat{\beta}_t^{\lambda} \in \underset{\beta \in \mathbb{R}^p}{\arg\min} \frac{1}{2} \|y_t - X_{t-1}\beta\|_2^2 + \lambda \|\beta\|_1 \quad \text{(coefficients update)}\\ X_t \in \underset{X \in \mathcal{N}}{\arg\min} \frac{1}{2} \|y_{t-1} - X\hat{\beta}_{t-1}^{\lambda}\|_2^2 \quad \text{(dictionary update)} \end{cases}$$

where \mathcal{N} is a set of **normalized dictionary** in $\mathbb{R}^{n \times p}$,

<u>Rem</u>: Applied to signals being small patches of images this reaches state-of-the-art on several image processing tasks *cf.* Mairal *et al.* (2010)

Model considered afterwards

We consider a simple model for theory and simulation validation

$$y = X\beta^* + \varepsilon$$

- ► additive white Gaussian noise $\varepsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$
- ► $X \in \mathbb{R}^{n \times p}$ has normalized columns *e.g.*, $\|\mathbf{x}_j\|_2 = \sqrt{n}$ for all j = 1, ..., p
- the true signal β^* is sparse
- ► the prediction error E ||Xβ* Xβ̂||²/n as a measure of performance for an estimator β̂

<u>Rem</u>: Note that p can be larger than n













Orthonormal dictionary (n = p)

Orthonormal case: $X^{\top}X = \mathrm{Id}_p$ or $\langle \mathbf{x}_i | \mathbf{x}_j \rangle = \delta_{i,j}$ for all i, j

$$\hat{\beta}^{\lambda} = \operatorname*{arg\,min}_{\beta \in \mathbb{R}^p} \left(\frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right)$$

Solution = **Soft-Thresholding**:

$$\hat{\beta}^{\lambda} = \begin{pmatrix} \eta_{\mathrm{ST},\lambda}(\langle \mathbf{x}_{1}, y \rangle) \\ \vdots \\ \eta_{\mathrm{ST},\lambda}(\langle \mathbf{x}_{p}, y \rangle) \end{pmatrix}$$

where

$$\eta_{\mathrm{ST},\lambda}(x) = \mathrm{sign}(x) \cdot (|x| - \lambda)_+$$



<u>Drawback</u>: it shrinks large coefficients toward zero by a factor λ !

Debiasing the Lasso the easy way

$$\hat{\beta}_{\text{Lasso}}^{\lambda} \in \operatorname*{arg\,min}_{\beta \in \mathbb{R}^{p}} \left(\frac{1}{2} \|y - X\beta\|_{2}^{2} + \lambda \|\beta\|_{1} \right)$$

Definition: active set

The **active set** or **support** of a vector β is the indexes of its non-zero coordinates

$$\operatorname{supp}[\beta] = \{ j \in [p] : \beta_j \neq 0 \}$$

<u>Rem</u>: The Lasso is sparse means its active set $\mathrm{supp}[\hat{\beta}_{\mathrm{Lasso}}^{\lambda}]$ is small

Definition: LSLasso

The Least Square Lasso (**LSLasso**): performs a least square fitting for the variables x_j that are activated

$$\hat{\beta}_{\text{LSLasso}}^{\lambda} \in \underset{\text{supp}[\beta] = \text{supp}[\hat{\beta}_{\text{Lasso}}^{\lambda}]}{\text{arg min}} \| Y - X\beta \|_{2}^{2}$$



Lasso

LSLasso



Lasso

LSLasso





Tuning parameters



Prediction error $||X\beta^* - X\hat{\beta}||_2^2/n$ (here r = 20)

Tuning parameters



Prediction error $||X\beta^* - X\hat{\beta}||_2^2/n$ (here r = 20)





$$k = 1$$

- 1. Compute over the training set the estimators for $\lambda_1 > \cdots > \lambda_r$, get $\hat{\beta}^{\lambda_1}, \dots, \hat{\beta}^{\lambda_r}$
- 2. Compute the prediction errors using the validation set and get $\operatorname{Error}_1^k, \ldots, \operatorname{Error}_r^k$



$$k = 2$$

- 1. Compute over the training set the estimators for $\lambda_1 > \cdots > \lambda_r$, get $\hat{\beta}^{\lambda_1}, \dots, \hat{\beta}^{\lambda_r}$
- 2. Compute the prediction errors using the validation set and get $\operatorname{Error}_1^k, \ldots, \operatorname{Error}_r^k$



$$k = 3$$

- 1. Compute over the training set the estimators for $\lambda_1 > \cdots > \lambda_r$, get $\hat{\beta}^{\lambda_1}, \dots, \hat{\beta}^{\lambda_r}$
- 2. Compute the prediction errors using the validation set and get $\operatorname{Error}_1^k, \ldots, \operatorname{Error}_r^k$



$$k = 4$$

- 1. Compute over the training set the estimators for $\lambda_1 > \cdots > \lambda_r$, get $\hat{\beta}^{\lambda_1}, \dots, \hat{\beta}^{\lambda_r}$
- 2. Compute the prediction errors using the validation set and get $\operatorname{Error}_1^k, \ldots, \operatorname{Error}_r^k$



$$k = 5$$

- 1. Compute over the training set the estimators for $\lambda_1 > \cdots > \lambda_r$, get $\hat{\beta}^{\lambda_1}, \dots, \hat{\beta}^{\lambda_r}$
- 2. Compute the prediction errors using the validation set and get $\operatorname{Error}_{1}^{k}, \ldots, \operatorname{Error}_{r}^{k}$



$$k = 6$$

- 1. Compute over the training set the estimators for $\lambda_1 > \cdots > \lambda_r$, get $\hat{\beta}^{\lambda_1}, \dots, \hat{\beta}^{\lambda_r}$
- 2. Compute the prediction errors using the validation set and get $\operatorname{Error}_{1}^{k}, \ldots, \operatorname{Error}_{r}^{k}$



$$k = 7$$

- 1. Compute over the training set the estimators for $\lambda_1 > \cdots > \lambda_r$, get $\hat{\beta}^{\lambda_1}, \dots, \hat{\beta}^{\lambda_r}$
- 2. Compute the prediction errors using the validation set and get $\operatorname{Error}_1^k, \ldots, \operatorname{Error}_r^k$



$$k = 8$$

- 1. Compute over the training set the estimators for $\lambda_1 > \cdots > \lambda_r$, get $\hat{\beta}^{\lambda_1}, \dots, \hat{\beta}^{\lambda_r}$
- 2. Compute the prediction errors using the validation set and get $\operatorname{Error}_{1}^{k}, \ldots, \operatorname{Error}_{r}^{k}$



$$k = 9$$

- 1. Compute over the training set the estimators for $\lambda_1 > \cdots > \lambda_r$, get $\hat{\beta}^{\lambda_1}, \dots, \hat{\beta}^{\lambda_r}$
- 2. Compute the prediction errors using the validation set and get $\operatorname{Error}_{1}^{k}, \ldots, \operatorname{Error}_{r}^{k}$



$$k = 10$$

- 1. Compute over the training set the estimators for $\lambda_1 > \cdots > \lambda_r$, get $\hat{\beta}^{\lambda_1}, \dots, \hat{\beta}^{\lambda_r}$
- 2. Compute the prediction errors using the validation set and get $\operatorname{Error}_1^k, \ldots, \operatorname{Error}_r^k$

Divide (X, y) in K fold along the samples



$$k = 10$$

- 1. Compute over the training set the estimators for $\lambda_1 > \cdots > \lambda_r$, get $\hat{\beta}^{\lambda_1}, \dots, \hat{\beta}^{\lambda_r}$
- 2. Compute the prediction errors using the validation set and get $\operatorname{Error}_1^k, \ldots, \operatorname{Error}_r^k$

Final Step: average the prediction errors get $\widetilde{\text{Error}}_1, \ldots, \widetilde{\text{Error}}_r$ and choose $\hat{i}^{\text{CV}} \in [r]$ achieving the smallest error

Divide (X, y) in K fold along the samples



$$k = 10$$

- 1. Compute over the training set the estimators for $\lambda_1 > \cdots > \lambda_r$, get $\hat{\beta}^{\lambda_1}, \dots, \hat{\beta}^{\lambda_r}$
- 2. Compute the prediction errors using the validation set and get $\operatorname{Error}_1^k, \ldots, \operatorname{Error}_r^k$

Final Step: average the prediction errors get $\widehat{\text{Error}}_1, \ldots, \widehat{\text{Error}}_r$ and choose $\hat{i}^{\text{CV}} \in [r]$ achieving the smallest error **Final Step (bis)**: compute $\hat{\beta}^{\lambda_i}$ over the whole (X, Y) for $i = \hat{i}^{\text{CV}}$

Drawbacks of *K*-fold Cross Validation

Say K = 10

Computational limits

Naive method:

- ► compute 10 times the estimator over datasets of size 90% of the original one
- \blacktriangleright compute 1 time the estimator over 100% of the dataset

Theoretical limits

- basic results suppose estimators are (almost) "independent"
- little is known on the performance of cross-validated Lasso / LSLasso type methods
- cf. Arlot and Celisse (2010)

... but a practical interest



... but a practical interest



Parameters: $\lambda_1 \quad \cdots \quad \lambda_r$









For simplicity, denote:

► $\overline{\beta}^i = \hat{\beta}_{\text{LSLasso}}^{\lambda_i}$ (*i.e.*, the least square estimator over the variables in \hat{S}^i)

• $\overline{\beta}^{i,j}$ the least square estimator over the variable in $\hat{S}^i \cup \hat{S}^j$ <u>Rem</u>: there are at most r^2 such estimators; not all needed for our purpose

Adaptive Validation in Prediction (AV_p)

Assume the active sets are ordered such that: $|\hat{S}^1| \leq \ldots \leq |\hat{S}^r|$

Definition AV_p

Let $a \ge 0$. The AV_p is the estimator $\overline{\beta} := \overline{\beta}^{\hat{i}}$ with active set $\hat{S} = \hat{S}^{\hat{i}}$, where

$$\widehat{i} := \min\left\{i \in [r-1] \left| \frac{\|X\overline{\beta}^i - X\overline{\beta}^{i,j}\|_2^2}{|\widehat{S}^i| + |\widehat{S}^{i,j}|} \le a, \forall j \in [r] : |\widehat{S}^j| \ge |\widehat{S}^i| \right\}$$

when the minimum exists, and $\hat{i}=r$ otherwise.

- ▶ a is a tuning parameter, as was K for Cross-Validation (to be discuss later)
- ▶ $\|X\overline{\beta}^i X\overline{\beta}^{i,j}\|_2^2$ measures the variation in prediction obtained by adding the variable from the set \hat{S}^j to \hat{S}^i

What does it mean?



What does it mean?



What does it mean?



Performance in prediction



Performance in prediction



Performance in prediction



Time efficiency



AV_p: the algorithm

```
Data: Y, X, \hat{S}^1, \ldots, \hat{S}^r, a
Result: \hat{i} \in [r], \overline{\beta} \in \mathbb{R}^p
 Initialize index: i \leftarrow 1
while i \leq r-1 do
        Initialize stopping TestFailure \leftarrow False and compute \overline{\beta}^i
        Initialize comparisons: j \leftarrow \min\{k \in [i] | |\hat{S}^k| > |\hat{S}^i|\}
        while (j \leq r) \& (TestFailure = False) do
               Compute \hat{S}^{i,j} and \overline{\beta}^{i,j}
              if \|X\overline{\beta}^i-X\overline{\beta}^{i,j}\|_2^2\leq a|\hat{S}^i|+a|\hat{S}^{i,j}| then
               i'' i \leftarrow i + 1
              else
                     TestFailure \leftarrow True
               end
        end
        if TestFailure == True then
             i \leftarrow i + 1
        else
               break
        end
end
\hat{i} \leftarrow i \text{ and } \overline{\beta} \leftarrow \overline{\beta}^i
```

When σ^2 (noise level) is known

In theory: a should be chosen proportionally to σ^2 In practice: $a = \sigma^2$ works fine ... when σ^2 is known

When σ^2 (noise level) is known

Need to estimate σ^2 with often p > n

- ▶ ℓ₁-penalized maximum likelihood over (β, σ²) Stadler et al. 2010
- Square root Lasso / Scaled Lasso Antoniadis (2010) , Belloni *et al.* (2011) , Sun and Zhang (2012)

When σ^2 (noise level) is known

In theory: a should be chosen proportionally to σ^2 In practice: $a = \sigma^2$ works fine ... when σ^2 is known

When σ^2 (noise level) is known

Need to estimate σ^2 with often p > n

- ▶ ℓ₁-penalized maximum likelihood over (β, σ²) Stadler et al. 2010
- Square root Lasso / Scaled Lasso Antoniadis (2010) , Belloni *et al.* (2011) , Sun and Zhang (2012)

.asso
$$\begin{cases} \hat{\beta}^{\lambda} \in \operatorname*{arg\,min}_{\beta \in \mathbb{R}^{p}} \left(\frac{1}{2} \|y - X\beta\|_{2}^{2} + \lambda \|\beta\|_{1} \right) \\ \hat{\sigma}^{2} = \|\hat{\beta}^{\lambda} - y\|_{2}^{2}/n \end{cases}$$

When σ^2 (noise level) is known

In theory: a should be chosen proportionally to σ^2 In practice: $a = \sigma^2$ works fine ... when σ^2 is known

When σ^2 (noise level) is known

Need to estimate σ^2 with often p > n

- ▶ ℓ₁-penalized maximum likelihood over (β, σ²) Stadler et al. 2010
- Square root Lasso / Scaled Lasso Antoniadis (2010) , Belloni *et al.* (2011) , Sun and Zhang (2012)

Square root Lasso

$$\begin{cases} \hat{\beta}^{\lambda} \in \operatorname*{arg\,min}_{\beta \in \mathbb{R}^{p}} \left(\boxed{\frac{1}{2} \|y - X\beta\|_{2}} + \lambda \|\beta\|_{1} \right) \\ \hat{\sigma}^{2} = \|\hat{\beta}^{\lambda} - y\|_{2}^{2}/n \end{cases}$$

When σ^2 (noise level) is known

In theory: a should be chosen proportionally to σ^2 In practice: $a = \sigma^2$ works fine ... when σ^2 is known

When σ^2 (noise level) is known

Need to estimate σ^2 with often p > n

- ▶ ℓ₁-penalized maximum likelihood over (β, σ²) Stadler et al. 2010
- Square root Lasso / Scaled Lasso Antoniadis (2010) , Belloni *et al.* (2011) , Sun and Zhang (2012)

Square root Lasso
$$\begin{cases} \hat{\beta}^{\lambda} \in \operatorname*{arg\,min}_{\beta \in \mathbb{R}^{p}} \left(\boxed{\frac{1}{2} \|y - X\beta\|_{2}} + \lambda \|\beta\|_{1} \right) \\ \hat{\sigma}^{2} = \|\hat{\beta}^{\lambda} - y\|_{2}^{2}/n \end{cases}$$

<u>Rem</u>: Reasonable rough estimation of $\hat{\sigma}^2$ for $\lambda = \sqrt{2n \log p}$

AV_p: theory behind

- ► Non-parametric statistics: AV_p idea ≈ the Lepski method Lepski (1990), Lepski, Mammen and Spokoiny (1997), Chichignoud, Lederer and Wainwright (2014)
- Image processing: popularized as the <u>ICI</u> for Intersection of Intervals of Confidence Katkovnik (1999)

Theorem for the AV_p applied to the Lasso: $\overline{\beta} = \overline{\beta}^i$

For a large enough, under some technical assumption over the estimators $\overline{\beta}^{i,j}$, then with high probability it holds that

$$\|X\overline{\beta} - X\beta\|_2^2 \leq \!\! 8a|\widetilde{S}|$$

where \widetilde{S} being the smallest active set among $\hat{S}^1,\ldots,\hat{S}^r$ containing the true support ${\rm supp}(\beta^*)$

Conclusion

Take home message

- New technique for tuning parameters of Lasso but not only: Sqrt-root Lasso, Thresholded Ridge/Tikohnov Regression.
- (Partial) theoretical guarantees
- Computationally more efficient than CV

Future work

- Improve the noise estimation step
- Encompass more method in the framework
- Generalized to other noise model

More info

Website: http://josephsalmon.eu

- Article on ArXiV next week
- Python Code on demand (soon on authors webpage)
- Slides online after the talk



Powered with MooseTeX

Références I

S. Arlot and A. Celisse.

A survey of cross-validation procedures for model selection. *Statistics surveys*, 4:40–79, 2010.

A. Antoniadis.

Comments on: ℓ_1 -penalization for mixture regression models. TEST, 19(2):257–258, 2010.

A. Belloni, V. Chernozhukov, and L. Wang.
 Square-root Lasso: Pivotal recovery of sparse signals via conic programming.

Biometrika, 98(4):791-806, 2011.

- S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. SIAM J. Sci. Comput., 20(1):33–61 (electronic), 1998.
- M. Chichignoud, J. Lederer, and M. Wainwright. Tuning lasso for sup-norm optimality. *ArXiv e-prints*, 2014.

Références II

V. Katkovnik.

A new method for varying adaptive bandwidth selection. *IEEE Trans. Image Process.*, 47(9):2567–2571, 1999.

O. V. Lepski.

On a problem of adaptive estimation in gaussian white noise. *Theory of Probability and its Applications*, 35(3):454–466, 1990.

• O. V. Lepski, E. Mammen, and V. G. Spokoiny.

Optimal spatial adaptation to inhomogeneous smoothness: an approach based on kernel estimates with variable bandwidth selectors.

Ann. Statist., 25(3):929-947, 1997.

- J. Mairal, F. Bach, J. Ponce, and G. Sapiro.
 Online learning for matrix factorization and sparse coding.
 J. Mach. Learn. Res., pages 19–60, 2010.
- N. Städler, P. Bühlmann, and Sara s van de Geer.
 *l*₁-penalization for mixture regression models.
 TEST, 19(2):209–256, 2010.

Références III

► T. Sun and C.-H. Zhang.

Scaled sparse linear regression. *Biometrika*, 99(4):879–898, 2012.

R. Tibshirani.

Regression shrinkage and selection via the lasso.

J. Roy. Statist. Soc. Ser. B, 58(1):267-288, 1996.