

Sparsity in Learning with the LASSO: Least Absolute Shrinkage and Selection Operator

Binh Nguyen – Telecom Paris

M2DS Alternants Research Seminar Course; 07/04/2022

Outline

Introduction

Lasso: Sparsity in Learning

Lasso Estimation: How do we fit the model

Hyper-Parameter Selection with the Lasso

Some Practical Info for this Course

- ▶ Each Thursday from now (April) to end of June
- ▶ First two sessions for each theme is lecture + having fun with practical Python (technically Jupyter) notebook for coding
- ▶ For interaction during classes: it's better for you to think as well, so prepare for some derivations/questions
- ▶ Four presentations as grading, final note is the average
- ▶ Advice: you should start working with the assigned paper as early as possible, because you might have some questions related to the paper and can check with me in the second session for each of the topics

Some Practical Info for this Course

Notice that:

- ▶ Group ordering will change for each topic
- ▶ 7 people - 3 groups (2-2-3) – expect that the group with 3 people will have to present more
- ▶ For now the plan for topic will be: Sparse regression with Lasso, Optimal Transport with Gromov-Wasserstein distance, and Flow-based Model (Normalizing Flows)

Overview

- ▶ Supervised learning with Linear Model
- ▶ The Lasso
- ▶ Solving Lasso with non-smooth optimization: proximal operators

Supervised Learning

▶ $y \in \mathbb{R}, \mathbf{x} \in \mathbb{R}^p, \varepsilon \sim \mathcal{N}(0, \sigma^2), f : \mathbb{R}^p \rightarrow \mathbb{R}$

▶ General relationship

$$y = f(\mathbf{x}) + \varepsilon$$

▶ Goal: learn f from some realizations of (\mathbf{x}, y)

Linear Regression

$$y = f(\mathbf{x}) + \varepsilon$$

Goal: learn f from some realizations of (\mathbf{x}, y)

- ▶ Assumption: f is linear:

$$f(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\beta}$$

- ▶ Question: what should we target now for learning f ?

Linear Regression

$$y = f(\mathbf{x}) + \varepsilon$$

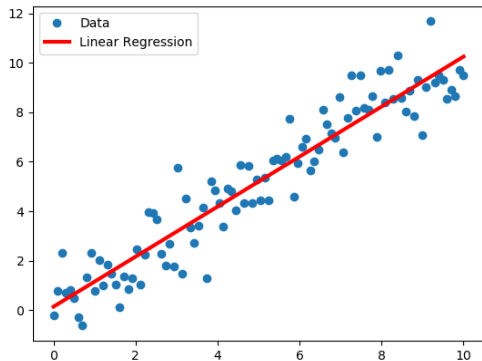
Goal: learn f from some realizations of (\mathbf{x}, y)

- ▶ Assumption: f is linear:

$$f(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\beta}$$

- ▶ Question: what should we target now for learning f ?
- ▶ Answer: Estimate $\boldsymbol{\beta}$ from some realizations of (\mathbf{x}, y)
- ▶ Question: how can we estimate $\boldsymbol{\beta}$ given data $(\mathbf{x}_i, y_i)_{i=1}^n$?

Linear Least Squares



$$\hat{\beta} = \operatorname{argmin}_{\beta} \sum_{i=1}^n \frac{1}{2} (y - \mathbf{x}^{\top} \beta)^2$$

Maximum (Log) Likelihood?

$$y = f(\mathbf{x}) + \varepsilon$$

- ▶ For one sample, define the likelihood:

$$\mathbb{P}(y = y_{obs} \mid \mathbf{x} = \mathbf{x}_{obs}) = \mathbb{P}(\varepsilon = y_{obs} - \mathbf{x}_{obs}^\top \boldsymbol{\beta})$$

Maximum (Log) Likelihood?

$$y = f(\mathbf{x}) + \varepsilon$$

- ▶ For one sample, define the likelihood:

$$\begin{aligned}\mathbb{P}(y = y_{obs} \mid \mathbf{x} = \mathbf{x}_{obs}) &= \mathbb{P}(\varepsilon = y_{obs} - \mathbf{x}_{obs}^\top \boldsymbol{\beta}) \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_{obs} - \mathbf{x}_{obs}^\top \boldsymbol{\beta})^2}{\sigma^2}\right)\end{aligned}$$

Maximum (Log-)Likelihood?

- ▶ For n *i.i.d.* samples $(\mathbf{x}_i, \mathbf{y}_i)_{i=1}^n$, the likelihood:

$$\begin{aligned}\mathcal{L}(\boldsymbol{\beta}) &\stackrel{\text{def.}}{=} \mathbb{P}(\mathbf{y}_1, \dots, \mathbf{y}_n \mid \mathbf{x}_1, \dots, \mathbf{x}_n) \\ &= \mathbb{P}(\mathbf{y}_1 \mid \mathbf{x}_1) \times \dots \times \mathbb{P}(\mathbf{y}_n \mid \mathbf{x}_n)\end{aligned}$$

Maximum Log-likelihood Estimator:

$$\boldsymbol{\beta}_{MLE} = \underset{\boldsymbol{\beta}}{\operatorname{argmax}} \log \mathcal{L}(\boldsymbol{\beta})$$

Maximum (Log-)Likelihood?

- For n *i.i.d.* samples $(\mathbf{x}_i, y_i)_{i=1}^n$, the likelihood:

$$\begin{aligned}\mathcal{L}(\boldsymbol{\beta}) &\stackrel{\text{def.}}{=} \mathbb{P}(y_1, \dots, y_n \mid \mathbf{x}_1, \dots, \mathbf{x}_n) \\ &= \mathbb{P}(y_1 \mid \mathbf{x}_1) \times \dots \times \mathbb{P}(y_n \mid \mathbf{x}_n)\end{aligned}$$

Maximum Log-likelihood Estimator:

$$\boldsymbol{\beta}_{MLE} = \underset{\boldsymbol{\beta}}{\operatorname{argmax}} \log \mathcal{L}(\boldsymbol{\beta}) = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \sum_{i=1}^n \frac{1}{2} (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2$$

Matrix Notations

- ▶ $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times p}, \mathbf{y} \in \mathbb{R}^n$
- ▶ $\sum_{i=1}^n \frac{1}{2} (\mathbf{y} - \mathbf{x}_i^\top \boldsymbol{\beta})^2 = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$
- ▶ ℓ_2 norm of a vector: $\|\mathbf{v}\| = \sqrt{\sum_{i=1}^n v_i^2}$

Closed-form Solution in low-dimension

Objective:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \|y - X\beta\|^2$$

- ▶ Question: low-dimension: $n > p$ – what is the closed-form?

Closed-form Solution in low-dimension

Objective:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \|y - X\beta\|^2$$

- ▶ Question: low-dimension: $n > p$ – what is the closed-form?

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

Maximum Likelihood Estimator/Linear Least Square

Objective:

$$\hat{\beta}_{MLE} = \underset{\beta}{\operatorname{argmin}} \|y - X\beta\|^2$$

▶ Pros:

- ▶ Consistent: $\lim_{n \rightarrow \infty} \hat{\beta}_{MLE} = \beta^*$.
- ▶ Clear statistical framework.

▶ Cons:

- ▶ Can behave badly when f is misspecified (*i.e.* f is not linear).
- ▶ Ill-posed when $n < p$: $X^\top X$ is not invertible.

So... what do we do when $n < p$?

Ockham's razor: only a few coefficients in β are important

- ▶ Or: induce some sparsity on $\hat{\beta}$
- ▶ ℓ_0 (pseudo)-norm: number of non-zero coefficients in β

$$\ell_0(\beta) \stackrel{\text{def.}}{=} \sum_{i=1}^n \mathbb{1}_{\beta_i \neq 0} = |\text{supp}(\beta)|$$

Subset Selection/Matching Pursuit

$$\beta_{subset} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 \quad \text{subject to} \quad \|\beta\|_0 \leq t$$

where $t > 0$ is an integer controlling the sparsity of the solution

Subset Selection/Matching Pursuit

$$\beta_{subset} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 \quad \text{subject to} \quad \|\beta\|_0 \leq t$$

where $t > 0$ is an integer controlling the sparsity of the solution

Problems:

- ▶ Non-convexity
- ▶ Instability: adding a single sample may completely change β and hence its support
- ▶ NP-Hard to solve (*i.e.* very long time to solve)

Least Absolute Shrinkage and Selection Operator

- ▶ Idea: relax non-convexity from the ℓ_0 norm to a convex problem

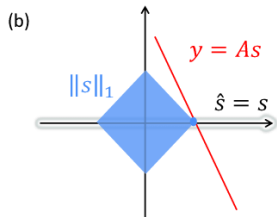
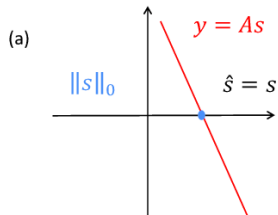
$$\|\beta\|_1 \stackrel{\text{def.}}{=} \sum_{i=1}^n |\beta_i|$$

- ▶ Which makes

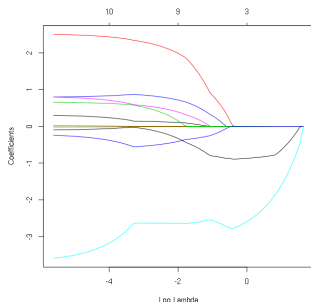
$$\beta_{\text{lasso}} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2} \|y - X\beta\|^2$$

$$\text{subject to } \|\beta\|_1 \leq t$$

where $t > 0$ controls the sparsity of the solution



Lagrangian Formulation



$$\beta_{lasso} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|_1$$

where $\lambda > 0$ controls the sparsity of the solution

- ▶ Promote sparsity: there is a threshold λ_{max} such that $\lambda > \lambda_{max}$ implies $\beta_{lasso} = 0$
- ▶ Question: how to find this λ_{max} ?

Finding λ_{max}

Reminder: the sub-gradient ∂f : a vector $g \in \mathbb{R}^p$ is a subgradient of $f : \mathbb{R}^p \rightarrow \mathbb{R}$ at x if

Finding λ_{max}

Reminder: the sub-gradient ∂f : a vector $g \in \mathbb{R}^p$ is a subgradient of $f : \mathbb{R}^p \rightarrow \mathbb{R}$ at \mathbf{x} if

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{y} - \mathbf{x})$$

► Idea: $\beta = 0$ is the solution to the lasso iff:

$$0 \in \partial_{\beta=0} \left(\frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|_1 \right)$$

Finding λ_{max}

Reminder: the sub-gradient ∂f : a vector $g \in \mathbb{R}^p$ is a subgradient of $f : \mathbb{R}^p \rightarrow \mathbb{R}$ at \mathbf{x} if

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{y} - \mathbf{x})$$

- Idea: $\beta = 0$ is the solution to the lasso iff:

$$0 \in \partial_{\beta=0} \left(\frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|_1 \right)$$

- We know that

$$\partial \left(\frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 \right) = \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\beta)$$

$$\partial_{\beta=0} (\lambda \|\beta\|_1) = [-\lambda, \lambda]^p$$

- So: for every $i \in [p]$: $(\mathbf{X}^\top \mathbf{y})_i \in [-\lambda, \lambda]$

Finding λ_{max}

Reminder: the sub-gradient ∂f : a vector $g \in \mathbb{R}^p$ is a subgradient of $f : \mathbb{R}^p \rightarrow \mathbb{R}$ at \mathbf{x} if

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{y} - \mathbf{x})$$

- Idea: $\beta = 0$ is the solution to the lasso iff:

$$0 \in \partial_{\beta=0} \left(\frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|_1 \right)$$

- We know that

$$\partial \left(\frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 \right) = \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\beta)$$

$$\partial_{\beta=0} (\lambda \|\beta\|_1) = [-\lambda, \lambda]^p$$

- So: for every $i \in [p]$: $(\mathbf{X}^\top \mathbf{y})_i \in [-\lambda, \lambda]$

$$\lambda_{max} = \|\mathbf{X}^\top \mathbf{y}\|_\infty$$

Why sparsity?

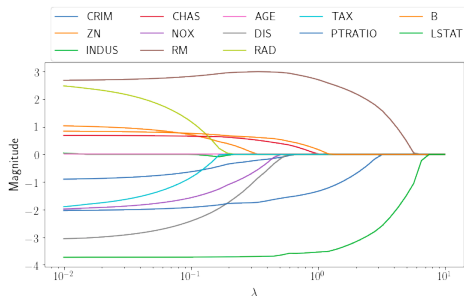


Figure: Lasso Path on Boston dataset

- ▶ Perform *model selection* and *estimation* at the same time: which variables in x are important
- ▶ Sparsity = faster computation and solvers

Why sparsity?

- ▶ Perform *model selection* and *estimation* at the same time: which variables in x are important
- ▶ Sparsity = faster computation and solvers
 - ▶ Example: $p = 1000$, $n = 1000$, computing $X\beta$ takes approx. $n \times p = 10^6$ operations
 - ▶ But: if we know only 10 of the coefficients in β is non-zero, it takes only 10^4 operations – or 100 times faster

The Simplest Lasso

- ▶ $n = 1, p = 1, x = 1$
- ▶ The problem reduces to

$$\min_{\beta} \frac{1}{2}(\mathbf{y} - \beta)^2 + \lambda|\beta|$$

- ▶ Proximity operator:

$$\text{prox}_{|\cdot|}(\mathbf{y}, \lambda) \stackrel{\text{def.}}{=} \underset{\beta}{\text{argmin}} \frac{1}{2}(\mathbf{y} - \beta)^2 + \lambda|\beta|$$

- ▶ Solution: soft-thresholding

$$\text{st}(\mathbf{y}, \lambda) \stackrel{\text{def.}}{=} \text{prox}_{|\cdot|}(\mathbf{y}, \lambda) = \begin{cases} 0, & \text{if } |\mathbf{y}| \leq \lambda \\ \mathbf{y} - \lambda & \text{if } \mathbf{y} > \lambda \\ \mathbf{y} + \lambda & \text{if } \mathbf{y} < -\lambda \end{cases}$$

Soft-Thresholding

Nice property: separability, if $\mathbf{x}, \beta \in \mathbb{R}^p$

$$\text{prox}_{\|\cdot\|_1}(\mathbf{x}, \lambda) = \underset{\beta}{\text{argmin}} \frac{1}{2}(\mathbf{x} - \beta)^2 + \lambda|\beta| = (\text{st}(\mathbf{x}_1, \lambda), \dots, \text{st}(\mathbf{x}_p, \lambda))$$

But How About p -dimension problem?

ISTA: Iterative Soft-Thresholding Algorithm

▶ Recall: $\beta_{lasso} = \operatorname{argmin}_{\beta} \frac{1}{2}(\mathbf{y} - \mathbf{X}\beta)^2 + \lambda|\beta|$

▶ Gradient of the smooth term:

$$\nabla_{\beta} \left(\frac{1}{2}(\mathbf{y} - \mathbf{X}\beta)^2 \right) = \mathbf{X}^{\top}(\mathbf{X}\beta - \mathbf{y})$$

▶ ISTA:

$$\beta^{t+1} = \operatorname{prox}_{\|\cdot\|_1}(\beta^t - \eta\mathbf{X}^{\top}(\mathbf{X}\beta - \mathbf{y}), \lambda\eta)$$

▶ Question: Where does this come from?

But How About p -dimension problem?

ISTA: Iterative Soft-Thresholding Algorithm

▶ Recall: $\beta_{lasso} = \operatorname{argmin}_{\beta} \frac{1}{2}(\mathbf{y} - \mathbf{X}\beta)^2 + \lambda|\beta|$

▶ Gradient of the smooth term:

$$\nabla_{\beta} \left(\frac{1}{2}(\mathbf{y} - \mathbf{X}\beta)^2 \right) = \mathbf{X}^{\top}(\mathbf{X}\beta - \mathbf{y})$$

▶ ISTA:

$$\beta^{t+1} = \operatorname{prox}_{\|\cdot\|_1}(\beta^t - \eta\mathbf{X}^{\top}(\mathbf{X}\beta - \mathbf{y}), \lambda\eta)$$

▶ Question: Where does this come from?

A Reminder: Convex Optimization 101

- ▶ Recall: Gradient descent, minimize $f(x)$ for a smooth f :

$$\text{Iterate } \mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t)$$

- ▶ Its *quadratic surrogate*:

$$\mathbf{x}^{t+1} = \underset{\mathbf{x}}{\operatorname{argmin}} f(\mathbf{x}^t) + \nabla f(\mathbf{x}^t)^\top (\mathbf{x} - \mathbf{x}^t) + \frac{1}{2\eta} \|\mathbf{x} - \mathbf{x}^t\|_2^2$$

Back to ISTA...

- ▶ The quadratic surrogate + ℓ_1 regularization term:

$$\beta^{t+1} = \underset{\beta}{\operatorname{argmin}} f(\beta^t) + \nabla f(\beta^t)^\top (\beta - \beta^t) + \frac{1}{2\eta} \|\beta - \beta^t\|_2^2 + \lambda \|\beta\|_1$$

- ▶ Take $f(\beta) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\beta)^2$

Back to ISTA...

- ▶ The quadratic surrogate + ℓ_1 regularization term:

$$\begin{aligned}\beta^{t+1} &= \underset{\beta}{\operatorname{argmin}} f(\beta^t) + \nabla f(\beta^t)^\top (\beta - \beta^t) + \frac{1}{2\eta} \|\beta - \beta^t\|_2^2 + \lambda \|\beta\|_1 \\ &= \underset{\beta}{\operatorname{argmin}} \nabla f(\beta^t)^\top (\beta - \beta^t) + \frac{1}{2\eta} \|\beta - \beta^t\|_2^2 + \lambda \|\beta\|_1\end{aligned}$$

- ▶ Take $f(\beta) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\beta)^2$

Back to ISTA...

- ▶ The quadratic surrogate + ℓ_1 regularization term:

$$\begin{aligned}\beta^{t+1} &= \operatorname{argmin}_{\beta} f(\beta^t) + \nabla f(\beta^t)^\top (\beta - \beta^t) + \frac{1}{2\eta} \|\beta - \beta^t\|_2^2 + \lambda \|\beta\|_1 \\ &= \operatorname{argmin}_{\beta} \nabla f(\beta^t)^\top (\beta - \beta^t) + \frac{1}{2\eta} \|\beta - \beta^t\|_2^2 + \lambda \|\beta\|_1 \\ &= \operatorname{argmin}_{\beta} \frac{1}{2\eta} \|\beta - (\beta^t - \eta \nabla f(\beta^t))\|_2^2 + \lambda \eta \|\beta\|_1\end{aligned}$$

- ▶ Take $f(\beta) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\beta)^2$

Back to ISTA...

- ▶ The quadratic surrogate + ℓ_1 regularization term:

$$\begin{aligned}\beta^{t+1} &= \operatorname{argmin}_{\beta} f(\beta^t) + \nabla f(\beta^t)^\top (\beta - \beta^t) + \frac{1}{2\eta} \|\beta - \beta^t\|_2^2 + \lambda \|\beta\|_1 \\ &= \operatorname{argmin}_{\beta} \nabla f(\beta^t)^\top (\beta - \beta^t) + \frac{1}{2\eta} \|\beta - \beta^t\|_2^2 + \lambda \|\beta\|_1 \\ &= \operatorname{argmin}_{\beta} \frac{1}{2\eta} \|\beta - (\beta^t - \eta \nabla f(\beta^t))\|_2^2 + \lambda \eta \|\beta\|_1 \\ &\stackrel{\text{def.}}{=} \operatorname{prox}_{\|\cdot\|_1}(\beta^t - \mathbf{X}^\top (\mathbf{X}\beta - \mathbf{y}), \lambda\eta)\end{aligned}$$

- ▶ Take $f(\beta) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\beta)^2$

Back to ISTA...

$$\beta^{t+1} = \text{prox}_{\|\cdot\|_1}(\beta^t - \eta X^\top (X\beta - y), \lambda\eta)$$

- ▶ Due to the prox operator many coefficients are 0
- ▶ One iteration takes $\mathcal{O}(\min(n, p) \times p)$
→ problematic for large p

Coordinate Descent

Idea: freeze others, update only one coefficient β_j at each iteration,
i.e.

$$\beta_j^{t+1} = \text{prox}_{\|\cdot\|_1}(\beta_j^t - \eta \mathbf{X}_{*,j}^\top (\mathbf{X}\boldsymbol{\beta} - \mathbf{y}), \lambda\eta)$$

$$\text{and } \beta_k^{t+1} = \beta_k^t \text{ for all } k \neq j$$

- ▶ Iteration cost $\mathcal{O}(n)$ if we know the *residual* $\mathbf{r} = \mathbf{X}\boldsymbol{\beta}^t - \mathbf{y}$
- ▶ Update residual: Once we update the coordinate j :

$$\mathbf{r} \leftarrow \mathbf{r} + (\beta_j^{t+1} - \beta_j^t) \mathbf{X}_{*,j}$$

A Bit More Advanced: FISTA - Fast(er) ISTA

Initialize $\beta^1 = z^1 = 0, \gamma^1 = 1$, then for each update

$$\blacktriangleright \beta_j^{t+1} = \text{prox}_{\|\cdot\|_1}(z^t - \eta X^\top (Xz^t - y), \lambda\eta)$$

$$\blacktriangleright \gamma^{t+1} = \frac{1 + \sqrt{1 + 4(\gamma^t)^2}}{2}$$

$$\blacktriangleright z^{t+1} = \beta^{t+1} + \frac{\gamma^t - 1}{\gamma^{t+1}}(\beta^{t+1} - \beta^t)$$

Some form of acceleration, which make convergence rate of fista $\mathcal{O}(1/T^2)$ instead of $\mathcal{O}(1/T)$ of ISTA

Parameter Selection for the Lasso

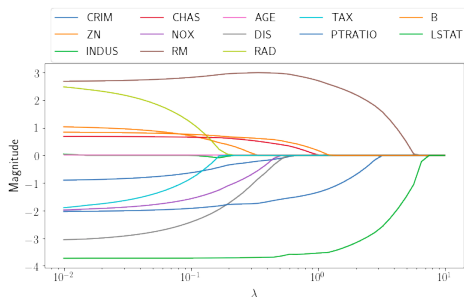


Figure: Lasso Path on Boston dataset

$$\beta_{lasso} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|_1$$

- ▶ We have talked about the role of λ as to control the sparsity level – but which sparsity level is optimal?

Selection of sparsity parameter λ

Caveat: to simplify things, definitions are taken from sklearn page – for a specific theoretical formula for the Lasso, check Zou et al. (2007)

- ▶ Equivalent to model selection: we can do that via a criterion
- ▶ Akaike's Information Criterion (AIC):

$$\text{AIC}(\mathcal{L}(\hat{\beta})) \stackrel{\text{def.}}{=} -2 \log(\mathcal{L}(\hat{\beta})) + 2d$$

where d is the degrees of freedom (number of parameters of models)

Zou, Hui, Trevor Hastie, and Robert Tibshirani. "On the degrees of freedom of the lasso." *The Annals of Statistics* 35.5 (2007): 2173-2192.

Selection of sparsity parameter λ

Caveat: to simplify things, definitions are taken from sklearn page – for a specific theoretical formula for the Lasso, check Zou et al. (2007)

- ▶ Equivalent to model selection: we can do that via a criterion
- ▶ Akaike's Information Criterion (AIC):

$$\text{AIC}(\mathcal{L}(\hat{\beta})) \stackrel{\text{def.}}{=} -2 \log(\mathcal{L}(\hat{\beta})) + 2d$$

where d is the degrees of freedom (number of parameters of models)

- ▶ Bayesian Information Criterion (BIC):

$$\text{BIC}(\mathcal{L}(\hat{\beta})) \stackrel{\text{def.}}{=} -2 \log(\mathcal{L}(\hat{\beta})) + \log(n)2d$$

- ▶ Question: What can we say about these two criteria as a function of the likelihood?

Zou, Hui, Trevor Hastie, and Robert Tibshirani. "On the degrees of freedom of the lasso." *The Annals of Statistics* 35.5 (2007): 2173-2192.

Selection of sparsity parameter λ

- ▶ Question: What's the formula for AIC and BIC in our settings?
(*i.e.* linear regression with with least squares loss and Gaussian noise $\varepsilon \sim \mathcal{N}(0, \sigma^2)$?)

Selection of sparsity parameter λ

- ▶ Question: What's the formula for AIC and BIC in our settings?
(*i.e.* linear regression with with least squares loss and Gaussian noise $\varepsilon \sim \mathcal{N}(0, \sigma^2)$?)

$$\text{AIC}(\hat{\beta}_\lambda) = n \log(2\pi\sigma^2) + \frac{1}{\sigma^2} \|\mathbf{y} - \mathbf{X}\hat{\beta}_\lambda\|^2 + 2d$$

$$\text{BIC}(\hat{\beta}_\lambda) = n \log(2\pi\sigma^2) + \frac{1}{\sigma^2} \|\mathbf{y} - \mathbf{X}\hat{\beta}_\lambda\|^2 + \log(n)d$$

Selection of sparsity parameter λ

- ▶ Question: What's the formula for AIC and BIC in our settings?
(*i.e.* linear regression with with least squares loss and Gaussian noise $\varepsilon \sim \mathcal{N}(0, \sigma^2)$?)

$$\text{AIC}(\hat{\beta}_\lambda) = n \log(2\pi\sigma^2) + \frac{1}{\sigma^2} \|\mathbf{y} - \mathbf{X}\hat{\beta}_\lambda\|^2 + 2d$$

$$\text{BIC}(\hat{\beta}_\lambda) = n \log(2\pi\sigma^2) + \frac{1}{\sigma^2} \|\mathbf{y} - \mathbf{X}\hat{\beta}_\lambda\|^2 + \log(n)d$$

- ▶ The first term on the RHS is a constant, so:

$$\lambda_{AIC} \stackrel{\text{def.}}{=} \underset{\lambda}{\operatorname{argmin}} \frac{1}{\sigma^2} \|\mathbf{y} - \mathbf{X}\hat{\beta}_\lambda\|^2 + 2d$$

$$\lambda_{BIC} \stackrel{\text{def.}}{=} \underset{\lambda}{\operatorname{argmin}} \frac{1}{\sigma^2} \|\mathbf{y} - \mathbf{X}\hat{\beta}_\lambda\|^2 + \log(n)d$$

- ▶ Note: we also have to estimate σ^2 from the residual

Selection of sparsity parameter λ – with AIC/BIC

Pros:

- ▶ Relatively fast as the regularization path only computed once
- ▶ Guarantee to be optimal theoretically

Selection of sparsity parameter λ – with AIC/BIC

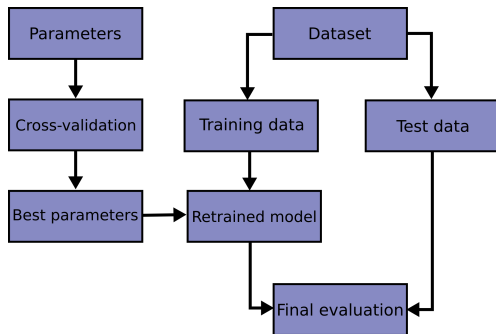
Pros:

- ▶ Relatively fast as the regularization path only computed once
- ▶ Guarantee to be optimal theoretically

Cons:

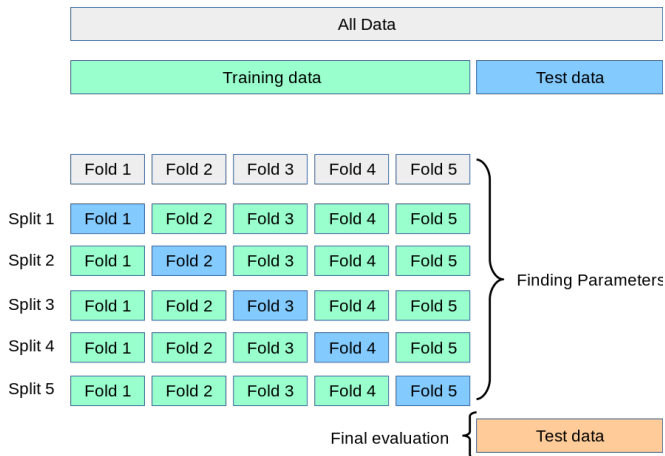
- ▶ Theoretical guarantees only with large n (asymptotic)
- ▶ Estimation of the term d (degrees of freedom) is tricky when $n \ll p$

In Practice: K-Folds Cross-Validation



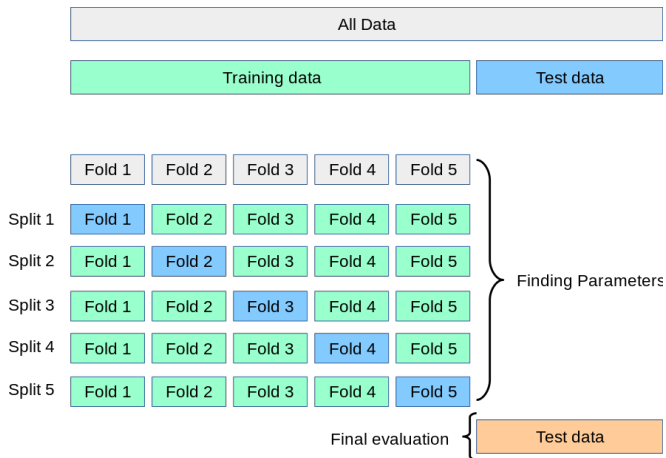
Cross-validation: based on sklearn page:
scikit-learn.org/stable/modules/cross_validation.html

In Practice: K-Folds Cross-Validation



- ▶ Performing model fitting on each (K-1) folds of the training data, evaluate on the remaining fold
- ▶ Which performance measure? (accuracy, MSE – mean squared errors, etc.)

In Practice: K-Folds Cross-Validation



- ▶ Easy with sklearn: `sklearn.model_selection` and even better: `sklearn.linear_model.LassoCV`

Next session

- ▶ Brief intro on a bit more advanced: hyper-parameter selection for Lasso with bi-level optimization
- ▶ Practical session with Jupyter notebook: playing with `sklearn.linear_model.Lasso`, etc. , and some optimization scheme