

# Sparsity in Learning with the LASSO 2

Binh Nguyen – Telecom Paris

M2DS Alternants Research Seminar Course; 14/04/2022

# Outline

Reminder

Variants of Lasso

Hyperparameter Optimization

# Outline

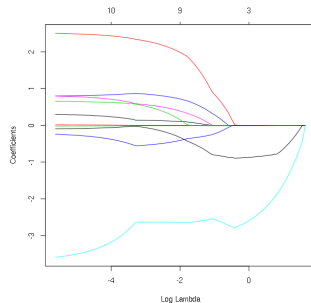
Reminder

Variants of Lasso

Hyperparameter Optimization

## Previously...

### Lasso: Least Absolute Shrinkage and Selection Operator



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

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

→ Promote sparsity: there is a threshold  $\lambda_{max}$  such that  $\lambda > \lambda_{max}$  implies  $\beta_{lasso} = 0$

# Outline

Reminder

Variants of Lasso

Hyperparameter Optimization

## Adaptive (weighted) Lasso

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

- ▶ Only one  $\lambda$  that dictates sparsity degree of all  $\beta_j$
- ▶ What if we want a scheme that is adaptive: coefficients with large magnitude (absolute value) receive smaller sparse penalty?

---

Zou H. (2006), 'The adaptive lasso and its oracle properties', Journal of the American Statistical Association 101(476), 1418–1429.

## Adaptive (weighted) Lasso

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

- ▶ Only one  $\lambda$  that dictates sparsity degree of all  $\beta_j$
- ▶ What if we want a scheme that is adaptive: coefficients with large magnitude (absolute value) receive smaller sparse penalty?  
→ Lasso with adaptive weights on  $\ell_1$ -regularization

$$\beta_{lasso} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \sum_{j=1}^p w_j |\beta_j|$$

where  $w_j \in [0, 1)$ .

---

Zou H. (2006), 'The adaptive lasso and its oracle properties', Journal of the American Statistical Association 101(476), 1418–1429.

## Adaptive (weighted) Lasso

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

- ▶ Only one  $\lambda$  that dictates sparsity degree of all  $\beta_j$
- ▶ What if we want a scheme that is adaptive: coefficients with large magnitude (absolute value) receive smaller sparse penalty?  
→ Lasso with adaptive weights on  $\ell_1$ -regularization

$$\beta_{lasso} \stackrel{\text{def.}}{=} \operatorname{argmin}_{\beta} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \sum_{j=1}^p w_j |\beta_j|$$

where  $w_j \in [0, 1)$ .

→ Optimization problem is still convex in  $\beta$

---

Zou H. (2006), 'The adaptive lasso and its oracle properties', Journal of the American Statistical Association 101(476), 1418–1429.



## Adaptive (weighted) Lasso

$$\beta_{lasso} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \sum_{j=1}^p w_j |\beta_j|$$

► Typically  $w_j$  are initialized as

$$w_j = \begin{cases} \frac{1}{|\beta_j^{\text{init}}|} & \text{if } \beta_j^{\text{init}} \neq 0 \\ 0 & \text{if } \beta_j^{\text{init}} = 0 \end{cases}$$

## Adaptive (weighted) Lasso

$$\beta_{lasso} \stackrel{\text{def.}}{=} \operatorname{argmin}_{\beta} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \sum_{j=1}^p w_j |\beta_j|$$

- ▶ Typically  $w_j$  are initialized as

$$w_j = \begin{cases} \frac{1}{|\beta_j^{\text{init}}|} & \text{if } \beta_j^{\text{init}} \neq 0 \\ 0 & \text{if } \beta_j^{\text{init}} = 0 \end{cases}$$

- ▶ But what is this  $\beta^{\text{init}}$ ?

## Adaptive (weighted) Lasso

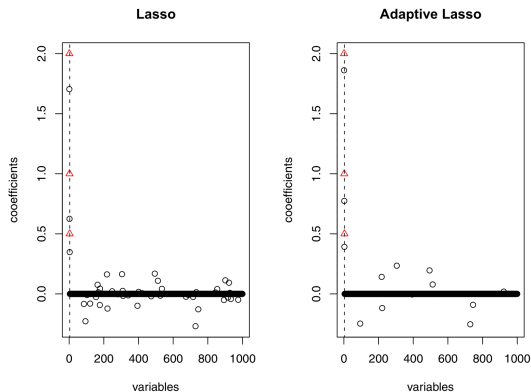
$$\beta_{lasso} \stackrel{\text{def.}}{=} \operatorname{argmin}_{\beta} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \sum_{j=1}^p w_j |\beta_j|$$

- ▶ Typically  $w_j$  are initialized as

$$w_j = \begin{cases} \frac{1}{|\beta_j^{\text{init}}|} & \text{if } \beta_j^{\text{init}} \neq 0 \\ 0 & \text{if } \beta_j^{\text{init}} = 0 \end{cases}$$

- ▶ But what is this  $\beta^{\text{init}}$ ?
- ▶ Just put a standard lasso for finding  $\beta^{\text{init}}$  (called screening operation)

# Adaptive (weighted) Lasso



**Fig. 2.4** Estimated regression coefficients in the linear model with  $p = 1000$  and  $n = 50$ . Left: Lasso. Right: Adaptive Lasso with Lasso as initial estimator. The 3 true regression coefficients are indicated with triangles. Both methods used with tuning parameters selected from 10-fold cross-validation.

---

Bühlmann, P., & Geer, S. A. van de. (2011). *Statistics for high-dimensional data: Methods, theory and applications*. Springer.

## Adaptive (weighted) Lasso

$$\beta_{lasso} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \sum_{j=1}^p |w_j \beta_j|$$

- ▶ Optimization problem is still convex in  $\beta$ , but how to solve now that there is multiple value of  $\lambda$  possible?

## Adaptive (weighted) Lasso

$$\beta_{lasso} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \sum_{j=1}^p |w_j \beta_j|$$

- ▶ Optimization problem is still convex in  $\beta$ , but how to solve now that there is multiple value of  $\lambda$  possible?
- ▶ Question: can we reformulate the adaptive lasso back to the original lasso?

## Adaptive (weighted) Lasso

$$\beta_{lasso} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \sum_{j=1}^p |w_j \beta_j|$$

- ▶ Optimization problem is still convex in  $\beta$ , but how to solve now that there is multiple value of  $\lambda$  possible?
- ▶ Question: can we reformulate the adaptive lasso back to the original lasso?
- ▶ Hint: use the change of variable  $\tilde{\beta}$  as some form of  $w$  and  $\beta$
- ▶ To the whiteboard...

## Adaptive (weighted) Lasso

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

and this means

$$\tilde{\beta}_{lasso} = \mathbf{W}^{-1}\beta_{lasso}$$

*i.e.* the solution of the adaptive Lasso is just a rescaling of the solution of original Lasso

→ enjoys theoretical guarantee (consistency, asymptotic normality) from the Lasso without additional assumptions



# Adaptive (weighted) Lasso

In `sklearn.linear_model.Lasso`

```
fit(X, y, sample_weight=None, check_input=True)
```

[\[source\]](#)

Fit model with coordinate descent.

**Parameters:**

- X** : *{ndarray, sparse matrix} of (n\_samples, n\_features)*  
Data.
- y** : *{ndarray, sparse matrix} of shape (n\_samples,) or (n\_samples, n\_targets)*  
Target. Will be cast to X's dtype if necessary.
- sample\_weight** : *float or array-like of shape (n\_samples,)*, *default=None*  
Sample weights. Internally, the `sample_weight` vector will be rescaled to sum to `n_samples`.  
*New in version 0.23.*
- check\_input** : *bool, default=True*  
Allow to bypass several input checking. Don't use this parameter unless you know what you do.

---

[https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.Lasso.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.Lasso.html)

## A cousin of Lasso: Elastic-Net

- ▶ A problem with Lasso: when there are high-correlations between variables, *e.g.*  $X_{*,i}$  and  $X_{*,j}$  empirically Lasso select one but not both...
- ▶ At most  $n$  variables will be selected by the lasso, so problematic when  $n \ll p$
- ▶ A solution: adding  $\ell_2$  norm to the lasso optimization problem: elastic net

---

Zou, Hui; Hastie, Trevor (2005). "Regularization and Variable Selection via the Elastic Net". *Journal of the Royal Statistical Society, Series B.* 67 (2): 301–320.

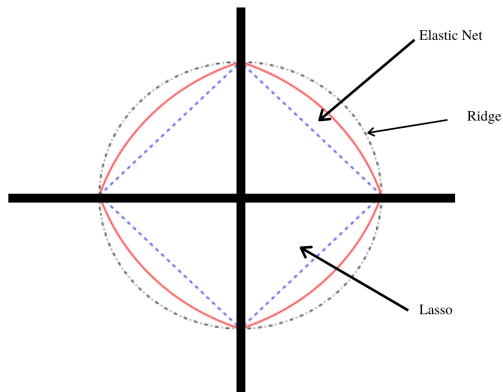
## A cousin of Lasso: Elastic-Net

- ▶ A problem with Lasso: when there are high-correlations between variables, *e.g.*  $X_{*,i}$  and  $X_{*,j}$  empirically Lasso select one but not both...
- ▶ At most  $n$  variables will be selected by the lasso, so problematic when  $n \ll p$
- ▶ A solution: adding  $\ell_2$  norm to the lasso optimization problem: elastic net

---

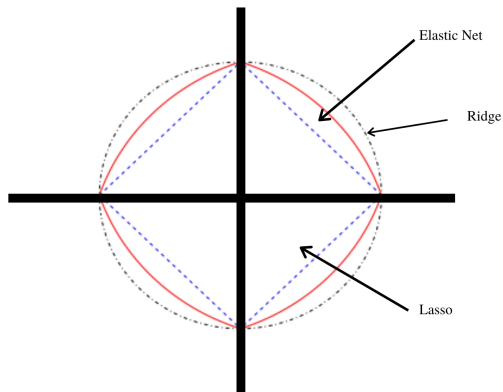
Zou, Hui; Hastie, Trevor (2005). "Regularization and Variable Selection via the Elastic Net". *Journal of the Royal Statistical Society, Series B.* 67 (2): 301–320.

# Elastic-Net



$$\beta_{enet} \stackrel{\text{def.}}{=} \frac{1}{2n} \underset{\beta}{\operatorname{argmin}} \|y - X\beta\|^2 + \lambda_1 \|\beta\|_1 + \frac{\lambda_2}{2} \|\beta\|_2^2$$

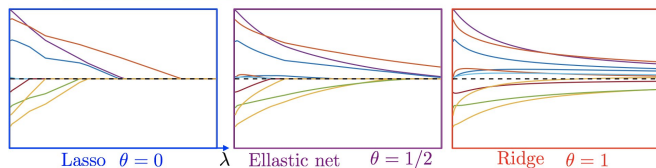
# Elastic-Net



$$\beta_{enet} \stackrel{\text{def.}}{=} \frac{1}{2n} \underset{\beta}{\operatorname{argmin}} \|y - X\beta\|^2 + \lambda_1 \|\beta\|_1 + \frac{\lambda_2}{2} \|\beta\|_2^2$$

but now we have two hyper-parameters  $\lambda_1$  and  $\lambda_2$ ?

# Elastic-Net



we can just set  $\theta = \frac{\lambda_2}{\lambda_1 + \lambda_2} \in [0, 1]$ , then the equivalent problem is

$$\beta_{enet} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|^2 + (1 - \theta) \|\beta\|_1 + \frac{\theta}{2} \|\beta\|_2^2$$

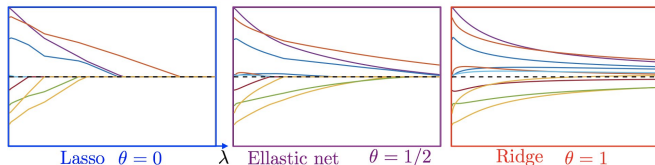
→ enet-path interpolates between Lasso and Ridge regression path

---

Image from Gabriel Peyré's twitter:

<https://twitter.com/gabrielpeyre/status/1318054267685621761>

# Elastic-Net



- ▶ Elastic-net solutions: interpolates between Lasso and Ridge regression solutions
- ▶ Question: this gives hint on finding the solution of Enet? (remember how we find solution for Lasso and for Ridge?)

## Elastic-Net, in Orthogonal Design settings

$$\beta_{enet} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \|y - X\beta\|_2^2 + \lambda_1 \|\beta\|_1 + \frac{\lambda_2}{2} \|\beta\|_2^2$$

in the case  $\frac{1}{n} X^\top X = \text{Id}$ , then  $\hat{\beta}^{LS} = 1/n (X^\top X)^{-1} X^\top y = X^\top y/n$

so for the first term:

$$\underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \|y - X\beta\|_2^2$$



## Elastic-Net, in Orthogonal Design settings

$$\boldsymbol{\beta}_{enet} \stackrel{\text{def.}}{=} \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda_1 \|\boldsymbol{\beta}\|_1 + \frac{\lambda_2}{2} \|\boldsymbol{\beta}\|_2^2$$

in the case  $\frac{1}{n} \mathbf{X}^\top \mathbf{X} = \text{Id}$ , then  $\hat{\boldsymbol{\beta}}^{LS} = 1/n (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top \mathbf{y} / n$

so for the first term:

$$\begin{aligned} & \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 \\ &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \frac{1}{2n} \{ \mathbf{y}^\top \mathbf{y} + \boldsymbol{\beta}^\top \boldsymbol{\beta} - 2\mathbf{y}^\top \mathbf{X}\boldsymbol{\beta} \} \end{aligned}$$

## Elastic-Net, in Orthogonal Design settings

$$\boldsymbol{\beta}_{enet} \stackrel{\text{def.}}{=} \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda_1 \|\boldsymbol{\beta}\|_1 + \frac{\lambda_2}{2} \|\boldsymbol{\beta}\|_2^2$$

in the case  $\frac{1}{n} \mathbf{X}^\top \mathbf{X} = \text{Id}$ , then  $\hat{\boldsymbol{\beta}}^{LS} = 1/n(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top \mathbf{y}/n$

so for the first term:

$$\underset{\boldsymbol{\beta}}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$$

$$= \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \frac{1}{2n} \{ \mathbf{y}^\top \mathbf{y} + \boldsymbol{\beta}^\top \boldsymbol{\beta} - 2\mathbf{y}^\top \mathbf{X}\boldsymbol{\beta} \}$$

$$= \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \frac{1}{2n} \{ \mathbf{y}^\top \mathbf{y} + n\boldsymbol{\beta}^\top \boldsymbol{\beta} + 2n\boldsymbol{\beta}^\top \hat{\boldsymbol{\beta}}^{LS} + n(\hat{\boldsymbol{\beta}}^{LS})^\top \hat{\boldsymbol{\beta}}^{LS} - n(\hat{\boldsymbol{\beta}}^{LS})^\top \boldsymbol{\beta} \}$$

## Elastic-Net, in Orthogonal Design settings

$$\beta_{enet} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda_1 \|\beta\|_1 + \frac{\lambda_2}{2} \|\beta\|_2^2$$

in the case  $\frac{1}{n} \mathbf{X}^\top \mathbf{X} = \text{Id}$ , then  $\hat{\beta}^{LS} = 1/n(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top \mathbf{y}/n$

so for the first term:

$$\begin{aligned} & \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \\ &= \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \{ \mathbf{y}^\top \mathbf{y} + \beta^\top \beta - 2\mathbf{y}^\top \mathbf{X}\beta \} \\ &= \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \left\{ \mathbf{y}^\top \mathbf{y} + n\beta^\top \beta + 2n\beta^\top \hat{\beta}^{LS} + n(\hat{\beta}^{LS})^\top \hat{\beta}^{LS} - n(\hat{\beta}^{LS})^\top \beta \right\} \\ &= \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \left\{ (\hat{\beta}^{LS} - \beta)^\top (n\hat{\beta}^{LS} - \beta) + \mathbf{y}^\top (\text{Id} - \mathbf{X}^\top \mathbf{X}) \mathbf{y} \right\} \end{aligned}$$

## Elastic-Net, in Orthogonal Design settings

$$\beta_{enet} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda_1 \|\beta\|_1 + \frac{\lambda_2}{2} \|\beta\|_2^2$$

in the case  $\frac{1}{n} \mathbf{X}^\top \mathbf{X} = \text{Id}$ , then  $\hat{\beta}^{LS} = 1/n(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top \mathbf{y}/n$

so for the first term:

$$\begin{aligned} & \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \\ &= \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \{ \mathbf{y}^\top \mathbf{y} + \beta^\top \beta - 2\mathbf{y}^\top \mathbf{X}\beta \} \\ &= \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \{ \mathbf{y}^\top \mathbf{y} + n\beta^\top \beta + 2n\beta^\top \hat{\beta}^{LS} + n(\hat{\beta}^{LS})^\top \hat{\beta}^{LS} - n(\hat{\beta}^{LS})^\top \beta \} \\ &= \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \{ (\hat{\beta}^{LS} - \beta)^\top (n\hat{\beta}^{LS} - \beta) + \mathbf{y}^\top (\text{Id} - \mathbf{X}^\top \mathbf{X}) \mathbf{y} \} \\ &= \underset{\beta}{\operatorname{argmin}} \frac{1}{2} \|\hat{\beta}^{LS} - \beta\|_2^2 \end{aligned}$$

## Elastic-Net, in Orthogonal Design settings

$$\beta_{enet} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda_1 \|\beta\|_1 + \frac{\lambda_2}{2} \|\beta\|_2^2$$

in the case  $\frac{1}{n} \mathbf{X}^\top \mathbf{X} = \text{Id}$

## Elastic-Net, in Orthogonal Design settings

$$\boldsymbol{\beta}_{enet} \stackrel{\text{def.}}{=} \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda_1 \|\boldsymbol{\beta}\|_1 + \frac{\lambda_2}{2} \|\boldsymbol{\beta}\|_2^2$$

in the case  $\frac{1}{n} \mathbf{X}^\top \mathbf{X} = \text{Id}$

► This means

$$\boldsymbol{\beta}_{enet} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \frac{1}{2} \sum_{j=1}^p (\hat{\beta}_j^{LS} - \beta_j)^2 + \lambda_1 \sum_{j=1}^p |\beta_j| + \frac{\lambda_2}{2} \sum_{j=1}^p \beta_j^2$$

## Elastic-Net, in Orthogonal Design settings

$$\beta_{enet} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \|y - X\beta\|_2^2 + \lambda_1 \|\beta\|_1 + \frac{\lambda_2}{2} \|\beta\|_2^2$$

in the case  $\frac{1}{n} X^\top X = \text{Id}$

- ▶ This means

$$\beta_{enet} = \underset{\beta}{\operatorname{argmin}} \frac{1}{2} \sum_{j=1}^p (\hat{\beta}_j^{LS} - \beta_j)^2 + \lambda_1 \sum_{j=1}^p |\beta_j| + \frac{\lambda_2}{2} \sum_{j=1}^p \beta_j^2$$

- ▶ The problem is separable: for each  $j$

$$\beta_j^{enet} = \underset{\beta}{\operatorname{argmin}} \frac{1}{2} (\hat{\beta}_j^{LS} - \beta_j)^2 + \lambda_1 |\beta_j| + \frac{\lambda_2}{2} \beta_j^2$$

## Elastic-Net, in Orthogonal Design settings

$$\beta_{enet} \stackrel{\text{def.}}{=} \underset{\beta}{\operatorname{argmin}} \frac{1}{2n} \|y - X\beta\|_2^2 + \lambda_1 \|\beta\|_1 + \frac{\lambda_2}{2} \|\beta\|_2^2$$

in the case  $\frac{1}{n} X^\top X = \text{Id}$

- ▶ This means

$$\beta_{enet} = \underset{\beta}{\operatorname{argmin}} \frac{1}{2} \sum_{j=1}^p (\hat{\beta}_j^{LS} - \beta_j)^2 + \lambda_1 \sum_{j=1}^p |\beta_j| + \frac{\lambda_2}{2} \sum_{j=1}^p \beta_j^2$$

- ▶ The problem is separable: for each  $j$

$$\begin{aligned} \beta_j^{enet} &= \underset{\beta}{\operatorname{argmin}} \frac{1}{2} (\hat{\beta}_j^{LS} - \beta_j)^2 + \lambda_1 |\beta_j| + \frac{\lambda_2}{2} \beta_j^2 \\ &= \underset{\beta}{\operatorname{argmin}} \frac{1}{2} \left( \beta_j - \frac{\hat{\beta}_j^{LS}}{1 + \lambda_2} \right)^2 + \frac{\lambda_1}{1 + \lambda_2} |\beta_j| \\ &\stackrel{\text{def.}}{=} \operatorname{prox}_{\|\cdot\|_1} \left( \beta_j - \frac{\hat{\beta}_j^{LS}}{1 + \lambda_2}, \frac{\lambda_1}{1 + \lambda_2} \right) \end{aligned}$$



# Elastic-Net

This means: in general settings, we can find solution of Enet with iterative optimization algorithm (from last session):

- ▶ ISTA, FISTA
- ▶ Coordinate descent (implemented in sklearn)

## Other Variants

- ▶ Group lasso
- ▶ Lasso for data matrix with missing elements
- ▶ Debiased Lasso

## Other Variants

- ▶ Group lasso
- ▶ Lasso for data matrix with missing elements
- ▶ Debiased Lasso

...which we will wait for presentations next week :-)

Reminder

Variants of Lasso

Hyperparameter Optimization

## Previously...

Lasso: Least Absolute Shrinkage and Selection Operator

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

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

- ▶ Choose  $\lambda$  based  $\lambda_{max} = \|\mathbf{X}^T \mathbf{y}\|_\infty$
- ▶ Reminder: when  $\lambda > \lambda_{max}$  all  $\beta_j$  will shrink to zero
- ▶ But  $\lambda$  to select? – cross-validation/Information Criterion

## Hyperparameter selection, the popular way

- ▶ Cross validation
- ▶ Criterion (AIC/BIC) that control model complexity

## Hyperparameter selection, the popular way

- ▶ Cross validation
- ▶ Criterion (AIC/BIC) that control model complexity
- ▶ Formalization: for Lasso

$$\hat{\beta}^{(\lambda)} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|\mathbf{y}^{\text{train}} - \mathbf{X}^{\text{train}}\beta\|^2 + \lambda \|\beta\|_1$$

- ▶ Subject to:

$$\mathcal{L}(\lambda) = \min_{\lambda} \|\mathbf{y}^{\text{val}} - \mathbf{X}^{\text{val}}\hat{\beta}^{(\lambda)}\|^2$$

# Hyperparameter selection, the popular way

- ▶ Cross validation
- ▶ Criterion (AIC/BIC) that control model complexity
- ▶ Formalization: for Lasso

$$\hat{\beta}^{(\lambda)} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|\mathbf{y}^{\text{train}} - \mathbf{X}^{\text{train}}\beta\|^2 + \lambda \|\beta\|_1$$

- ▶ Subject to:

$$\mathcal{L}(\lambda) = \min_{\lambda} \|\mathbf{y}^{\text{val}} - \mathbf{X}^{\text{val}}\hat{\beta}^{(\lambda)}\|^2$$

→ Today: hyper-parameter selection with bi-level optimization



## Hyperparameter Selection: Bilevel Optimization?

$$\mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) = \underbrace{\min_{\lambda} \|y^{\text{val}} - X^{\text{val}} \hat{\beta}^{(\lambda)}\|^2}_{\text{outer optimization problem}} \text{ s.t. } \underbrace{\hat{\beta}^{(\lambda)} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} h(\beta, \lambda)}_{\text{inner optimization problem}}$$

Caveat: for the moment we deviate from Lasso, and assume the case  $h$  is at least twice-differentiable

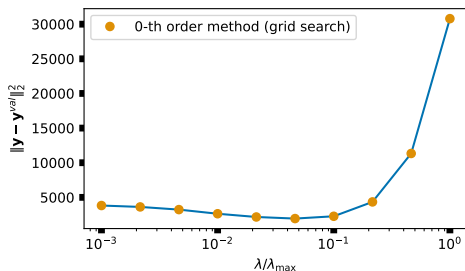
## Grid-search as a zero-order optimization method

$$\mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) = \underbrace{\min_{\lambda} \|y^{\text{val}} - X^{\text{val}} \hat{\beta}^{(\lambda)}\|^2}_{\text{outer optimization problem}} \text{ s.t. } \underbrace{\hat{\beta}^{(\lambda)} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} h(\beta, \lambda)}_{\text{inner optimization problem}}$$

Grid-search with cross-validation (assume 1-fold CV):

- ▶ Defines a range of values for  $\lambda$
- ▶ For each  $\lambda$ , solves the inner problem, then calculate the outer loss
- ▶ Choose  $\lambda \in \text{grid}(\lambda)$  that that minimizes the outer loss

# Grid-search as a zero-order optimization method



Grid-search with cross-validation (assume 1-fold CV):

- ▶ Defines a range of values for  $\lambda$
- ▶ For each  $\lambda$ , solves the inner problem, then calculate the outer loss
- ▶ Choose  $\lambda \in \text{grid}(\lambda)$  that that minimizes the outer loss

---

Example from: <https://qb3.github.io/sparse-ho/index.html>

# First-order hyperparameter-optimization?

$$\mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) = \underbrace{\min_{\lambda} \|\mathbf{y}^{\text{val}} - \mathbf{X}^{\text{val}} \hat{\beta}^{(\lambda)}\|^2}_{\text{outer optimization problem}} \text{ s.t. } \underbrace{\hat{\beta}^{(\lambda)} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} h(\beta, \lambda)}_{\text{inner optimization problem}}$$

- Idea: gradient descent?

$$\lambda^{(t+1)} = \lambda^{(t)} - \eta \nabla \mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda^{(t)})$$

# First-order hyperparameter-optimization?

$$\mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) = \underbrace{\min_{\lambda} \|\mathbf{y}^{\text{val}} - \mathbf{X}^{\text{val}} \hat{\beta}^{(\lambda)}\|^2}_{\text{outer optimization problem}} \text{ s.t. } \underbrace{\hat{\beta}^{(\lambda)} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} h(\beta, \lambda)}_{\text{inner optimization problem}}$$

- Idea: gradient descent?

$$\lambda^{(t+1)} = \lambda^{(t)} - \eta \nabla \mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda^{(t)})$$

# First-order hyperparameter-optimization

$$\mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) = \underbrace{\min_{\lambda} \|\mathbf{y}^{\text{val}} - \mathbf{X}^{\text{val}} \hat{\beta}^{(\lambda)}\|^2}_{\text{outer optimization problem}} \text{ s.t. } \underbrace{\hat{\beta}^{(\lambda)} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} h(\beta, \lambda)}_{\text{inner optimization problem}}$$

- ▶ Previous calculus classes tell us that

$$\nabla \mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) = \partial_{\lambda} \hat{\beta}^{(\lambda)\top} \nabla_1 \mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) + \nabla_2 \mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda)$$

# First-order hyperparameter-optimization

$$\mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) = \underbrace{\min_{\lambda} \|\mathbf{y}^{\text{val}} - \mathbf{X}^{\text{val}} \hat{\beta}^{(\lambda)}\|^2}_{\text{outer optimization problem}} \text{ s.t. } \underbrace{\hat{\beta}^{(\lambda)} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} h(\beta, \lambda)}_{\text{inner optimization problem}}$$

- ▶ Previous calculus classes tell us that

$$\nabla \mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) = \partial_{\lambda} \hat{\beta}^{(\lambda)\top} \nabla_1 \mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) + \nabla_2 \mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda)$$

- ▶ Question: what is problematic in computation of this gradient?

# First-order hyperparameter-optimization

$$\mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) = \underbrace{\min_{\lambda} \|\mathbf{y}^{\text{val}} - \mathbf{X}^{\text{val}} \hat{\beta}^{(\lambda)}\|^2}_{\text{outer optimization problem}} \text{ s.t. } \underbrace{\hat{\beta}^{(\lambda)} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} h(\beta, \lambda)}_{\text{inner optimization problem}}$$

- ▶ Previous calculus classes tell us that

$$\nabla \mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) = \partial_{\lambda} \hat{\beta}^{(\lambda)\top} \nabla_1 \mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) + \nabla_2 \mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda)$$

- ▶ Question: what is problematic in computation of this gradient?
- ▶  $\hat{\beta}^{(\lambda)}$  is the solution of another optimization problem...



# Implicit Function Theorem to the rescue

Remember the inner problem:

$$\hat{\beta}(\lambda) \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} h(\beta, \lambda)$$

# Implicit Function Theorem to the rescue

Remember the inner problem:

$$\hat{\beta}^{(\lambda)} \in \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} h(\beta, \lambda)$$

►  $\hat{\beta}^{(\lambda)}$  is an implicit function of  $\lambda$ , characterized by

$$\nabla_1 h(\hat{\beta}^{(\lambda)}, \lambda) = 0$$

# Implicit Function Theorem to the rescue

Remember the inner problem:

$$\hat{\beta}^{(\lambda)} \in \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} h(\beta, \lambda)$$

- ▶  $\hat{\beta}^{(\lambda)}$  is an implicit function of  $\lambda$ , characterized by

$$\nabla_1 h(\hat{\beta}^{(\lambda)}, \lambda) = 0$$

- ▶ Implicit Function Theorem: if  $\mathcal{L}$  and  $h$  are continuously differentiable, then there exists a unique  $\hat{\beta}^{(\lambda)}$ , and we have

$$\begin{aligned}\partial_\lambda \hat{\beta}^{(\lambda)} &= -[\nabla_1^2 h(\hat{\beta}^{(\lambda)}, \lambda)]^{-1} \nabla_{1,2}^2 h(\hat{\beta}^{(\lambda)}, \lambda) \\ &= -[H_{\beta,h}]^{-1} \nabla_{1,2}^2 h(\hat{\beta}^{(\lambda)}, \lambda)\end{aligned}$$

# Implicit Function Theorem to the rescue

Remember the inner problem:

$$\hat{\beta}^{(\lambda)} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} h(\beta, \lambda)$$

- ▶  $\hat{\beta}^{(\lambda)}$  is an implicit function of  $\lambda$ , characterized by

$$\nabla_1 h(\hat{\beta}^{(\lambda)}, \lambda) = 0$$

- ▶ Implicit Function Theorem: if  $\mathcal{L}$  and  $h$  are continuously differentiable, then there exists a unique  $\hat{\beta}^{(\lambda)}$ , and we have

$$\begin{aligned}\partial_\lambda \hat{\beta}^{(\lambda)} &= -[\nabla_1^2 h(\hat{\beta}^{(\lambda)}, \lambda)]^{-1} \nabla_{1,2}^2 h(\hat{\beta}^{(\lambda)}, \lambda) \\ &= -[H_{\beta,h}]^{-1} \nabla_{1,2}^2 h(\hat{\beta}^{(\lambda)}, \lambda)\end{aligned}$$

- ▶ Question: where does this equation come from?

# First-order hyperparameter-optimization

$$\mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) = \underbrace{\min_{\lambda} \|y^{\text{val}} - X^{\text{val}} \hat{\beta}^{(\lambda)}\|^2}_{\text{outer optimization problem}} \text{ s.t. } \underbrace{\hat{\beta}^{(\lambda)} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} h(\beta, \lambda)}_{\text{inner optimization problem}}$$

► So:

$$\begin{aligned} \nabla \mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) &= \partial_{\lambda} \hat{\beta}^{(\lambda)\top} \nabla_1 \mathcal{L} + \nabla_2 \mathcal{L} \\ &= -[\nabla_{1,2}^2 h]^\top [H_{\beta, h}]^{-1} \nabla_1 \mathcal{L} + \nabla_2 \mathcal{L} \end{aligned}$$

---

Y. Bengio. Gradient-based optimization of hyperparameters. *Neural computation*, 12(8):1889–1900, 2000.

# First-order hyperparameter-optimization

$$\mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) = \underbrace{\min_{\lambda} \|\mathbf{y}^{\text{val}} - \mathbf{X}^{\text{val}} \hat{\beta}^{(\lambda)}\|^2}_{\text{outer optimization problem}} \text{ s.t. } \underbrace{\hat{\beta}^{(\lambda)} \in \underset{\beta \in \mathbb{R}^p}{\text{argmin}} h(\beta, \lambda)}_{\text{inner optimization problem}}$$

► So:

$$\begin{aligned} \nabla \mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) &= \partial_{\lambda} \hat{\beta}^{(\lambda)\top} \nabla_1 \mathcal{L} + \nabla_2 \mathcal{L} \\ &= -[\nabla_{1,2}^2 h]^\top [H_{\beta, h}]^{-1} \nabla_1 \mathcal{L} + \nabla_2 \mathcal{L} \end{aligned}$$

► But: any problem remains?

---

Y. Bengio. Gradient-based optimization of hyperparameters. *Neural computation*, 12(8):1889–1900, 2000.

# First-order hyperparameter-optimization

$$\mathcal{L}(\hat{\beta}(\lambda), \lambda) = \underbrace{\min_{\lambda} \|y^{\text{val}} - X^{\text{val}} \hat{\beta}(\lambda)\|^2}_{\text{outer optimization problem}} \text{ s.t. } \underbrace{\hat{\beta}(\lambda) \in \underset{\beta \in \mathbb{R}^p}{\text{argmin}} h(\beta, \lambda)}_{\text{inner optimization problem}}$$

► So:

$$\begin{aligned}\nabla \mathcal{L}(\hat{\beta}(\lambda), \lambda) &= \partial_{\lambda} \hat{\beta}(\lambda)^{\top} \nabla_1 \mathcal{L} + \nabla_2 \mathcal{L} \\ &= -[\nabla_{1,2}^2 h]^{\top} [H_{\beta,h}]^{-1} \nabla_1 \mathcal{L} + \nabla_2 \mathcal{L}\end{aligned}$$

- But: any problem remains?
- Inverting Hessian is generally very costly, and not possible when  $n < p$ ...

---

Y. Bengio. Gradient-based optimization of hyperparameters. Neural computation, 12(8):1889–1900, 2000.

# First-order hyperparameter-optimization

$$\nabla \mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda) = -[\nabla_{1,2}^2 h]^\top [H_{\beta,h}]^{-1} \nabla_1 \mathcal{L} + \nabla_2 \mathcal{L}$$

Pedregosa (2016): at iteration  $k$  we have a tolerance  $\epsilon_k$  small enough

1. With  $\lambda_k$ , solve the inner optimization problem, obtain  $\hat{\beta}^{\lambda_k}$
2. Approximate  $[H_{\beta,h}]^{-1} \nabla_1 \mathcal{L}$  by solving for  $q_k$  s.t

$$\|H_{\hat{\beta}^{\lambda_k}, h} q_k - \nabla_1 \mathcal{L}\| \leq \epsilon_k$$

3. Approximate  $\nabla \mathcal{L}(\hat{\beta}^{(\lambda)}, \lambda)$  with

$$p_k = -[\nabla_{1,2}^2 h]^\top q_k + \nabla_2 \mathcal{L}(\hat{\beta}^{\lambda_k}, \lambda_k)$$

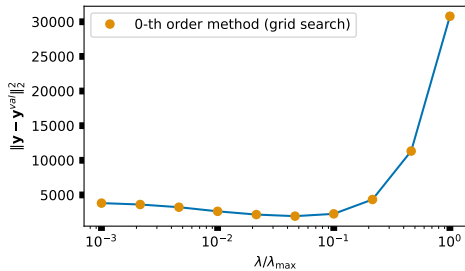
4. Update  $\lambda_{k+1} = \text{ProjGD}(\lambda_k, p_k, \eta)$

—→ no inversion of the Hessian

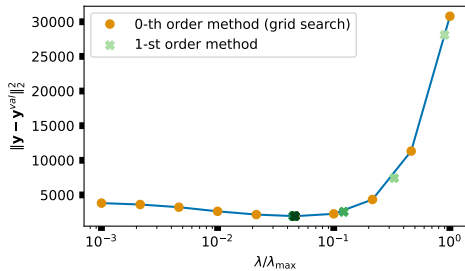
Pedregosa, F. (2016). Hyperparameter optimization with approximate gradient. In International conference on machine learning (pp. 737-746). PMLR.



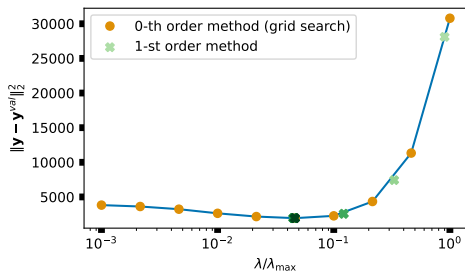
# First-order hyperparameter-optimization



# First-order hyperparameter-optimization



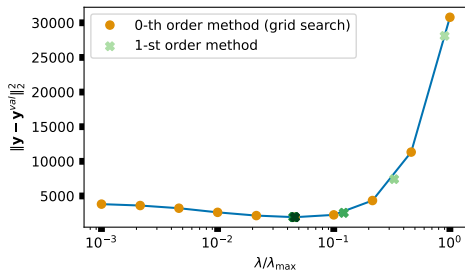
# First-order hyperparameter-optimization



- ▶ Still: we requires  $h$  to be smooth
- ▶ But what about the case for Lasso?

$$h(\boldsymbol{\beta}, \lambda) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|_1$$

# First-order hyperparameter-optimization



→ Check the work of Bertrand et al. (2020)

- ▶ Also leverage the sparsity induced by the Lasso for the computation
- ▶ Faster than implicit forward differentiation methods

---

Bertrand, Q., Klopfenstein, Q., et al. (2020). Implicit differentiation of Lasso-type models for hyperparameter optimization. Proceedings of the 37th International Conference on Machine Learning