

# Information and Model Selection

Vahid Partovi Nia

Advanced Machine Learning: Lecture 03

February 25, 2018



POLYTECHNIQUE  
MONTREAL



Least Angle

Information

① Least Angle

② Information



# Multiple Regression

Least Angle

Information

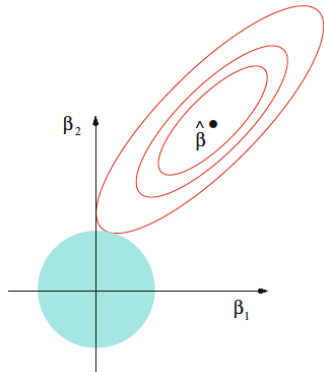
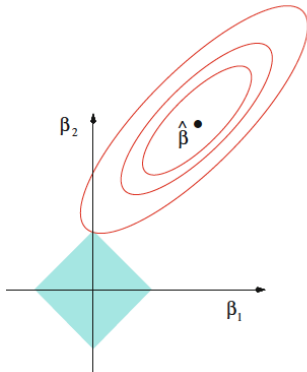
- ① initialize  $\mathbf{r}_0 = \mathbf{1}$
- ②  $j = 1$
- ③ regress  $\mathbf{x}_j$  on  $(\mathbf{r}_0, \dots, \mathbf{r}_{j-1})$
- ④  $\hat{\gamma}_{lj} = \frac{\mathbf{r}_l^\top \mathbf{x}_j}{\mathbf{r}_l^\top \mathbf{r}_l}$
- ⑤ orthogonalize  $\mathbf{r}_j = \mathbf{x}_j - \sum_{k=1}^{j-1} \hat{\gamma}_{kj} \mathbf{r}_k$
- ⑥  $j = j + 1$  go to 3
- ⑦  $\hat{\beta}_p = \frac{\mathbf{y}^\top \mathbf{r}_p}{\mathbf{r}_p^\top \mathbf{r}_p}$



# Visual penalization

Least Angle

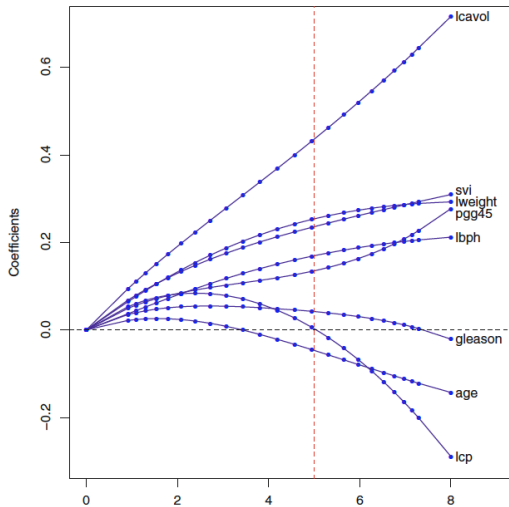
Information



# Ridge coefficients

Least Angle

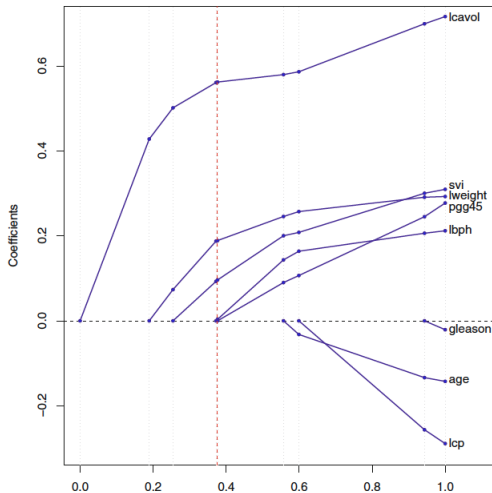
Information



# Lasso coefficients

Least Angle

Information



# Penalization and Selection

Least Angle

Information

$$S(\beta) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) + P_\lambda(\beta)$$

- $P_\lambda(\beta)$  must be non-differentiable on the axes to select.
- Lasso has partially linear path. This helps to develop the path algorithm.
- Linear selection path appears while  $S(\beta)$  is partially quadratic with non-differentiable penalization on axes.



- ① standardize  $\mathbf{x}_j$
- ② set  $\beta_j = 0, j = 1, \dots, p,$
- ③ initialize  $\mathbf{r} = \mathbf{y} - \bar{\mathbf{y}}$
- ④ Find the most correlated  $\mathbf{x}_j$  with  $\hat{\mathbf{r}}_j = \operatorname{argmax} \mathbf{r}^\top \mathbf{x}_j$
- ⑤ Move  $\beta_j$  towards its least squares,  $\beta_j = \delta \frac{\mathbf{r}^\top \mathbf{x}_j}{\mathbf{x}_j^\top \mathbf{x}_j}$
- ⑥ Update residual  $\mathbf{r} = \mathbf{y} - \beta_j \mathbf{x}_j$  until  $\mathbf{x}_k$  have more correlation.
- ⑦ Continue until  $p$  predictors are in.

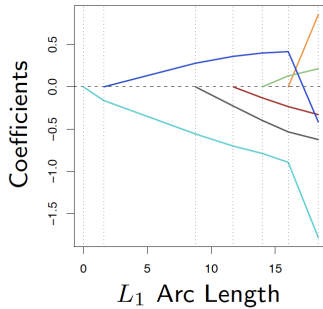




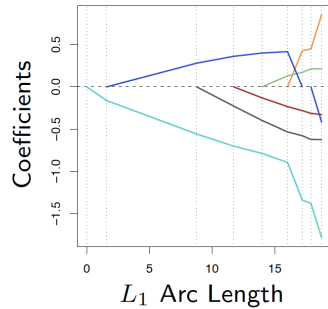
- ① standardize  $\mathbf{x}_j$
- ② set  $\beta_j = 0, j = 1, \dots, p, k = 0, A_k = \emptyset$
- ③ initialize  $\mathbf{r}_k = \mathbf{y} - \bar{\mathbf{y}}$
- ④ Add the most correlated predictor to  $A_k$ .
- ⑤  $\boldsymbol{\beta}_{A_k} = (\mathbf{X}_{A_k}^\top \mathbf{X}_{A_k})^{-1} \mathbf{X}_{A_k}^\top \mathbf{y}$
- ⑥  $\mathbf{r}_{A_k} = \mathbf{y} - \mathbf{X}_{A_k} \boldsymbol{\beta}_{A_k}$
- ⑦  $\delta_{A_k} = (\mathbf{X}_{A_k}^\top \mathbf{X}_{A_k})^{-1} \mathbf{X}_{A_k}^\top \mathbf{r}_{A_k}$
- ⑧  $\boldsymbol{\beta}_{A_k}(\delta) = \boldsymbol{\beta}_{A_k} + \delta \times \delta_{A_k}$ , increase  $\delta$  until another predictor (out of  $A_k$ ) is more correlated with  $\mathbf{y}$ . Go to 4.



Least Angle Regression



Lasso



- If a nonzero coefficient hits zero, put variables out of  $A_k$
- Recompute least squares.
- Go back to LAR algorithm.



Lasso:

$$\hat{\beta}^{\text{lasso}} = \operatorname{argmin} (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) + \lambda \sum_{j=1}^p |\beta_j|$$

Dantzig:

$$\hat{\beta}^{\text{DS}} = \operatorname{argmin} \|\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\beta)\|_\infty + \lambda \sum_{j=1}^p |\beta_j|$$



$$\int_0^1 h(x)dx$$

$$\int_{-\infty}^{\infty} h(x)dx$$



- Why do we need parametric models?
- Why do we use likelihood?
- Why maximum likelihood is good?
- What information means?
- How information is related to data?



KL divergence between the assumed class  $f(x | \theta)$  from true data distribution  $f(x | \theta_0)$  is

$$\begin{aligned}\text{KL}(\theta_0, \theta) &= \int \log \left\{ \frac{f(x | \theta_0)}{f(x | \theta)} \right\} f(x | \theta_0) \\ &= \mathbb{E}_{\theta_0} \left\{ \frac{f(x | \theta_0)}{f(x | \theta)} \right\}\end{aligned}$$



$\mathbb{KL}$  divergence between the assumed class  $f(x | \theta)$  from true data distribution  $f(x | \theta_0)$  is

$$\begin{aligned}\mathbb{KL}(\theta_0, \theta) &= \int \log \left\{ \frac{f(x | \theta_0)}{f(x | \theta)} \right\} f(x | \theta_0) \\ &= \mathbb{E}_{\theta_0} \left\{ \frac{f(x | \theta_0)}{f(x | \theta)} \right\}\end{aligned}$$

$$\mathbb{KL}(\theta_0, \theta) \neq \mathbb{KL}(\theta, \theta_0)$$

Cross entropy of the assumed class  $f(x | \theta)$  from true data distribution  $f(x | \theta_0)$  is

$$\mathbb{H}(\theta, \theta_0) = \int \log f(x | \theta) f(x | \theta_0) dx$$





$$\text{KL}(\theta_0, \theta) = \mathbb{H}(\theta_0, \theta_0) - \mathbb{H}(\theta, \theta_0)$$



- $\text{KL}(\theta_0, \theta) > 0$  iff  $f(x | \theta_0) \neq f(x | \theta)$  on a set of  $x$  with positive measure.
- $\text{KL}(\theta_0, \theta) = 0$  iff  $f(x | \theta_0) = f(x | \theta)$  almost everywhere.
- $\text{KL}_n(\theta_0, \theta) = n\text{KL}(\theta_0, \theta)$  for a set of i.i.d observations  $(x_1, \dots, x_n)$ .
- $\frac{\partial \mathbb{H}(\theta, \theta_0)}{\partial \theta} \big|_{\theta=\theta_0} = 0$
- $\frac{\partial^2 \mathbb{H}(\theta, \theta_0)}{\partial \theta \partial \theta^\top} \big|_{\theta=\theta_0} = -J(\theta_0)$  where  $J(\cdot)$  is the observed information.



Suppose  $A = \{A_1, \dots, A_k\}$  with probabilities  $p_1, \dots, p_k$ . Define  $A'$  to be an  $A$ -similar event as  $A' = \{A_1, \dots, A_k, A_{k+1}\}$  with probabilities  $p_1, \dots, p_k, p_{k+1} = 0$ .

- If two sets  $A$  and  $B$  are independent  
 $\mathbb{H}(A \times B) = \mathbb{H}(A) + \mathbb{H}(B)$ .
- $\mathbb{H}(A) = \mathbb{H}(A')$  .

The only function that satisfies the above two properties is  $\mathbb{H}(A) = \lambda \sum_i p_i \log p_i$ .  
Why this result is important?



# More about entropy

Least Angle

Information

| $A_1$ | $A_2$ |
|-------|-------|
| 0.1   | 0.9   |
| 0.49  | 0.51  |
| <hr/> |       |
| 0.69  | 0.325 |

