

# STAT406- Methods of Statistical Learning Lecture 6

Matias Salibian-Barrera

UBC - Sep / Dec 2017

“Far better an approximate answer to the right question, which is often vague, than an exact answer to the wrong question, which can always be made precise.”

John Tukey. The future of data analysis. *Annals of Mathematical Statistics*, 33(1), (1962), p. 13.

# Effective degrees of freedom

- How many “effective” parameters are we using?
- In linear regression, we have  $p$  parameters
- A more general definition is as follows. For a fitting method producing  $\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n$ ,

$$\text{edf} = \frac{1}{\sigma^2} \sum_{i=1}^n \text{cov}(\hat{y}_i, y_i)$$

Efron, B. (1986). How biased is the apparent error rate of a prediction rule? Journal of the American Statistical Association, 81(394):461-470.

# Effective degrees of freedom

- It is easy to see that for least squares predictors, we have

$$\hat{\mathbf{y}} = \mathbf{H} \mathbf{y}$$

with

$$\mathbf{H} = \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'$$

and

$$\text{edf} = \frac{1}{\sigma^2} \sum_{i=1}^n \text{cov}(\hat{y}_i, y_i) = \text{trace}(\mathbf{H}) = p$$

# Effective degrees of freedom

- More in general, for any linear predictor

$$\hat{\mathbf{y}} = \mathbf{S} \mathbf{y}$$

we have

$$\text{edf} = \text{trace}(\mathbf{S}) = \sum_{i=1}^n \mathbf{s}_{i,i}$$

# Effective degrees of freedom

- The ridge regression fit satisfies

$$\hat{\mathbf{y}} = \mathbf{S}_\lambda \mathbf{y}$$

where

$$\mathbf{S}_\lambda = \mathbf{X} (\mathbf{X}'\mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}'$$

$$\text{trace}(\mathbf{S}) = ?$$

# Effective degrees of freedom

- Using the singular value decomposition (SVD) of  $\mathbf{X}$

$$\mathbf{X} = \mathbf{U} \mathbf{\Lambda} \mathbf{V}'$$

where  $\mathbf{U} \in \mathbb{R}^{n \times p}$ ,  $\mathbf{V} \in \mathbb{R}^{p \times p}$  with

$$\mathbf{U}'\mathbf{U} = \mathbf{I}_p = \mathbf{V}'\mathbf{V}$$

and

$$\mathbf{\Lambda} = \text{diag}(d_1, \dots, d_p),$$

we have

$$\text{trace}(\mathbf{S}) = \sum_{i=1}^p \left( \frac{d_i^2}{d_i^2 + \lambda} \right)$$

# Effective degrees of freedom

- For example, in the Air Pollution data example, if we use

$$\lambda = \exp(6)$$

we get

$$\text{edf} = 9.9$$



# Model / feature selection - LASSO

- Another regularized method is given by LASSO

$$\min_{\alpha, \beta} \sum_{i=1}^n (y_i - \alpha - \beta' \mathbf{x}_i)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

$$\min_{\alpha, \beta} \sum_{i=1}^n (y_i - \alpha - \beta' \mathbf{x}_i)^2 + \lambda \|\beta\|_1$$

for some  $\lambda > 0$

# Model / feature selection - LASSO

- The above is equivalent to

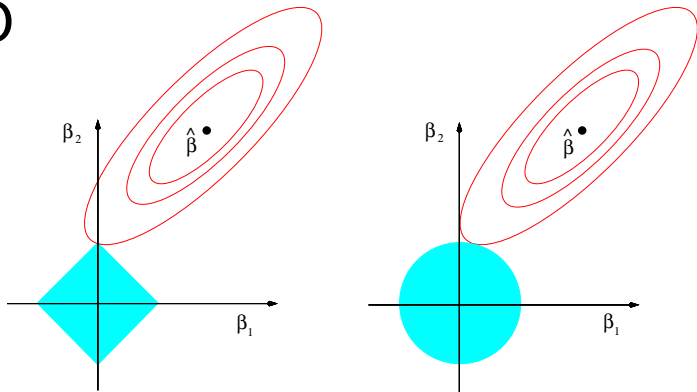
$$\min_{\alpha, \beta} \sum_{i=1}^n (y_i - \alpha - \beta' \mathbf{x}_i)^2$$

subject to

$$\sum_{j=1}^p |\beta_j| \leq K$$

for some  $K > 0$

# LASSO



**FIGURE 3.11.** *Estimation picture for the lasso (left) and ridge regression (right). Shown are contours of the error and constraint functions. The solid blue areas are the constraint regions  $|\beta_1| + |\beta_2| \leq t$  and  $\beta_1^2 + \beta_2^2 \leq t^2$ , respectively, while the red ellipses are the contours of the least squares error function.*

# Credit data - glmnet output

