

# Regression Analysis

## Model Selection

**Nicoleta Serban, Ph.D.**

*Professor*

Stewart School of Industrial and Systems Engineering

Regularized Regression:  
Penalties



## About This Lesson



# Bias-Variance Tradeoff

**Prediction Risk:** Measure of the Bias-Variance Tradeoff

$$R(S) = \frac{1}{n} \sum_{i=1}^n E(\hat{Y}_i(S) - Y_i^*)^2$$

Irreducible error

Mean Square Error

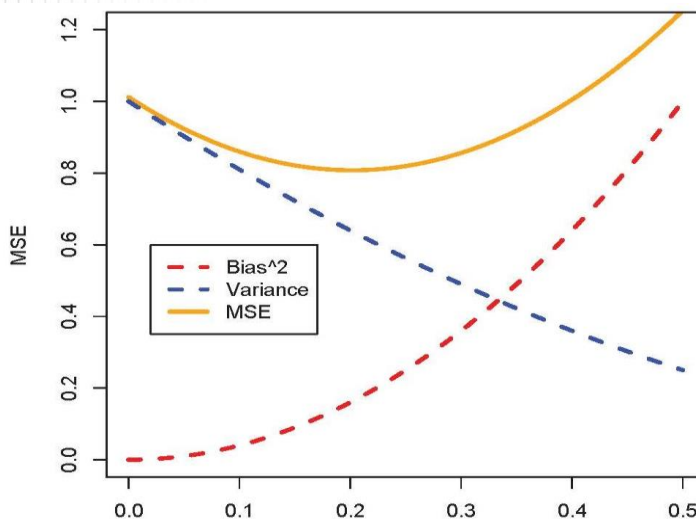
$$= V(Y_i^*) + \text{Bias}^2(\hat{Y}_i(S)) + V(\hat{Y}_i(S))$$

for a submodel  $S$ , with  $\hat{Y}_i(S)$  the fitted response for model  $S$  and  $Y_i^*$  the future observation.

- It is possible to find a model with lower MSE than the full model!
- It is “generic” in statistics: introducing some bias often yields in a decrease in MSE.

Georgia  
Tech

# Bias-Variance Tradeoff



Georgia  
Tech

# Biased Regression: Penalties

Not all biased models are better.

**We need a way to find “good” biased models!**

- Penalize large values of  $\beta$ s jointly
  - Should lead to “multivariate” shrinkage of the vector  $\beta$
- Goal is really to penalize “complex” models
  - Heuristically, “large” is interpreted as “complex model”
    - If truth really is complex, this may not work!
      - It will then be hard to build a good model anyways



# Regularized Regression

## Without Penalization

Estimate  $(\beta_0, \beta_1, \dots, \beta_p)$  by minimizing the sum of squared errors

$$\sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}))^2$$

## With Penalization

Estimate  $(\beta_0, \beta_1, \dots, \beta_p)$  by minimizing the penalized sum of squared errors

$$\sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}))^2 + \lambda \text{Penalty}(\beta_1, \dots, \beta_p)$$

The bigger  $\lambda$ , the bigger the penalty for model complexity.



# Regularized Regression (cont'd)

The penalized sum of squared errors:

$$Q(\beta_1, \dots, \beta_p) = \sum_{i=1}^n \left( y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \right)^2 + \lambda \text{Penalty}(\beta_1, \dots, \beta_p)$$

We consider three choices for the penalty:

## $L_0$ penalty

$\|\beta\|_0 = \#\{j: \beta_j \neq 0\} \Rightarrow$  Minimizing  $Q$  means searching through all submodels

## $L_1$ penalty (LASSO Regression)

$\|\beta\|_1 = \sum_{j=1}^p |\beta_j| \Rightarrow$  Minimizing  $Q$  forces many  $\beta_j$ s to be zeros

## $L_2$ penalty (Ridge Regression)

$\|\beta\|_2 = \sqrt{\sum_{j=1}^p \beta_j^2} \Rightarrow$  Minimizing  $Q$  accounts for multicollinearity



# Comparing Penalties

- $L_0$  penalty
  - Provides best model given a selection criterion
  - Requires fitting all submodels
- $L_1$  penalty
  - Measures sparsity
- $L_2$  penalty
  - Easy to implement
  - Does not do variable selection

**Example:** Consider vectors  $\mathbf{u} = (1, 0, \dots, 0)$  and  $\mathbf{v} = (\frac{1}{\sqrt{p}}, \dots, \frac{1}{\sqrt{p}})$ , both of length  $p$ .

Vector  $\mathbf{u}$  is sparse, because it contains mostly zeros.

Using the  $L_1$  norm, we have  $\|\mathbf{u}\|_1 = \sum_{i=1}^p |u_i| = 1$  and  $\|\mathbf{v}\|_1 = \sum_{i=1}^p |v_i| = \sqrt{p}$ .

Using the  $L_2$  norm, we have  $\|\mathbf{u}\|_2 = \sqrt{\sum_{i=1}^p u_i^2} = 1$  and  $\|\mathbf{v}\|_2 = \sqrt{\sum_{i=1}^p v_i^2} = 1$ .

The  $L_1$  penalty rewards the sparsity of  $\mathbf{u}$ ; the  $L_2$  penalty makes no distinction.



# Summary

