



NYC DATA SCIENCE  
**ACADEMY**

# Regularization and Cross Validation

---

Data Science Bootcamp

---

# Outline

---

- ❖ **Part 1: Ridge & Lasso Regression**
- ❖ **Part 2: Cross-Validation**
- ❖ **Part 3: Review**

*PART 1*

# Ridge & Lasso Regression

# Subset Selection

---

- ❖ In subset selection, we identify a **subset of  $p$  predictors** that we believe are related to the response. We then fit a model using least squares regression on the reduced set of variables.
- ❖ Subset selection can be performed in a myriad of ways:
  - Best subset selection
  - Forward stepwise regression
  - Backward stepwise regression
  - Both stepwise regression
- ❖ Select a model based on selected criteria:
  - AIC
  - BIC
  - $R^2_{\text{Adj.}}$

# Shrinkage/Regularization

---

- ❖ In shrinkage/regularization, we fit a model involving all predictors; however, the estimated coefficients are **shrunken towards 0** relative to the least squares estimates. As a result:
  - Estimate variance is reduced.
  - Variable selection can be performed.
  
- ❖ In order to regularize the coefficient estimates, we must further **constrain** the original least squares minimization problem. How do we do this?
  - Ridge regression
  - Lasso regression

## Multiple Linear Regression: Mathematically

---

- ❖ Recall that in multiple linear regression we wish to quantify the relationship between X and Y as follows:

$$Y \approx \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$

- ❖ Our original task was to find the estimates of  $\beta_0, \beta_1, \beta_2, \dots, \beta_p$  that **reduce the sum of the squared vertical distances** from the observations to the regression surface (i.e., the RSS) as much as possible.
  - Solved the minimization problem using basic calculus and linear algebra.

$$RSS = \sum_{i=1}^n e_i^2$$

## Ridge Regression: Mathematically

---

- ❖ Ridge regression is an [extension of the minimization problem](#) posed by multiple linear regression; rather than attempting to simply reduce the RSS, ridge regression attempts to minimize the following:

$$RSS + \lambda \sum_{j=1}^p \beta_j^2$$

- ❖ **NB:** Here,  $\lambda$  is a [tuning parameter](#) that essentially determines how the ridge regression will operate.

## Ridge Regression: Mathematically

---

- ❖ As in the multiple linear regression setting, ridge regression attempts to find coefficient estimates that **render the RSS as small as possible**, thus fitting the data well.
- ❖ Additionally, there is an added **shrinkage penalty** (the extra term containing the tuning parameter  $\lambda$ ).
  - Notice, the shrinkage penalty is simply the sum of the squared coefficient estimates. This penalty will be small when the estimates are close to 0. Therefore, it has the effect of **shrinking the coefficient estimates as a group!**
- ❖ The value of  $\lambda$  **determines the relative impact** of the RSS and shrinkage penalty terms on the resulting coefficient estimates.
  - A small  $\lambda$  penalizes the RSS more than the shrinkage penalty.
  - A large  $\lambda$  penalizes the shrinkage penalty more than the RSS.

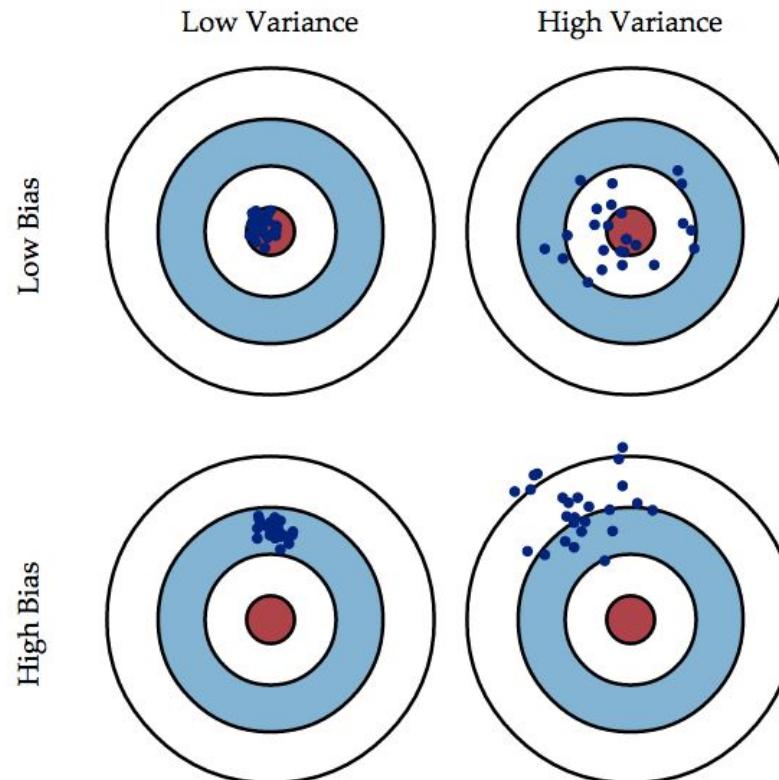
## Ridge Regression: Mathematically

---

- ❖ Recall that the standard least squares coefficient estimates are **scale equivariant**: if we multiply a predictor variable by a constant, the corresponding least squares coefficient estimate will be scaled down by the same constant.
  - Regardless of how a predictor variable is scaled, the resulting product with the corresponding estimated coefficient will **remain the same**.
- ❖ In ridge regression, this is not the case. Coefficient estimates can **change dramatically** when multiplying a given predictor by a constant due to the shrinkage penalty; it depends on the sum of the squared coefficients!
- ❖ To avoid the issue of overvaluing or undervaluing certain predictor variables simply based on their magnitudes, we must **standardize the variables** prior to performing ridge regression.

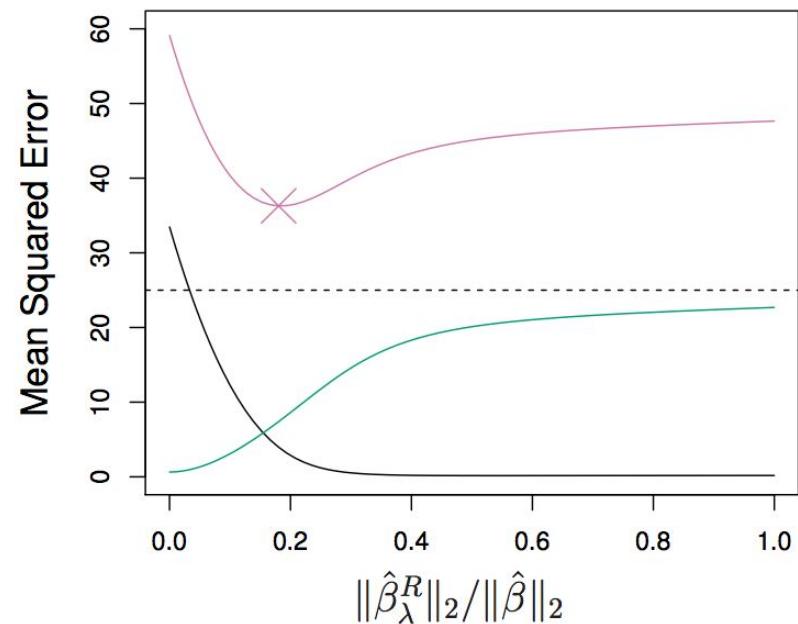
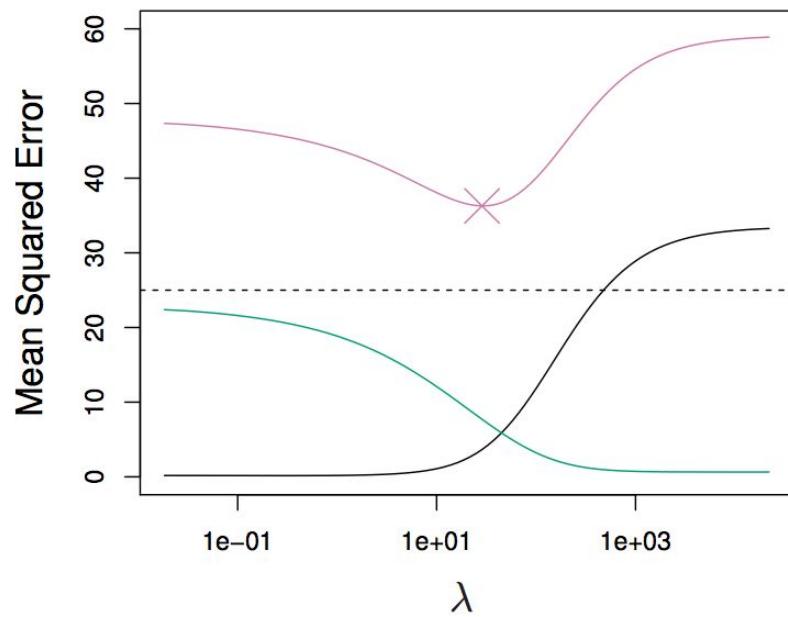
# Ridge Regression: Mathematically

- ❖ Ideally, in any scenario, we hope to uncover a model that has:
  - Low bias (i.e., high accuracy)
  - Low variance (i.e., high precision)



# Ridge Regression: Mathematically

- ❖ By shrinking the coefficient estimates towards 0 by increasing  $\lambda$ , we see that:
  - The bias (black) increases slightly but remains relatively small.
  - The variance (green) reduces substantially.
  - The mean squared error (red) of the predictions drops.



## Lasso Regression: Mathematically

---

- ❖ The main disadvantage of ridge regression is that, while parameter estimates are shrunken, they **only asymptotically approach 0** as we increase the value of  $\lambda$ .
  - Thus, the resulting model still includes estimates for all parameters.
- ❖ Lasso regression is another **extension of the minimization problem** posed by multiple linear regression; rather than attempting to simply reduce the RSS, lasso regression attempts to minimize the following:

$$RSS + \lambda \sum_{j=1}^p |\beta_j|$$

- ❖ **NB:** The only difference between the lasso and ridge regressions is that lasso implements the  **$l_1$  penalty** (norm) rather than the  **$l_2$  penalty**.

## Lasso Regression: Mathematically

---

- ❖ While both ridge and lasso regression have the effect of shrinking coefficient estimates towards 0, the lasso **necessarily forces** some coefficient estimates to be exactly 0 (when  $\lambda$  is sufficiently large).
- ❖ Lasso regression has the added advantage of essentially performing variable selection, yielding models that are **both accurate and parsimonious**.
- ❖ Once again, the value of  $\lambda$  **determines the relative impact** of the RSS and shrinkage penalty terms on the resulting coefficient estimates.
  - In both ridge and lasso regression, it is important to select an appropriate value of  $\lambda$  by means of **cross-validation**.

## Lasso VS Ridge

---

- ❖ Why is it the case that lasso regression results in coefficient estimates that are **exactly 0**? Why does ridge regression only end up **shrinking** the coefficients?
- ❖ We can restate the optimization problems of lasso and ridge regression, respectively, as the following **Lagrangian multiplier** scenarios:

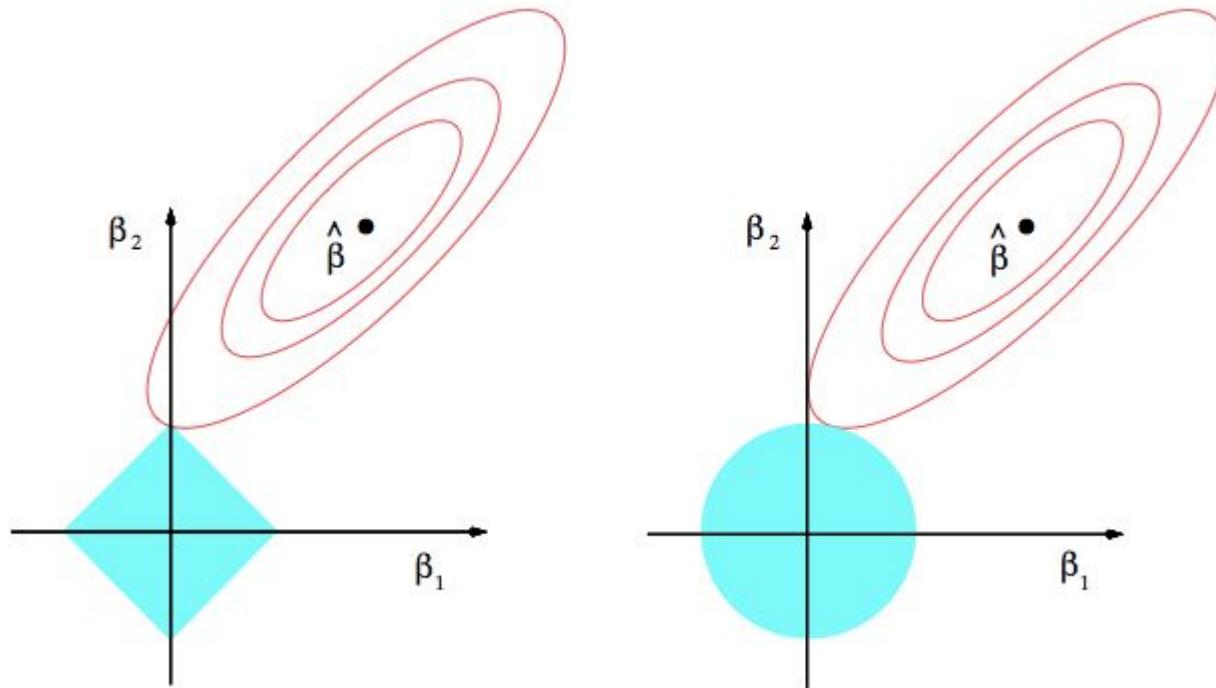
$$\underset{\beta}{\text{minimize}} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \quad \text{subject to} \quad \sum_{j=1}^p |\beta_j| \leq s$$

and

$$\underset{\beta}{\text{minimize}} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \quad \text{subject to} \quad \sum_{j=1}^p \beta_j^2 \leq s,$$

## Lasso VS Ridge: Visually

- ❖ The **red ellipses** represent the contours of the least squares error function.
- ❖ The **blue regions** represent the constrained regions for lasso and ridge, respectively.



## Selecting the Tuning Parameter: $\lambda$

---

- ❖ How do we choose the best value of  $\lambda$ ?
- ❖ Let's take a look at the method of [cross-validation](#).

*PART 2*

# Cross-Validation

# The Training & Test Sets

---

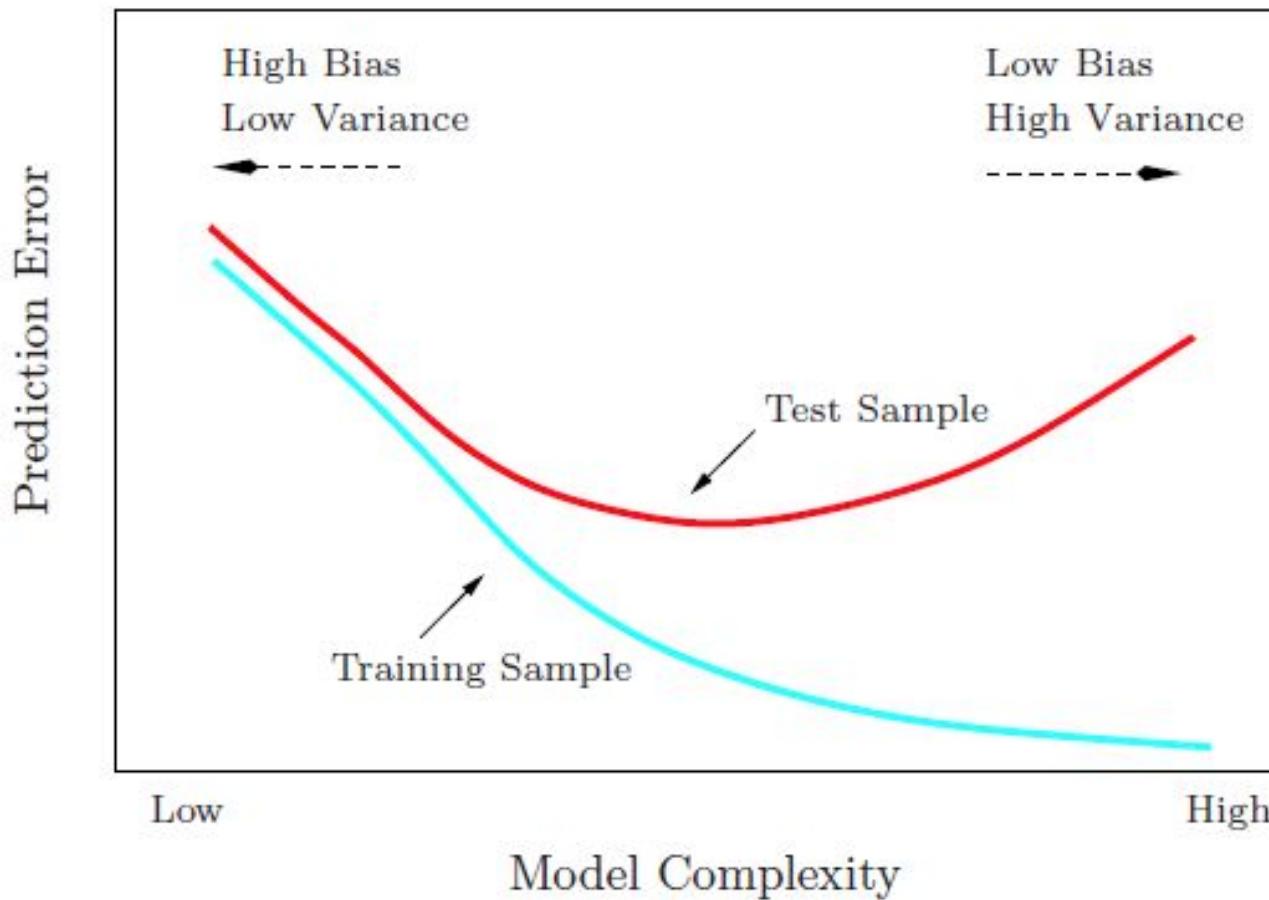
- ❖ When conducting cross-validation, we generally split the observations in our data into different (non-overlapping) sections:
  - The **training** set
  - The **test** set
- ❖ We use the observations **in the training set** to fit an initial model, but then can calculate the error in different ways:
  - The **training error** is calculated by applying the fitted model to the observations used to create the initial model (the training set).
  - The **test error** is calculated by applying the fitted model to the observations not used to create the initial model (the test set); this is like assessing the predictions of new observations.
- ❖ Why does the training error **underestimate** the test error?

# The Bias-Variance Tradeoff

---

- ❖ Recall that:
  - Bias is **how far off** on the average the model is from the truth.
  - Variance is how much the estimate **varies** about its average.
  
- ❖ With low model complexity:
  - **Bias is high** because predictions are more likely to stray from the truth with an inflexible model.
  - **Variance is low** because there are only few parameters being fit.
  
- ❖ With high model complexity:
  - **Bias is low** because the model can adapt to more subtleties in the data.
  - **Variance is high** because we have more parameters to estimate from the same amount of data.

# The Bias-Variance Tradeoff



# K-Fold Cross-Validation

---

- ❖ The results of **K-fold cross-validation** can help determine the best model at hand by estimating the test error among ultimate models.
  
- ❖ The process boils down to the following steps:
  - **Divide** your data into  $K$  (relatively) equal parts.
  - Leave out one of the  $K$  parts (call it part  $k$ ), and **put it to the side**.
  - Fit the model to the remaining ( $K - 1$ ) parts all together as your **training set**.
  - Use part  $k$  as your **test set** to estimate the prediction error.
  - **Repeat** this process  $K$  times, once each for the different splits of your data.

## K-Fold Cross-Validation

---

- ❖ The test error can be estimated from the results of this process by essentially computing a **weighted average** of the  $K$  folds as follows:

$$CV_K = \sum_{k=1}^K \frac{n_k}{N} MSE_k$$

- ❖ Here:
  - $K$  is the number of groups (folds).
  - $n_k$  is the number of observations in fold  $k$  out of the total  $N$  observations.
  - $MSE_k$  is the mean squared error obtained by using fold  $k$  as the test set, and the remaining data as the training set.
- ❖ **NB:** For classification problems, the MSE is simply replaced by the error rate.

## Yet Again, How do we Choose $K$ ?

---

- ❖ 5- or 10-fold cross-validation is typically used because these values have been empirically shown to yield estimates of the test error rate that tend to neither suffer from extreme bias nor high variance.

## Selecting the Tuning Parameter: $\lambda$

---

- ❖ In order to best implement the ridge and lasso regression methods, we need to have a way of determining **the best value of  $\lambda$**  (or, equivalently, the constraint  $s$  in the Lagrange multiplier formulation of the problem).
- ❖ As mentioned earlier, **cross-validation** helps us check by iterating across a slew of  $\lambda$  values and computing the cross-validation error rate for each.
  - Split the data into training and test sets (10-fold CV).
  - Select the  $\lambda$  for which the cross-validation error is the smallest.
- ❖ Lastly, refit the model **using all available observations**, this time with the best selected value of the tuning parameter.

*PART 3*

# Review

# Review

---

- ❖ Part 1: Alternatives to PCA: Ridge & Lasso Regression
  - Alternatives to PCA
  - Subset Selection
  - Shrinkage/Regularization
  - Mathematically
    - Multiple Linear Regression
    - Ridge Regression
    - Lasso Regression
  - Lasso VS Ridge
    - Visually
- ❖ Part 2: Cross-Validation
  - The Training & Test Sets
  - The Bias-Variance Tradeoff
  - $K$ -Fold Cross-Validation
  - Yet Again, How do we Choose  $K$ ?
  - Selecting the Tuning Parameter:  $\lambda$