



Resampling Methods

ddebarr@uw.edu

2017-02-02

“Oh sure, going in that direction will totally minimize the objective function“

-- Sarcastic Gradient Descent

-- John Urschel, Baltimore Ravens Offensive Lineman, MIT PhD Candidate (Math)



Course Outline

1. Introduction to Statistical Learning
2. Linear Regression
3. Classification
4. Resampling Methods
5. Linear Model Selection and Regularization
6. Moving Beyond Linearity
7. Tree-Based Methods
8. Support Vector Machines
9. Unsupervised Learning
10. Neural Networks and Genetic Algorithms



Agenda

	5 Resampling Methods	175
	5.1 Cross-Validation	176
	5.1.1 The Validation Set Approach	176
	5.1.2 Leave-One-Out Cross-Validation	178
Discriminant Analysis (from last week)	5.1.3 k -Fold Cross-Validation	181
	5.1.4 Bias-Variance Trade-Off for k -Fold Cross-Validation	183
Resampling Methods	5.1.5 Cross-Validation on Classification Problems	184
Hands-On Labs (including caret)	5.2 The Bootstrap	187
	5.3 Lab: Cross-Validation and the Bootstrap	190
	5.3.1 The Validation Set Approach	191
	5.3.2 Leave-One-Out Cross-Validation	192
	5.3.3 k -Fold Cross-Validation	193
	5.3.4 The Bootstrap	194
	5.4 Exercises	197

Validation Set Approach

- Observations are randomly divided into training and validation sets
- 50/50 split can be used; but 80/20 may be more common
- Training set appears in blue; Validation set appears in beige

A diagram illustrating the validation set approach. At the top, a long horizontal bar with a black border contains the numbers '1 2 3' on the left and 'n' on the right. A large black arrow points downwards from the center of this bar to a split bar below. The split bar consists of two adjacent rectangular segments. The left segment is light blue and contains the numbers '7 22 13'. The right segment is light beige and contains the number '91'.

1 2 3

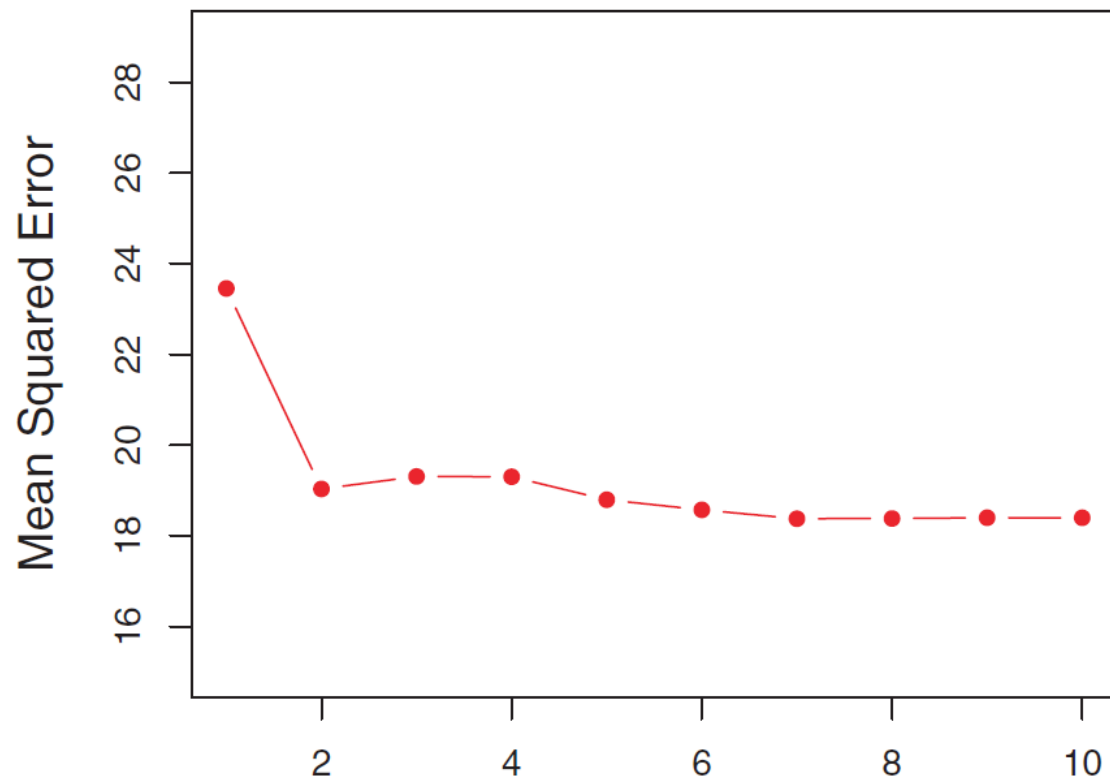
n

7 22 13

91

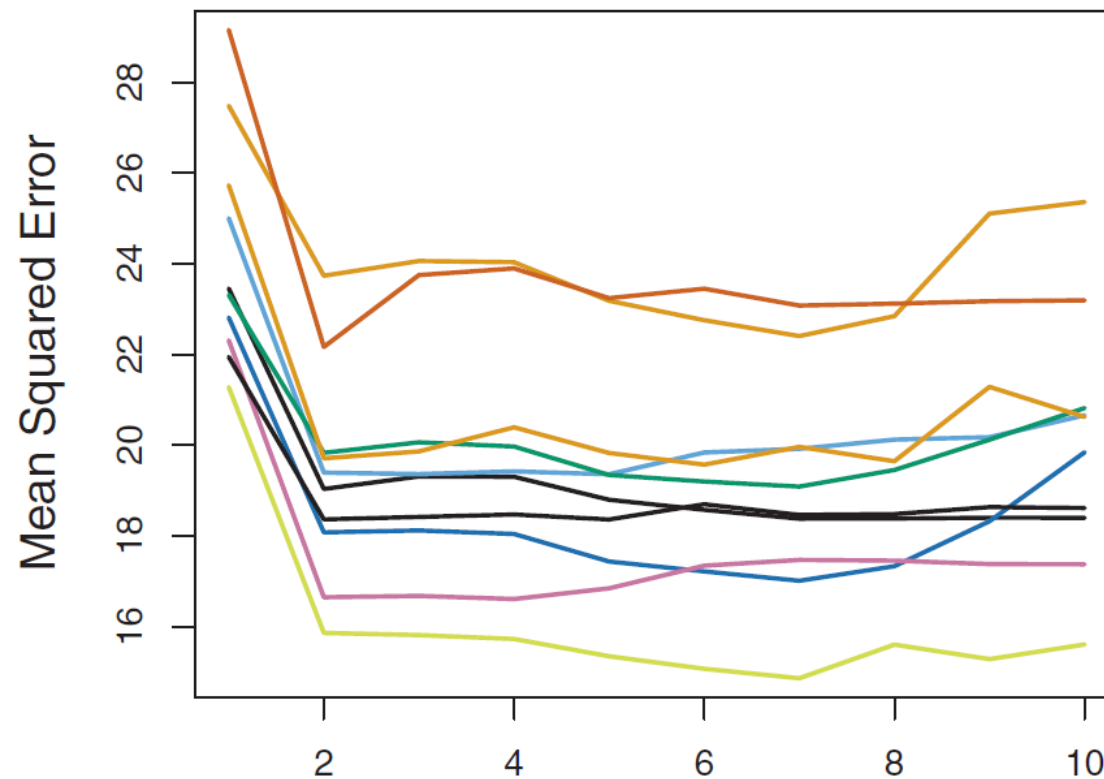


Repeated Validation Set Approach



Degree of Polynomial

One Train/Validation Set



Degree of Polynomial

Ten Train/Validation Sets

MSE Reported for Validation Set

Leave One Out Cross Validation (LOOCV)

- Each observation takes its turn as the validation set
- 'n' models are constructed; one for each validation set



$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \text{MSE}_i$$



LOOCV Short Cut for Linear Regression

- We can use the leverage statistic to turn the error estimates for the training set into a LOOCV estimate
- For multiple regression, we use the entries of the diagonal of the “projection” matrix (sometimes called the hat matrix, because it is used to derive \hat{y})

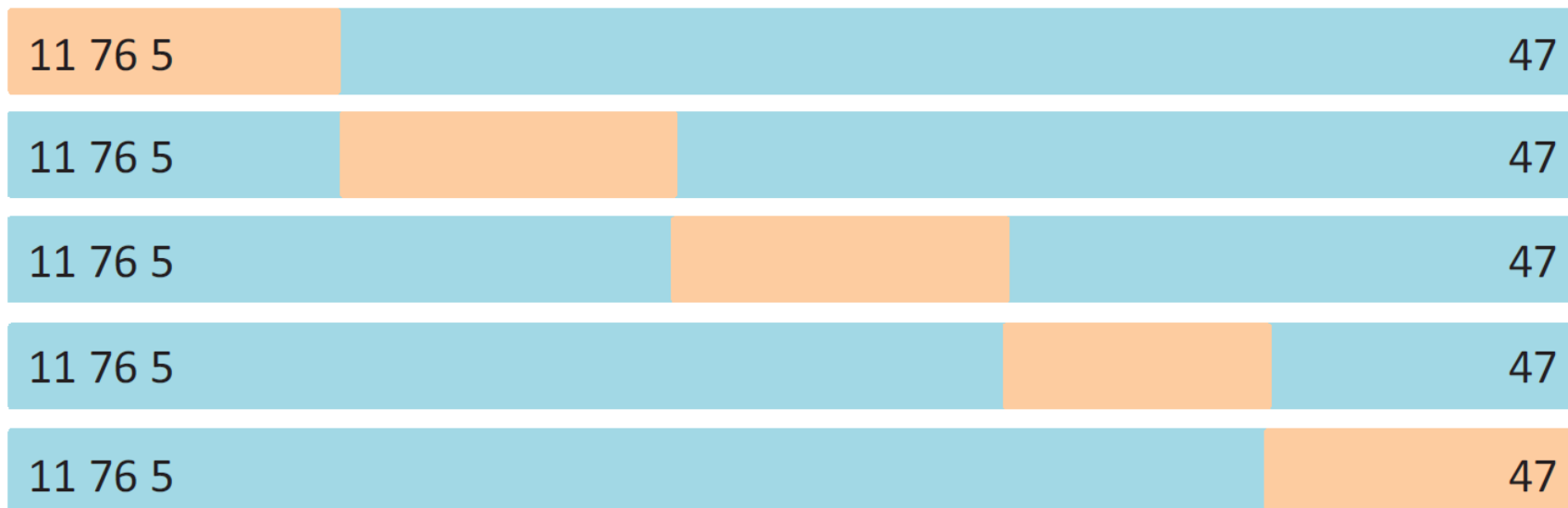
$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{y}_i}{1 - h_i} \right)^2$$

$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{i'=1}^n (x_{i'} - \bar{x})^2}$$

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

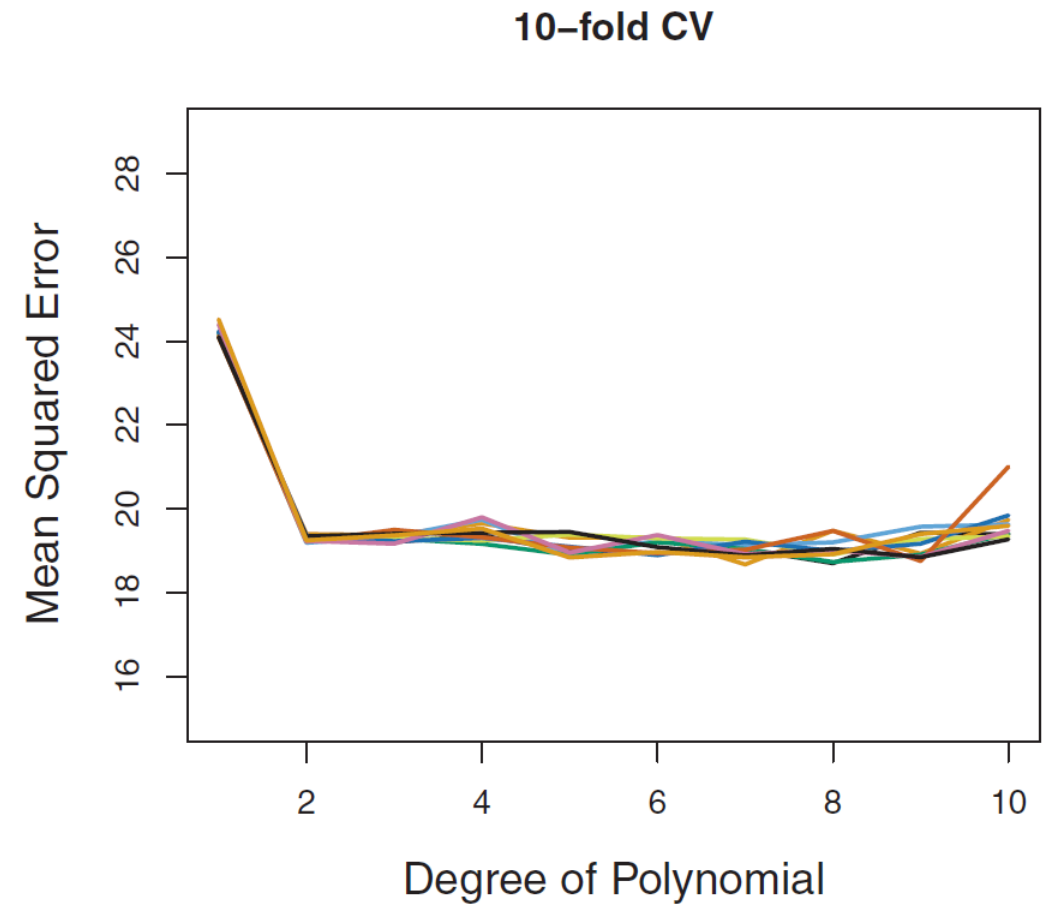
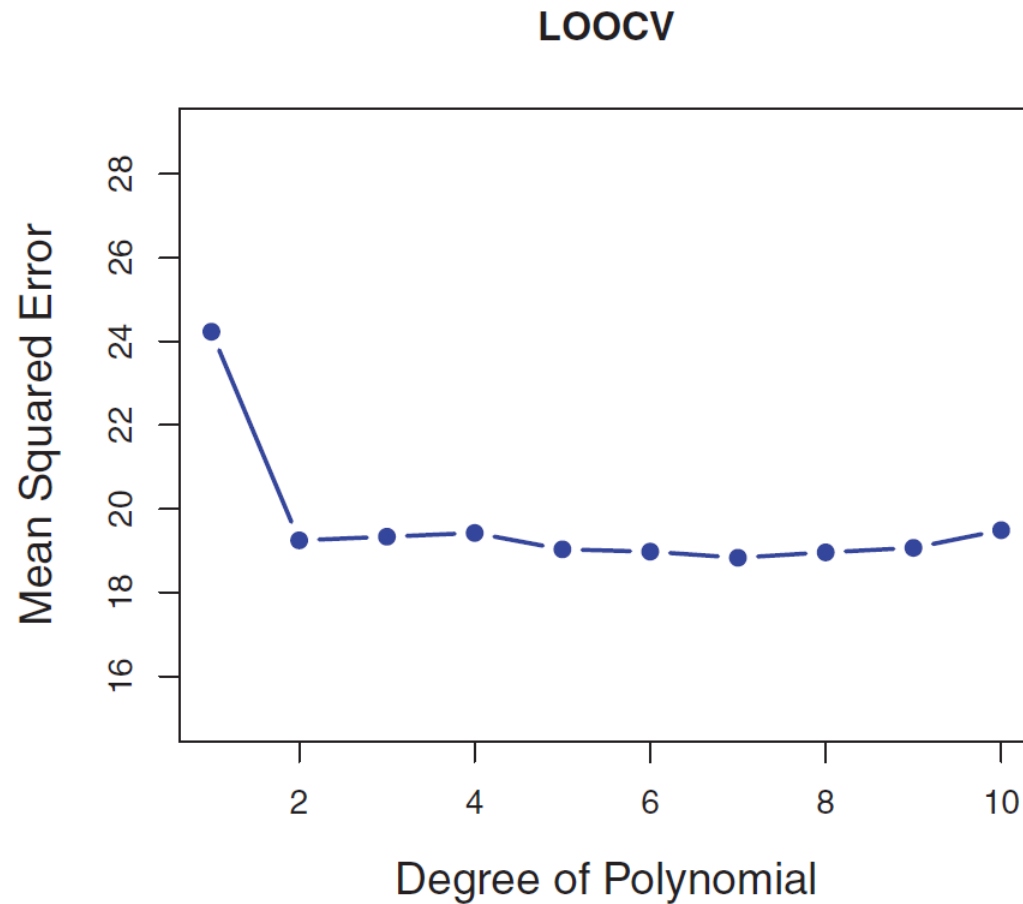
K-Fold Cross Validation

- Training set appears in blue; Validation set appears in beige
- 'k' models are constructed; one for each validation set



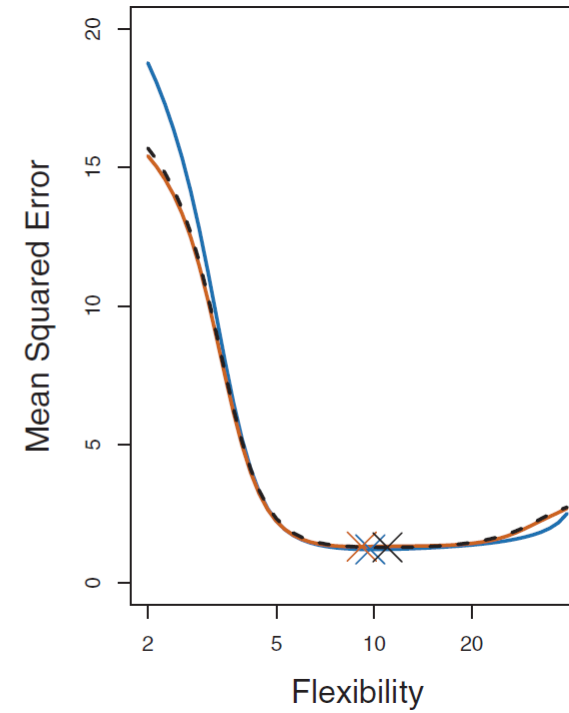
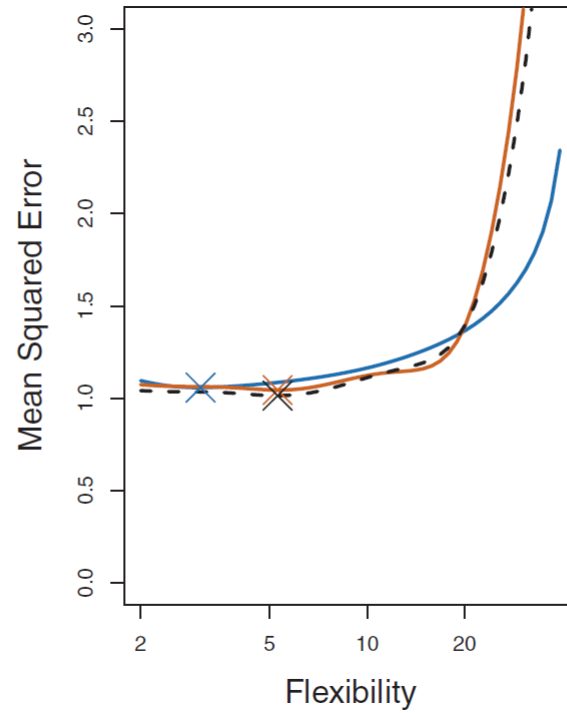
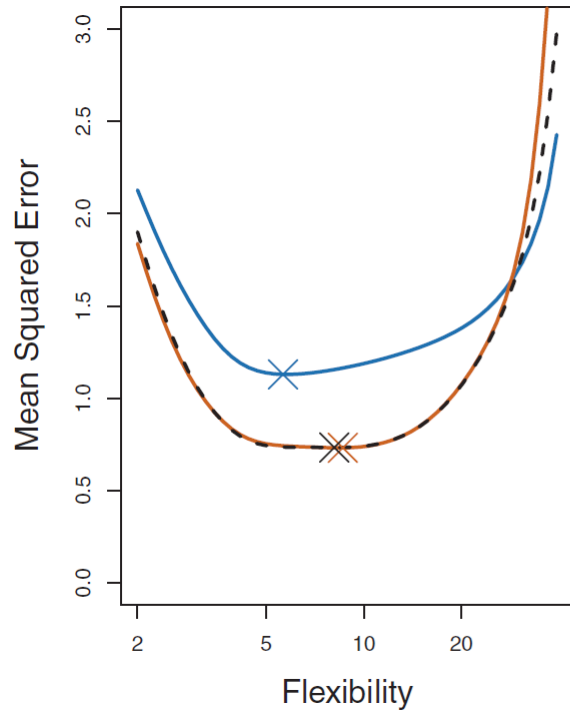
$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^k \text{MSE}_i$$

LOOCV versus Repeated 10-fold CV

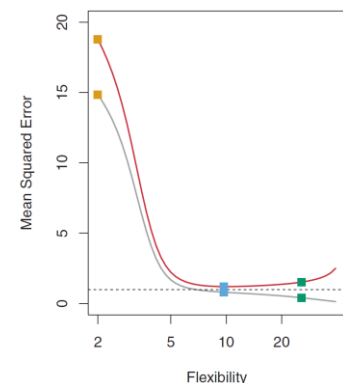
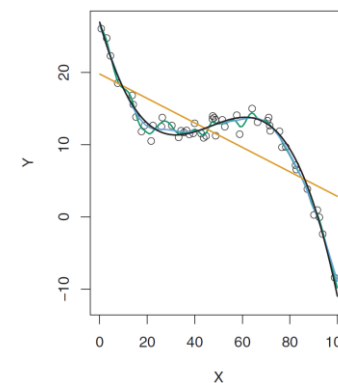
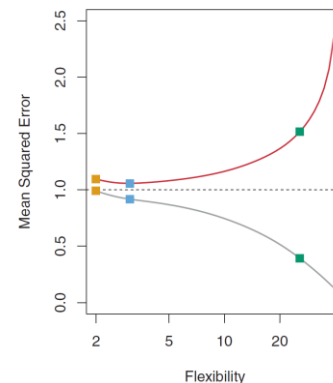
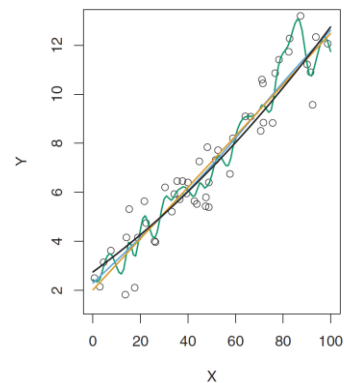
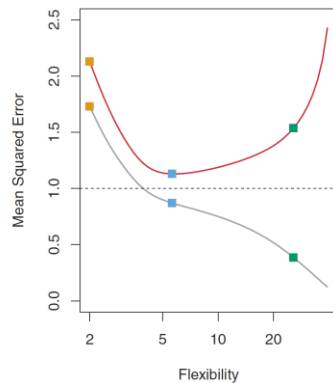
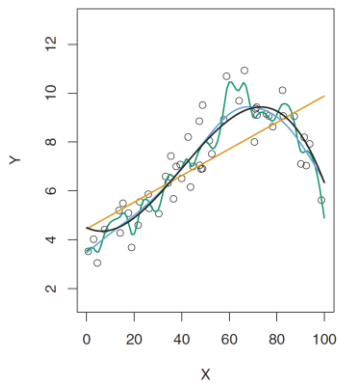




Cross Validation Applied to Simulated Data



True Error: blue
 LOOCV: black
 10-fold CV: beige





Cross Validation for Classification Problems

- Analogous definitions apply
- Example: LOOCV

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \text{Err}_i$$

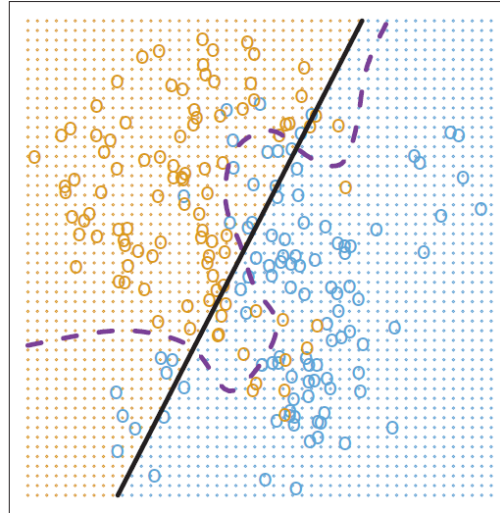
$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \beta_3 X_2 + \beta_4 X_2^2$$



Example Quadratic Logistic Regression Model

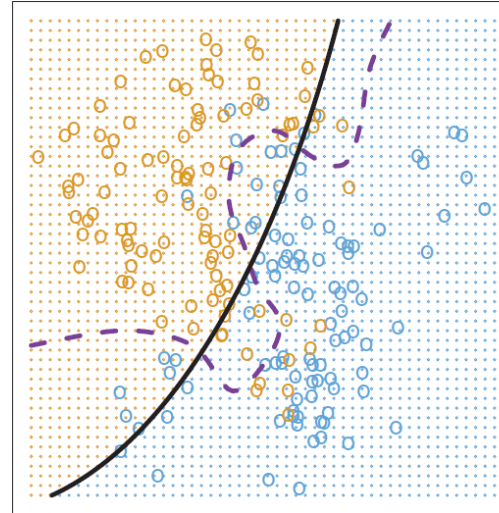
Error = 0.201

Degree=1



Error = 0.197

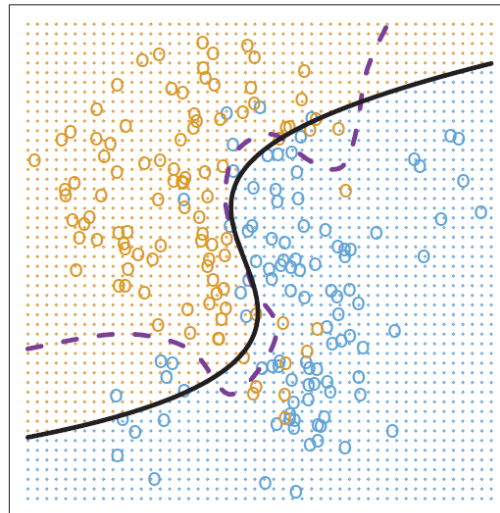
Degree=2



Bayes Error (dashed line) = 0.133

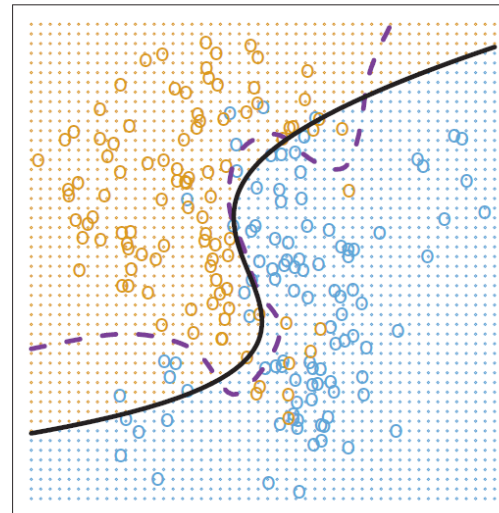
Error = 0.160

Degree=3



Error = 0.161

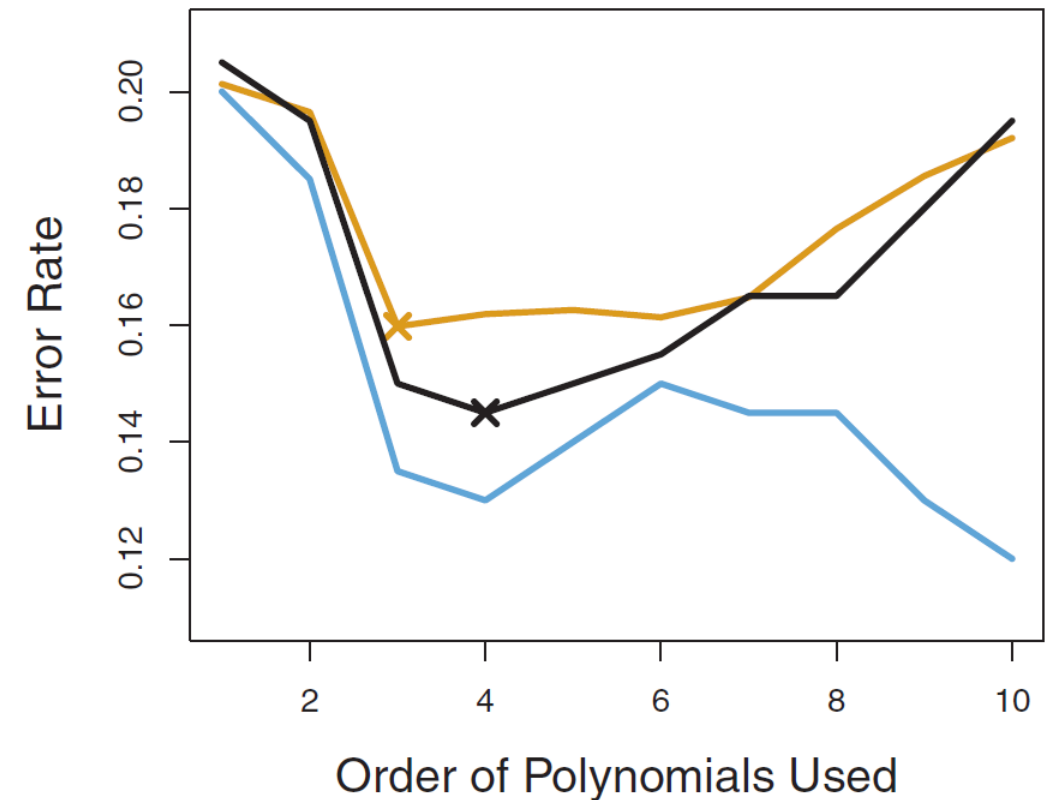
Degree=4



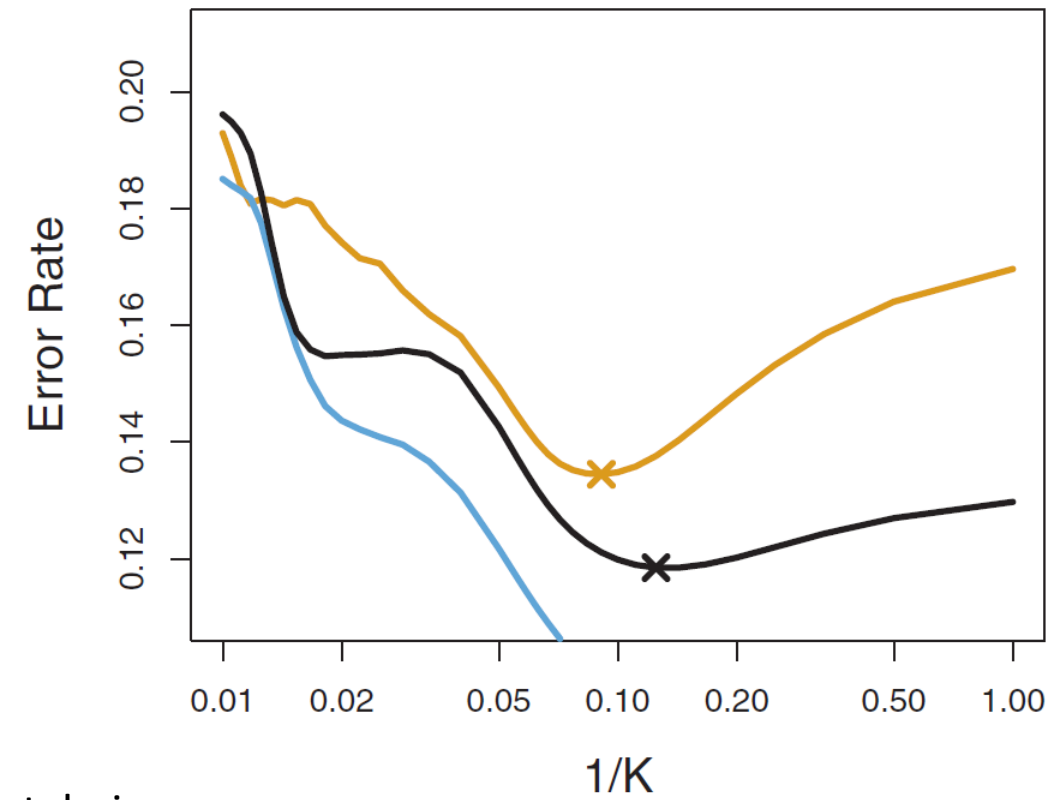


Test versus Cross Validation versus Train Error

Error rates for simulated data from previous slide



Logistic regression



K Nearest Neighbor

Test: beige
CV: black
Train: blue

Minimizing the Variance of a Sum of Possibly Correlated Variables

$$\text{Var}(\alpha X + (1-\alpha)Y) = \text{Var}(\alpha X) + \text{Var}((1-\alpha)Y) + 2\text{Cov}(\alpha X, (1-\alpha)Y)$$

$$= \alpha^2 \text{Var}(X) + (1-\alpha)^2 \text{Var}(Y) + 2\alpha(1-\alpha)\text{Cov}(X, Y)$$

... so the gradient is ...

$$\frac{\partial}{\partial \alpha} \left(\alpha^2 \text{Var}(X) + (1-\alpha)^2 \text{Var}(Y) + 2\alpha(1-\alpha)\text{Cov}(X, Y) \right) = 2\alpha \text{Var}(X) + 2(1-\alpha)(-1)\text{Var}(Y) + (2-4\alpha)\text{Cov}(X, Y)$$

$$= 2\alpha \text{Var}(X) - (2-2\alpha)\text{Var}(Y) + (2-4\alpha)\text{Cov}(X, Y)$$

... solving for α when gradient equals 0 ...

$$2\alpha \text{Var}(X) - (2-2\alpha)\text{Var}(Y) + (2-4\alpha)\text{Cov}(X, Y) = 0$$

$$2\alpha \text{Var}(X) - 2\text{Var}(Y) + 2\alpha \text{Var}(Y) + 2\text{Cov}(X, Y) - 4\alpha \text{Cov}(X, Y) = 0$$

$$2\alpha \text{Var}(X) + 2\alpha \text{Var}(Y) - 4\alpha \text{Cov}(X, Y) = 2\text{Var}(Y) - 2\text{Cov}(X, Y)$$

$$\alpha \text{Var}(X) + \alpha \text{Var}(Y) - 2\alpha \text{Cov}(X, Y) = \text{Var}(Y) - \text{Cov}(X, Y)$$

$$\alpha (\text{Var}(X) + \text{Var}(Y) - 2\text{Cov}(X, Y)) = \text{Var}(Y) - \text{Cov}(X, Y)$$

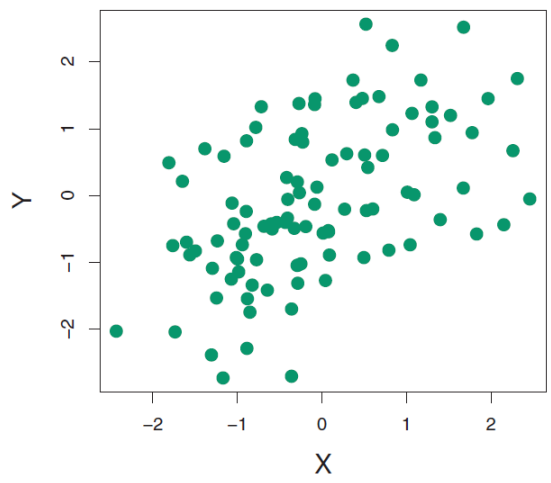
$$\alpha = \frac{\text{Var}(Y) - \text{Cov}(X, Y)}{\text{Var}(X) + \text{Var}(Y) - 2\text{Cov}(X, Y)}$$



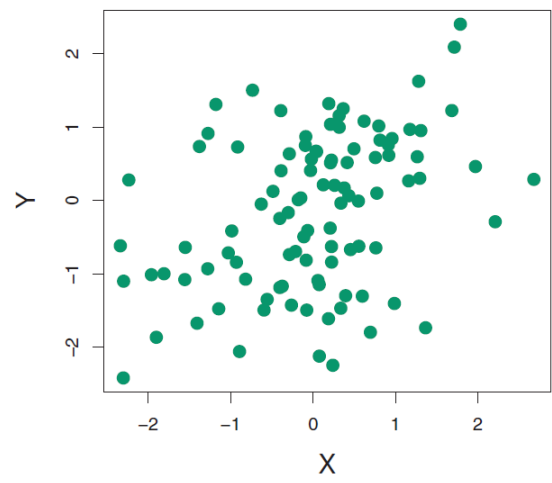
Simulated Data Sets

α is the proportion of money to be invested in X to minimize the variance of the return (risk)

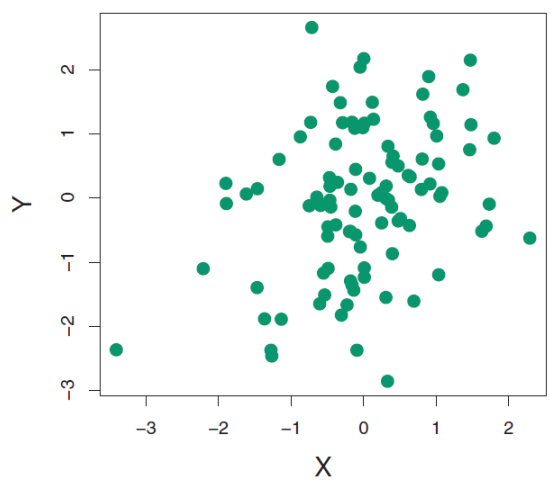
$\alpha = 0.576$



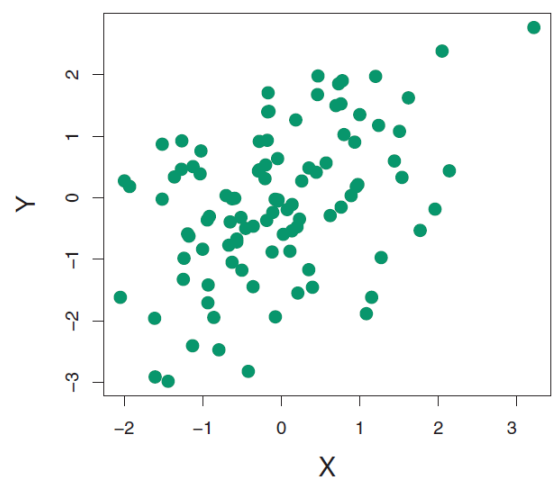
$\alpha = 0.532$



$\alpha = 0.657$



$\alpha = 0.651$



Mean and Standard Deviation of α Estimates for 1000 Samples of 100 Simulated Returns

$$\bar{\alpha} = \frac{1}{1,000} \sum_{r=1}^{1,000} \hat{\alpha}_r = 0.5996$$

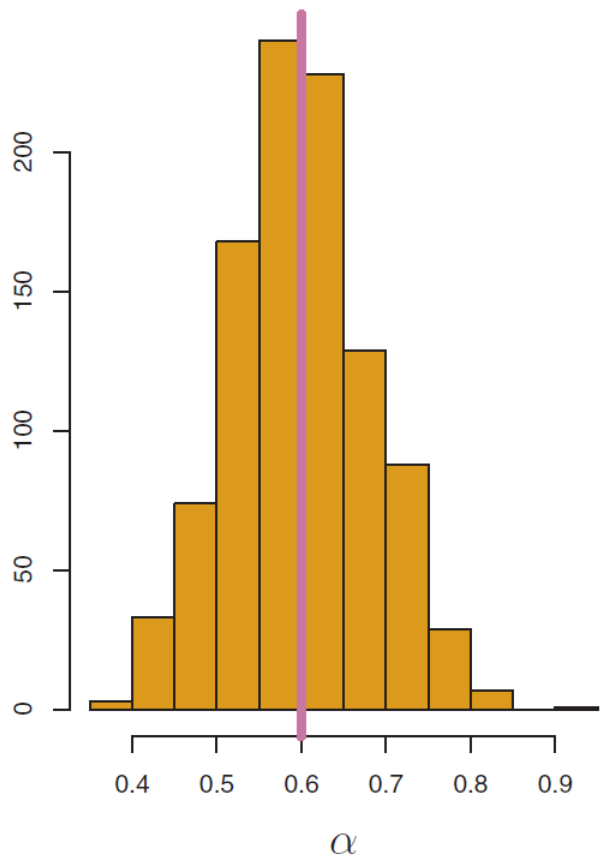
$$\sqrt{\frac{1}{1,000 - 1} \sum_{r=1}^{1,000} (\hat{\alpha}_r - \bar{\alpha})^2} = 0.083$$

Unfortunately, they did *not* share the parameters used to generate the samples. ☹️

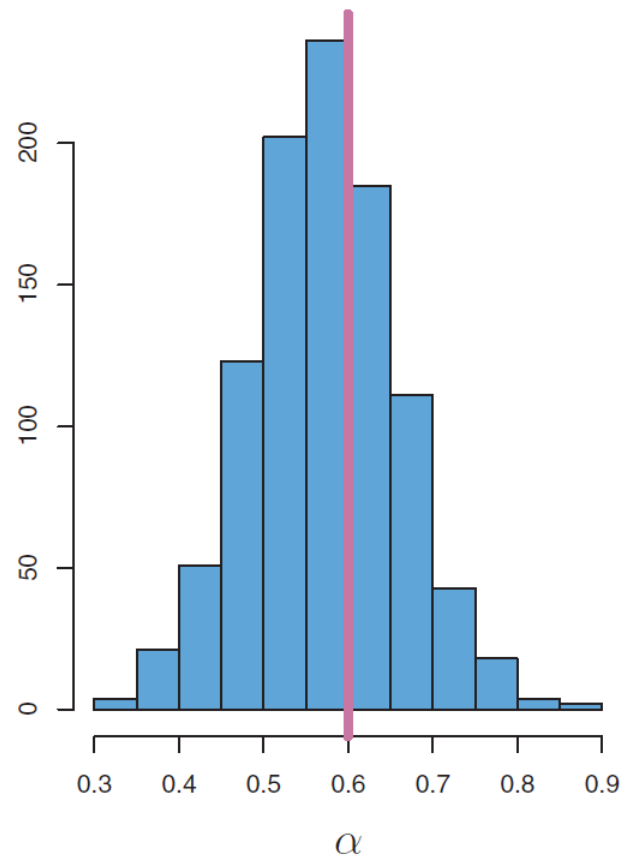
I say we grab torches and pitch forks, then head over to the statistics department. Who's with me? 😊

Maybe they're just trying to prep us: we can't handle the truth [we won't know the truth for real data].

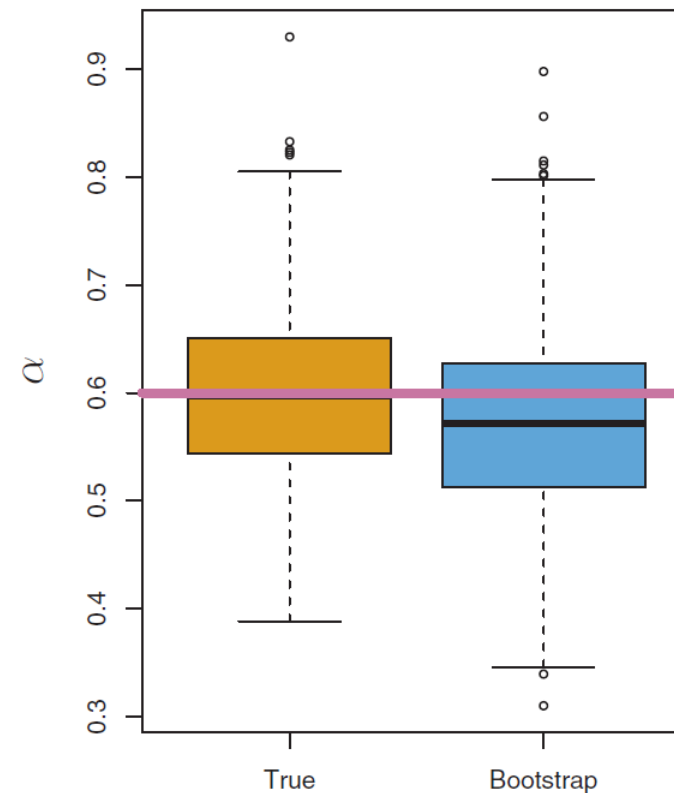
Simulated Values versus Bootstrap Values



Simulated Data
(based on True parameters)



Bootstrap Samples



The Bootstrap: Quantifying Uncertainty via Standard Error

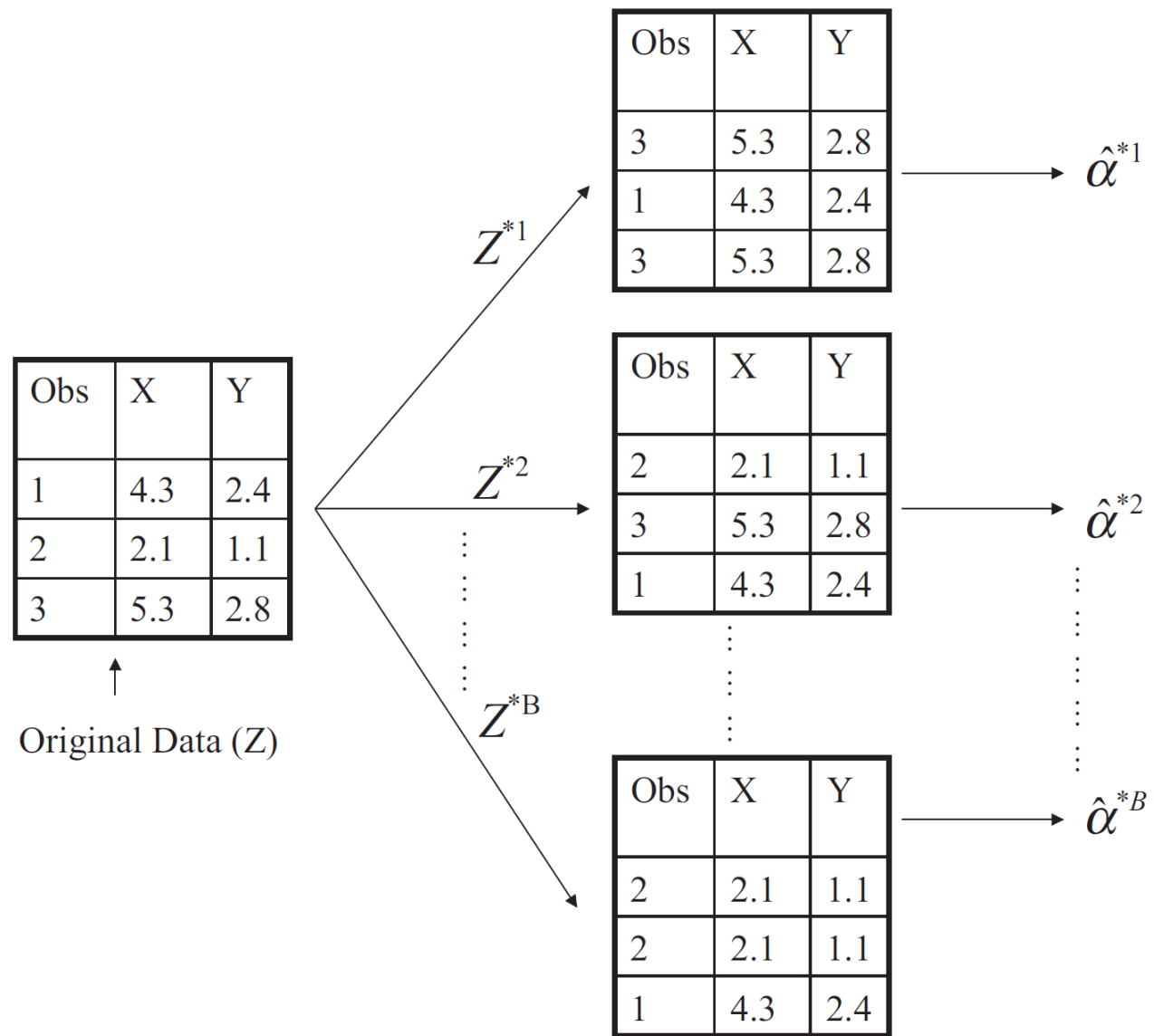
Notice that the “Standard Error” is simply the standard deviation of the alpha values from the bootstrap samples!

$$SE_B(\hat{\alpha}) = \sqrt{\frac{1}{B-1} \sum_{r=1}^B \left(\hat{\alpha}^{*r} - \frac{1}{B} \sum_{r'=1}^B \hat{\alpha}^{*r'} \right)^2}$$



Bootstrap Sample

A bootstrap sample of a data set of 'n' observations is created by drawing 'n' random samples from the data set *with* replacement



Yet Another Example

How do we compute a 95% confidence interval for correlation?

One way ...

```
library(MASS)
X = mvrnorm(1000, mu = c(0, 0), Sigma = matrix(c(9, -4.5, -4.5, 4), nrow=2))

Pearson.Correlation.Confidence.Interval = function(vector1, vector2, confidence = 0.95) {
  z = qnorm(1 - (1 - confidence) / 2)
  n = length(vector1)
  r = cov(vector1, vector2) / (sd(vector1) * sd(vector2))
  return(tanh(atanh(r) + z * c(-1, 0, 1) * sqrt(1 / (n - 3))))
}
```

Is there an easier way?

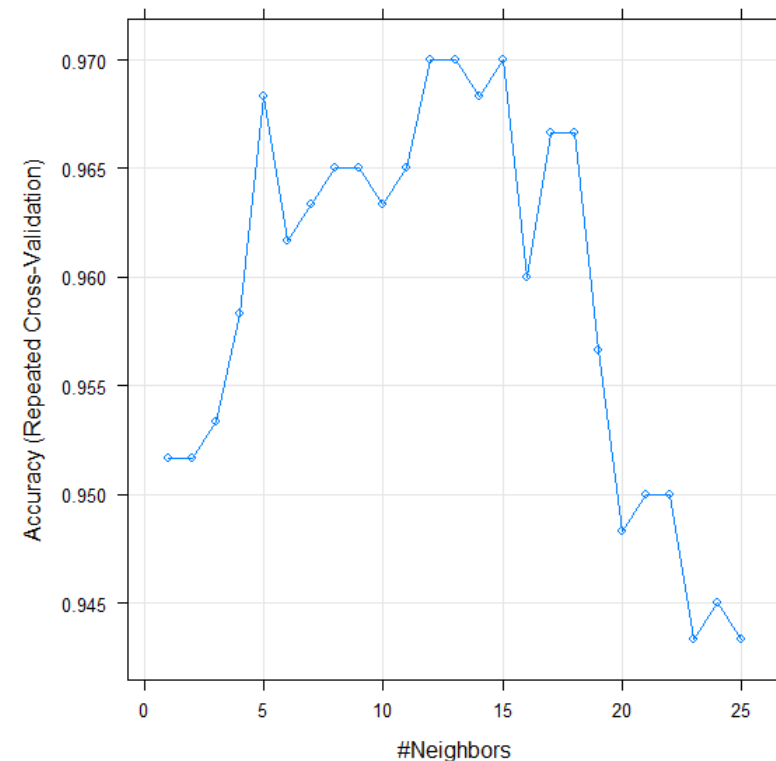
[maybe you need to understand uncertainty for some business metric]

```
library(boot)
Pearson.Correlation = function(X, index) { return(cov(X[index,1], X[index,2]) / (sd(X[index,1]) * sd(X[index,2]))) }
boot.ci(boot(X, statistic = Pearson.Correlation, R = 10000), conf = 0.95, type="bca")
```



Repeated 5 Fold Cross Validation via CARET

```
set.seed(2^17-1)
library(caret)
library(class)
input = iris # flower classification
summary(input)
indices = tapply(1:nrow(input), input[,5], sample)
trn = rbind(input[indices$setosa[1:40],],
            input[indices$versicolor[1:40],],
            input[indices$virginica[1:40],])
tst = rbind(input[indices$setosa[41:50],],
            input[indices$versicolor[41:50],],
            input[indices$virginica[41:50],])
selection = train(trn[,1:4], trn[,5], method = "knn", metric = "Accuracy", maximize = T,
                 trControl = trainControl(method = "repeatedcv", number = 5, repeats = 5),
                 tuneGrid = data.frame(k = 1:25))
table(tst[,5], knn(trn[,1:4], tst[,1:4], trn[,5], k = selection$bestTune))
plot(selection)
```





Agenda

	5 Resampling Methods	175
	5.1 Cross-Validation	176
	5.1.1 The Validation Set Approach	176
	5.1.2 Leave-One-Out Cross-Validation	178
Discriminant Analysis (from last week)	5.1.3 k -Fold Cross-Validation	181
	5.1.4 Bias-Variance Trade-Off for k -Fold Cross-Validation	183
Resampling Methods	5.1.5 Cross-Validation on Classification Problems	184
Hands-On Labs (including caret)	5.2 The Bootstrap	187
	5.3 Lab: Cross-Validation and the Bootstrap	190
	5.3.1 The Validation Set Approach	191
	5.3.2 Leave-One-Out Cross-Validation	192
	5.3.3 k -Fold Cross-Validation	193
	5.3.4 The Bootstrap	194
	5.4 Exercises	197