

# Ch 5.1.1-2: Leave One Out Cross-validation

## Lecture 10 - CMSE 381

Prof. Elizabeth Munch

Michigan State University

::

Dept of Computational Mathematics, Science & Engineering

Mon, Oct 3, 2022

## Announcements:

- Fourth homework due next monday
- Office hours
- Drops
- Grade conversion

## Last time:

- Exam

Percent	Convert
≥ 90%	4.0
≥ 85%	3.5
≥ 80%	3
≥ 75%	2.5
≥ 70%	2
≥ 65%	1.5
≥ 60%	1
< 60%	0

# Covered in this lecture

- LOO CV
- Outliers
- Leverage statistic

## Section 1

Validation set

# What's the problem?

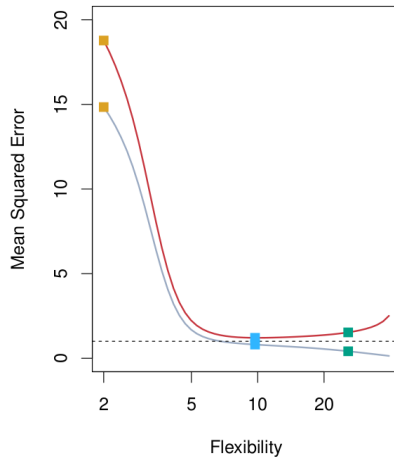
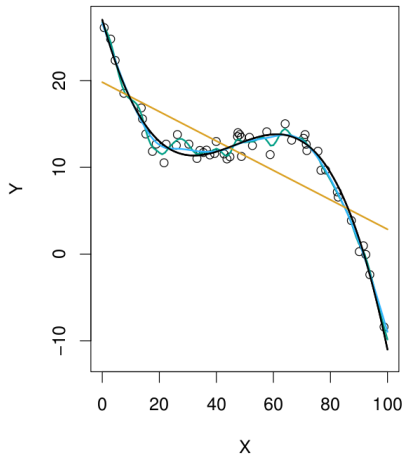
- How well is my ML method doing? *Model Assessment*
- Which method is best for our data?
- How many features should I use? Which ones? *Model selection*
- What is the uncertainty in the learned parameters?

# Training Error vs Testing Error

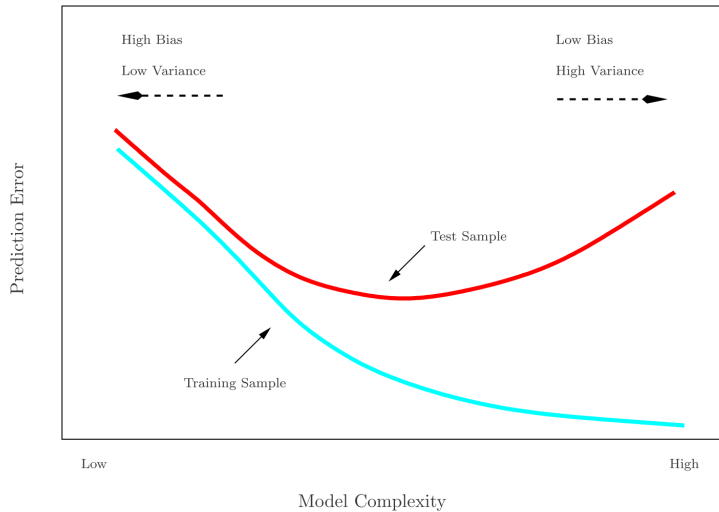
**Training Error**

**Testing Error**

# Throw-back Wednesday

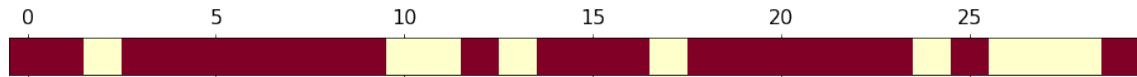


# Model tradeoffs





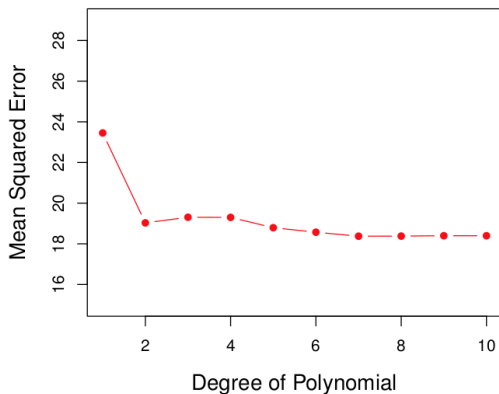
## Validation set approach



- Divide randomly into two parts:
  - ▶ Training set
  - ▶ Validation/Hold-out/Testing set
- Fit model on training set
- Use fitted model to predict response for observations in the test set
- Evaluate quality (e.g. MSE)

# Coding example in jupyter notebook

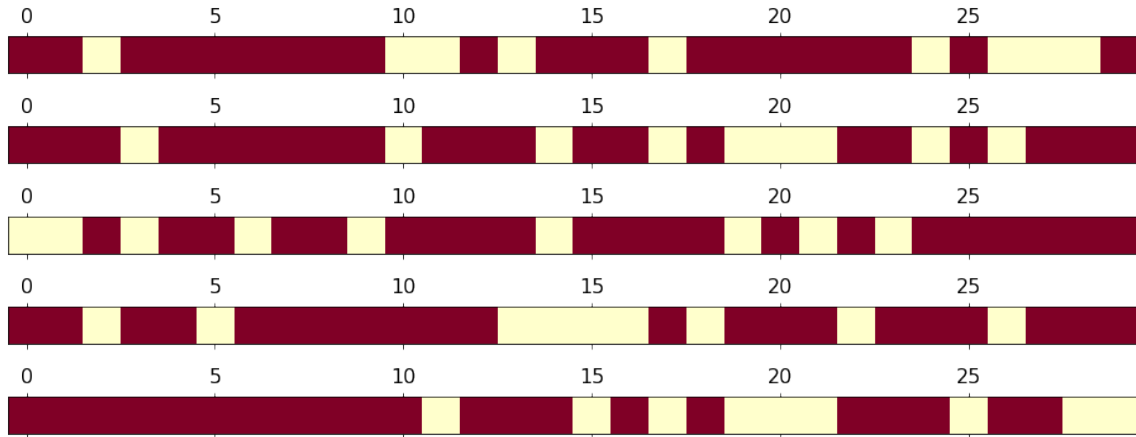
## Example with the auto data



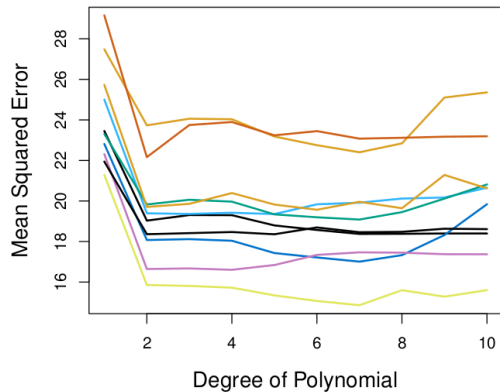
Predicting mpg using horsepower:

$$\text{mpg} = \beta_0 + \beta_1 \text{hp} + \beta_2 \text{hp}^2 + \cdots + \beta_p \text{hp}^p$$

## Rinse and repeat



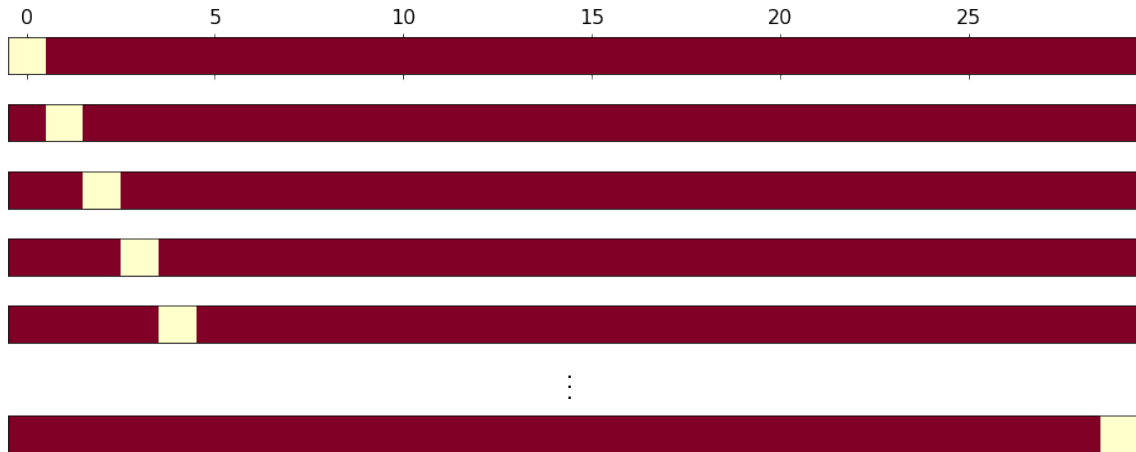
## Again example with auto data



## Section 2

### Leave-One-Out Cross-Validation (LOOCV)

# The idea



# The idea in mathy words

- Remove  $(x_1, y_1)$  for testing.
- Train the model on  $n - 1$  points:  
 $\{(x_2, y_2), \dots, (x_n, y_n)\}$
- Calculate  $\text{MSE}_1 = (y_1 - \hat{y}_1^2)$

Return the score:

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

- Remove  $(x_2, y_2)$  for testing.
- Train the model on  $n - 1$  points:  
 $\{(x_1, y_1), (x_3, y_3), \dots, (x_n, y_n)\}$
- Calculate  $\text{MSE}_2 = (y_2 - \hat{y}_2^2)$
  
- Rinse and repeat



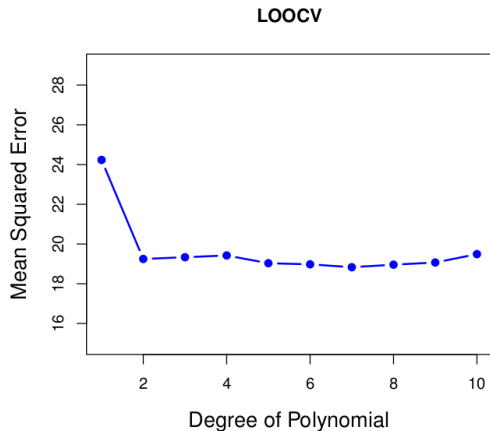
Do the LOOCV coding section

# LOOCV Pros and Cons

**Advantages:**

**Disadvantages:**

## Again example with auto data

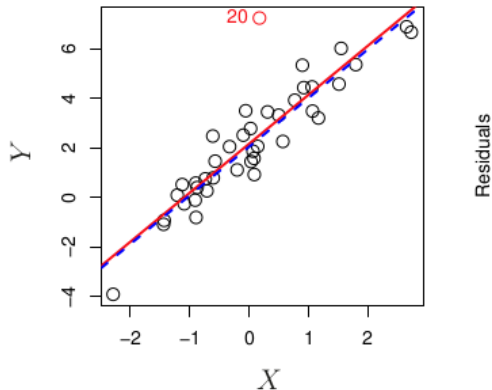


## Section 3

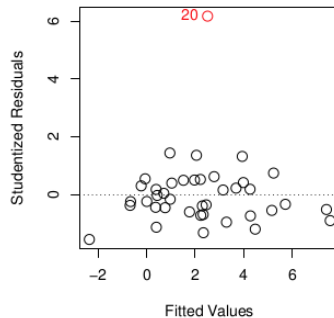
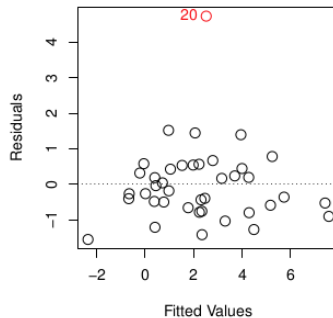
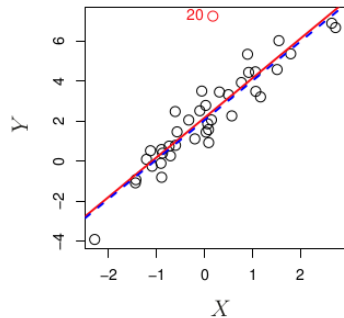
The one time you can cheat (by not computing every model fit)

# Outliers

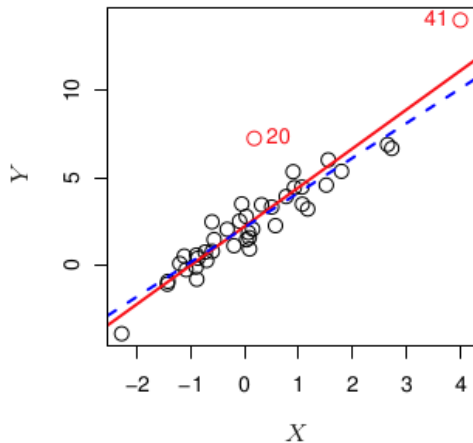
An *outlier* is a point for which  $y_i$  is far from the value predicted by the model.



# Residuals

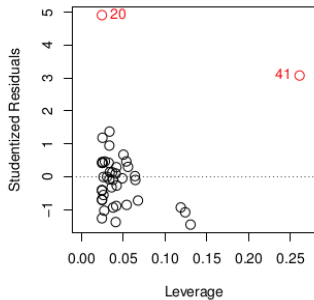
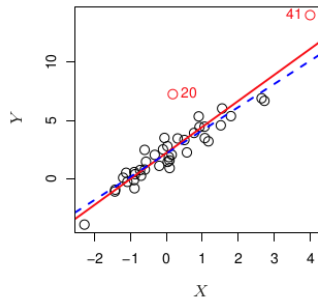


# High Leverage



Observations with *high leverage* have an unusual value for  $x_i$ .

# Leverage statistic



Version for  $p = 1$

$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2}$$



## Leverage statistic properties

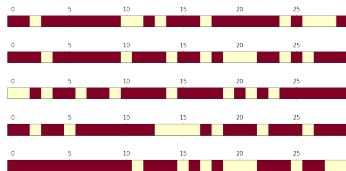
$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2}$$

# Speeding up LOOCV

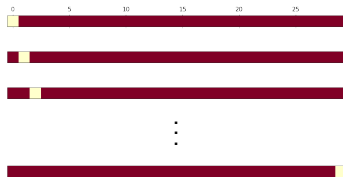
**Warning:** This only works for least squares linear or polynomial regression.

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

## Validation set



## LOO-CV



## LOO-CV Score

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

## Cheap trick for regression

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

# Next time

14	M	Oct 3	Leave one out CV	5.1.1, 5.1.2	
15	W	Oct 5	k-fold CV	5.1.3	
16	F	Oct 7	More k-fold CV	5.1.4	
17	M	Oct 10	CV for classification	5.1.5	HW #4 Due
18	W	Oct 12	Resampling methods: Bootstrap	5.2	
19	F	Oct 14	Subset selection	6.1	
20	M	Oct 17	Shrinkage: Ridge	6.2.1	HW #5 Due
21	W	Oct 19	Shrinkage: Lasso	6.2.2	
22	F	Oct 21	Dimension Reduction	6.3	
	M	Oct 24	No class - Fall break		
21	W	Oct 26	More dimension reduction; High dimensions	6.4	
22	F	Oct 28	Polynomial & Step Functions.	7.1,7.2	HW #6 Due
23	M	Oct 31	Review		
24	W	Nov 2	<b>Midterm #2</b>		