

Ch 5.1.3-4: k -Fold Cross-Validation

Lecture 13 - CMSE 381

Prof. Elizabeth Munch

Michigan State University

::

Dept of Computational Mathematics, Science & Engineering

Wed, Oct 2, 2024

Announcements

Last time:

- Validation Set
- LOOCV

Announcements:

- Exam 1 grades
- HW #4 Posted.
 - ▶ Changed Deadline! Due Wednesday Oct 9.

Lec #	Date			Reading	HW
12	Mon	9/30	Leave one out CV	5.1.1, 5.1.2	
13	Wed	10/2	k-fold CV	5.1.3	
14	Fri	10/4	More k-fold CV,	5.1.4-5	
15	Mon	10/7	k-fold CV for classification	5.1.5	
16	Wed	10/9	Resampling methods: Bootstrap	5.2	HW #4 Due Weds 10/9
17	Fri	10/11	Subset selection	6.1	
18	Mon	10/14	Shrinkage: Ridge	6.2.1	
19	Wed	10/16	Shrinkage: Lasso	6.2.2	
20	Fri	10/18	Dimension Reduction	6.3	HW #5 Due Fri 10/18
	Mon	10/21	No class - Fall break		
	Wed	10/23	Review		
	Fri	10/25	Midterm #2		

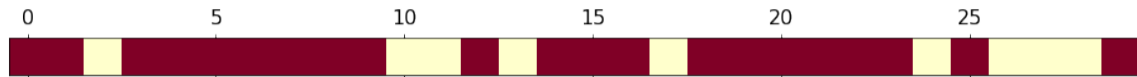
Covered in this lecture

- k -fold CV

Section 1

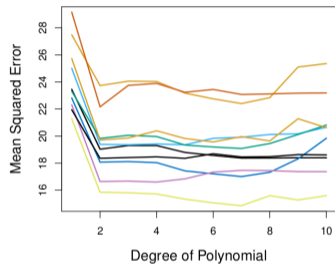
Last time

Validation set approach

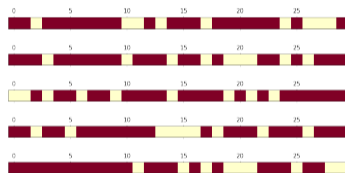


- Divide randomly into two parts:
 - ▶ Training set
 - ▶ Validation/Hold-out/Testing set
- Fit model on training set
- Use fitted model to predict response for observations in the test set
- Evaluate quality (e.g. MSE)

Problems



Ex. Predict mpg using horsepower



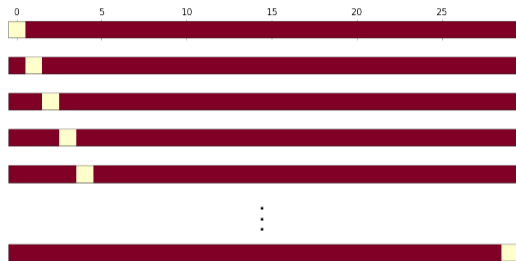
- Highly variable results, no consensus about the error
- Tends to overestimate test error rate

Leave One Out CV (LOOCV)

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

- 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 $\text{MSE}_2 = (y_2 - \hat{y}_2)^2$

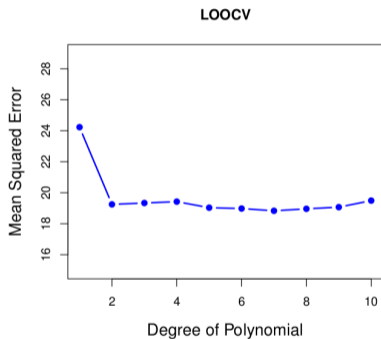
- Rinse and repeat



Return the score:

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

Pros and Cons

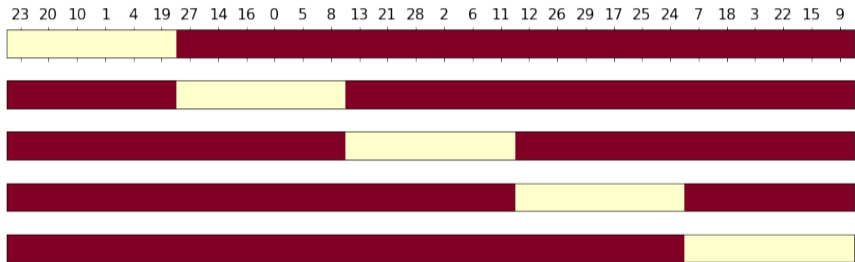


- No variance
- Higher computation cost

Section 2

k-Fold CV

The idea



Mathy version

- Randomly split data into k -groups (folds)
- Approximately equal sized. For the sake of notation, say each set has ℓ points
- Remove i th fold U_i and reserve for testing.
- Train the model on remaining points
- Calculate
$$\text{MSE}_i = \frac{1}{\ell} \sum_{(x_j, y_j) \in U_i} (y_j - \hat{y}_j)^2$$
- Rinse and repeat

Return

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

By hand first!

There are 10 students in the class, and we have data points for each. They have already been randomly permuted below. Write down the training/testing sets for a 3-fold CV

- Damien
- Alice
- Greta
- Jasmin
- Benji
- Inigo
- Firas
- Carina
- Enrique
- Hubert

Fold 1

Fold 2

Fold 3

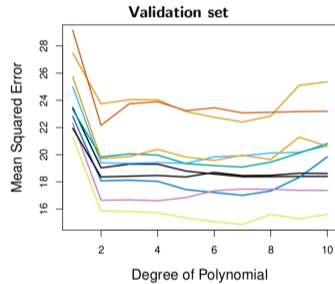
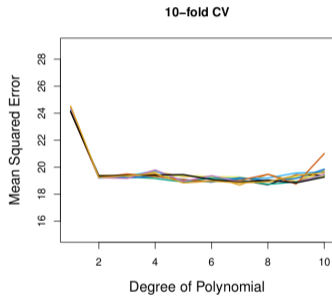
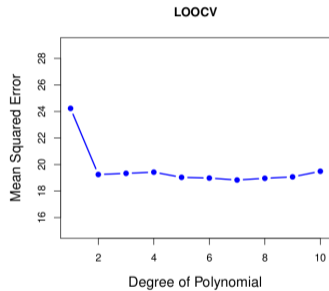
Coding - Building k -fold CV

Pros and Cons

Pros:

Cons:

Comparison



Next time

Lec #	Date		Reading	HW	
12	Mon	9/30	Leave one out CV	5.1.1, 5.1.2	
13	Wed	10/2	k-fold CV	5.1.3	
14	Fri	10/4	More k-fold CV,	5.1.4-5	
15	Mon	10/7	k-fold CV for classification	5.1.5	
16	Wed	10/9	Resampling methods: Bootstrap	5.2	HW #4 Due Weds 10/9
17	Fri	10/11	Subset selection	6.1	
18	Mon	10/14	Shrinkage: Ridge	6.2.1	
19	Wed	10/16	Shrinkage: Lasso	6.2.2	
20	Fri	10/18	Dimension Reduction	6.3	HW #5 Due Fri 10/18
	Mon	10/21	No class - Fall break		
	Wed	10/23	Review		
	Fri	10/25	Midterm #2		