# Ch 5.1.3-4: *k*-Fold Cross-Validation
## Lecture 13 - CMSE 381

Prof. Lianzhang Bao

Michigan State University
::
Dept of Computational Mathematics, Science & Engineering

Mon, Feb 17, 2025

**Last time:**

- Validation Set
- LOOCV

**Announcements:**

- Exam 1 grades.... hopefully soon
- HW #4 will be posted soon.
  - Due Sunday 3/2.

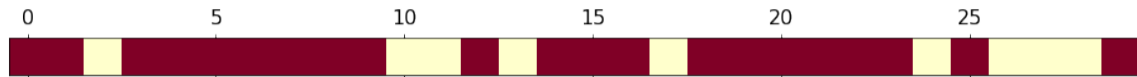| 12 | F | 2/14 | Leave one out CV | 5.1.1, 5.1.2 | | |
|---|---|---|---|---|---|---|
| 13 | M | 2/17 | k-fold CV | 5.1.3 | | |
| 14 | W | 2/19 | More k-fold CV | 5.1.4-5 | | Q5 |
| 15 | F | 2/21 | k-fold CV for classification | 5.1.5 | | |
| 16 | M | 2/24 | Subset selection | 6.1 | | |
| 17 | W | 2/26 | Shrinkage: Ridge | 6.2.1 | | |
| 18 | F | 2/28 | Shrinkage: Lasso | 6.2.2 | HW #4 Due Sun 3/2 | |
| | M | 3/3 | Spring Break | | | |
| | W | 3/5 | Spring Break | | | |
| | F | 3/7 | Spring Break | | | |
| 19 | M | 3/10 | PCA | 6.3 | | |
| 20 | W | 3/12 | PCR | 6.3 | | Q6 |
| | F | 3/14 | *Review* | | HW #5 Due Sun 3/16 | |
| | M | 3/17 | **Midterm #2** | | | |

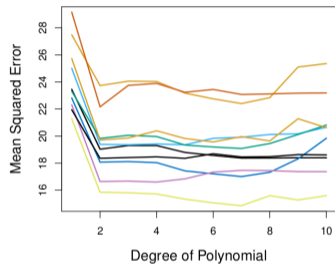# Covered in this lecture

- $k$-fold CV

# Section 1

## Last time

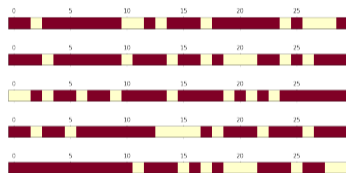# Validation set approach



- 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)

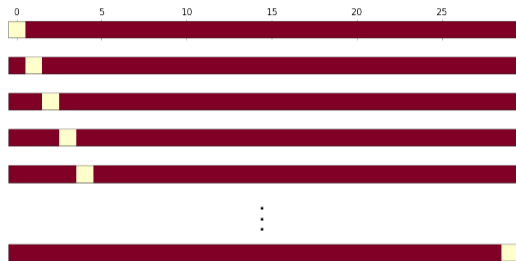Ex. Predict `mpg` using
`horsepower`

- 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: $\{(x_2, y_2), \cdots, (x_n, y_n)\}$
- Calculate $\mathrm{MSE}_1 = (y_1 - \hat{y}_1)^2$

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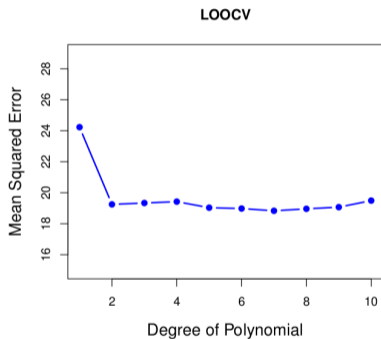- Calculate $\mathrm{MSE}_2 = (y_2 - \hat{y}_2)^2$

- Rinse and repeat



Return the score:

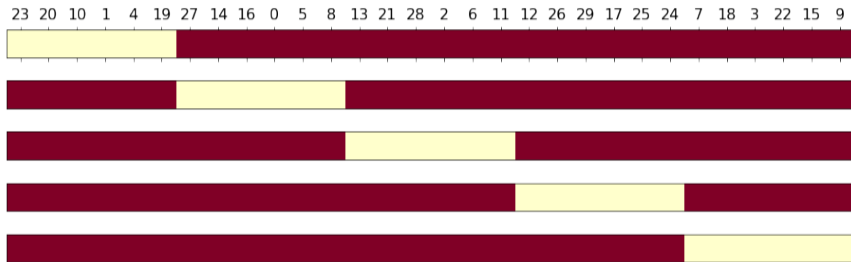$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^{n} \mathrm{MSE}_i$$

# Pros and Cons



- No variance
- Higher computation cost

Section 2

$k$-Fold CV

# The idea



23 20 10 1 4 19 27 14 16 0 5 8 13 21 28 2 6 11 12 26 29 17 25 24 7 18 3 22 15 9

## Mathy version

- Randomly split data into $k$-groups (folds)

- Approximately equal sized. For the sake of notation, say each set has $\ell$ points

- Remove $i$th fold $U_i$ and reserve for testing.

- Train the model on remaining points

- Calculate
  $\mathrm{MSE}_i = \frac{1}{\ell} \sum_{(x_j, y_j) \in U_i} (y_j - \hat{y}_j)^2$

- Rinse and repeat

Return

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

## 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

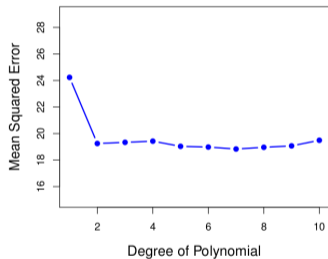| | **Fold 1** | **Fold 2** | **Fold 3** |
|---|---|---|---|
| Damien | | | |
| Alice | | | |
| Greta | | | |
| Jasmin | | | |
| Benji | | | |
| Inigo | | | |
| Firas | | | |
| Carina | | | |
| Enrique | | | |
| Hubert | | | |

# Pros and Cons

**Pros:**

**Cons:**

# Comparison

# Next time

| | | | | | | |
|---|---|---|---|---|---|---|
| 12 | F | 2/14 | Leave one out CV | 5.1.1, 5.1.2 | | |
| 13 | M | 2/17 | k-fold CV | 5.1.3 | | |
| 14 | W | 2/19 | More k-fold CV | 5.1.4-5 | | Q5 |
| 15 | F | 2/21 | k-fold CV for classification | 5.1.5 | | |
| 16 | M | 2/24 | Subset selection | 6.1 | | |
| 17 | W | 2/26 | Shrinkage: Ridge | 6.2.1 | | |
| 18 | F | 2/28 | Shrinkage: Lasso | 6.2.2 | HW #4 Due Sun 3/2 | |
| | M | 3/3 | Spring Break | | | |
| | W | 3/5 | Spring Break | | | |
| | F | 3/7 | Spring Break | | | |
| 19 | M | 3/10 | PCA | 6.3 | | |
| 20 | W | 3/12 | PCR | 6.3 | | Q6 |
| | F | 3/14 | *Review* | | HW #5 Due Sun 3/16 | |
| | M | 3/17 | **Midterm #2** | | | |