

Cross-Validation

Wednesday, September 22, 2021 6:07 PM

- Hold-Out / Validation
 - Training Error vs. Test Error \leftarrow out of sample

$$E[e_{oos}] = ?$$

$$e_{\text{Train}} \\ e_{\text{Test-set}} \quad (\text{Test/Train Split}) \quad [\text{e.g., } 80/20] \\ e_{\text{CV}}$$

- Test/Train Split
 - Holdout set: Randomly dividing available sample data (Validation) into training and testing sets

* Stratified Sampling \rightarrow Equal ratios of classes/ranges

$$E[e_{oos}] = e_{\text{Test}}$$

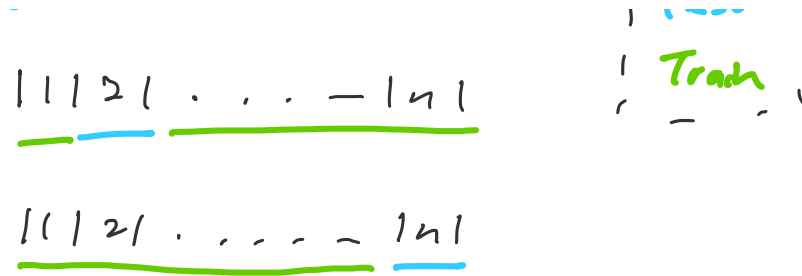
* Issues:

1. High variability in Test Set on random splits
2. We end up with smaller Training set
 - \hookrightarrow Higher Error Rate in Training (Overestimate)

- Cross-Validation

- LOOCV

$$\underbrace{1 \mid 1 \mid 2 \mid \dots \mid n \mid}_{\text{Test ?}}$$



* Distribution estimate for $e_{oos} \rightarrow E[e_{oos}] = \mu_{LOOCV}$
 median
 ...

$$\therefore CV(n) = \frac{1}{n} \sum_{i=1}^n MSE_i$$

$$MSE_i = (y_i - \hat{y}_i)^2$$

↳ Less bias \Rightarrow we fit on $n-1$ points!

↳ No randomness on test-train split since all n points are used as test sets

• K-Fold CV

- k th fold is used as test
- $1 \dots k-1$ fold is used as training

$$CV(k) = \frac{1}{k} \sum_{i=1}^k MSE_i$$

* LOOCV is k -fold with $k=n$

* Bias-Variance Issues in k -Fold/LOOCV

- We are taking mean error of k models to estimate e_{oos}

all k models as $k \rightarrow n$ get more correlated!

- Mean of many correlated values has a higher variance than the sets of many uncorrelated values!