

Cross Validation

DS 301

Iowa State University

Cross-validation

An alternative to the approaches we discussed is to directly estimate the test error using cross-validation.

CMSE7

c1) validation set approach

data



train

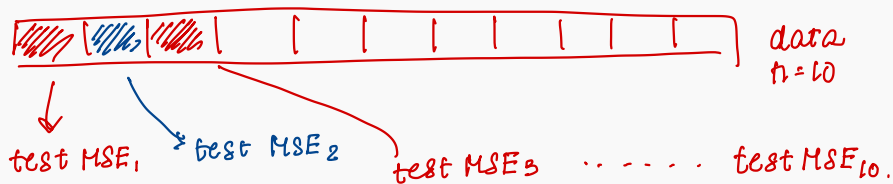
\hat{f}

test

$$\text{test MSE: } \frac{\sum_{i=1}^n (\hat{f}(x_i) - y_i)^2}{n}$$

Leave-one-out cross-validation (LOOCV)

- leave out a single observation
↳ test set.
- $(n-1)$ observations remaining → training set.



⇒ average of these:

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

- utilize as much data as possible
- every observation gets the chance to be the test set.
- deterministic.

Drawbacks of LOOCV

- computationally expensive
 - ↳ need to fit model n times.
- theoretical problems w/ LOOCV.

k-fold cross-validation

data

k_1
k_2
k_3
k_4
k_5

$k=5$

each fold gets the chance to be the test set.

- $k_2, k_3, k_4, k_5 \rightarrow$ training set
 $k_1 \rightarrow$ test set.

test MSE_{k_1}

- $k_1, k_3, k_4, k_5 \rightarrow$ training set
 $k_2 \rightarrow$ test set.

best MSE_{k_2}

\vdots

test MSE_{k_5}

$$CV_{k=5} = \frac{1}{5} \sum_{i=1}^k \text{best } MSE_{k_i}$$

$k=5, k=10, \dots$

See R script: `cross_validation.R` and `cv_subset.R`