

- Recall  $y = f(x) + \varepsilon$ ,  $E(\varepsilon) = 0 \Rightarrow E(y|x) = f(x)$ .  
We construct/train an estimator  $\hat{f}$  for  $f$ , based on  $n$  iid observations  $(x_i, y_i)$ .  
How well does  $\hat{f}$  perform?

- We want to know the expected test MSE  $\theta = E[(y_{\text{new}} - \hat{f}(x_{\text{new}}))^2]$

(Types of estimators for validation)

- ① We want to estimate  $\theta$  by  $\hat{\theta}_{\text{val}} = (\text{of validation})$   $\left\{ \begin{array}{l} \text{training data } x_{\text{new}}, y_{\text{new}} \\ \text{tends to be too large/pessimistic on average (biased).} \end{array} \right.$   
VALIDATION as arithmetic mean of all:  $\hat{\theta}_{\text{val}} = \frac{1}{|\text{Test set}|} \sum_{i \in \text{Test set}} (y_i - \hat{f}(x_i))^2$ , where  $\hat{f}$  was trained on the training set.

- ② LOOCV:  $\hat{\theta}_{\text{loocv}} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}^{(-i)}(x_i))^2$  - less biased.
- ③ Kfold CV:  $\hat{\theta}_{\text{Kfold CV}} = \frac{1}{K} \sum_{k=1}^K \left[ \frac{1}{|F_k|} \sum_{i \in F_k} (y_i - \hat{f}_k^{(-F_k)}(x_i))^2 \right]$  which has less variance? (debatable).  
$$= \frac{1}{K} \sum_{k=1}^K \text{MSE}_{\text{Kfold}}$$

- Motivation ex. gene expression data.

$n \ll p$ , where there are many genes ( $p$ ) and few observations ( $n$ ).  
observations ex. binary response variable indicating disease status.

- Goal: Predict disease status based on gene expression

- Consider a KNN estimator. If response variables are binary, then you can just do a majority vote of neighboring  $y$ 's.  
ex. 1-NN classifier  $\hat{f}(x_0) = y_0 = y_i$  for  $i$  s.t.  $x_i$  is closest to  $x_0$ .

Preprocessing and preselection. We take 20 genes with largest marginal correlation with  $y$ .  
Do 1-NN classification based on these 20 genes (ie fit a classifier)  
Do CV to evaluate performance. For any  $\hat{f}$ , we know

$$\theta = E(\text{test error rate of } \hat{f}) = \frac{1}{2} \quad (\text{ie if } \exists \text{ no info on } x, y, \text{ then it is random chance})$$

- Rcode line 151: `stopifnot()` to check data  
38: `cv.knn` function

- It is good to look @, investigate how folds are (randomly) chosen. (Need to be careful)  
Preselect features based on just the training data (?) (fold?) not the whole set!  
(Then, Apply the procedure to evaluate other data (see Rcode line 38))  
see Rcode lines 63-66 for how preselection differs for different num of folds  $K$ .

- As we choose  $K$ , assess the amount of variability and bias / center pt of the:  
CV error rate (should be at / close to 0.5).

→ output of line 116 (is incorrect procedure)

- We should choose folds once per estimate calculation. Then, choose (all folds again) for another estimate.   
→ ie choose folds for an estimate

- We want to have the distributions of CV error rates and that they are unbiased; and get info on variance. If each fold is like a different observation, recalling  $\text{MSE}_k$ ,  $\text{Var}(\hat{\theta}_{\text{Kfold CV}}) = \text{Var}\left(\frac{1}{K} \sum_{k=1}^K \text{MSE}_k\right)$

$$\Rightarrow \text{Then } \text{Var}(\hat{\theta}_{\text{Kfold CV}}) = \frac{1}{K^2} \sum_{k=1}^K \text{Var}(\text{MSE}_k) = \frac{1}{K} \text{Var}(\text{MSE}_k)$$

→ This is good...

This is only true if the  $K$  folds are independent