

Tentative solution to RecEx Module 5: Resampling

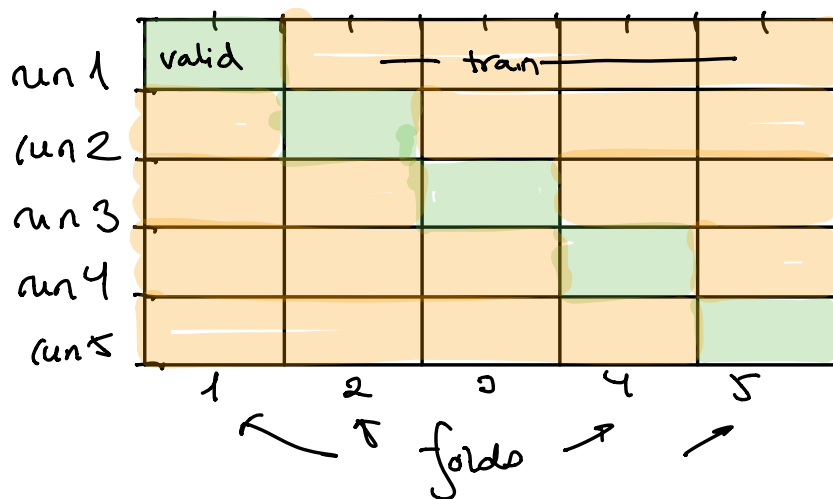
Cross-validation

1) Explain how k -fold cross-validation is implemented.

Drawing: $k=5$ (for simplicity)

First - shuffle the data: indices $[1, 2, 17, \dots, 28]$


Then partition into k groups = folds.



in general

$$\frac{n!}{k! \left[\left(\frac{n}{k} \right)! \right]^k} \text{ possible}$$

ways to do this
(multinomial with a twist)
see below

train = 

validate = 

In run 1 fold 1 is kept aside and folds 2-5 are used to train the method (maybe many times, once for every model complexity). Then error is calculated on the validation fold.

Repeat k times:

$$Q_n = \frac{1}{n} \sum_{j=1}^k \text{MSE}_j \cdot q_j$$

$$MSE_j = \frac{1}{n_j} \sum_{i \in G} (y_i - \hat{y}_i)^2 \text{ for MSE}$$

\nearrow obs in validation fold
 \nwarrow prediction of x_i in validation fold using fitted model from folds "j"
 \nwarrow not j.

other loss functions may be O/L loss.

Regression: find the ^{optimal} number of neighbors in kNN-regression.

Classification: choose between QDA or LDA in classification

n obj. delar inn i k grupper med m i kvar, $n = km$

(ant. perm. av de n) = $n!$

$$\begin{aligned}
 (\text{ant. perm. av de } n) &= (\underbrace{\text{ant. m\u00e4ter \u00e5 dela inn i k grupper}}_B) \\
 &\cdot (\underbrace{\text{ant. perm. av grupperna}}_{= k!}) \cdot (\underbrace{\text{ant. m\u00e4ter \u00e5 permutera i varje grupp}}_{m!})^k \\
 &= B k! (m!)^k
 \end{aligned}$$

$$B = \frac{n!}{k! (m!)^k} = \binom{n}{k} \frac{(n-k)!}{(m!)^k}$$

2) Advantages & disadvantages of k -fold CV relative to

a) the validation set

D: computational complexity

A: bias = generally larger sample size for each k -fold then validation set, which means "more data \rightarrow better fit" and therefore not overestimate the test set error

Bias?: compared to using the full data set for model fit.

A: different validation sets may give very different test error, so the results are variable - much more than for k -fold.

b) LOOCV \Leftarrow no randomness in splits:

A: less computational efforts for k than n , unless nice formula as for multiple linear regression.

A: less bias - in the sense that LOOCV use a larger set to fit to data ($n-1$ obs) which gives a less biased version of the test set error.

D: higher variance: [According to our textbook] we are averaging the output from n fitted models that are trained on nearly the same data \Rightarrow correlated positively. ³

This happens to a less degree with k -fold, since the k models more different data end are thus less variable.

$$W_n = ((y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + \dots + (y_n - \hat{y}_n)^2) \frac{1}{n}$$

$$\text{Var}(W_n) = \begin{array}{c} \text{sum variances} \\ \text{of terms} \end{array} + 2 \cdot \begin{array}{c} \text{covariances} \\ \text{of terms} \end{array}$$

\nearrow
 this part tend to be
 larger for LOOCV than k -fold
 when correlation higher between
 models

c) choice of k in k -fold. We just know that

$k=n = \text{LOOCV}$ (small bias - high variance) \swarrow of estimator for test
 + generally high comp. demand \swarrow set error

$k=2$ (larger bias - lower variance)
 + less comp. challenging

$\underbrace{\hspace{10em}}_{\text{simulations}}$
 and empirical research has found $k=5$ or $k=10$
 to be good choices!

3) Case (as in R-code) : classification setup with two classes

- $n=50$ observations of $p=5000$ predictors

a) choose to use only $d=25$ predictors, but choose the top d from absolute correlation coeff between the p preds. and the class label.

b) then use logistic regression with the d predictors.

⇒ How to do CV? On $\underbrace{a+b}_{\text{right}}$ or only on $\underbrace{b}_{\text{wrong}}$?

Wrong: if only b , then all data used to find the predictors → gives overoptimistic result (miscl. rate of 0% can be found)

right: both $a+b \Rightarrow$ all is good

See R-code in problem and run to see what the misclassification rate is.

Bootstrepping

1)

a) $P(\text{draw } x_i) = \frac{1}{n}$, $P(\text{not draw } x_i) = 1 - \frac{1}{n}$

b) $P(\text{not any } x_i\text{'s}) = (1 - \frac{1}{n})^n$

$P(\text{at least one } x_i) = 1 - (1 - \frac{1}{n})^n$

c) $P(x_i \text{ in boot sample}) = 1 - (1 - \frac{1}{n})^n \approx 1 - \exp(-1)$
 $= 0.632$

d) R-code to check result and see how fast
 $1 - (1 - \frac{1}{n})^n \rightarrow 0.632$ (in n).

2) Bootstrap to estimate $SD(\hat{\beta})$:

for (b in $1:B$) {

Draw with replacement from data to get

$(X, Y)_b^*$: bootstrap sample b , $b=1, \dots, B$.

fit $Y = X\beta + \varepsilon$ and keep $\hat{\beta}_b$

}

Calculate $\hat{SD}(\hat{\beta}) = \sqrt{\frac{1}{B-1} \sum_{b=1}^B \left(\hat{\beta}_b - \frac{1}{B} \sum \hat{\beta}_b \right)^2}$

Why do we want to do this - when we really know that $\hat{SD}(\hat{\beta}) = \hat{\sigma} \cdot \text{diag}((X^T X)^{-1})$?

And we might also do $\hat{Cov}(\hat{\beta})$ but then use

$$\hat{Cov}(\hat{\beta}) = \frac{1}{n-1} \sum_{b=1}^B (\hat{\beta}_b - \bar{\hat{\beta}})(\hat{\beta}_b - \bar{\hat{\beta}})^T$$

$$\bar{\hat{\beta}} = \frac{1}{B} \sum_{b=1}^B \hat{\beta}_b$$

3) is covered on page 195 of the ISL book