

Lecture 11: Cross validation

Reading: Chapter 5

STATS 202: Data mining and analysis

Sergio Bacallado
October 16, 2013

Validation set approach

Goal: Estimate the test error for a supervised learning method.

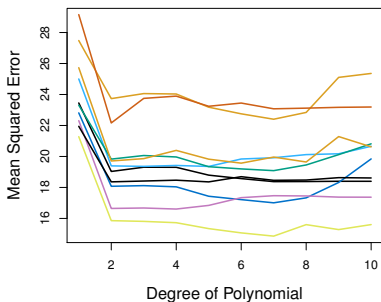
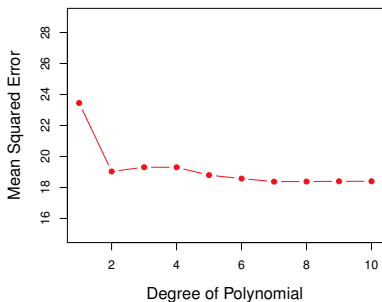
Strategy:

- ▶ Split the data in two parts.
- ▶ Train the method in the first part.
- ▶ Compute the error on the second part.



Validation set approach

Polynomial regression to estimate mpg from horsepower in the Auto data.



Problem: Every split yields a different estimate of the error.

Leave one out cross-validation

- ▶ For every $i = 1, \dots, n$:
 - ▶ train the model on every point except i ,
 - ▶ compute the test error on the held out point.
- ▶ Average the test errors.



Leave one out cross-validation

- ▶ For every $i = 1, \dots, n$:
 - ▶ train the model on every point except i ,
 - ▶ compute the test error on the held out point.
- ▶ Average the test errors.

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i^{(-i)})^2$$

Prediction for the i sample without using the i th sample.

Leave one out cross-validation

- ▶ For every $i = 1, \dots, n$:
 - ▶ train the model on every point except i ,
 - ▶ compute the test error on the held out point.
- ▶ Average the test errors.

$$\text{CV}_{(n)} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(y_i \neq \hat{y}_i^{(-i)})$$

... for a classification problem.

Leave one out cross-validation

Computing $CV_{(n)}$ can be computationally expensive, since it involves fitting the model n times.

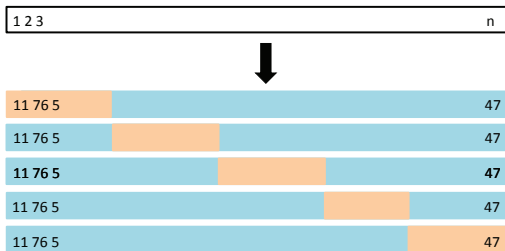
For linear regression, there is a shortcut:

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{y}_i}{1 - H_{ii}} \right)^2$$

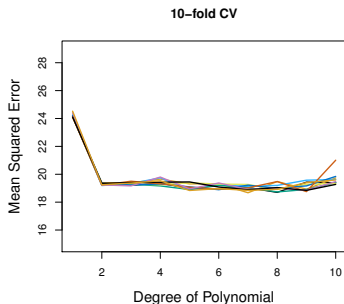
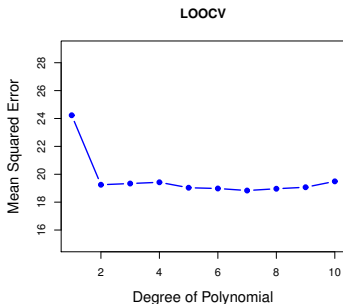
where H_{ii} is the leverage statistic.

k -fold cross-validation

- ▶ Split the data into k subsets or *folds*.
- ▶ For every $i = 1, \dots, k$:
 - ▶ train the model on every fold except the i th fold,
 - ▶ compute the test error on the i th fold.
- ▶ Average the test errors.

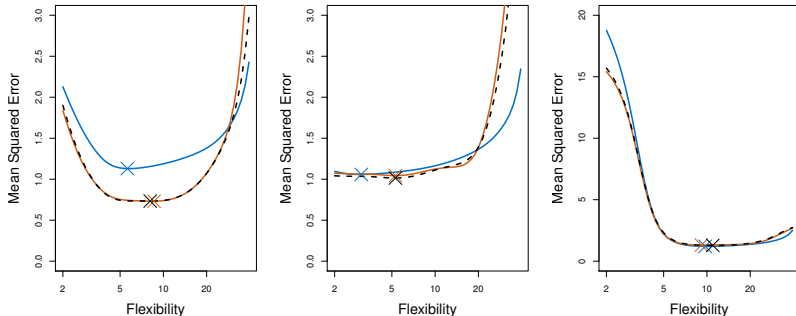


LOOCV vs. k -fold cross-validation



- ▶ k -fold CV depends on the chosen split.
- ▶ In k -fold CV, we train the model on less data than what is available. This introduces **bias** into the estimates of test error.
- ▶ In LOOCV, the training samples highly resemble each other. This increases the **variance** of the test error estimate.

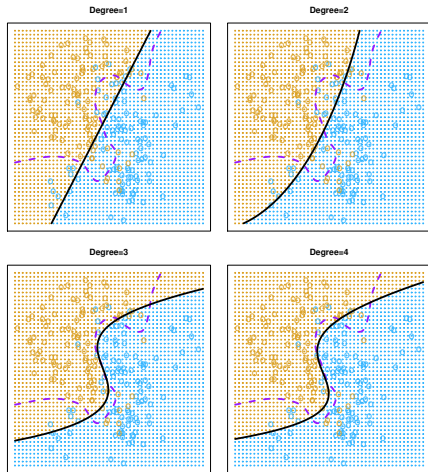
Choosing an optimal model



Even if the error estimates are off, choosing the model with the minimum cross validation error often leads to the method with minimum test error.

Choosing an optimal model

In a classification problem, things look similar.

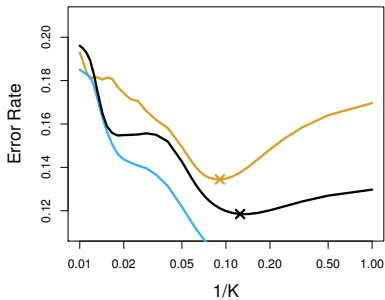
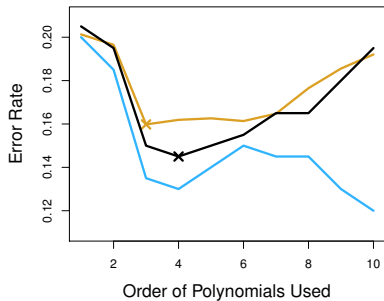


- - - Bayes boundary

— Logistic regression
with polynomial predictors
of increasing degree.

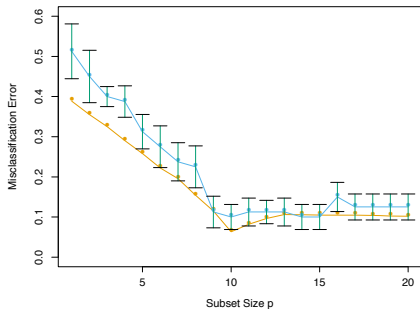
Choosing an optimal model

In a classification problem, things look similar.



The one standard error rule

Forward stepwise selection



Blue: 10-fold cross validation

Yellow: True test error

- ▶ A number of models with $10 \leq p \leq 15$ have the same CV error.
- ▶ The vertical bars represent 1 standard deviation in the test error from the 10 folds.
- ▶ **Rule of thumb:** Choose the simplest model whose CV error is no more than one standard deviation above the model with the lowest CV error.

The wrong way to do cross validation

Reading: Section 7.10.2 of The Elements of Statistical Learning.

We want to classify 200 individuals according to whether they have cancer or not. We use logistic regression onto 1000 measurements of gene expression.

Proposed strategy:

- ▶ Using all the data, select the 20 most significant genes using z -tests.
- ▶ Estimate the standard error of logistic regression with these 20 predictors via 10-fold cross validation.

The wrong way to do cross validation

To see how that works, let's use the following simulated data:

- ▶ Each gene expression is standard normal and independent of all others.
- ▶ The response (cancer or not) is sampled from a coin flip — no correlation to any of the “genes”.

What should the misclassification rate be for any classification method using these predictors?

Roughly 50%.

The wrong way to do cross validation

We run this simulation, and obtain a CV error rate of 3%!

Why is this?

- ▶ Since we only have 200 individuals in total, among 1000 variables, at least some will be correlated with the response.
- ▶ We do variable selection using *all the data*, so the variables we select have some correlation with the response in every subset or fold in the cross validation.

The **right** way to do cross validation

- ▶ Divide the data into 10 folds.
- ▶ For $i = 1, \dots, 10$:
 - ▶ Using every fold except i , perform the variable selection and fit the model with the selected variables.
 - ▶ Compute the error on fold i .
- ▶ Average the 10 test errors obtained.

In our simulation, this produces an error estimate of close to 50%.

Moral of the story: Every aspect of the model that involves using the data — variable selection, for example — must be cross-validated.