

## Overfitting, cross-validation and Nearest Neighbor with R

**Objective:** The objective of this exercise is to understand how cross-validation can be used to avoid overfitting as well as the  $k$ -nearest neighbor method.

**Piazza discussion forum:** You can get help by asking questions on Piazza: [piazza.com/dtu.dk/fall2019/october2019](https://piazza.com/dtu.dk/fall2019/october2019)

**Software installation:** Extract the R toolbox from the Dropbox folder . Start R and go to the `<base-dir>/02450Toolbox_R/` directory by setting the working directory through `setwd(<base-dir>/02450Toolbox_R/)` and run `source('setup.R')`. Remember the purpose of the exercises is not to re-write the code from scratch but to work with the scripts provided in the directory `<base-dir>/02450Toolbox_R/Scripts/`

**Representation of data in R:**

	R var.	Type	Size	Description
	<b>X</b>	Numeric	$N \times M$	Data matrix: The rows correspond to $N$ data objects, each of which contains $M$ attributes.
	<b>attributeNames</b>	Cell array	$M \times 1$	Attribute names: Name (string) for each of the $M$ attributes.
	<b>N</b>	Numeric	Scalar	Number of data objects.
	<b>M</b>	Numeric	Scalar	Number of attributes.
	<b>y</b>	Numeric	$N \times 1$	Dependent variable (output): For each data object, <b>y</b> contains an output value that we wish to predict.
Regression	<b>y</b>	Numeric	$N \times 1$	Class index: For each data object, <b>y</b> contains a class index, $y_n \in \{0, 1, \dots, C - 1\}$ , where $C$ is the total number of classes.
	<b>classNames</b>	Cell array	$C \times 1$	Class names: Name (string) for each of the $C$ classes.
	<b>C</b>	Numeric	Scalar	Number of classes.
Classification				All variables mentioned above appended with <b>_train</b> or <b>_test</b> represent the corresponding variable for the training or test set.
	<b>*_train</b>	—	—	Training data.
	<b>*_test</b>	—	—	Test data.

### 6.1 Decision tree pruning using cross-validation

In this exercise we will use cross-validation to prune a decision tree. When applying cross-validation the observed data is split into training and test sets, i.e., **X\_train**,

`y_train` and `X_test` and `y_test`. We train the model on the training data and evaluate the performance of the trained model on the test data.

- 6.1.1 Inspect and run the script `ex6_1_1.R`. The script load the `wine2.mat` file with wine data by the command `library(R.matlab)` and `dat <- readMat('Data/wine2.ma`. In this version of the wine data, outliers have already been removed. Notice how the script divides the data into a training and a test data set. Now, we want to find optimally pruned decision tree, by modifying its maximum depth. For different values of parameter (depth from 2 to 20) explain how the script fits the decision tree, and compute the classification error on the training and test set (holdout cross-validation). Notice how the script plot the training and test classification error as a function of the pruning level. What does this plot tell you?

Script details:

- *Take a look at the function `cvFolds` in the package `cvTools` and see how it can be used to partition the data into a training and a test set ( $K=2$ ).*
- *To classify the training set using the classification tree, you can use the `predict` function. Type `?predict.rpart` to get help.*
- *To compare two vectors of strings, you can use `==` to test for equality or `!=` to test for inequality for each pair of strings in the two vectors of equal length.*

What appears to be the optimal tree depth? Do you get the same result when you run your code again, generating a new random split between training and test data? What other parameters of the tree could you optimize in cross-validation?

- 6.1.2 Inspect the script `ex6_1_2.R`. The script repeat the exercise above, using 10-fold cross-validation. To do this, the data set is divided into 10 random training and test folds. For each fold, a decision tree is fitted on the training set and it's performance is evaluated on the test set. Finally, the average classification error is computed across the 10 cross-validation folds.

Script details:

- *As before, `cvFolds` can be used to partition the data into the 10 training and test partitions.*

What appears to be the optimal tree depth? Do you get the same result when you run your code again, generating a new random split between training and test data? How about 100-fold cross-validation or leave-one-out cross-validation?

## 6.2 Variable selection in linear regression

In this exercise we consider cross-validation for variable selection and model performance evaluation in linear regression. We will try to predict the body-weight of

a person based on a number of body measurements using linear regression with feature subset selection. The data is a subset of the data available at <http://www.sci.usq.edu.au/courses/STA3301/resources/Data/> described in [1]. To measure how well we can predict the body-weight, we will use the squared error between the true and estimated body-weight.

In our estimation we will use two levels of cross-validation: 1) On the outer level, we use 5-fold cross-validation to estimate the performance of our model, i.e., we compute the squared error averaged over 5 test sets. 2) On the inner level, we use 10-fold cross-validation to perform sequential feature selection (see figure 1).

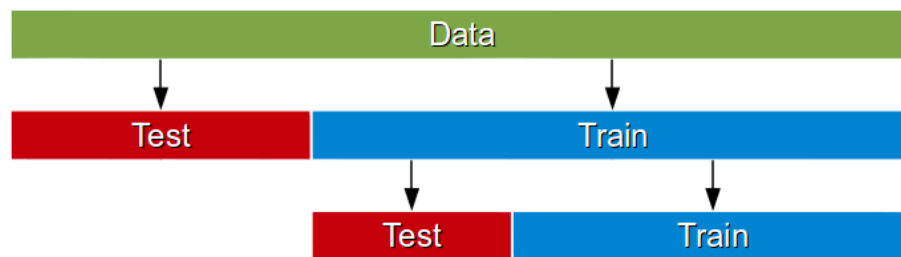


Figure 1: Multi-level cross validation

6.2.1 You can load the body data into R with the commands `library(R.matlab)` and `dat <- readMat('Data/body.mat')`. The data set contains data for the 23 attributes in the matrix `X` and the body-weight in `y`.

Inspect and run the script `ex6_2_1.R`. The script applies 5-fold cross-validation to the problem of fitting a linear regression model to estimate the body-weight based on the attributes. Explain how the script, when fitting the models, compares two methods: 1) using all 23 attributes, and 2) using 10-fold cross-validation to perform sequential feature selection, thus choosing a subset of the 23 attributes.

Explain how the script computes the 5-fold cross-validated training and test error with and without sequential feature selection. Explain how it can be seen that without feature selection, the model overfits. Explain how it can be seen the feature selection tends to choose features such as height and waist girth, and disregard features such as the wrist diameter, which seems reasonable when predicting body-weight.

Script details:

- *Again, you may use `cvFolds` to set up the cross-validation partitions needed.*
- *To fit a linear regression model, use the functions `lm` and `predict` as in the previous exercises.*

- To perform sequential feature selection, you can use the function `forwardSelection`. One of the inputs required by this function is a function that implements a criterion function:  
`criterion = function(XTRAIN,YTRAIN,XTEST,YTEST){...}`.
- By default, `forwardSelection` uses 10-fold crossvalidation for the feature selection, but you can use the argument `cvK` if you want to use another number of folds.
- One way to supply the criterion function is to write a function, e.g., called `funLinreg`. The function may either be written in the same file, or in a separate file `funLinreg.R` that must then be run before the function can be used. This function must take the four arguments shown above, fit a linear model using the training data, evaluate the fitted model on the test data, and compute the squared test error as an output. Using this approach, the first argument to `forwardSelection` should be `funLinreg`.

**Optional:** Try modifying the solution to use backward feature subset selection. Does it give the same result?

### 6.3 K-nearest neighbor classification

In this exercise we will use the k-nearest neighbors (KNN) method for classification. First, we will consider 4 different synthetic datasets, that can be loaded into R using the commands `library(R.matlab)` and `dat <- readMat('Data/synth1.mat')` ... `dat <- readMat('D`

- 6.3.1 Consider the script `ex6_3_1.R`. For each of the four synthetic datasets, do the following. Load the dataset into R and examine it by making a scatter plot. Classify the test data `X_test` using a k-nearest neighbor classifier. Choose a suitable number of neighbors. Examine the accuracy and error rate.

Script details:

- The function `knn` in the R package `FNN` can be used to perform k-nearest neighbors classification.
- To plot a confusion matrix, you can use the function `confmatplot` in the course toolbox. This function also displays the accuracy and error rate.

Which distance measures worked best for the four problems? Can you explain why? How many neighbors were needed for the four problems? Can you give an example of when it would be good to use a large/small number of neighbors? Consider e.g. when clusters are well separated versus when they are overlapping.

In general we can use cross-validation to select the optimal distance metric and number of nearest neighbors  $k$  although this can be computationally expensive. We will return to the Iris data we have considered in previous exercises, and attempt to classify the Iris flowers using KNN.

- 6.3.2 Consider the script `ex6_3_2.R`. The script loads the Iris data into R. Explain how the script uses leave-one-out crossvalidation to estimate the number of neighbors,  $k$ , for the  $k$ -nearest neighbors classifier and plots the crossvalidated average classification error as a function of  $k$  for  $k = 1, \dots, 40$ .

Script details:

- *To load the Iris data, you can run your solution to exercise 4.1.1.*
- *Use `cvFolds` to set up the crossvalidation partitions needed.*
- *As before, use the `knn` function for  $k$ -nearest neighbors classification.*

- 6.3.3 Discussion: What are the benefits and drawbacks of K-nearest neighbor classification and regression compared to logistic regression, decision trees and linear regression? (Hint: There are two important aspects of classification and regression methods, how well the methods can *predict* unlabeled data and how well the method *describe* what aspects in the data causes the data to be classified a certain way .)

#### 6.4 Task for the report

The report will make use of cross-validation, but in conjunction with methods we have not seen yet. Please see report description for more information.

## References

- [1] Grete Heinz, Louis J Peterson, Roger W Johnson, and Carter J Kerk. Exploring relationships in body dimensions. *Journal of Statistics Education*, 11(2), 2003.