

# Lab 04C: Week 13

---

## Contents

### 1 Questions

- 1.1 Who has diabetes?
- 1.2 Violent crime rates by US states
- 1.3 Who has diabetes? [Revisited]

### 2 For practice after the computer lab

The **specific aims** of this lab are:

- estimate logistic regression models
- evaluate out-of-sample performance using cross validation
- perform classification using k-nearest neighbours and evaluate out-of-sample performance using cross validation
- perform dimension reduction using PCA/SVD
- perform cluster analysis using k-means and hierarchical approaches

The unit **learning outcomes** addressed are:

- LO1 Formulate domain/context specific questions and identify appropriate statistical analysis.
- LO3 Construct, interpret and compare numerical and graphical summaries of different data types including large and/or complex data sets.
- LO7 Perform statistical machine learning using a given classifier, and create a cross-validation scheme to calculate the prediction accuracy.

*You can use time in this lab to ask your tutor for further guidance on your approach for the group project modelling and report writing. For example, you may want to ask for clarification on any of the feedback that the marker or your peers provided on your presentation.*

## 1 Questions

### 1.1 Who has diabetes?

---

This dataset is originally from the National Institute of Diabetes and Digestive and Kidney Diseases

(Johannes 1988). The objective is to predict whether or not a patient has diabetes based on certain clinical measurements. The larger database of patients has been subset such that all observations here are females at least 21 years old of Pima Indian heritage.

The data sets consists of several medical predictor variables and one target variable,  $y$  which equals 1 if an individual is diabetic and 0 otherwise. Predictor variables includes the number of pregnancies the patient has had (`npreg`), their BMI, insulin level (`serum`), age, triceps skin fold thickness `skin`, diastolic blood pressure (`bp`), plasma glucose concentration (`glu`) and diabetes pedigree function (`ped`).

It is available from various places, including [Kaggle](#) and the **reglogit** R package and I've uploaded a copy to the [GitHub server](#).

```
library(tidyverse)
pima = readr::read_csv("https://raw.githubusercontent.com/DATA2002/data/master/pima.csv")
glimpse(pima)
```

Rows: 768

Columns: 9

```
$ npreg <dbl> 6, 1, 8, 1, 0, 5, 3, 10, 2, 8, 4, 10, 10, 1, 5, 7, 0, ...
$ glu   <dbl> 148, 85, 183, 89, 137, 116, 78, 115, 197, 125, 110, 16...
$ bp    <dbl> 72, 66, 64, 66, 40, 74, 50, 0, 70, 96, 92, 74, 80, 60,...
$ skin  <dbl> 35, 29, 0, 23, 35, 0, 32, 0, 45, 0, 0, 0, 0, 23, 19, 0...
$ serum <dbl> 0, 0, 0, 94, 168, 0, 88, 0, 543, 0, 0, 0, 0, 846, 175,...
$ bmi   <dbl> 33.6, 26.6, 23.3, 28.1, 43.1, 25.6, 31.0, 35.3, 30.5, ...
$ ped   <dbl> 0.627, 0.351, 0.672, 0.167, 2.288, 0.201, 0.248, 0.134...
$ age   <dbl> 50, 31, 32, 21, 33, 30, 26, 29, 53, 54, 30, 34, 57, 59...
$ y     <dbl> 1, 0, 1, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, ...
```

### 1.1.1 k-nearest neighbours

1. Perform k-nearest neighbours on the data with  $k = 5$ . How does this perform in-sample? How does this perform out-of-sample? Don't use the **caret** package, write your own CV method.
2. Let's establish out of sample performance over a range of  $k$  values. Put an additional `for` loop around your CV loop above to get an estimate of out-of-sample performance for  $k = 1, 2, \dots, 50$ . Visualise your results.
3. Use the **caret** package to perform *repeated* CV to identify the optimal value for  $k$ .

### 1.1.2 Comparison

4. Compare the out of sample accuracies for logistic regression, decision tree, random forest and k-nearest neighbours the different methods using 5 fold CV with 10 repeats.
5. Which model do you prefer?

6. Are our results generalisable to other populations?

## 1.2 Violent crime rates by US states

---

This is the same data (`USArrests`) we considered for PCA in the lecture. It is built into R, so can be accessed just by typing `USArrests`.

### 1.2.1 Hierarchical clustering

The `hclust()` function can be used for hierarchical clustering after we have calculated a set of *distances* (a measure of dissimilarity) between each of the observation. The distance calculation can be done using the `dist()` function, with the default distance metric being the Euclidean distance.

```
dplyr::glimpse(USArrests)
```

1. Calculate the pairwise distances using the `dist()` function applied to the scaled `USArrests` data frame.
2. Use the `hclust()` function applied to the pairwise distances to perform a hierarchical cluster analysis. Visualise it using the `plot()` function.
3. Use the `cutree()` function to cut the tree to make 4 groups. Add the identified clusters as a new column in the `USArrests` data frame. Plot `Murder` against `UrbanPop` and colour the points by the identified cluster groups.

### 1.2.2 k-means

Use the `kmeans()` function to perform k-means clustering on the scaled `USArrests` data frame and specify 4 clusters using the `centers` parameter. Plot `Murder` against `UrbanPop` and colour the points by the identified cluster groups.

## 1.3 Who has diabetes? [Revisited]

---

### 1.3.1 Clustering and dimension reduction

Let's ignore the labels for now and let the clustering algorithm try to discover structure in the data.

1. Perform k-means clustering on the Pima Indians data (without the `y` variable) for  $k = 2$ . Create a confusion matrix and compare the clusters discovered by the k-means algorithm with the observed outcomes (`pima$y`).

2. Use PCA to perform dimension reduction on the data (again without the  $y$  variable). How much of the variation in the data can be explained by the first two principal components?
3. Visualise the first two principal components on a biplot. On separate plots, overlay
  - a. the clusters identified by the k-means clustering, and
  - b. the original data labels.

## 2 For practice after the computer lab

As additional practice questions, I recommend these two DataCamp chapters:

- [`</>` Logistic Regression](#)
  - [`</>` Decision trees](#)
  - [`</>` k-nearest neighbors](#)
  - [`</>` Cluster analysis in R](#)
  - [`</>` Dimension reduction with PCA](#)
- 

## References

- Jed Wing, Max Kuhn. Contributions from, Steve Weston, Andre Williams, Chris Keefer, Allan Engelhardt, Tony Cooper, Zachary Mayer, et al. 2018. *Caret: Classification and Regression Training*. <https://CRAN.R-project.org/package=caret>.
- Johannes, R S. 1988. "Using the ADAP Learning Algorithm to Forecast the Onset of Diabetes Mellitus." *Johns Hopkins APL Technical Digest* 10: 262–66.
- Liaw, Andy, and Matthew Wiener. 2002. "Classification and Regression by randomForest." *R News* 2 (3): 18–22. <https://CRAN.R-project.org/doc/Rnews/>.
- Therneau, Terry, and Beth Atkinson. 2018. *Rpart: Recursive Partitioning and Regression Trees*. <https://CRAN.R-project.org/package=rpart>.