# Exercise Lesson03: KNN

## Machine Learning using R

## Exercise 1: Crabs data set

The `crabs` data set describes different morphological measurements of two different crab-species (see https://stat.ethz.ch/R-manual/R-devel/library/MASS/html/crabs.html for a detailed description of the data). We want to train a KNN-classifier on the crabs data set.

a) The data set is contained in the r package `MASS`. Load the package in order to access the data frame `crabs`.

b) Get a first impression of the data and make sure that everything is coded correctly.

c) Fit a knn classifier to the data using k=3. We try to model the species (`sp`) of the crabs based on the five morphological measurements (if unsure, check the data description to make sure you pick the correct columns). (**Hint**: `knn()`, you need to load the package `class` for the `knn()` function)

d) Display the confusion matrix of the predictions (on the training-data itself) and calculate the training error (see slides).

## Exercise 2: NYC taxi trip records

The taxi trip data set is a collection of information regarding taxi trips in New York City. The data set includes information regarding (amongst others) the pick-up and drop-off location, trip distances, fare prices and passenger count. See https://www1.nyc.gov/site/tlc/about/tlc-trip-record-data.page for more information. We work with a reduced version of the data set, which only includes a sample of the total records.

a) Load the data with the command `read.csv('taxi.csv', stringsAsFactors=TRUE)` and get a first impression of the data. There should only be one factor in the data which is the colour of the taxi (there are green and yellow cabs in NY). How many rows does the data contain? How many observations are there for green and yellow taxi trips? (**Note**: the pick-up and drop-off location IDs shouldn't ideally be coded as numbers, however, for this exercises we will assume that location IDs with similar numbers are also closer to each other and will therefore keep the two variables as numericals)

b) We want to compare the training and test error when using a k-nearest neighbor classification. The goal is to predict the colour of the taxi based on the trip information. Remove a random sample from the taxi data (size=1000), we will use this sample as a test data set. The remaining data will be used as training data. Calculate the training and test error of the KNN algorithm with k=4. (**Hint**: `sample()`, fix the random seed with `set.seed()`)

c) **Extra:** We now want to compare the training and test error for different k. Write a loop in which you fit a knn classifier to the training data (from previous exercise) with k ranging from 1 to 30. Collect for each k the training error and the test error in a table. Also: Standardize the data (train and test set) before applying the knn classifier.

d) **Extra:** Plot the test and training error rates against k. Which k should be chosen based on this quick analysis? Why is this the best k?