

# COMP10200: Assignment 1, part b

© Sam Scott, Mohawk College, 2020

## Get Some Machine Learning Data

Go to the UCI Machine Learning Repository at <https://archive.ics.uci.edu/ml/> and click the “View All Data Sets” link in the header.

In the selectors at right, choose “Classification”, “Numerical”, and “10 to 100” attributes. Find a data set with a reasonable number of attributes (say 15 or so, although if you find an interesting data set with more features than that, it’s fine).

Read over the description of the data set and look at the data files. Make sure that the data is the right kind of data. What you’re looking for is all numerical attributes except for the class label which can be a string or a set of discrete values (0, 1, 2, etc.). Make sure the data set makes sense for a classification task – if you’re not sure, ask the instructor.

Some data sets have a single file, others separate training and test data, still others separate labels into a separate file from the data. For this exercise, we want a single CSV file with data and labels combined and with a header row that contains the name of each feature. You can stitch it together in Excel if you need to. In some cases, you will have to enter the feature names yourself. Here’s an example of the type of file you want to end up with. Of course, yours will have more features and more examples.

```
ID, Age, Height, Weight, Diagnosis
1, 22, 180, 50, Sick
2, 45, 175, 55, Healthy
3, 99, 130, 35, Healthy
```

## Read & Split the Data

Write a Python script that reads the CSV file and splits the data set randomly into testing and training data. You’ll need 5 Numpy arrays for this: feature names, training data, training labels, testing data, and testing labels. Make sure the training set contains about 75% of the data and the testing set has the other 25%. Each time the script is run, it should produce a different split. There are a number of ways to do this with tools built in to **numpy**. There are also special tools in **sklearn** for creating testing and training splits. Feel free to explore those on your own and use them.

## Summarize It

Print summary information for the training and testing data: The name of each feature, followed by it’s minimum, maximum, mean, and median in the training set and in the testing set. Make proper use of the NumPy array type for this (i.e. no loops!)

## Graph It

Explore your training data by creating 2D scatter plots of pairs of features. Make sure the different class labels are represented with different colors. Find 3 pairs that seem promising for separating the classes and add code to your Python script to produce those 3 scatter plots. (Hint: Call `plt.figure(1)` and

draw the first plot, then call `plt.figure(2)` and draw the second plot, etc.). Make sure each graph has a useful title and has the axes labelled.

Then produce a bar chart showing the overall frequency of each label. You can find information on how to do bar charts this in Chapter 3 of *Data Science from Scratch* (on eLearn).

Make proper use of the NumPy array type for this (i.e. no loops!)

## Commenting

Follow the documentation conventions for the course. You'll find these in a python file in the Student Resources section of eLearn.

Within the code, put a headline on each section (Read It, Summarize It, Graph It, etc.) using a `##` comment. Use `#` comments on lines where you think it might not be obvious what you're doing.

## Handing In

Zip together your Python scripts plus the original or modified file that you got from the UCI Machine Learning Repository (i.e. the one that your script reads), and hand it in to the drop box.

## Evaluation

Your code will be evaluated for Performance, Structure, and Documentation using the Rubric attached to the drop box.