## Assignment 1 - Introduction to Machine Learning

For this assignment, you will be using the Breast Cancer Wisconsin (Diagnostic) Database to create a classifier that can help diagnose patients. First, read through the description of the dataset (below).

```
In [1]: import numpy as np
        import pandas as pd
        from sklearn.datasets import load_breast_cancer

        cancer = load_breast_cancer()

        #print(cancer.DESCR) # Print the data set description
```

The object returned by `load_breast_cancer()` is a scikit-learn Bunch object, which is similar to a dictionary.

### Question 0 (Example)

How many features does the breast cancer dataset have?

*This function should return an integer.*

### Question 1

Scikit-learn works with lists, numpy arrays, scipy-sparse matrices, and pandas DataFrames, so converting the dataset to a DataFrame is not necessary for training this model. Using a DataFrame does however help make many things easier such as munging data, so let's practice creating a classifier with a pandas DataFrame.

Convert the sklearn.dataset cancer to a DataFrame.

*This function should return a (569, 31) DataFrame with*

*columns =*

```
['mean radius', 'mean texture', 'mean perimeter', 'mean area',
'mean smoothness', 'mean compactness', 'mean concavity',
'mean concave points', 'mean symmetry', 'mean fractal dimension',
'radius error', 'texture error', 'perimeter error', 'area error',
'smoothness error', 'compactness error', 'concavity error',
'concave points error', 'symmetry error', 'fractal dimension error',
'worst radius', 'worst texture', 'worst perimeter', 'worst area',
'worst smoothness', 'worst compactness', 'worst concavity',
'worst concave points', 'worst symmetry', 'worst fractal dimension',
'target']
```

*and index =*

```
RangeIndex(start=0, stop=569, step=1)
```

### Question 2

What is the class distribution? (i.e. how many instances of `malignant` (encoded 0) and how many benign (encoded 1)?)

*This function should return a Series named target of length 2 with integer values and index = ['malignant', 'benign']*

## Question 3

Split the DataFrame into X (the data) and y (the labels).

*This function should return a tuple of length 2: (X, y), where*

- *X, a pandas DataFrame, has shape (569, 30)*
- *y, a pandas Series, has shape (569,).*

## Question 4

Using `train_test_split`, split X and y into training and test sets (X_train, X_test, y_train, and y_test).

**Set the random number generator state to 0 using random_state=0 to make sure your results match the autograder!**

*This function should return a tuple of length 4: (X_train, X_test, y_train, y_test), where*

- *X_train has shape (426, 30)*
- *X_test has shape (143, 30)*
- *y_train has shape (426,)*
- *y_test has shape (143,)*

## Question 5

Using KNeighborsClassifier, fit a k-nearest neighbors (knn) classifier with X_train, y_train and using one nearest neighbor (n_neighbors = 1).

*This function should return a sklearn.neighbors.classification.KNeighborsClassifier.*

## Question 6

Using your knn classifier, predict the class label using the mean value for each feature.

Hint: You can use `cancerdf.mean()[:-1].values.reshape(1, -1)` which gets the mean value for each feature, ignores the target column, and reshapes the data from 1 dimension to 2 (necessary for the precict method of KNeighborsClassifier).

*This function should return a numpy array either array([ 0.]) or array([ 1.])*

## Question 7

Using your knn classifier, predict the class labels for the test set X_test.

*This function should return a numpy array with shape (143,) and values either 0.0 or 1.0.*

## Question 8

Find the score (mean accuracy) of your knn classifier using X_test and y_test.

*This function should return a float between 0 and 1*