# Module 1

September 16, 2020

---

*You are currently looking at **version 1.0** of this notebook. To download notebooks and datafiles, as well as get help on Jupyter notebooks in the Coursera platform, visit the Jupyter Notebook FAQ course resource.*

---

## 0.1 Applied Machine Learning, Module 1: A simple classification task

### 0.1.1 Import required modules and load data file

```
In [ ]: %matplotlib notebook
        import numpy as np
        import matplotlib.pyplot as plt
        import pandas as pd
        from sklearn.model_selection import train_test_split

        fruits = pd.read_table('readonly/fruit_data_with_colors.txt')

In [ ]: fruits.head()

In [ ]: # create a mapping from fruit label value to fruit name to make results eas
        lookup_fruit_name = dict(zip(fruits.fruit_label.unique(), fruits.fruit_name
        lookup_fruit_name
```

The file contains the mass, height, and width of a selection of oranges, lemons and apples. The heights were measured along the core of the fruit. The widths were the widest width perpendicular to the height.

### 0.1.2 Examining the data

```
In [ ]: # plotting a scatter matrix
        from matplotlib import cm

        X = fruits[['height', 'width', 'mass', 'color_score']]
        y = fruits['fruit_label']
        X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)

        cmap = cm.get_cmap('gnuplot')
        scatter = pd.scatter_matrix(X_train, c= y_train, marker = 'o', s=40, hist_k
```

```
In [ ]: # plotting a 3D scatter plot
        from mpl_toolkits.mplot3d import Axes3D

        fig = plt.figure()
        ax = fig.add_subplot(111, projection = '3d')
        ax.scatter(X_train['width'], X_train['height'], X_train['color_score'], c =
        ax.set_xlabel('width')
        ax.set_ylabel('height')
        ax.set_zlabel('color_score')
        plt.show()
```

### 0.1.3 Create train-test split

```
In [ ]: # For this example, we use the mass, width, and height features of each fru
        X = fruits[['mass', 'width', 'height']]
        y = fruits['fruit_label']

        # default is 75% / 25% train-test split
        X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)
```

### 0.1.4 Create classifier object

```
In [ ]: from sklearn.neighbors import KNeighborsClassifier

        knn = KNeighborsClassifier(n_neighbors = 5)
```

### 0.1.5 Train the classifier (fit the estimator) using the training data

```
In [ ]: knn.fit(X_train, y_train)
```

### 0.1.6 Estimate the accuracy of the classifier on future data, using the test data

```
In [ ]: knn.score(X_test, y_test)
```

### 0.1.7 Use the trained k-NN classifier model to classify new, previously unseen objects

```
In [ ]: # first example: a small fruit with mass 20g, width 4.3 cm, height 5.5 cm
        fruit_prediction = knn.predict([[20, 4.3, 5.5]])
        lookup_fruit_name[fruit_prediction[0]]

In [ ]: # second example: a larger, elongated fruit with mass 100g, width 6.3 cm, h
        fruit_prediction = knn.predict([[100, 6.3, 8.5]])
        lookup_fruit_name[fruit_prediction[0]]
```

### 0.1.8 Plot the decision boundaries of the k-NN classifier

```
In [ ]: from adspy_shared_utilities import plot_fruit_knn

        plot_fruit_knn(X_train, y_train, 5, 'uniform')   # we choose 5 nearest neig
```

### 0.1.9 How sensitive is k-NN classification accuracy to the choice of the 'k' parameter?

```
In [ ]: k_range = range(1,20)
        scores = []

        for k in k_range:
            knn = KNeighborsClassifier(n_neighbors = k)
            knn.fit(X_train, y_train)
            scores.append(knn.score(X_test, y_test))

        plt.figure()
        plt.xlabel('k')
        plt.ylabel('accuracy')
        plt.scatter(k_range, scores)
        plt.xticks([0,5,10,15,20]);
```

### 0.1.10 How sensitive is k-NN classification accuracy to the train/test split proportion?

```
In [ ]: t = [0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2]

        knn = KNeighborsClassifier(n_neighbors = 5)

        plt.figure()

        for s in t:

            scores = []
            for i in range(1,1000):
                X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
                knn.fit(X_train, y_train)
                scores.append(knn.score(X_test, y_test))
            plt.plot(s, np.mean(scores), 'bo')

        plt.xlabel('Training set proportion (%)')
        plt.ylabel('accuracy');

In [ ]:
```