# Supervised Learning with scikit-learn

## 1). Classification

**a). k-Nearest Neighbors: Fit**

```python
# Import KNeighborsClassifier from sklearn.neighbors
from sklearn.neighbors import KNeighborsClassifier

# Create arrays for the features and the response variable
y = df['party'].values
X = df.drop('party', axis=1).values

# Create a k-NN classifier with 6 neighbors
knn = KNeighborsClassifier(6)

# Fit the classifier to the data
knn.fit(X,y)
```

**b). k-Nearest Neighbors: Predict**

```
# Import KNeighborsClassifier from sklearn.neighbors
from sklearn.neighbors import KNeighborsClassifier


# Create arrays for the features and the response variable
y = df['party'].values
X = df.drop('party',axis=1).values


# Create a k-NN classifier with 6 neighbors: knn
knn = KNeighborsClassifier(6)


# Fit the classifier to the data
knn.fit(X,y)


# Predict the labels for the training data X
y_pred = knn.predict(X)


# Predict and print the label for the new data point X_new
new_prediction = knn.predict(X_new)
print("Prediction: {}".format(new_prediction))
```

\<script.py\> output:

   Prediction: ['democrat']
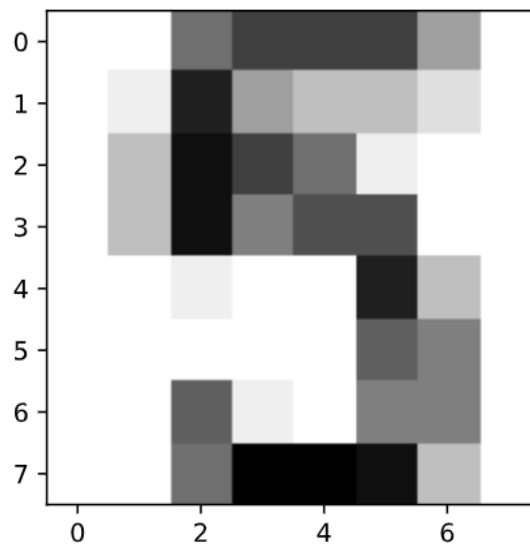
**c). <u>The digits recognition dataset</u>**

**# Import necessary modules**

**from sklearn import datasets**

**import matplotlib.pyplot as plt**

**# Load the digits dataset: digits**

**digits = datasets.load_digits()**

**# Print the keys and DESCR of the dataset**

**print(digits.keys())**

**print(digits.DESCR)**

**# Print the shape of the images and data keys**

**print(digits.images.shape)**

**print(digits.data.shape)**

**# Display digit 1010**

**plt.imshow(digits.images[1010], cmap=plt.cm.gray_r, interpolation='nearest')**

**plt.show()**

**d). <u>Train/Test Split + Fit/Predict/Accuracy</u>**

**# Import necessary modules**

**from sklearn.neighbors import KNeighborsClassifier**

**from sklearn.model_selection import train_test_split**

**# Create feature and target arrays**

**X = digits.data**

**y = digits.target**

**# Split into training and test set**

**X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state=42, stratify=y)**

**# Create a k-NN classifier with 7 neighbors: knn**

**knn = KNeighborsClassifier(n_neighbors=7)**

**# Fit the classifier to the training data**

**knn.fit(X_train, y_train)**

**# Print the accuracy**

**print(knn.score(X_test, y_test))**

<script.py> output:

   0.983333333333

**e). <u>Overfitting & Underfitting</u>**

**# Setup arrays to store train and test accuracies**

**neighbors = np.arange(1, 9)**

**train_accuracy = np.empty(len(neighbors))**

**test_accuracy = np.empty(len(neighbors))**

**# Loop over different values of k**

**for i, k in enumerate(neighbors):**

   **# Setup a k-NN Classifier with k neighbors: knn**

   **knn = KNeighborsClassifier(n_neighbors=k)**


   **# Fit the classifier to the training data**

   **knn.fit(X_train, y_train)**

   **#Compute accuracy on the training set**

   **train_accuracy[i] = knn.score(X_train, y_train)**

   **#Compute accuracy on the testing set**

   **test_accuracy[i] = knn.score(X_test, y_test)**

**# Generate plot**

**plt.title('k-NN: Varying Number of Neighbors')**

**plt.plot(neighbors, test_accuracy, label = 'Testing Accuracy')**

**plt.plot(neighbors, train_accuracy, label = 'Training Accuracy')**

**plt.legend()**

**plt.xlabel('Number of Neighbors')**

**plt.ylabel('Accuracy')**

**plt.show()**