

Programming assignment 1: k-Nearest Neighbors classification

```
In [11]: import numpy as np
from sklearn import datasets, model_selection
import matplotlib.pyplot as plt
%matplotlib inline
```

Introduction

For those of you new to Python, there are lots of tutorials online, just pick whichever you like best)

If you never worked with Numpy or Jupyter before, you can check out these guides

- <https://docs.scipy.org/doc/numpy-dev/user/quickstart.html>
- <http://jupyter.readthedocs.io/en/latest/>

Your task

In this notebook code to perform k-NN classification is provided. However, some functions are incomplete. Your task is to fill in the missing code and run the entire notebook.

In the beginning of every function there is docstring, which specifies the format of input and output. Write your code in a way that adheres to it. You may only use plain python and `numpy` functions (i.e. no scikit-learn classifiers).

Exporting the results to PDF

Once you complete the assignments, export the entire notebook as PDF and attach it to your homework solutions. The best way of doing that is

1. Run all the cells of the notebook.
2. Download the notebook in HTML (click File > Download as > .html)
3. Convert the HTML to PDF using e.g. <https://www.spjda.com/html-to-pdf> or `wkhtmltopdf` for Linux ([tutorial](#))
4. Concatenate your solutions for other tasks with the output of Step 3. On a Linux machine you can simply use `pdflatex`, there are similar tools for other platforms too. You can only upload a single PDF file to Moodle.

This way is preferred to using `nbsconvert`, since `nbsconvert` clips lines that exceed page width and makes your code harder to grade.

Load dataset

The iris data set (https://en.wikipedia.org/wiki/Iris_flower_data_set) is loaded and split into train and test parts by the function `load_dataset`.

```
In [12]: def load_dataset(split):
    """Load and split the dataset into training and test parts.

    Parameters
    -----
    split : float in range (0, 1)
        Fraction of the data used for training.

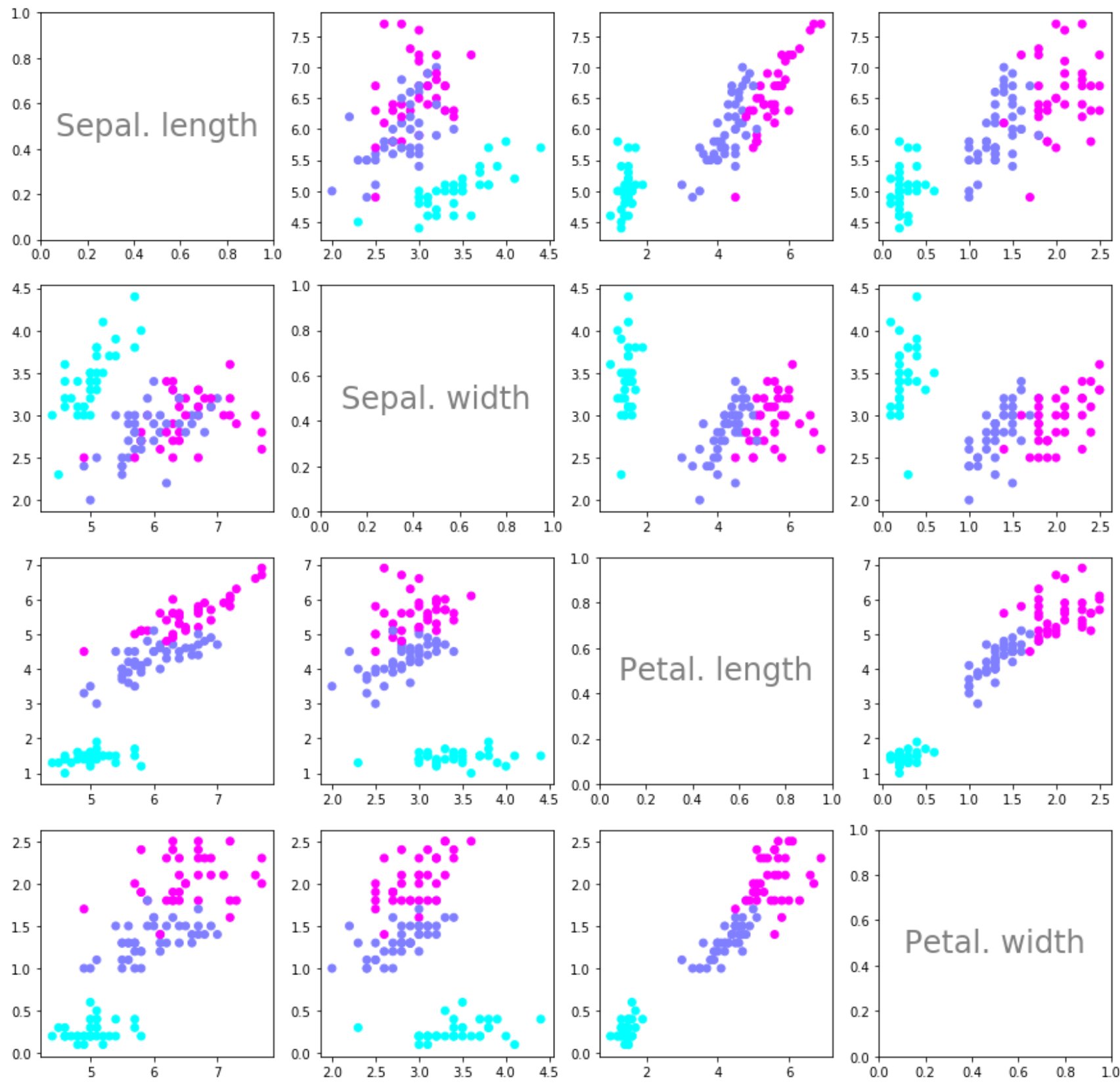
    Returns
    -----
    X_train : array, shape (N_train, 4)
        Training features.
    y_train : array, shape (N_train)
        Training labels.
    X_test : array, shape (N_test, 4)
        Test features.
    y_test : array, shape (N_test)
        Test labels.
    """
    dataset = datasets.load_iris()
    X, y = dataset['data'], dataset['target']
    X_train, X_test, y_train, y_test = model_selection.train_test_split(X, y, random_s
tate=123, test_size=(1 - split))
    return X_train, X_test, y_train, y_test

In [13]: # prepare data
split = 0.75
X_train, X_test, y_train, y_test = load_dataset(split)
```

Plot dataset

Since the data has 4 features, 16 scatterplots (4x4) are plotted showing the dependencies between each pair of features.

```
In [14]: f, axes = plt.subplots(4, 4, figsize=(15, 15))
for i in range(4):
    for j in range(4):
        if j == 0 and i == 0:
            axes[i,j].text(0.5, 0.5, 'Sepal. length', ha='center', va='center', size=24, alpha=.5)
        elif j == 1 and i == 1:
            axes[i,j].text(0.5, 0.5, 'Sepal. width', ha='center', va='center', size=24, alpha=.5)
        elif j == 2 and i == 2:
            axes[i,j].text(0.5, 0.5, 'Petal. length', ha='center', va='center', size=24, alpha=.5)
        elif j == 3 and i == 3:
            axes[i,j].text(0.5, 0.5, 'Petal. width', ha='center', va='center', size=24, alpha=.5)
        else:
            axes[i,j].scatter(X_train[:,j], X_train[:,i], c=y_train, cmap=plt.cm.cool)
```



Task 1: Euclidean distance

Compute Euclidean distance between two data points.

```
In [15]: def euclidean_distance(x1, x2):
    """Compute Euclidean distance between two data points.

    Parameters
    -----
    x1 : array, shape (4)
        First data point.
    x2 : array, shape (4)
        Second data point.

    Returns
    -----
    distance : float
        Euclidean distance between x1 and x2.
    """
    # If the instances are lists or tuples
    x1 = np.array(x1)
    x2 = np.array(x2)

    # Euclidean formula
    return np.linalg.norm(x1-x2)
```

Task 2: get k nearest neighbors' labels

Get the labels of the *k* nearest neighbors of the datapoint *x_{new}*.

```
In [16]: def get_neighbors_labels(X_train, y_train, x_new, k):
    """Get the labels of the k nearest neighbors of the datapoint x_new.

    Parameters
    -----
    X_train : array, shape (N_train, 4)
        Training features.
    y_train : array, shape (N_train)
        Training labels.
    x_new : array, shape (4)
        Data point for which the neighbors have to be found.
    k : int
        Number of neighbors to return.

    Returns
    -----
    neighbors_labels : array, shape (k)
        Array containing the labels of the k nearest neighbors.
    """

    distances = []

    """
    Appending distances list with euclidian distances between X_test observation
    and all X_train observations.
    Format is as follows (X_train observation feature values, distance, label of X_train observation)
    """

    for index in range(len(X_train)):
        dist = euclidean_distance(x_new, X_train[index])
        distances.append((X_train[index], dist, y_train[index]))

    # Sorting distances in descending order
    distances.sort(key=lambda x: x[1])

    neighbors = []

    # Appending distances in range of k nearest neighbors
    for x in range(k):
        neighbors.append(distances[x])

    return neighbors
```

Task 3: get the majority label

For the previously computed labels of the *k* nearest neighbors, compute the actual response, i.e. give back the class of the majority of nearest neighbors. In case of a tie, choose the "lowest" label (i.e. the order of tie resolutions is 0 > 1 > 2).

```
In [17]: def get_response(neighbors_labels, num_classes=3):
    """Predict label given the set of neighbors.

    Parameters
    -----
    neighbors_labels : array, shape (k)
        Array containing the labels of the k nearest neighbors.
    num_classes : int
        Number of classes in the dataset.

    Returns
    -----
    y : int
        Majority class among the neighbors.
    """
    counter = {}

    # Adding counter dictionary the labels as keys and occurrences as values
    for label in range(len(neighbors_labels)):
        response = neighbors_labels[label][-1]
        if response in counter:
            counter[response] += 1
        else:
            counter[response] = 1

    # Getting the highest occurrence number in counter dictionary
    highest = max(counter.values())

    # Making a list of the corresponding keys of the highest occurrence number in dictionary
    class_votes = [k for k, v in counter.items() if v == highest]

    # If class_votes has a single value in it, it stays the same. If it has 2 or more keys in it,
    # we choose the smallest.
    class_votes = min(class_votes)
    return class_votes
```

Task 4: compute accuracy

Compute the accuracy of the generated predictions.

```
In [18]: def compute_accuracy(y_pred, y_test):
    """Compute accuracy of prediction.

    Parameters
    -----
    y_pred : array, shape (N_test)
        Predicted labels.
    y_test : array, shape (N_test)
        True labels.
    """
    correct = 0
    for i in range(len(y_pred)):
        if y_pred[i] == y_test[i]:
            correct += 1
    correct = correct / len(y_pred)
    return correct

In [19]: # This function is given, nothing to do here.
def predict(X_train, y_train, X_test, k):
    """Generate predictions for all points in the test set.

    Parameters
    -----
    X_train : array, shape (N_train, 4)
        Training features.
    y_train : array, shape (N_train)
        Training labels.
    X_test : array, shape (N_test, 4)
        Test features.
    k : int
        Number of neighbors to consider.

    Returns
    -----
    y_pred : array, shape (N_test)
        Predictions for the test data.
    """
    y_pred = []
    for x_new in X_test:
        neighbors = get_neighbors_labels(X_train, y_train, x_new, k)
        y_pred.append(get_response(neighbors))
    return y_pred
```

Testing

Should output an accuracy of 0.9473684210526315.

```
In [20]: # prepare data
split = 0.75
X_train, X_test, y_train, y_test = load_dataset(split)
print('Training set: {} samples'.format(X_train.shape[0]))
print('Test set: {} samples'.format(X_test.shape[0]))

# generate predictions
k = 3
y_pred = predict(X_train, y_train, X_test, k)
accuracy = compute_accuracy(y_pred, y_test)
print('Accuracy = {}'.format(accuracy))

Training set: 112 samples
Test set: 38 samples
Accuracy = 0.9473684210526315
```

```
In [ ]:
```