

COMPUTER LAB 1 - k-NN classifier

Duration: 6 hours

Exercise 1 - Synthetic dataset

In this exercise, you will employ a synthetic dataset (file Lab1_Ex_1_Synthtetic.hdf5), containing labelled training data and test data for two classes. For each example the first two columns represent the features, while the last column represents the label.

Task: your task is to implement a k-NN classifier, which calculates the probability that a given test example belongs to each class, and outputs a class label as the class with the highest probability. You will evaluate the classifier performance computing the average classification accuracy (i.e. the fraction of test examples that have been classified correctly in respect to the full test set).

In particular, you should perform the following:

- Train a k-NN classifier for different values of k.
- Compare accuracy on the training set and the test set. Calculating accuracy of the training set means that you will have to classify each sample in the training set as if it were a test sample; one expects that classification of training samples will perform well, and this may also be used to validate your implementation. Accuracy is defined as the ratio between the number of test samples that are correctly classified, and the total number of test samples. Create a graph using the matplotlib library showing the evolution of the accuracy for different values of k over the test set. Create a second graph to show the evolution of the accuracy for different values of k over the train set and compare the two.
- Identifying overfitting and underfitting in the obtained results.

Note that, for this computer lab, you do not need to employ a validation set.

Other indications:

- The student is required to implement the k-NN algorithm from scratch. Only the numpy library is allowed, while other libraries such as scikit_learn are forbidden.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import h5py
4 import math
5 import pandas as pd
6
7 #Implement a function to compute the Euclidean distance between two vectors, and one to implemer
8 #   - Taking a sample
9 #   - Computing all the distances between the sample element and the elements of the training se
10 #   - sort the the training set based on the distances to the element (the use of functions like
11 #   - select the top k elements in terms of distance
12 #   - evaluate to which class the majority of these k elements belongs to (e.g., it is possible
```

```

1 #Change the path to match the position of your file
2 #The Dataset can be loaded using the file option in Google Colab (the directory icon on the left)
3 Dataset1 = h5py.File('/content/Lab1_Ex_1_Synthtic.hdf5')
4 Data = np.array(Dataset1.get('Dataset'))
5
6
7 Data.ndim
8
9 Train_Set = Data[:200,:]
10 Test_Set = Data[200:,:]
11
12 #To be completed by the student
13 #Function to calculate Euclidean distance
14
15 def eucdist (sample, data):    # - Taking a sample
16     dist = np.sqrt((sample[0]-data[0])**2 + (sample[1]-data[1])**2)
17     return dist
18
19
20 #Defining K-NN
21
22 def k_nn (sample,n):
23     knn_mat = Train_Set
24     dist=[]
25     for j in range (len(knn_mat)):
26         dist.append(eucdist(sample,knn_mat[j]))
27
28     knn_mat = pd.DataFrame(knn_mat)
29     knn_mat.insert(loc=3, column=3, value=dist)
30     knn_mat = knn_mat.sort_values(3)
31     knn = knn_mat.values[:n]
32
33     class1 = 0
34     class2 = 0
35
36     for k in range(len(knn)):
37         if knn[k,2] == 0 :
38             class1 += 1
39         else:
40             class2 += 1
41
42
43     if class1>class2 :
44         return 0
45     else :
46         return 1
47
48
49 k = 10
50 a = k_nn(Test_Set[199] ,k)
51 print('The class is : ' , a)
52

```

The class is : 1

So far we have created the KNN classifier, Now we calculate the Accuracy

```

1 def acc (test , k):
2     true = 0
3     test_predicted = 0
4     for i in range(len(test)):
5         test_predicted = (k_nn(test[i] ,k))
6         if test [i][2] == test_predicted:
7             true += 1
8
9     return (true / len(test ))* 100
10
11 acc_test =  acc(Test_Set, 10)
12 print('The Accuracy for the test dataset is:', acc_test)
13 acc_train =  acc(Train_Set, 10)
14 print('The Accuracy for the train dataset is:',acc_train)
15

```

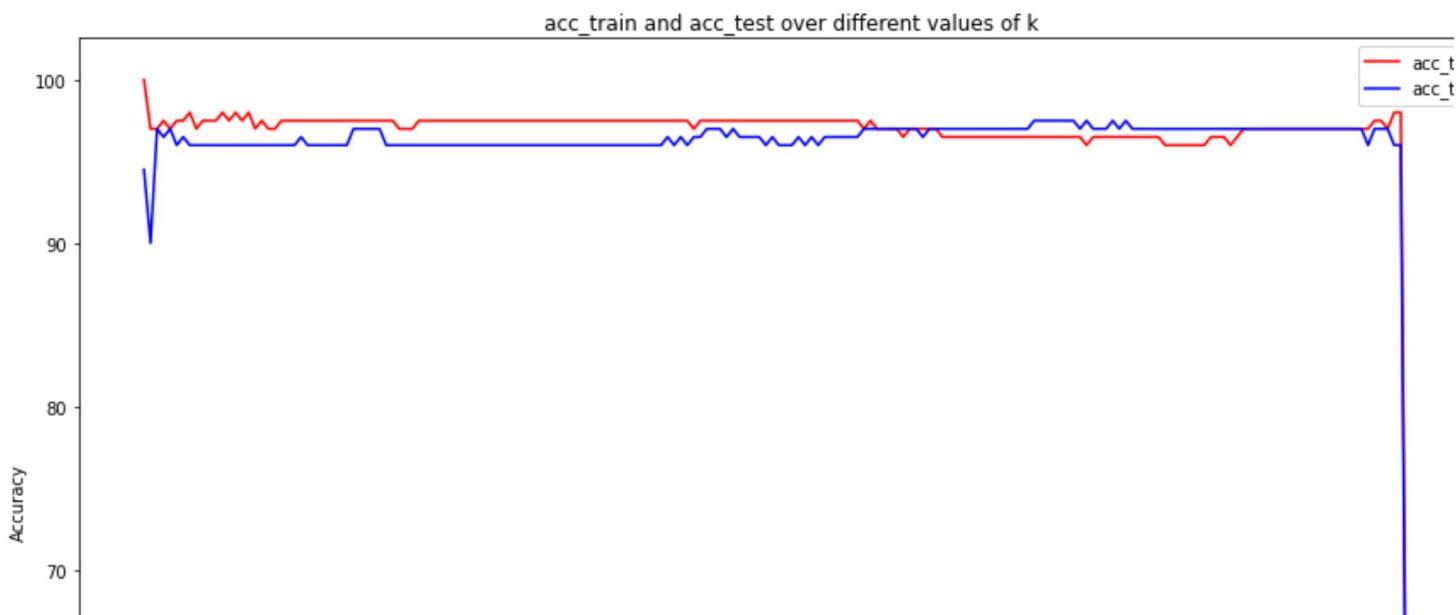
The Accuracy for the test dataset is: 96.0
The Accuracy for the train dataset is: 97.5

By calculating the accuracy we can see that the the training data set has a bit higher accuracy and hence there's no overfitting . Now , we plot the accuracy for different values of k

```

1 k = []
2 acc_train_mat = []
3 acc_test_mat = []
4 for i in range(1,len(Train_Set)):
5     k.append(i)
6     acc_train_mat.append(acc(Train_Set , i))
7     acc_test_mat.append(acc(Test_Set , i))
8
9 f = plt.figure()
10 f.set_figwidth(15)
11 f.set_figheight(10)
12 plt.plot(k, acc_train_mat , color = 'red')
13 plt.plot(k, acc_test_mat , color = 'blue')
14 plt.legend(['acc_train', 'acc_test'])
15 plt.title('acc_train and acc_test over different values of k')
16 plt.xlabel ('k')
17 plt.ylabel('Accuracy')
18 plt.show()

```



Student's comments to exercise 1

The k-nearest neighbors algorithm, also known as KNN or k-NN, is a non-parametric, supervised (Train data have label) learning classifier, which uses proximity (We choose euc distance) to make classifications or predictions about the grouping of an individual data point. KNN sometimes called a lazy learner algorithm as it just store data rather than learning from data for optimizing some parameter.

1. The k value in the k-NN algorithm defines how many neighbors will be checked to determine the classification of a specific query point.
2. It is therefore very intensive to the value of k and as k increases the accuracy is decreasing.(The initial increase is for k = 0 which is not important)
3. The figure depicts best results achieved with the values of k < 120. In this period there is no overfitting and the accuracy is acceptable and above 90% for both of the training and testing data set. with k increasing , the underfitting happens as acc_test > acc_train . and the worst scenario happens near k = 200 with around 50% accuracy which is almost terrible for a binary classification problem.

Overall speaking, One the one hand, a small value of k, means that noise will have a higher influence on the result. On the other hand, if k is selected to be too large, the model becomes too generalized and fails to accurately predict the data points in both train and test sets.

Exercise 2 - Wine dataset

Part 1

In this exercise, a real problem will be examined. The dataset used in this exercise was derived from wine quality dataset from the work "*Modeling wine preferences by data mining from physicochemical properties*" by P. Cortez, A. Cerdeira, F. Almeida, T. Matos and J. Reis.

For each element of the dataset 11 features are provided, representing different wine characteristics, such as density, pH and alcholic content, and the final column consists of a quality evaluation on a scale from 1 to 10. More information can be found at <https://archive.ics.uci.edu/ml/datasets/wine+quality>.

A subset of the dataset containing 400 elements is provided. Create a training set and a test set of 200 samples each. The objective is to:

- Predict the wine quality over the test set using the k-NN algorithm and evaluating the prediction accuracy for different values of k. Create a graph using the matplotlib library showing the evolution of the accuracy for different values of k over the test set.
- Identifying overfitting and underfitting in the obtained results.

Part 2

The prediction of the wine quality could also be framed as a regression. Estimate the accuracy and the Mean Square Error achieved using linear regression. For this task is possible to use the library sklearn and the function `linear_model.LinearRegression()`

```

1 #Defining Eucdist calculator function
2
3 Dataset2 = h5py.File('/content/Lab1_Ex_2_wine.hdf5')
4
5 Data = np.array(Dataset2.get('Dataset'))
6
7 # Normalizing data
8 Train_Set = Data[:200,:]
9 Train_Set = pd.DataFrame(Train_Set)
10 y = Train_Set[11]
11 Train_Set = Train_Set.drop(columns=11)
12 Train_Set=(Train_Set-Train_Set.mean())/Train_Set.std()
13 Train_Set.insert(loc=11, column=11 ,value= y)
14 Train_Set = Train_Set.values
15
16 # Normalizing data
17 Test_Set = Data[200:,:]
18 Test_Set = pd.DataFrame(Test_Set)
19 y = Test_Set[11]
20 Test_Set = Test_Set.drop(columns=11)
21 Test_Set = (Test_Set-Test_Set.mean())/Test_Set.std()
22 Test_Set.insert(loc=11, column=11 ,value= y)
23 Test_Set = Test_Set.values
24
25
26 def eucdist (sample, data):      #    - Taking a sample
27     dist = 0
28     for i in range(len(data)-1):
29         dist += (sample[i]-data[i])**2
30     return np.sqrt(dist)
31
32 # Defining knn
33
34 def k_nn (sample,k):
35     knn_mat = Train_Set
36     dist=[]
37     for j in range (len(knn_mat)):
38         dist.append(eucdist(sample,knn_mat[j]))
39
40     knn_mat = pd.DataFrame(knn_mat)
41     knn_mat.insert(loc=12, column=12, value=dist)
42     knn_mat = knn_mat.sort_values(12)
43     knn = knn_mat.values[:,k]
44
45 # class prediction section
46

```

```

47     classmat = np.array(range(10)) + 1
48     predictmat = np.zeros(10)
49
50     for k in range(len(knn)):
51         for j in classmat:
52             if knn[k,11] == j :
53                 predictmat[j-1] += 1
54
55     return np.argmax(predictmat) + 1
56
57
58 a = k_nn(Test_Set[12] ,10)
59 print(a)
60
61
62
63
64

```

6

Double-click (or enter) to edit

```

1 # Calculating Acc
2 def acc (test , k):
3     true = 0
4     test_predicted = 0
5     for i in range(len(test)):
6         test_predicted = (k_nn(test[i] ,k))
7         if test [i][11] == test_predicted:
8             true += 1
9
10    return (true / len(test ))* 100
11
12 acc_test =  acc(Test_Set, 5)
13 print(acc_test)
14 acc_train =  acc(Train_Set, 5)
15 print(acc_train)
16

```

43.5
66.0

```

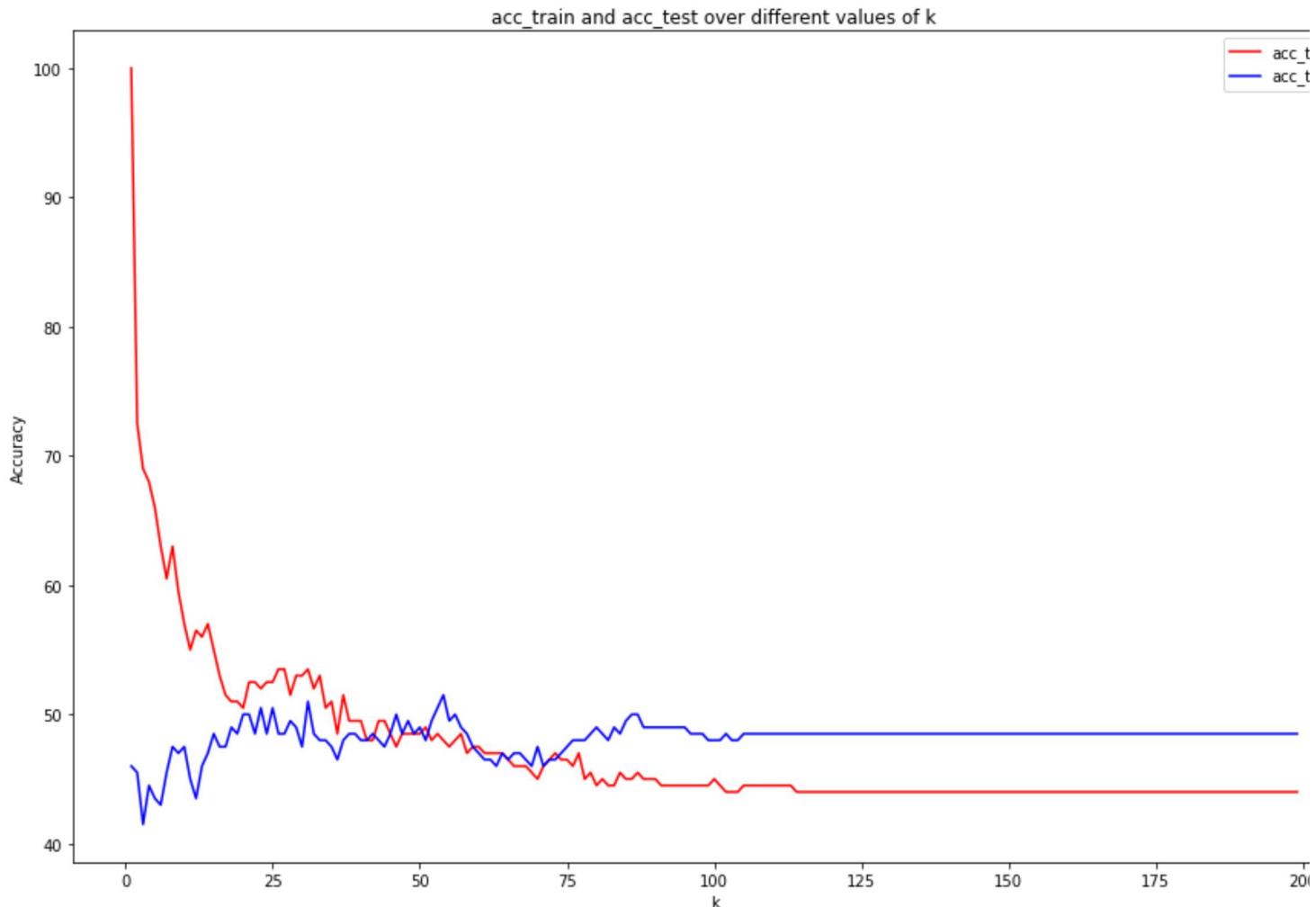
1 k = []
2 acc_train_mat = []
3 acc_test_mat = []
4 for i in range(1,len(Train_Set)):
5     k.append(i)
6     acc_train_mat.append(acc(Train_Set , i))
7     acc_test_mat.append(acc(Test_Set , i))
8
9 f = plt.figure()
10 f.set_figwidth(15)
11 f.set_figheight(10)
12 plt.plot(k, acc_train_mat , color = 'red')
13 plt.plot(k, acc_test_mat , color = 'blue')
14 plt.legend(['acc_train', 'acc_test'])

```

```

15 plt.title('acc_train and acc_test over different values of k')
16 plt.xlabel ('k')
17 plt.ylabel('Accuracy')
18 plt.show()

```



```

1 #Part 2
2 Dataset2 = h5py.File('/content/Lab1_Ex_2_wine.hdf5')
3 Data = np.array(Dataset2.get('Dataset'))
4
5 # Normalizing and saving train data labels
6
7 Train_Set = Data[:200,:]
8 Train_Set = pd.DataFrame(Train_Set)
9 y_train = Train_Set[11]
10 y_train = np.array(y_train.values)
11 Train_Set = Train_Set.drop(columns=11)
12 Train_Set=(Train_Set-Train_Set.mean())/Train_Set.std()
13 Train_Set = np.array(Train_Set.values)
14
15 # Normalizing and saving test data labels
16
17 Test_Set = Data[200:,:]
18 Test_Set = pd.DataFrame(Test_Set)
19 y_test = Test_Set[11]
20 y_test = np.array(y_test.values)

```

```

21 Test_Set = Test_Set.drop(columns=11)
22 Test_Set = (Test_Set - Test_Set.mean()) / Test_Set.std()
23 Test_Set = np.array(Test_Set.values)
24
25 from sklearn import linear_model
26 clf = linear_model.LinearRegression()
27
28 # Training the model
29 clf.fit(Train_Set , y_train)
30 Test_predicted = np.print(clf.predict(Test_Set))
31 Train_predicted = np.print(clf.predict(Train_Set))
32
33 #To be completed by the student

1 # Calculating Acc
2 def accuracy (test , predict):
3     true = 0
4     for i in range(len(test)):
5         if test[i] == predict[i]:
6             true += 1
7
8     return (true / len(test))* 100
9
10 acc_test = accuracy(Test_predicted, y_test)
11 print(acc_test)
12 acc_train = accuracy(Train_predicted, y_train)
13 print(acc_train)

45.0
47.0

```

Student's comments to exercise 2

- KNN

1. We normalize data in order to avoid KNN assigning higher weights to the features with a higher amplitude.
2. The k value in the k-NN algorithm defines how many neighbors will be checked to determine the classification of a specific query point and it is therefore very intensive to the value of k and as k increases the accuracy is decreasing.(The initial increase is for k = 0 which is not important)
3. The figure shows that there's a huge overfitting for the k < 25 while underfitting happens for most of the values of k > 50. The total accuracy is not acceptable. The reason probably is Unbalanced class distribution, meaning that some classes are underrepresented, hence KNN suffering.

- Linear Regression

1. We basically use Linear regression to regress some continuous parameter based on a series of other parameters; hence, linear regression (and descretizing its result using rounding) is not the best solution for classification (discrete).

Exercise 3: Phoneme Dataset

In this exercise the Phoneme dataset is examined <https://catalog.ldc.upenn.edu/LDC93s1>. Each line represents 256 samples gathered at a 16 kHz of different speech signals. The objective is to classify whether the sound emitted is a "sh", "iy", "dcl", "aa", "ao" phoneme.

Again, a subset of the dataset containing 400 elements is provided. Create a training set and a test set of 200 samples each.

- Classify the samples which compose the test set using the k-NN algorithm and evaluate the prediction accuracy for different values of k. Create a graph using the matplotlib library showing the evolution of the accuracy for different values of k over the test set.
- Identifying overfitting and underfitting in the obtained results.

```
1 Dataset3 = h5py.File('/content/Lab1_Ex_3.hdf5')
2
3 Data = np.array(Dataset3.get('Dataset'))
4
5 # Normalizing data
6 Train_Set = Data[:200,:]
7
8 Train_Set = pd.DataFrame(Train_Set)
9 y = Train_Set[256]
10 Train_Set = Train_Set.drop(columns=256)
11 Train_Set=(Train_Set-Train_Set.mean())/Train_Set.std()
12 Train_Set.insert(loc=256, column=256 ,value= y)
13 Train_Set = Train_Set.values
14
15 # Normalizing data
16 Test_Set = Data[200:,:]
17 Test_Set = pd.DataFrame(Test_Set)
18 y = Test_Set[256]
19 Test_Set = Test_Set.drop(columns=256)
20 Test_Set = (Test_Set-Test_Set.mean())/Test_Set.std()
21 Test_Set.insert(loc=256, column=256 ,value= y)
22 Test_Set = Test_Set.values
23
24
25 # Ecuc dist
26
27 def eucdist (sample, data):      #   - Taking a sample
28     dist = 0
29     for i in range(len(data)-1):
30         dist += (sample[i]-data[i])**2
31     return np.sqrt(dist)
32
33 # Defining knn
34
35 def k_nn (sample,k):
36     knn_mat = Train_Set
37     dist=[]
38     for j in range (len(knn_mat)):
39         dist.append(eucdist(sample,knn_mat[j]))
40
41     knn_mat = pd.DataFrame(knn_mat)
42     knn_mat.insert(loc=257, column=257, value=dist)
43     knn_mat = knn_mat.sort_values(257)
44     knn = knn_mat.values[:k]
```

```

45
46 # class prediction section
47
48     classmat = np.array(range(5))
49     predictmat = np.zeros(5)
50
51     for k in range(len(knn)):
52         for j in classmat:
53             if knn[k,256] == j :
54                 predictmat[j] += 1
55
56     return np.argmax(predictmat)
57
58
59
60
61

```

1. We normalize data in order to avoid KNN assigning higher weights to the features with a higher amplitude
2. We can use the same euclid function as the previous exercise
3. We also can use the same KNN function but with the following changes:

- Changing the label column (256 instead of 11)
- Modifying the class prediction section with respect to the number of the classes (5)

```

1 # Calculating Acc
2 def acc (test , k):
3     true = 0
4     test_predicted = 0
5     for i in range(len(test)):
6         test_predicted = (k_nn(test[i] ,k))
7         if test [i][256] == test_predicted:
8             true += 1
9
10    return (true / len(test ))* 100
11
12 acc_test =  acc(Test_Set, 5)
13 print('Test Accuracy:' ,acc_test)
14 acc_train =  acc(Train_Set, 5)
15 print('Train Accuracy:' ,acc_train)
16

```

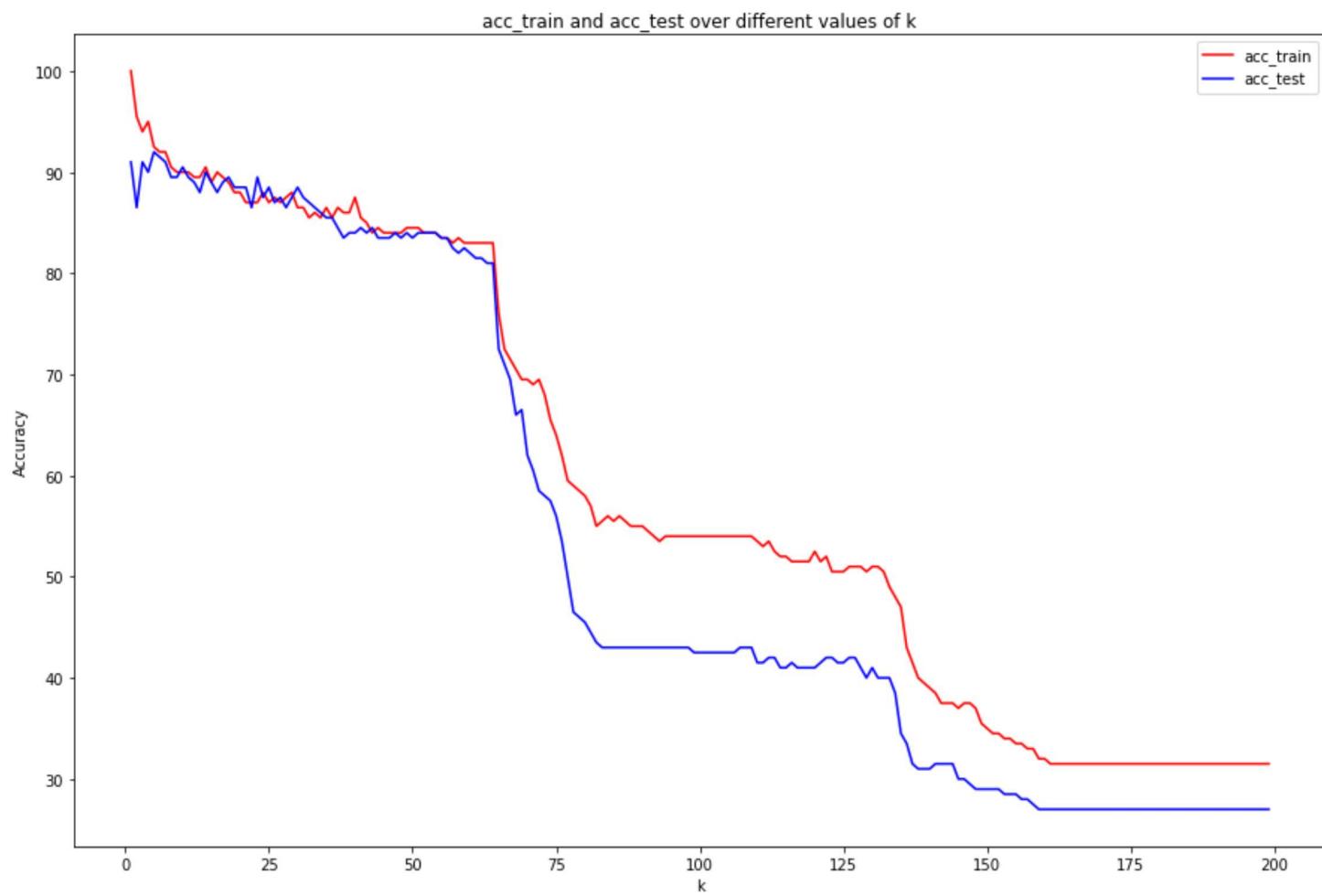
Test Accuracy: 92.0
 Train Accuracy: 92.5

- We measure the accuracy by comparing the KNN prediction with the true label of the data
-
1. The accuracy values are acceptable
 2. The small difference between the testing and training datasets' accuracy doesn't demonstrate any overfitting

```

1 k = []
2 acc_train_mat = []
3 acc_test_mat = []
4 for i in range(1,len(Train_Set)):
5     k.append(i)
6     acc_train_mat.append(acc(Train_Set , i))
7     acc_test_mat.append(acc(Test_Set , i))
8
9 f = plt.figure()
10 f.set_figwidth(15)
11 f.set_figheight(10)
12 plt.plot(k, acc_train_mat , color = 'red')
13 plt.plot(k, acc_test_mat , color = 'blue')
14 plt.legend(['acc_train', 'acc_test'])
15 plt.title('acc_train and acc_test over different values of k')
16 plt.xlabel ('k')
17 plt.ylabel('Accuracy')
18 plt.show()

```



Student's comments to exercise 3

1. KNN is a very slow algorithm (Since It must calculate the distance of a sample point from all of the other points in the training dataset) - Specially in cases like this excersise with many features(256), the calculation time can be problematic. (48 minutes execution time for this block)

2. The k value in the k-NN algorithm defines how many neighbors will be checked to determine the classification of a specific query point.
3. It is therefore very intensive to the value of k and as k increases the accuracy is decreasing.
4. On the one hand, a small value of k, means that noise will have a higher influence on the result. On the other hand, if k is selected to be too large, the model becomes too generalized and fails to accurately predict the data points in both train and test sets. This situation is known as underfitting.

Exercise 1 – Model fitting for continuous distributions: Multivariate Gaussian

In this exercise, you will employ a dataset based on the classic dataset *Iris Plants Database* <https://archive.ics.uci.edu/ml/datasets/iris>. You will be provided a subset of this dataset comprising only two classes (*Iris Setosa* and *Iris Versicolour*), and only two features per class (*petal length* in cm and *petal width* in cm). The objective is to determine the kind of iris based on the content of the features.

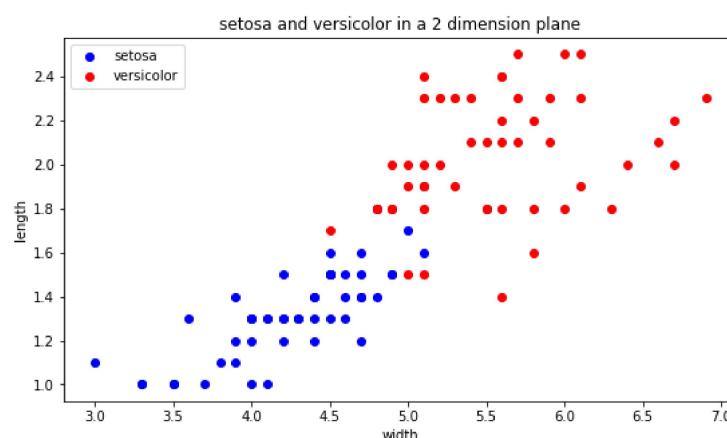
Task: you have to fit class-conditional Gaussian multivariate distributions to the data, and visualize the probability density functions. In particular, you should perform the following:

- Divide the dataset in two parts (*Iris Setosa* which corresponds to class zero and *Iris Versicolour* which correspond to class 1). Then work only on one class at a time.
- Plot the data of each class (use the `plt.scatter()` function)
- Visualize the histogram of petal length and petal width (use e.g. the `plt.hist()` function)
- Calculate the maximum likelihood estimate of the mean and covariance matrix under a multivariate Gaussian model, independently for each class (these are the parameters of the class-conditional distributions). Note: is the Gaussian model good for these data?
- Visualize the 2-D joint pdf of petal length and width for the two classes.

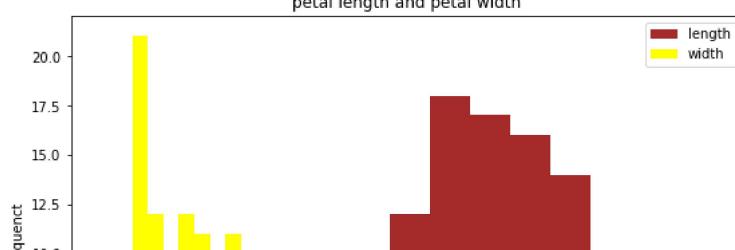
```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import h5py
4 Dataset1 = h5py.File("/content/Lab2_Ex_1_Iris.hdf5")
5 data = np.array(Dataset1.get('Dataset'))
```

New section

```
1 #Separate the dataset in the two classes, you can use the numpy function argsort and unique to do this.
2 indx0 = np.where(data[:,2] == 0) # Returns a list of indices of the rows which their third element is equal to zero(creates a filter)
3 indx1 = np.where(data[:,2] == 1) #
4 caption0 = 'Setosa'
5 caption1 = 'Versicolour'
6 setosa = data[indx0][:,:2] # uses the index0 list to separate (filter) the rows which their third element is equal to zero and slice the first 2 columns
7 versicolour = data[indx1][:,:2] #
8 #print (20 *'*_',f'{caption0}', 20 *'*_','\n', setosa)
9 #print (20 *'*_',f'{caption1}', 20 *'*_','\n', versicolour)
10 %%# Draw the scatter plot of the two classes on the same image
11 plt.figure(figsize=(9,5))
12 setosa_petal_length = setosa[:,0]
13 setosa_petal_width = setosa[:,1]
14 versicolour_petal_length = versicolour[:,0]
15 versicolour_petal_width = versicolour[:,1]
16 plt.scatter(setosa_petal_length, setosa_petal_width , c = 'blue' )
17 plt.scatter(versicolour_petal_length, versicolour_petal_width, c = 'red')
18 labels = ['setosa' , 'versicolor']
19 plt.legend(labels , loc='upper left')
20 plt.xlabel('width')
21 plt.ylabel('length')
22 plt.title("setosa and versicolor in a 2 dimension plane")
23 plt.show()
24 #Draw the scatter plot of the two classes on the same image
```



```
1 #Visualize the histogram of petal length and petal width (use e.g. the plt.hist() function)
2 plt.figure(figsize=(9,6))
3 total_petal_length = np.hstack((setosa_petal_length, versicolour_petal_length))
4 total_petal_width = np.hstack((setosa_petal_width, versicolour_petal_width))
5 plt.hist(total_petal_length , color= 'brown', histtype='bar', label='length')
6 plt.hist(total_petal_width, color = 'yellow', histtype='bar', label= 'width')
7 plt.legend()
8 plt.xlabel('amplitude')
9 plt.ylabel('frequency')
10 plt.title("petal length and petal width")
11 plt.show()
```



The figure demonstrate sum of two random variables in a gaussian-shape figure

```

1 #Calculate mean and covariance matrix under a multivariate Gaussian model. Scalar products can be computed with the function np.matmul()
2
3 #The transpose can be obtained with the function np.transpose
4
5 #To make a scalar product between arrays in the form [Mx1]x[1xM] starting from 1D array A it may be necessary to add a new axis using
6 #the function A[:,np.newaxis]
7 def mean_calculator (x):
8     x_mean = [0,0]
9     for j in range(x.shape[1]):
10        for i in range(x.shape[0]):
11            x_mean[j] += x[i,j]
12        x_mean[j] /= x.shape[0]
13    return x_mean
14 ###
15 def cov_calculator(y):
16     y_mean = np.array(mean(y))
17     y_cov = np.empty([len(y.shape), len(y.shape)])
18     for k in range(y.shape[0]):
19         y_cov += (np.subtract( y[k][:, np.newaxis] @ y[k][:, np.newaxis].T , y_mean[:, np.newaxis] @ y_mean[:, np.newaxis].T))
20     return np.divide(y_cov,y.shape[0])
21 ###
22 setosa_mean = mean_calculator(setosa)
23 setosa_cov = cov_calculator(setosa)
24 versicolour_mean = mean_calculator(versicolour)
25 versicolour_cov = cov_calculator(versicolour)
26
27 print('      MLE estimate:\n\n')
28
29 print("* setosa:\n\n", 'mean: ', "\n", mean_calculator(setosa), "  ", "\n\n", 'cov: ', "\n", cov_calculator(setosa))
30 print("  ")
31 print("\n* versicolour:\n\n", 'mean: ', "\n", mean_calculator(versicolour), "  ", "\n\n", 'cov: ', "\n", cov_calculator(versicolour))

      MLE estimate:

* setosa:
mean:
[4.26, 1.3259999999999998]

cov:
[[0.73574521 0.19258053]
 [0.19258053 0.15136837]]
```

```

* versicolour:
mean:
[5.552, 2.025999999999994]

cov:
[[0.3132109 0.05169961]
 [0.05169961 0.07695137]]
```

```

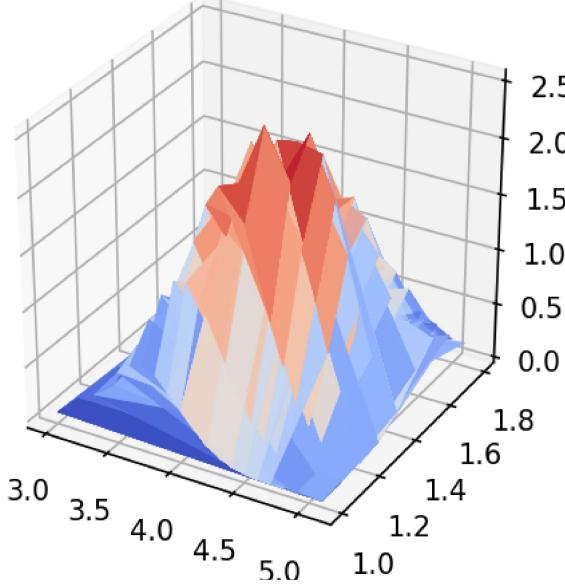
1
2 #Visualize the 2-D joint pdf of petal length and width, a pdf function can be initialized by providing mean and covariance matrix
3 from scipy.stats import multivariate_normal
4 from matplotlib import cm
5 from matplotlib.ticker import LinearLocator
6
7
8 pdf_class0 = multivariate_normal(mean=mean_calculator(setosa), cov= cov_calculator(setosa)) #class0 = setosa
9 #pdf_class1 = multivariate_normal(mean=mean(versicolour), cov=cov(versicolour)) #class1 = versicolour
10 #Create a grid of x and y values on which to sample the pdf, this is done by providing a list of x-y of coordinates to the function pdf_class0.pdf(...)
11 #A 3D view of the pdf can be obtained using the function ax.plot_surface
12
13 #Code Example:
14
15 X = setosa_petal_length
16 Y = setosa_petal_width
17
18 X, Y = np.meshgrid(X, Y)
19 X_flat = X.flatten()
20 Y_flat = Y.flatten()
21 XY_list = np.concatenate((X_flat[:,np.newaxis],Y_flat[:,np.newaxis]),axis=1)
22 PDF_values = np.reshape(pdf_class0.pdf(XY_list), np.shape(X))
23
24 fig1, ax1 = plt.subplots(subplot_kw={"projection": "3d"}, figsize=(4, 4), dpi=150)
25 ax1.plot_surface(X, Y, PDF_values, cmap=cm.coolwarm, alpha=0.7, linewidth=1)
26 plt.title('setosa 2-D joint pdf of petal length and width', fontsize=7)
27 plt.show
28
29
30
31 pdf_class1 = multivariate_normal(mean=mean_calculator(versicolour), cov = cov_calculator(versicolour))
32
33 XX = versicolour_petal_length
34 YY = versicolour_petal_width
35
36 XX, YY = np.meshgrid(XX, YY)
37 XX_flat = XX.flatten()
38 YY_flat = YY.flatten()
```

```

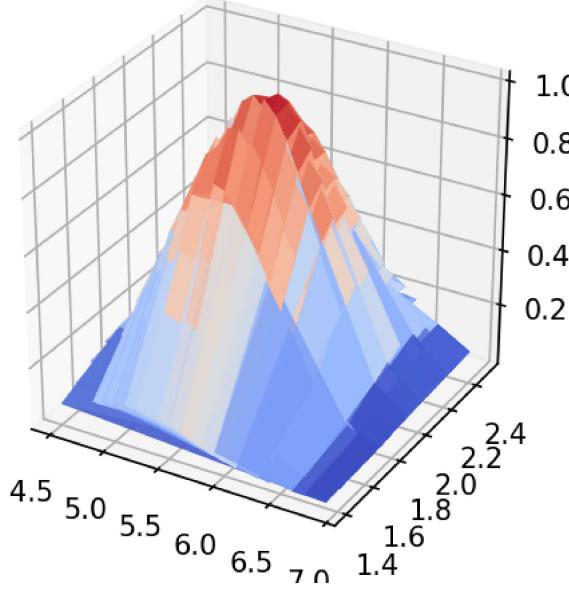
39 XYXY_list = np.concatenate((XX_flat[:,np.newaxis],YY_flat[:,np.newaxis]),axis=1)
40 PDF_values = np.reshape(pdf_class1.pdf(XYXY_list), np.shape(XX))
41
42 fig2, ax2 = plt.subplots(subplot_kw={"projection": "3d"}, figsize=(4, 4), dpi=150)
43 ax2.plot_surface(XX, YY, PDF_values, cmap=cm.coolwarm, alpha=0.7, linewidth=1)
44 plt.title('versicolor 2-D joint pdf of petal length and width', fontsize=7)
45 plt.show()
46 #To change the orientation of the 3D plot the function ax.view_init(), for a view from above select ax.view_init( 90, 0)
47
48 #After visualizing the pdf you can plot the points of the dataset on the estimated pdf using ax.scatter3D()
49 #For a better visualization of the points we suggest to make the pdf plot semi-transparent using the alpha parameter
50
51 #Code Example:
52 #PDF_points_class0 = pdf_class0.pdf(setosa_petal_length)
53 #plt.axes.scatter3D(setosa_petal_length, setosa_petal_width, PDF_points_class0, s=10)
54
55 #Note: the sample code was written only for class 0 but two plots have to be done, one for class 0 and one for class 1
56 #Note 2: the content of variables like Features_Class0_Points_Class0_Feature0_Points_Class0_Feature1 should be substituted with variables created by t
<function matplotlib.pyplot.show(*args, **kw)

```

setosa 2-D joint pdf of petal length and width



versicolor 2-D joint pdf of petal length and width



Student's comments to Exercise 1

1. The joint probability density functions of the two classes are depicted
2. Setosa shape has a problem in coding since the probability is more than 1!!!!
3. The figure is a 3D bellshape plane since it is the combination of 2 different bell-shape gaussian MLE
4. We can clearly see that the probability is maximum for the mean value of the data in both of the datasets

Exercise 2 - Model fitting for discrete distributions: Bag of Words

In this exercise, you will employ a real dataset (file *SMSSpamCollection*). The SMS Spam Collection v.1 (<https://www.kaggle.com/datasets/uciml/sms-spam-collection-dataset>) is a set of SMS tagged messages that have been collected for SMS Spam research. It contains one set of SMS messages in English of 5,574 messages, tagged according being ham (legitimate) or spam. Task: you have to fit the parameters employed by a Naïve Bayes Classifier, using a Bernoulli model. Under this model, the parameters are:

- π_c , the prior probabilities of each class.
- θ_{jc} , the probability that feature j is equal to 1 in class c.

Model fitting can be done using the pseudocode at the end of the Lecture 3 slides.

Display the class-conditional densities θ_{jc1} and θ_{jc2} . Try to identify “uninformative” features (i.e., features j such that $\theta_{jc1} \approx \theta_{jc2}$).

```

1 import numpy as np
2 import sklearn
3 from sklearn.feature_extraction.text import CountVectorizer
4 from sklearn import metrics
5 import matplotlib.pyplot as plt
6 import pandas as pd
7 import seaborn as sns
8
9 #from NLP_base_library import text_splitting, print_dtm
10 import sys

```

```

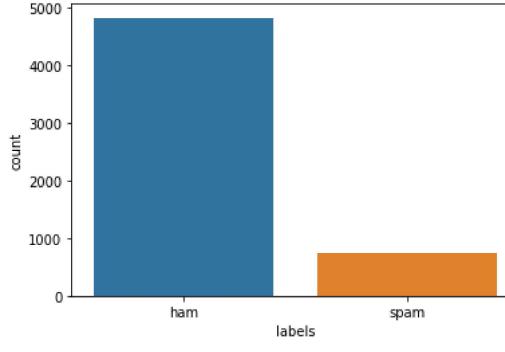
11 import matplotlib.pyplot as plt
12
13 data=pd.read_csv('/content/SMSSpamCollection', encoding = "ISO-8859-1",sep='\t',header=None)
14 data.rename(columns = {0:'labels', 1:'text'}, inplace=True)
15
16 sns.countplot(data['labels'])
17
18 #Dataset Preprocessing: conversion from a collection of sentences into a Bag_of_words representation
19 from nltk.corpus import stopwords
20 from nltk.stem.porter import PorterStemmer
21 import re
22 import nltk
23
24 ps = PorterStemmer()
25 nltk.download('stopwords')
26
27 from nltk.corpus import stopwords
28 from nltk.stem.porter import PorterStemmer
29 ps = PorterStemmer()
30
31 corpus = []
32 for i in range(0, len(data)):
33     review = re.sub('[^a-zA-Z]', ' ', data['text'][i]) #remove all characters which are not letters
34     review = review.lower() #set everything to lower case
35     review = review.split() #split into words
36     #remove stopwords (articles and other uninformative words), replace words with their "stem" (basic form without conjugations and pluralizations )
37     review = [ps.stem(word) for word in review if not word in stopwords.words('english')]
38     review = ' '.join(review)
39     corpus.append(review)
40
41 from sklearn.feature_extraction.text import CountVectorizer
42 cv = CountVectorizer(max_features=2500)
43
44
45
46 X = cv.fit_transform(corpus).toarray() #Convert dataset into Bag of words using the most informative 2500 words
47 feature_names = cv.get_feature_names_out()
48 print(f"Some examples of words which are used as features: {feature_names}") #This vector of strings contain the kind of words that are used as features
49
50 y = pd.get_dummies(data['labels'])
51 y = y.iloc[:,1].values #convert the name of the two classes ("ham" and "spam") into numeric values "0" and "1"
52
53 #Partition Training and Testing set
54 #The training set is composed of 2000 SMS with 2500 features each, which identify whether a word from dictionary is present or not following the BoW representation
55 X_train = X[:2000,:]
56
57 #y_train is composed of the 2000 labels for the training examples, indicating whether a message is real (class 0) or spam (class 1)
58 y_train = y[:2000]
59
60 X_test = X[2000:3000,:] #1000 SMS are used in the testing set
61 y_test = y[2000:3000]

```

```

/usr/local/lib/python3.8/dist-packages/seaborn/_decorators.py:36: FutureWarning: Pass the following
warnings.warn(
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
Some examples of words which are used as features: ['aah' 'aathi' 'abi' ... 'zed' 'zindgi' 'zoe']
(2000, 2500)

```



```

1 from google.colab import drive
2 drive.mount('/content/drive')

Mounted at /content/drive

1 #Evaluate the probabilities of the two classes, and the class conditional densities.
2 Nc_spam = 0 #Number of the spam messages (NC for spam)
3 Nc_ham = 0 #Number of the legit messages (NC for legit)
4 for i in range (len(y_train)):
5     if y_train[i] == 1 : Nc_spam += 1
6     else: Nc_ham += 1
7
8 ntotal = len(y_train) #Total number of the messages
9 prior_spam = Nc_spam/ntotal #Calculationg prior spam
10 prior_ham = Nc_ham/ntotal #Calculationg prior ham
11
12 # X (bag of word) contains the total number of the words (features)
13 # theta jc for the calss "spam" :
14
15
16 #claculating Njc of each class (Njc is a number for each feature and a vector of these numbers for a class)
17 Njc_class0 = []
18 Njc_class1 = []
19
20 for i in range(X_train.shape[1]): #Looping through columns(words/features) of X (word's BoW)
21     Njc0 = 0
22     Njc1 = 0
23     for j in range (len(X_train[:,i])): #Looping through the elements of column i (word i) of matrix X
24         if X_train[:,i][j] == 1 and y_train[j] == 0 :
25             Njc0 += 1
26             if X_train[:,i][j] == 1 and y_train[j] == 1 :

```

```

27     Njc1 += 1
28     Njc_class0.append(Njc0)
29     Njc_class1.append(Njc1)
30
31 #calculating Nc of each class (Nc is a int number) (summing all the NJc for each class)
32 Nc_class0 = 0
33 Nc_class1 = 0
34
35 for i in range (len(Njc_class0)):
36     Nc_class0 += Njc_class0[i]
37 for i in range (len(Njc_class1)):
38     Nc_class1 += Njc_class1[i]
39
40 #####print (Nc_class0 , Nc_class1)
41
42 #####print (len(Njc_class0), Njc_class0)
43 #####print (len(Njc_class1), Njc_class1)
44
45 #Calculating tetajc matrix by dividing each elemen of Njc by Nc
46 tetaj_class0 = np.divide(Njc_class0, Nc_class0)
47 tetaj_class1 = np.divide(Njc_class1, Nc_class1)
48
49 #####print (tetaj_class0)
50 # print (tetaj_class1)
51 # Finding uninformative tetajc by setting a minimum distance condition between tetaj_class 0 and tetaj_class1
52
53 uninformative = []
54 a = len(tetaj_class0)
55 i = 0
56 while i < a:
57     if np.abs(tetaj_class0[i] - tetaj_class1[i]) < (tetaj_class1[i]/10):
58         uninformative.append(i)
59         a -= 1
60     i += 1
61
62
63 print('Total number of the uninformative features:', len(uninformative),
64      "\nThe uninformative features' indexes:", uninformative,
65      "\nTotal Number of informative features:", X_train.shape[1] - len(uninformative))
66
67
68 print('prior_spam:', prior_spam , 'prior_legit:', prior_ham) # printing the prior of each class
    Total number of the uninformative features: 32
    The uninformative features' indexes: [165, 167, 194, 228, 264, 604, 607, 657, 703, 707, 798, 816, 840, 901, 937, 971, 983, 1262, 1299, 1362, 1456, 1476]
    Total Number of informative features: 2468
    prior_spam: 0.14 prior_legit: 0.86

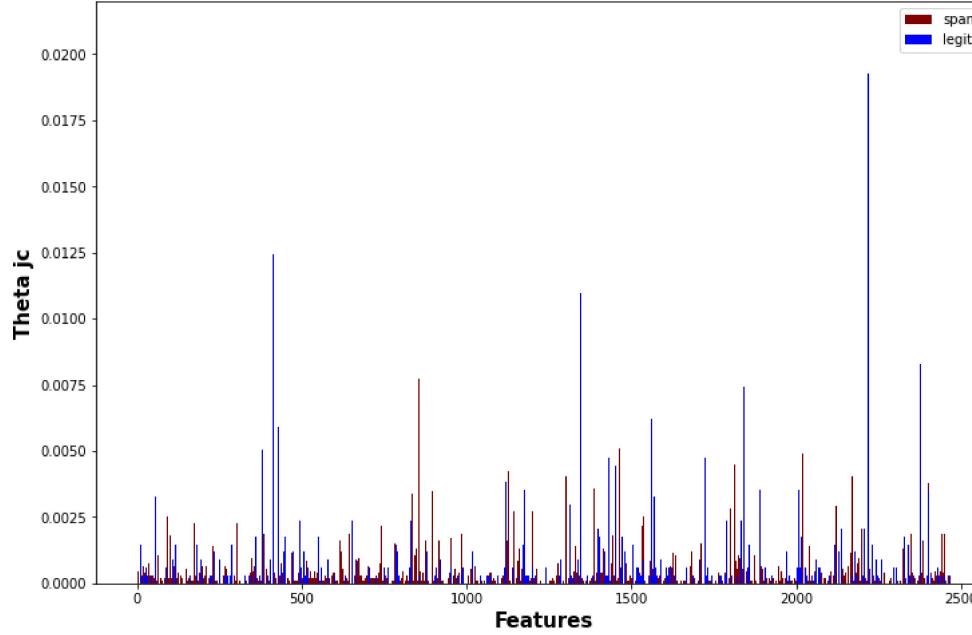
```

◀ | ▶

```

1 #Display the class-conditional densities θjc1 and θjc2.
2 barWidth = 1
3 fig = plt.subplots(figsize =(12, 8))
4 bin1 = np.arange(len(tetaj_class0))
5 bin2 = [x + barWidth for x in bin1]
6 plt.bar(bin1, tetaj_class0, color ='maroon', width = barWidth , label='spam')
7 plt.bar(bin2, tetaj_class1, color ='blue', width = barWidth, label = 'legit')
8 plt.xlabel('Features', fontweight = 'bold', fontsize = 15)
9 plt.ylabel('θjc', fontweight = 'bold', fontsize = 15)
10 plt.legend()
11 plt.ylim(0, 0.022)
12 plt.show()

```



Student's comments to Exercise 2

- Naive bayes classifiers are generative and supervised classifiers which are based on bayes theorem. Since NBC assumes that features in a particular class are totally independent ► *In text recognition problems, they can not understand the order of the words *
- In BoW problems we consider the total space of the n words we have as the features (each word is a feature) and the documents to be classified (here SMS) is a n dimensional vector
- Here we have defined 2 classes (legit and spam)
- To model the data using NBC, We use bernouli distribution for fitting into the data to find prior and theta jc
- The figure shows θ_{jc0} and θ_{jc1} Which corresponds to the probability of the words belonging to each of the classes , It is clear that most of the word have a much higher chance to appear in a legit sms
- To find **uninformative features** we can set aa minimum distance condition between tetaj_class 0 and tetaj_class1 : $|\theta_{jc0} - \theta_{jc1}| < \theta_{jc0}$

Exercise 3 - Classification – discrete data

In this exercise, you will design a Naïve Bayes Classifier (NBC) for the Bag of Words (BoW) features for document classification that have been prepared in Exercise 2. In particular, in exercise 2, you have already estimated the following parameters:

- The prior probabilities of each class, $\pi_c = p(y = c)$.
- The class-conditional probabilities of each feature, $\theta_{jc} = p(x_j = 1|y = c)$.

These parameters have been estimated from the training data. In this exercise, you will use the test data, and classify each test vector using an NBC whose model has been fitted in Exercise 2. In particular, you will do the following:

- For each test vector, calculate the MAP estimate of the class the test vector belongs to. Remember: the MAP classifier chooses the class that maximizes $\max_c \log p(y = c|x) \propto \log p(x|y = c) + \log p(c)$. In the NBC, the features (i.e. each entry of x) are assumed to be statistically independent, so $p(x|y = c) = \prod_{j=1}^D p(x_j|y = c)$. This formula allows you to calculate $p(x|y = c)$ for a given test vector x using the parameters θ_{jc} already calculated in Exercise 2. Note that, after the logarithm, the product becomes a summation. It is much better to use the logarithm in order to avoid underflow.
- See how the accuracy changes when the prior is not taken into account (e.g. by comparing the MLE and MAP estimate).
- After classifying a test vector using the NBC, the obtained class can be compared with the truth (vector y_{test}).
- The accuracy of the classifier can be computed as the percentage of times that the NBC provides the correct class for a test vector.
- Repeat the same operations using the training data as test data, and compare the accuracy of the classifier on the training and test data.
- Note: It is expected that students implement the Naive Bayes classifier from scratch without using pre-made functions such as `sklearn.naive_bayes`

Optional:

If you plot the class-conditional densities as done at the end of Exercise 2, you will see that many features are uninformative; e.g., words that appear very often (or very rarely) in documents belonging to either class are not very helpful to classify a document. The NBC can perform a lot better if these uninformative features are disregarded during the classification, i.e. only a subset of the features, chosen among the most informative ones, are retained. To rank the features by “significance”, one can employ the mutual information between feature x_j and class y (see Sec. 3.5.4 of the textbook):

$$I(X, Y) = \sum_{x_j} \sum_y p(x_j, y) \log \frac{p(x_j, y)}{p(x_j)p(y)}$$

For binary features, the mutual information of feature j can be written as:

$$I_j = \sum_c \left[\theta_{jc} \pi_c \log \frac{\theta_{jc}}{\theta_j} + (1 - \theta_{jc}) \pi_c \log \frac{1 - \theta_{jc}}{1 - \theta_j} \right]$$

with $\theta_j = p(x_j = 1) = \sum_c \pi_c \theta_{jc}$. For this part, you should:

- Calculate I_j for all features. Note: try to avoid divisions by zero adding the machine precision constant `eps` to the denominators.
- Rank the features by decreasing values of I_j , and keep only the K most important ones.
- Run the classifier employing only the K most important features, and calculate the accuracy.
- Plot the accuracy as a function of K .

```
1 # @title
2 # Evaluate the MAP on the test set
3 # MAP : Max log P(y=c|x) ~ log p(x|y = c) + log P(c)
4 # Calculating P(x|y=c) = n P(Xj|y=c)
5 #ham = 0      #spam = 1
6
7 MAP_sentences_0 = []
8 MAP_sentences_1 = []
9 MAP_sentence_i_0 = 0
10 MAP_sentence_i_1 = 0
11 ##### Test #####
12 for i in range(len(X_test)):
13     pxj0 = 1
14     pxj1 = 1
15     for j in range(len(X_test[i])):
16         if X_test[i][j] == 1:
17             if tetaj_class0[j] != 0:
18                 pxj0 *= tetaj_class0[j]
19             if tetaj_class1[j] != 0:
20                 pxj1 *= tetaj_class1[j]
21     MAP_sentence_i_0 = np.log(pxj0) + np.log(prior_ham)
22     MAP_sentences_0.append(MAP_sentence_i_0)
23     MAP_sentence_i_1 = np.log(pxj1) + np.log(prior_spam)
24     MAP_sentences_1.append(MAP_sentence_i_1)
25 y_hat_test = []
26 #print ('MAP0',MAP_sentences_0)
27 #print ('MAP1',MAP_sentences_1)
28
29 # Calculating MAP estimate of each test sentence (y_hat_test):
30
31 for i in range (len(MAP_sentences_0)):
32     if np.abs(MAP_sentences_0 [i]) > np.abs(MAP_sentences_1 [i]):
33         y_hat_test.append(0)
34     if np.abs(MAP_sentences_0 [i]) < np.abs(MAP_sentences_1 [i]):
35         y_hat_test.append(1)
36
37 # calculating the accuracy of estimation by comparing y_hat_test and y_test :
38 neq=0
39 for i in range (len(y_test)):
40     if y_test[i] == y_hat_test[i]:
41         neq +=1
42
43 print('The Accuracy for MAP on test is equal to:' ,neq/len(y_test)*100 ,'%')
44
45
46 ##### Train #####
47
48 for i in range(len(X_train)):
49     pxj0 = 1
50     pxj1 = 1
```

```

51     for j in range (len(X_train[i])):
52         if X_train[i][j] == 1:
53             if tetaj_class0[j] != 0 :
54                 pxj0 *= tetaj_class0[j]
55             if tetaj_class1[j] != 0:
56                 pxj1 *= tetaj_class1[j]
57             MAP_sentence_i_0 = np.log(pxj0) + np.log(prior_ham)
58             MAP_sentences_0.append(MAP_sentence_i_0)
59             MAP_sentence_i_1 = np.log(pxj1) + np.log(prior_spam)
60             MAP_sentences_1.append(MAP_sentence_i_1)
61 y_hat_train =[]
62 #print ('MAP0',MAP_sentences_0)
63 #print ('MAP1',MAP_sentences_1)
64
65 # Calculating MAP estimate of each test sentence (y_hat_train):
66
67 for i in range (len(MAP_sentences_0)):
68     if np.abs(MAP_sentences_0 [i]) > np.abs(MAP_sentences_1 [i]):
69         y_hat_train.append(0)
70     if np.abs(MAP_sentences_0 [i]) < np.abs(MAP_sentences_1 [i]):
71         y_hat_train.append(1)
72
73 # calculating the accuracy of estimation by comparing y_hat_test and y_test :
74 neq=0
75 for i in range (len(y_train)):
76     if y_train[i] == y_hat_train[i]:
77         neq +=1
78
79 print('The Accuracy for MAP on train is equal to:' ,neq/len(y_train)*100 ,'%')
80
    The Accuracy for MAP on test is equal to: 82.6 %
    The Accuracy for MAP on train is equal to: 70.25 %

1 #Evaluate the MLE on the test set
2
3 MLE_sentences_0 = []
4 MLE_sentences_1 = []
5 MLE_sentence_i_0 = 0
6 MLE_sentence_i_1 = 0
7 ##### Train #####
8 for i in range(len(X_test)):                                # X_test : all of the test sentences (1000)
9     pxj0 = 1
10    pxj1 = 1
11    for j in range (len(X_test[i])):
12        if X_test[i][j] == 1:                               # Looping through the words of sentence i
13            if tetaj_class0[j] != 0 :                      # Checking if the word is available in the sentence
14                pxj0 *= tetaj_class0[j]                      # Calculating P(x|y=c) = n (Xj|y=c) of the class 0 sentence
15            if tetaj_class1[j] != 0:                        # Calculating P(x|y=c) = n (Xj|y=c) of the class 1 sentence
16                pxj1 *= tetaj_class1[j]
17            MLE_sentence_i_0 = np.log(pxj0)              # Calculating MLE of the sentence i for class 0 (ham) : Max log P(y=c|x) ~ log p(x|y =c) ||| **Without |
18            MLE_sentences_0.append(MLE_sentence_i_0)       # Creating a vector of Class 0 Map estimate of each sentence
19            MLE_sentence_i_1 = np.log(pxj1)              # Calculating MLE of the sentence i for class 1 (spam) : Max log P(y=c|x) ~ log p(x|y =c) ||| **Without |
20            MLE_sentences_1.append(MLE_sentence_i_1)       # Creating a vector of Class 1 Map estimate of each sentence
21 y_hat_test =[]
22 #print ('MAP0',MAP_sentences_0)
23 #print ('MAP1',MAP_sentences_1)
24
25 # Calculating MAP estimate of each test sentence (y_hat_test):
26
27 for i in range (len(MLE_sentences_0)):
28     if np.abs(MLE_sentences_0 [i]) > np.abs(MLE_sentences_1 [i]):
29         y_hat_test.append(0)
30     if np.abs(MLE_sentences_0 [i]) < np.abs(MLE_sentences_1 [i]):
31         y_hat_test.append(1)
32     if np.abs(MLE_sentences_0 [i]) == np.abs(MLE_sentences_1 [i]): # The SMS with the same MLE for both classes are considered as legit
33         y_hat_test.append(0)
34
35
36
37 # calculating the accuracy of estimation by comparing y_hat_test and y_test :
38 neq=0
39 for i in range (len(y_test)):
40     if y_test[i] == y_hat_test[i]:
41         neq +=1
42
43 print('The Accuracy for MLE Test is equal to:' ,neq/len(y_test)*100 ,'%')
44
45 ##### Test#####
46 for i in range(len(X_train)):                                # X_train : all of the test sentences (1000)
47     pxj0 = 1
48     pxj1 = 1
49     for j in range (len(X_train[i])):
50         if X_train[i][j] == 1:                               # Looping through the words of sentence i
51             if tetaj_class0[j] != 0 :                      # Checking if the word is available in the sentence
52                 pxj0 *= tetaj_class0[j]                      # Calculating P(x|y=c) = n (Xj|y=c) of the class 0 sentence
53             if tetaj_class1[j] != 0:                        # Calculating P(x|y=c) = n (Xj|y=c) of the class 1 sentence
54                 pxj1 *= tetaj_class1[j]
55             MLE_sentence_i_0 = np.log(pxj0)              # Calculating MLE of the sentence i for class 0 (ham) : Max log P(y=c|x) ~ log p(x|y =c) ||| **Without |
56             MLE_sentences_0.append(MLE_sentence_i_0)       # Creating a vector of Class 0 Map estimate of each sentence
57             MLE_sentence_i_1 = np.log(pxj1)              # Calculating MLE of the sentence i for class 1 (spam) : Max log P(y=c|x) ~ log p(x|y =c) ||| **Without |
58             MLE_sentences_1.append(MLE_sentence_i_1)       # Creating a vector of Class 1 Map estimate of each sentence
59 y_hat_train =[]
60 #print ('MAP0',MAP_sentences_0)
61 #print ('MAP1',MAP_sentences_1)
62
63 # Calculating MAP estimate of each train sentence (y_hat_train):
64
65 for i in range (len(MLE_sentences_0)):
66     if np.abs(MLE_sentences_0 [i]) > np.abs(MLE_sentences_1 [i]):
67         y_hat_train.append(0)
68     if np.abs(MLE_sentences_0 [i]) < np.abs(MLE_sentences_1 [i]):
69         y_hat_train.append(1)
70     if np.abs(MLE_sentences_0 [i]) == np.abs(MLE_sentences_1 [i]): # The SMS with the same MLE for both classes are considered as legit
71         y_hat_train.append(0)
72

```

```

73
74 # calculating the accuracy of estimation by comparing y_hat_train and y_train :
75 neq=0
76 for i in range (len(y_train)):
77     if y_train[i] == y_hat_train[i]:
78         neq +=1
79
80 print('The Accuracy for MLE Train is equal to:' ,neq/len(y_train)*100 ,'%')
81
82 The Accuracy for MLE Test is equal to: 85.9 %
83 The Accuracy for MLE Train is equal to: 73.1 %

```

Student's comments to Exercise 3

Comparing accuracy of the Test and Train:

- The Accuracy of test dataset is more than the accuracy for the training data set, if the difference was small there was no problem, but in this situation with more than 10% difference, it seems that we have underfitting. Maybe happen due to not using all of the feature that is in the data!!!!

Comparing the accuracy of MLE and MAP:

- The Accuracy of MLE is a bit higher than MAP for both the test and training dataset ! We expected more accuracy for MAP since it also observes prior in addition to the likelihood but sometime in large dataset MLE can become becomes more robust and less sensitive to the choice of prior

Exercise 4 – Plotting the ROC curve

For the discrete data classification problem of Exercise 3, analyse the performance of the classifier plotting the complete ROC curve, instead of simply measuring the accuracy. This requires to do the following.

- Instead of classifying the documents choosing class 1 if $p(y = 1|x) > p(y = 2|x)$, now you can generalize this to choosing class 1 if $\frac{p(y=1|x)}{p(y=2|x)} > \tau \in [0, \infty)$ for some threshold τ that determines the compromise between true positive rate (TPR) and false positive rate (FPR).
- Choose a **reasonable** range of values for τ . For each value of τ , compute the TPR and FPR on the dataset (Hint: determine suitable minimum and maximum values for τ , and sample densely enough in that range).
- Plot a curve of TPR as a function of FPR – this is the ROC curve for this classifier.
- Determine an estimate of the Equal Error Rate (EER), i.e. the point of the ROC curve such that TPR+FPR=1.

```

1 from tkinter import N
2 #Use the posterior probabilities previously and classify using the formula above to estimate tau
3
4 # Calculating MAP estimate of each test sentence (y_hat_test):
5
6 def ROC_estimate (a,b,c):
7     mapclass0 = a
8     mapclass1 = b
9     taw=c
10
11    y_hat_test_ROC = []
12    for i in range (len(MAP_sentences_0)):
13        if np.abs(mapclass0 [i]) / np.abs(mapclass1[i]) > taw:
14            y_hat_test_ROC.append(0)
15        if np.abs(mapclass0 [i]) / np.abs(mapclass1[i]) < taw:
16            y_hat_test_ROC.append(1)
17
18    return y_hat_test_ROC
19
20 # calculating the accuracy of estimation by comparing y_hat_test and y_test :
21
22 def Accuracy (a,b):
23     y = a
24     yhat = b
25     neq=0
26     for i in range (len(y)):
27         if y[i] == yhat[i]:
28             neq +=1
29     acc = neq/len(y_test)*100
30     return acc
31
32
33
34
35 # Function to measure TPR
36 def TPR (a,b): #based on "spam as the true"
37     test = a # test is what the classifier says about the class of the test points in the test dataset
38     true = b # true is the correct classes of the test points in the test dataset
39     ntrue = 0
40     for i in range (len(true)): # Calculating the total number of y=1
41         if true [i] == 0:
42             ntrue += 1
43
44     npos = 0
45     for i in range (len(true)):
46         if true [i] == 0 and test [i] == 0:
47             npos += 1
48
49     tpr = npos / ntrue
50     return tpr
51
52 """
53 # Function to measure FPR
54 def FPR (a,b): #based on "spam as the true"
55     test = a # test is what the classifier says about the class of the test points in the test dataset
56     true = b # true is the correct classes of the test points in the test dataset
57     ntrue = 0
58     for i in range (len(true)): #calculating the total number of y=0

```

```

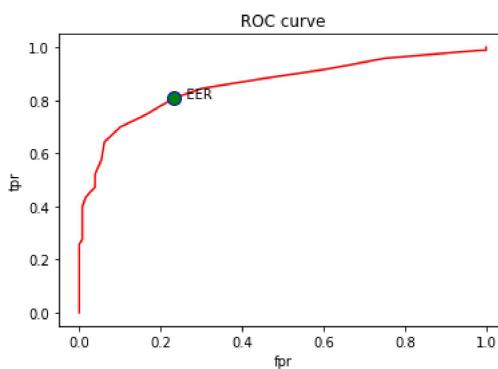
59     if true [i] == 1:
60         ntrue += 1
61
62     npos = 0
63     for i in range (len(true)):
64         if true [i] == 1 and test [i] == 0:
65             npos += 1
66
67     fpr = npos / ntrue
68     return fpr
69
70 tpr = TPR (ROC_estimate(MAP_sentences_0, MAP_sentences_1, 15) , y_test )
71
72 fpr = FPR (ROC_estimate(MAP_sentences_0, MAP_sentences_1, 15) , y_test )
73
74 # Generating a range of taw for calculating tpr and fpr based on the different amounts of taw
75 taw_range = np.arange(0,18,0.1)
76 # Calculating tpr and fpr for different values of taw
77 tpss = []
78 fpss = []
79 EER = []
80 EERpoint = []
81
82 for i in taw_range:
83     tpr = TPR (ROC_estimate(MAP_sentences_0, MAP_sentences_1, i) , y_test )
84     tpss.append(tpr)
85     fpr = FPR (ROC_estimate(MAP_sentences_0, MAP_sentences_1, i) , y_test )
86     fpss.append(fpr)
87     # Finding EER point where tpr + fpr = 1
88     if np.abs(TPR (ROC_estimate(MAP_sentences_0, MAP_sentences_1, i) , y_test ) + FPR (ROC_estimate(MAP_sentences_0, MAP_sentences_1, i) , y_test ) - 1) < 0.05:
89         EER.append(i) #saving the taw of EER
90         EERpoint.append (FPR (ROC_estimate(MAP_sentences_0, MAP_sentences_1, i) , y_test)) # saving the FPR in which EER happens
91         EERpoint.append (TPR (ROC_estimate(MAP_sentences_0, MAP_sentences_1, i) , y_test )) # saving the TPR in which EER happens
92
93
94 # printing the data
95 print ('\n')
96 print ('The taw in which EE happens: ', EER[0])
97 print ('\n')
98 print ('FPR and TPR of taw = EER is equal to :',EERpoint)
99 print ('\n')
100 print ('The Accuracy for taw = EER is equal to:' ,Accuracy(y_test , ROC_estimate(MAP_sentences_0, MAP_sentences_1, EER[0])) ,'%')
101
102 # plotting the ROC curve
103 print ('\n\n')
104 plt.plot(fpss,tpss, color = "r")
105 plt.xlabel("fpr")
106 plt.ylabel("tpr")
107 plt.title('ROC curve')
108 plt.plot(EERpoint[0], EERpoint[1], marker="o", markersize=10, markeredgecolor="blue", markerfacecolor="green")
109 plt.annotate(" EER", (EERpoint[0], EERpoint[1]))
110 plt.show()
111
112
113
114 #Note: To estimate the TPR you need to compute the number of cases where class 1 is correctly predicted,
115 #this value then has to be divided by the number of elements in the test set which belong to class 1
116
117 #The FPR is computed by selecting the number of cases where class 1 is predicted incorrectly,
118 #this value then has to be divided by the number of elements in the test set which do not belong to class 1
119
120 #The point corresponding to the EER can be found by plotting on the ROC curve the function y = 1 - x

```

The taw in which EE happens: 1.1

FPR and TPR of taw = EER is equal to : [0.23255813953488372, 0.8082663605051664]

The Accuracy for taw = EER is equal to: 80.3000000000001 %



Student's comments to Exercise 4

- Most of the comments related to the code are added between the code lines
- EER point indicates the threshold value where the classifier makes the same number of false positive and false negative predictions
- curve is almost acceptable and near to the top left of the figure (High AUC)
- The EER point is founded by putting an if condition $|TPR + FPR - 1| < 0.05$ (small range) in the loop of creating the ROC values

Exercise 5 – Classification – continuous data

This exercise employs the Iris dataset already employed in Exercise 1, and performs model fitting and classification using several versions of Gaussian discriminative analysis. However, for this exercise the available data have to be divided into two sets, namely *training* and *test* data.

You will have to 1) re-fit the training data to the specific model (see below), 2) classify each of the test samples, and 3) calculate the accuracy of each classifier.

Classifiers to be employed:

- Two-class quadratic discriminant analysis (fitting: both mean values and covariance matrices are class-specific – same as in exercise 1).
- Two-class linear discriminant analysis (fitting: class-specific mean values as in the previous case. Shared covariance matrix is calculated putting together the elements of both classes; the mean values should also be recalculated accordingly).

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import h5py
4
5 Dataset1 = h5py.File('/content/Lab2_Ex_1_Iris.hdf5')
6 Data = np.array(Dataset1.get('Dataset'))
7
8 Train = Data[:50,:]
9 Test = Data[50,:,:]
10
11 """***** Creating a training and a test for each of the classes *****
12
13 #Separate the training dataset into the two classes
14 #Create filters
15 indx0 = np.where(Train[:,2] == 0)
16 indx1 = np.where(Train[:,2] == 1)
17 # Apply filters
18 setosa_train = Train[indx0][:,0:2]
19 versicolour_train = Train[indx1][:,0:2]
20
21 #Separate the test dataset into the two classes
22 #Create filters
23 indx0 = np.where(Test[:,2] == 0)
24 indx1 = np.where(Test[:,2] == 1)
25 # Apply filters
26 setosa_test = Test[indx0][:,0:2]
27 versicolour_test = Test[indx1][:,0:2]
28
29 """***** Defining Functions to calculate mean and cov *****
30
31 def mean (x):
32     x_mean = [0,0]
33     for j in range(x.shape[1]):
34         for i in range(x.shape[0]):
35             x_mean[j] += x[i,j]
36         x_mean[j] /= x.shape[0]
37     return x_mean
38
39 def cov(y):
40     y_mean= mean(y)
41     y_mean = np.array(y_mean)
42     y_cov = np.empty([len(y.shape), len(y.shape)])
43     for i in range(y.shape[0]):
44         y_cov += (np.subtract( y[i][:, np.newaxis] @ y[i][:, np.newaxis].T , y_mean[:, np.newaxis] @ y_mean[:, np.newaxis].T))
45     return np.divide(y_cov,y.shape[0])
46
47 """***** Calculating mean and cov *****
48
49 # class specified means
50 mean_setosa_train = mean (setosa_train)
51 mean_versicolour_train = mean (versicolour_train)
52 #print ('mean_setosa_train:', mean_setosa_train)
53 #print ('mean_versicolour_train:', mean_versicolour_train)
54 # class specified covariance
55 cov_setosa_train = cov (setosa_train)
56 cov_versicolour_train = cov (versicolour_train)
57 #print ('cov_setosa_train:\n', cov_setosa_train)
58 #print ('cov_versicolour_train:\n', cov_versicolour_train)
59 # shared covariance
60 train_classless = Train[:,0:2]
61 cov_train = cov(train_classless)
62 #print ('cov_train:\n', cov_train)
63
64 """***** 1.Quadratic Discriminative Analysis *****
65
66 # Since the priors for setosa & versicolour is uniform we can use "Nearest Centroid" classifiers
67
68 # Defining nearest centroid classifier for QDA
69 def nearest_centroid (a):
70     D = 1
71     x = a
72
73     x_minus_mean = np.subtract(x,mean_setosa_train)
74     x_minus_mean_transpose = x_minus_mean.T
75     cov_inverse_setosa = np.linalg.inv(cov_setosa_train)
76     D_setosa = x_minus_mean_transpose @ cov_inverse_setosa @ x_minus_mean
77
78     x_minus_mean = np.subtract(x,mean_versicolour_train)
79     x_minus_mean_transpose = x_minus_mean.T
80     cov_inverse_versicolour = np.linalg.inv(cov_versicolour_train)
81     D_versicolour = x_minus_mean_transpose @ cov_inverse_versicolour @ x_minus_mean
82
83     if D_setosa > D_versicolour:
84         a=1
85     else:
86         a=0
87     return a
88
89 # Calculating the estimation vector
90 Test1 = Test[:,0:2]
91 estimated_calss_quadratic = []
92 for i in range(len(Test)):
93     estimated_calss_quadratic.append(nearest_centroid(Test1[i]))
94
95 # Calculating the accuracy
```

```

96 num = 0
97 for i in range (len(Test)):
98     if Test[i,2] == estimated_calss_quadratic[i]:
99         num +=1
100
101 print ("Accuracy of Quadratic Discriminative Analysis : ", num/len(Test))
102
103
104 """***** 2.Linear Discriminative Analysis ****"""
105
106 # We use discriminant function
107
108 def discriminant (a):
109     x = a
110     meansetosa_minus_meanversicolour = np.subtract(mean_setosa_train, mean_versicolour_train)
111     cov_total_inverse = np.linalg.inv(cov_train)
112     discriminant_x = meansetosa_minus_meanversicolour.T @ cov_total_inverse @ x
113
114     return discriminant_x
115
116
117 meansetosa_plus_meanversicolour = np.add(mean_setosa_train, mean_versicolour_train)
118
119
120 def linear_discriminant_classifier(a):
121     x = a
122     clas = 2
123     if discriminant(x) < 0.5* discriminant(meansetosa_plus_meanversicolour):
124         clas = 1
125     if discriminant(x) > 0.5* discriminant(meansetosa_plus_meanversicolour):
126         clas = 0
127
128     return clas
129
130
131 # Calculating the estimation vector
132 Test1 = Test[:,0:2]
133 estimated_calss_linear = []
134 for i in range(len(Test)):
135     estimated_calss_linear.append(linear_discriminant_classifier(Test1[i]))
136
137 # Calculating the accuracy
138 num = 0
139 for i in range (len(Test)):
140     if Test[i,2] == estimated_calss_linear[i]:
141         num +=1
142
143 print ("Accuracy of Linear Discriminative Analysis : ", num/len(Test))
144
145
146
Accuracy of Quadratic Discriminative Analysis :  0.92
Accuracy of Linear Discriminative Analysis :  0.9

```

Student's comments to Exercise 5

Quadratic Analysis

The classification function in QDA is:

$$D_i^2(y) = (y - \bar{y}_i)' S_i^{-1} (y - \bar{y}_i), \quad i = 1, 2, \dots, k$$

The observation y is assigned to the group for which $D_i^2(y)$ is smallest. **Note:** In QDA the covariance matrix is class specified

We define the function based on the formula and calculate the Accuracy

Linear Anayisus

In LDA, The classification procedure assigns the observation vector y to group 1 if its discriminant function $z = a'y$ is closer to z_1 or group 2 if its discriminant function is closer to z

$$z = a'y = (\bar{y}_1 - \bar{y}_2)' S_{p1}^{-1} y$$

$$z > \frac{1}{2}(\bar{z}_1 + \bar{z}_2)$$

We can now express the classification function regarding the observation vector y . $\frac{1}{2}(\bar{z}_1 + \bar{z}_2) = \frac{1}{2}(\bar{y}_1 - \bar{y}_2)' S_{p1}^{-1} (\bar{y}_1 + \bar{y}_2)$

Thus the observation vector y is assigned to a group determined by the following:

Assign y to group 1 if:

$$a'y = (\bar{y}_1 - \bar{y}_2)' S_{p1}^{-1} y > \frac{1}{2}(\bar{y}_1 - \bar{y}_2)' S_{p1}^{-1} (\bar{y}_1 + \bar{y}_2)$$

Or assign y to group 2 if:

$$a'y = (\bar{y}_1 - \bar{y}_2)' S_{p1}^{-1} y < \frac{1}{2}(\bar{y}_1 - \bar{y}_2)' S_{p1}^{-1} (\bar{y}_1 + \bar{y}_2)$$

S_{p1}^{-1} is the shared covariance matrix

Conclusion: Although QDA performed 2% better due to its flexibility but it might also overfit in some cases

Exercise 6 – Classification – continuous data

Classify the data in the phoneme dataset from Lab. 1 using quadratic discriminant analysis, linear discriminat analysis and a Naive Bayes classifier.

Compute the accuracy of each classifier and compare its performance with that of the k-nn classifier developed in Lab. 1.

Note: in Lab1 we had to employ a subset of the oiginal dataset due to the fact that k-nn has a quadratic complexity making it unfit for use on large datasets. The algorithms illustrated in this Lab have smaller complexity and thus it is possible to train on more data.

For this exercise you can use the sklearn library.

```

1 import numpy as np
2 import h5py
3 from scipy.stats import multivariate_normal

```

```

4
5
6 Dataset2 = h5py.File("/content/Lab2_Ex_6_phoneme.hdf5")
7 Data = np.array(Dataset2.get('Dataset'))
8
9 Train = Data[:4000,:]
10 Test = Data[4000:,:]
11 len_dat = np.shape(Test)[0]
12 Train_x = Train[:,255]
13 Train_y = Train[:,256]
14 Test_x = Test[:,255]
15 Test_y = Test[:,256]
16
17 from sklearn import metrics
18 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis, QuadraticDiscriminantAnalysis
19 from sklearn.naive_bayes import GaussianNB
20
21 #Part 1
22 clf = LinearDiscriminantAnalysis()
23 clf.fit(Train_x, Train_y)
24 linear_predict = clf.predict(Test_x)
25
26 #Complete here
27
28 #Part 2
29 clf2 = QuadraticDiscriminantAnalysis()
30 clf2.fit(Train_x, Train_y)
31 Quadratic_predict = clf2.predict(Test_x)
32 #Complete here
33
34 #Part 3
35 clf3 = GaussianNB()
36 clf3.fit(Train_x, Train_y)
37 Gaussian_predict = clf3.predict(Test_x)
38 #Complete here
39
40
41 # Accuracy Function
42 def Acc (a,b):
43     predicted = a
44     correct = b
45     num = 0
46     for i in range(len(predicted)):
47         if predicted[i] == correct[i] :
48             num += 1
49
50     acc = num / len(predicted)
51     return acc
52
53 # Calculating Accuracy
54 # LinearDiscriminantAnalysis
55 Linear_acc = Acc (linear_predict , Test_y )
56 print ("The accuracy of LinearDiscriminantAnalysis on the test data is equal to : (%) \n " , Linear_acc*100)
57
58 # QuadraticDiscriminantAnalysis
59 Quadratic_acc = Acc (Quadratic_predict , Test_y )
60 print ("The accuracy of QuadraticDiscriminantAnalysis on the test data is equal to : (%) \n " , Quadratic_acc*100)
61
62 # GaussianDiscriminantAnalysis
63 Gaussian_acc = Acc (Gaussian_predict , Test_y )
64 print ("The accuracy of GaussianDiscriminantAnalysis on the test data is equal to : (%) \n " , Gaussian_acc*100)
65
66

The accuracy of LinearDiscriminantAnalysis on the test data is equal to : (%)
94.10609037328095
The accuracy of QuadraticDiscriminantAnalysis on the test data is equal to : (%)
88.99803536345776
The accuracy of GaussianDiscriminantAnalysis on the test data is equal to : (%)
89.9803536345776

```

Student's comments to Exercise 6

Compared to KNN classifier in lab 1 with 96% Accuracy of the test data there's a lower accuracy !!! KNN was far better!

The LDA performed better and the probable reason is that the Gaussian assumptions holds, meaning that the class distributions are well-separated and approximately Gaussian, hence LDA can provide better accuracy than QDA.

COMPUTER LAB 3 - Principal component analysis

Duration: 3 hours

Introduction:

Hyperspectral images are scientific images of the Earth, acquired by satellites or aircrafts; rather than having three R/G/B color channels, these images have a lot more "color" components obtained through a fine sampling of the wavelength (hence the name "hyper"-spectral). The resulting 3-dimensional dataset has one image (spectral band or "color") for every sampled wavelength, which represents the measured radiance from each pixel at that specific wavelength. Hyperspectral images are very useful for image analysis. For every pixel at a given spatial position, it is possible to extract a so-called spectral vector, i.e. the 1-dimensional vector of values assumed by that pixel at all wavelengths. Assuming that each pixel is composed of just one substance, the spectral vector represents the radiance of that substance at all the wavelengths that have been sampled. Spectral vectors, therefore, can be used to infer which substance is contained in a given pixel – a typical classification problem that has a lot of practical applications in agriculture, analysis of land use / land cover, and other applications related to the study of the environment.

In this lab you will use a real hyperspectral image that has been acquired by the AVIRIS instrument, an airborne hyperspectral imager operated by the NASA. The image represents a scene of Indian Pines (Indiana, USA). It has a size of 145x145 pixels and 220 spectral bands. Along with the image, a ground truth is available, in terms of labels specifying which class (out of 16) each pixel belongs to. The classes are reported below; for more information, please see http://www.ehu.eus/ccwintco/index.php/Hyperspectral_Remote_Sensing_Scenes#Indian_Pines

#	Class	Samples
1	Alfalfa	46
2	Corn-notill	1428
3	Corn-mintill	830
4	Corn	237
5	Grass-pasture	483
6	Grass-trees	730
7	Grass-pasture-mowed	28
8	Hay-windrowed	478
9	Oats	20
10	Soybean-notill	972
11	Soybean-mintill	2455
12	Soybean-clean	593
13	Wheat	205
14	Woods	1265
15	Buildings-Grass-Trees-Drives	386
16	Stone-Steel-Towers	93

The purpose of this computer lab is twofold: To apply PCA to the spectral vectors in order to reduce their dimensionality.

- To apply PCA to the spectral vectors in order to reduce their dimensionality.
- To perform classification on the reduced data (optional)

Exercise 1 – PCA

In this exercise, you will employ the Indian Pines dataset. You will not do this for the entire dataset, but only for the spectral vectors belonging to two classes (as in the optional exercise you will perform 2-class classification on the PCA coefficients).

Reminder: the input to the PCA must always have zero mean: besides the sample covariance, you will have to compute the **mean value μ** over the training set and subtract it from each test vector before applying PCA.

Task: You have to reduce the dimensionality of the spectral vectors of the two classes you have chosen using PCA. In particular, you should perform the following:

- Extract spectral vectors of two classes, as described above (see sample code below).
- Estimate the sample covariance matrix of the dataset as a whole (i.e., considering together spectral vectors of the two classes)
- Perform the eigenvector decomposition of the sample covariance matrix. You can use the numpy linalg.eig function, which outputs the matrix containing the eigenvectors as columns, and a diagonal matrix containing the eigenvalues on the main diagonal.
 - Note: in the output matrix, eigenvectors/eigenvalues are not necessarily ordered by eigenvalue magnitude. You should sort them by yourself.
- Choose a number of dimensions K<=220.
- Construct the eigenvector matrix W for K components (i.e., select the last K columns)
- Using W, compute the PCA coefficients for each spectral vector in the data set
- Then from the PCA coefficients obtain an approximation of the corresponding vector and compute the error (mean square error - MSE)
- Plot the average MSE over the test set as a function of K.
- Plot the eigenvectors corresponding to the 3 largest eigenvalues – this will give you an idea of the basis functions employed by PCA

```

1 import numpy as np
2 import h5py
3 import scipy.io
4 import pandas as pd
5
6 mat = scipy.io.loadmat('Indian_pines.mat')
7 indian_pines = np.array(mat['indian_pines'])
8 print(indian_pines.shape, len(indian_pines[0]))
9
10 mat = scipy.io.loadmat('Indian_pines_gt.mat')
11 indian_pines_gt = np.array(mat['indian_pines_gt'])
12 print(indian_pines_gt.shape, indian_pines_gt)

(145, 145, 220) 145
(145, 145) [[3 3 3 ... 0 0 0]
 [3 3 3 ... 0 0 0]
 [3 3 3 ... 0 0 0]
 ...

```

```
[0 0 0 ... 0 0 0]
[0 0 0 ... 0 0 0]
[0 0 0 ... 0 0 0]
```

Extract spectral vectors of two classes, as described above (see sample code below).

```
1 # Coosing Class 4 , Corn , 237 samples
2 class1_value = 4 # corn
3
4 class1=np.zeros((1500,220)) # Create an empty 2d array to store spectral data of each of the classes
5 n=0
6 for i in range(145):
7     for j in range(145):
8         if indian_pines_gt[i,j]== class1_value:
9             class1[n,:]= indian_pines[i,j,:]/1.
10            n=n+1
11 class1=class1[:n,:]
12
13 # Choising Class 14 , Woods , 1265 samples
14 class2_value = 14 # woods
15
16 class2=np.zeros((1500,220)) # Create an empty 2d array to store spectral data of each of the classes
17 n=0
18 for i in range(145):
19     for j in range(145):
20         if indian_pines_gt[i,j]== class2_value:
21             class2[n,:]= indian_pines[i,j,:]/1.
22             n=n+1
23 class2=class2[:n,:]
24
25 print (class1.shape , class1)
26 print (class2.shape , class2)

(237, 220) [[3170. 4001. 4330. ... 1024. 1019. 1004.]
 [3167. 4005. 4423. ... 1038. 1015. 1000.]
 [3002. 4133. 4503. ... 1031. 1014. 1013.]
 ...
 [3705. 4497. 4785. ... 1035. 1013. 1024.]
 [3868. 4373. 4777. ... 1029. 1009. 1014.]
 [2592. 4369. 4599. ... 1038. 1013. 1019.]]
(1265, 220) [[2741. 4135. 4147. ... 998. 1001. 1010.]
 [2567. 4135. 4324. ... 1005. 1005. 1000.]
 [3167. 4128. 4320. ... 1001. 1000. 1004.]
 ...
 [2725. 3982. 4099. ... 1004. 1008. 1009.]
 [3320. 3862. 4015. ... 1004. 1000. 1009.]
 [3154. 3986. 4099. ... 996. 1003. 1000.]]
```

Estimate the sample covariance matrix of the dataset as a whole (i.e., considering together spectral vectors of the two classes)

```
1 # Insert code here
2 bothclass = np.concatenate((class1,class2) ,axis=0) # Concatenating the data of the classes
3
4 # Converting to DataFrame
5 bothclass_df = pd.DataFrame (bothclass) # Convert to a dataframe to be able to get mean
6
7 # Shuffling the data
8 bothclass_df = bothclass_df.sample(frac = 1)
9
10
11 # Creating the test and training datasets
12 bothclass_test_df = bothclass_df.iloc[1000:1502 , :] # Choose the last 502 points(pixels) as the test dataset
13 bothclass_test = bothclass_test_df.to_numpy() # ~502*220
14 bothclass_df = bothclass_df.iloc[:1000 , :] # Choose the first 1000 points as the training dataset
15 bothclass = bothclass_df.to_numpy()
16
17 # Normalizing the training data
18 mean_bothclass_df = bothclass_df.mean() # Calcualte the mean
19 mean_bothclass = mean_bothclass_df.to_numpy() # Convert back to an array ~(220,)
20 bothclass -= mean_bothclass # subtracting mean vector from all data (all rows)
21 bothclass = bothclass.T # Reshaping the data ~220*1000
22
23 # Normalizing the test data
24 bothclass_test -= mean_bothclass
25 bothclass_test = bothclass_test.T # ~220*502 # Normalized
26
27
28 # Calculating cov matrix based on the training matrix (bothclass)
29 cov_bothclass = np.cov(bothclass) # ~220*220
30
31
```

Double-click (or enter) to edit

Perform the eigenvector decomposition of the sample covariance matrix. You can use the numpy linalg.eig function, which outputs the matrix containing the eigenvectors as columns, and a diagonal matrix containing the eigenvalues on the main diagonal.

```
1 # Note: in the output matrix, eigenvectors/eigenvalues are not necessarily ordered by eigenvalue magnitude. You should sort them first.
2 from numpy import linalg as LA
3
4
5 w, v = LA.eig(cov_bothclass) # w : eigenvalues ~a 200 element list # v : eigen vectors matrix ~220*220
6 w = np.array([w]) # convert w from a list to an array to be concatenatable ~1*220
7
8 # Sorting
9 v_sorting = np.concatenate((v,w) , axis = 0) # append w to v (as the last row) to sort it based on the values of w ~221*220
10 v_sorting_df = pd.DataFrame(v_sorting) # convert array to a dataframe to be sortable
11 v_sorting_df_sorted = v_sorting_df.sort_values(by = 220, axis = 1 , ascending= False) # sort columns descending based on the values of w (last row)
12 v_sorting_sorted = v_sorting_df_sorted.to_numpy() # Convert back to an array
13 v_sorted = v_sorting_sorted[:220,:]. # Removing w from the last row of v ~220*220
```

Choose a number of dimensions K<=220. Construct the eigenvector matrix W for K components (i.e., select the last K columns)

```

1 # Considering k first eigenvectors
2
3 k = 50
4 v_new = v_sorted[:, :k]           # ~ 220 * k
5 v_new_transpose = v_new.T        # ~ k * 220
6

```

Using W, compute the PCA coefficients for each spectral vector in the data set. Then from the PCA coefficients obtain an approximation of the corresponding vector and compute the error (mean square error - MSE)

```

1 # Note: remember to remove the mean from the vectors of the dataset
2
3
4 z = v_new_transpose @ bothclass_test      # z = w.T @ x_test      # ~ k * 1000
5
6 #####
7 bothclass_hat = v_new @ z            # Calculating z_hat
8
9 #####
10 # Calculating MSE
11 SE = 0                                # square error
12 for i in range(bothclass_hat.shape[1]):
13     a = np.subtract(bothclass_hat[:, i], bothclass_test[:, i])
14     SE = a @ a.T
15     SE += SE
16 MSE = SE / bothclass_hat.shape[1]
17
18 print('The MSE is equal to:', MSE)
19
20

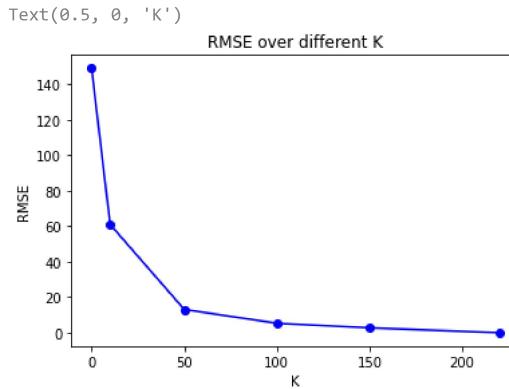
```

The MSE is equal to: 171.16245422275685

```

1 # plot MSE over K for k = 10 , 50 , 70 , 100
2 import matplotlib.pyplot as plt
3
4
5 def RMSE(a):
6     k = a
7     v_new = v_sorted[:, :k]           # ~ 220 * k
8     v_new_transpose = v_new.T        # ~ k * 220
9     z = v_new_transpose @ bothclass_test      # z = w.T @ x_test      # ~ k * 1000
10    bothclass_hat = v_new @ z          # Calculating z_hat
11
12    SE = 0                                # square error
13    for i in range(bothclass_hat.shape[1]):
14        a = np.subtract(bothclass_hat[:, i], bothclass_test[:, i])
15        SE = a.T @ a
16        SE += SE
17    MSE = SE / bothclass_hat.shape[1]
18    RMSE = MSE**0.5
19    return RMSE
20
21 a = RMSE(0)
22 b = RMSE(10)
23 c = RMSE(50)
24 d = RMSE(100)
25 e = RMSE(150)
26 f = RMSE(220)
27
28 RMSE_list = [a, b, c, d, e, f]
29 K_list = [0, 10, 50, 100, 150, 220]
30
31
32 plt.plot(K_list, RMSE_list, '-o', color='blue')
33 plt.title("RMSE over different K")
34 plt.ylabel('RMSE')
35 plt.xlabel('K')
36
37

```



Plot the average MSE over the test set as a function of K.

```

1 import matplotlib.pyplot as plt
2
3 #Insert code here
4
5 def RMSE_over_K(a):
6     k = a
7     v_new = v_sorted[:, :k]           # Slice the K first Eigen values
8
9

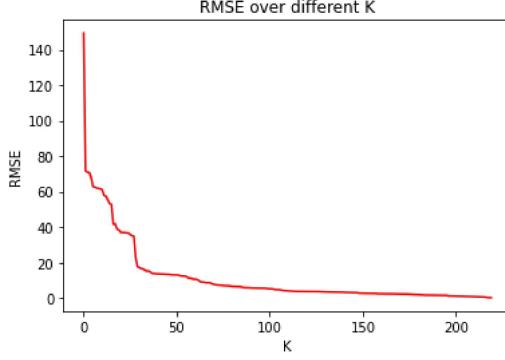
```

```

8
9     v_new_transpose = v_new.T
10    z = v_new_transpose @ bothclass_test
11    bothclass_hat = v_new @ z
12    SE = 0
13    for i in range (bothclass_hat.shape[1]):
14        a = np.subtract(bothclass_hat[:,i], bothclass_test[:,i])
15        SE = a.T @ a
16        SE +=SE
17    MSE = SE/bothclass_hat.shape[1]
18    RMSE = MSE ** 0.5
19    return RMSE
20
21 # Create a vector of MSE based on K
22 RMSE_k = []
23 for j in range(len(bothclass_hat)):
24     RMSE_k.append(RMSE_over_K(j))
25 print(len(RMSE_k),RMSE_k)
26 k = np.arange(0,220,1)
27
28 plt.plot(k, RMSE_k, color="red")
29 plt.title("RMSE over different K")
30 plt.ylabel('RMSE')
31 plt.xlabel('K')
32 plt.show()

```

220 [149.23858863352567, 71.73053702968461, 70.8877977661219, 70.88373667591365, 68.14855007021173.

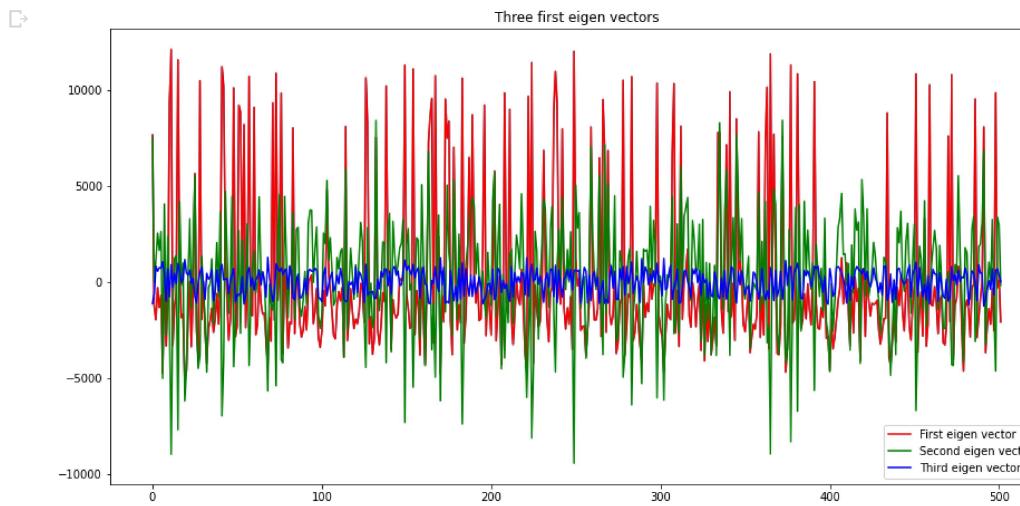


Plot the eigenvectors corresponding to the 3 largest eigenvalues – this will give you an idea of the basis functions employed by PCA

```

1 #Insert code here
2 a = plt.figure(figsize=(15,7.5))
3 first_eigenvector = z[0,:]
4 second_eigenvector = z[1,:]
5 third_eigenvector = z[2,:]
6
7 plt.plot(first_eigenvector)
8 plt.plot(first_eigenvector, color ="red" , label ='First eigen vector')
9 plt.plot(second_eigenvector, color ="green", label ='Second eigen vector')
10 plt.plot(third_eigenvector, color ="blue", label ='Third eigen vector')
11 plt.title('Three first eigen vectors')
12 plt.legend()
13 plt.show()

```



Student's comments to Exercise 1

- The comments related to the codes are embeddd into the code boxes
- The first 2 curvy plots clearly shows the high amount of the data carried by the first components and it is clear that by adding more component the reconstruction error tends to get lower and lower; hence it is possible to remove many of the latent variables without loosing much of the data.
- In total, we can infer the following from the plot of RMSE by k : Dimensionality reduction quality - Explanation of Data variance - Data compression quality
- The last figure clearly depicts that the highest variance value is for the first pricipal component the second for the second one and

Exercise 2 – Classification using dimensionality reduction and whitening (optional)

In this exercise you will apply a simple 2-class linear discriminant analysis (LDA) classifier to the data belonging to the two classes, before and after **whitening** (i.e. applying PCA plus rescaling each coefficient to unit variance: $y = \Lambda^{-\frac{1}{2}} W^T x$). Note that matrix Λ is obtained as one of the

outputs of *eig.m*). This classifier makes the Naïve Bayes Classifier assumption that the features are statistically independent, thus the shared covariance matrix is taken as $\Sigma = I$. We also assume that class 0 and class 1 are equiprobable. Thus, letting μ_0 and μ_1 be the mean of vectors in class 0 and 1, we define $x_0 = 0.5(\mu_0 + \mu_1)$ and $w = \mu_1 - \mu_0$. A test vector x is classified into class 1 or 0 depending on whether $\text{sign}(w^T(x - x_0))$ is equal to +1 or -1.

When the classifier is applied to the PCA coefficients, you simply replace μ_0 and μ_1 with their reduced versions (subtracting μ and applying the same PCA matrix computed over the training set), and recalculate x_0 and w accordingly.

Task: Divide the Indian Pines dataset into training and test sets (e.g. 75% of the data of each class to be used as training data, and 25% as test data). Train this classifier on the training data, and apply it to the test data. In particular, you should perform the following:

- Plot the mean vector of class 0 and 1 – this will give you a visual description of the differences among vectors of either class.
- Apply the classifier to the original data (without PCA) and compute its accuracy
- Apply the classifier to the PCA coefficients for different values of K, and compute its accuracy
- Apply the classifier to the original data where only the first K features have been retained, and compute its accuracy (this is a more brutal way to reduce dimensionality)
- Try to classify the data using a Support Vector Machine and compare the results

Plot the mean vector of class 0 and 1 – this will give you a visual description of the differences among vectors of either class.

```
1 #Insert code here
```

Apply the classifier to the original data (without PCA) and compute its accuracy

```
1 #For this section you can use the code you developed for the previous Lab or the sklearn function LinearDiscriminantAnalysis
```

Apply the classifier to the PCA coefficients for different values of K, and compute its accuracy

```
1 #Insert code Here
```

**Apply the classifier to the original data where only a subset of K features (selected randomly) have been retained, and compute its accuracy
(this is a more brutal way to reduce dimensionality)**

```
1 #Insert code here
```

Try to classify the data using a Support Vector Machine and compare the results

```
1 #For this section you can use the sklearn library
2
3 #from sklearn import svm
4
5 #svc = svm.SVC()
6 #svc.fit(Train, Train_Label)
7 #predict = svc.predict(Test)
8
9 #print(metrics.accuracy_score(predict, Test_Label))
```

Student's comments to Exercise 2

Add comments to the results of Exercise 2 here (may use LaTeX for formulas if needed).

Duration: 3 hours

Introduction:

In this lab, you are provided with the set of coordinates (x,y – horizontal and vertical) describing the trajectories of pedestrians moving across a scene. Your task is to simulate the observed positions of the pedestrians by adding observation noise, then to track the subjects using a Kalman filter. In other words, you need to estimate the next (x,y) positions, from the observations of the previous positions.

Simulating the observed coordinates

Choose one of the trajectories in the dataset. This data will be considered the real trajectory. Generate the observed directory by adding observation noise $\delta_t \sim \mathcal{N}(0, \sigma_R^2)$ to the (x,y) coordinates.

Designing the Kalman filter

Your task is to **design a Kalman filter** based on a constant velocity model, which tracks the next (x,y) position of the object, from the observation of the previous positions. The code must be based on the following model.

- The state vector contains coordinates and velocities: $z_t^T = (z_{1t}, z_{2t}, v_{1t}, v_{2t})$ (see slides). The object has initial coordinates (0,0) and velocity (Δ, Δ).
- Only the coordinates (but not the velocities) are observed. This leads to a linear dynamical system with:

$$A = \begin{pmatrix} 1 & 0 & \Delta & 0 \\ 0 & 1 & 0 & \Delta \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

$$C = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}$$

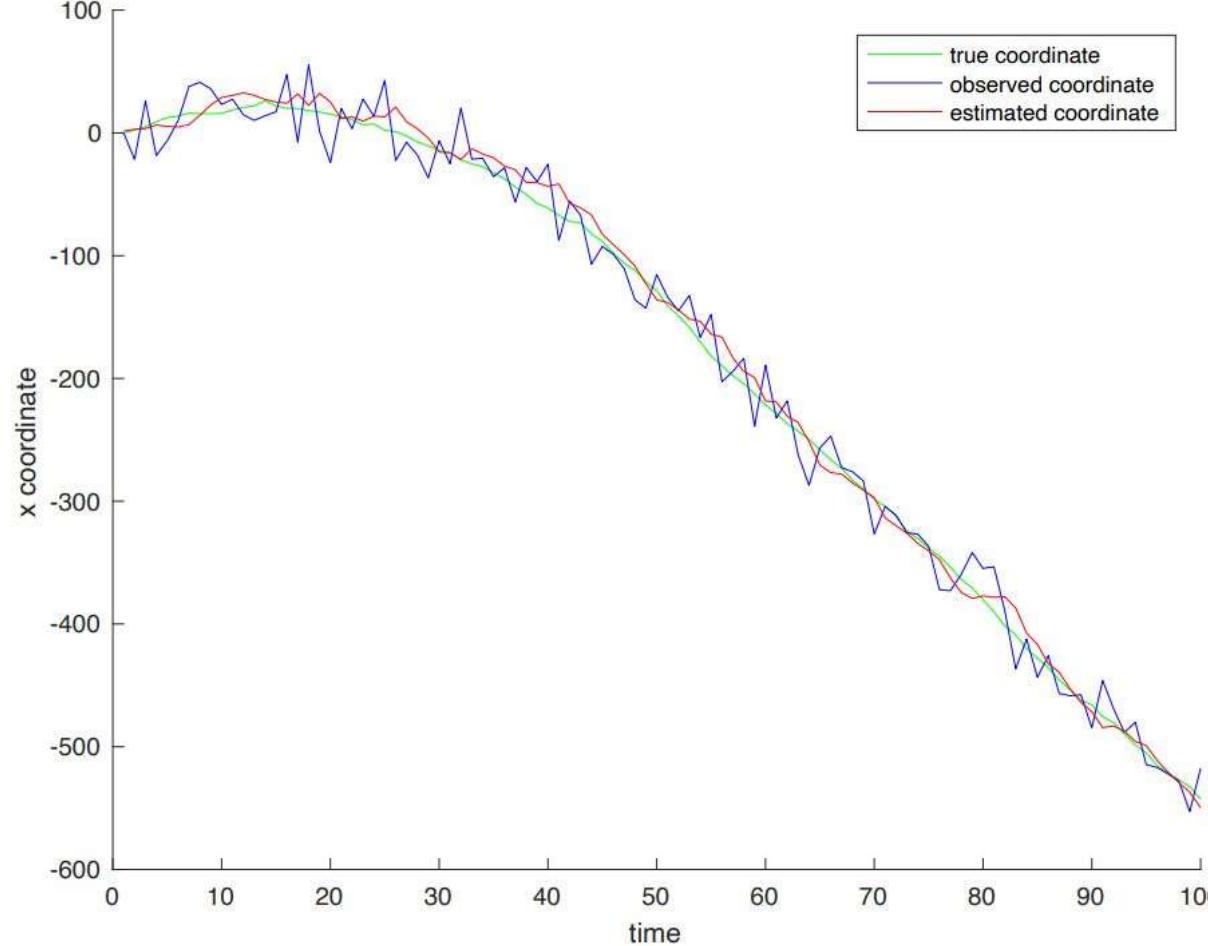
and $B = D = 0$.

- Σ_Q , and Σ_R should be set to:

$$\Sigma_Q = \sigma_Q^2 \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

$$\Sigma_R = \sigma_R^2 \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

To see if your Kalman filter is working well, you should plot the estimated position of the object over time with respect to the true position (i.e., the first two entries of the state vector) and the observed position. Depending on the chosen parameters, for each coordinate the graph may look something like this:



Suggestion: when implementing your Kalman filter, you will have to choose initial values for μ_t and Σ_t . Provided that you do not make very unreasonable assumptions, the Kalman filter will update those estimates from observed data, so the initial choices are not very critical.

Test your Kalman filter modifying the values of some of the parameters, including standard deviations σ_Q and σ_R , initial values for μ_t and Σ_t and the value of Δ .

New section

```

1 from google.colab import drive
2 drive.mount('/content/drive')
3
4 !pip install ndjson

```

Mounted at /content/drive
 Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public/simple/>
 Collecting ndjson
 Downloading ndjson-0.3.1-py2.py3-none-any.whl (5.3 kB)
 Installing collected packages: ndjson
 Successfully installed ndjson-0.3.1

```

1 import ndjson
2 import pandas as pd
3 import numpy as np
4
5 #Dataset origin: https://paperswithcode.com/dataset/trajnet-1
6 #The crowds_students001 file is loaded, and formatted as a list of numpy vectors
7
8 with open('/content/drive/MyDrive/Labs/Lab4/crowds_students001_trackonly.ndjson') as f:
9     data = ndjson.load(f)
10

```

```
11 p=-1'
12
13 person_dict = []
14
15 for ii in range(len(data)):
16     if(p!=data[ii]['track']['p']):
17         p=data[ii]['track']['p']
18         person_dict.append([])
19     person_dict[p].append((data[ii]['track']['x'],data[ii]['track']['y']))
20
21 person_dict_numpy = []
22 for ii in range(148):
23     person_dict_numpy.append(np.array(person_dict[ii], dtype=float))
24
25 print(data[1]['track']['p'])
26 print('data : ', len(data) , data)
27 print('person_dict 68: ', person_dict[68])
28 print('person_dict_numpy : ', len(person_dict_numpy))

0
data : 21335 [{'track': {'f': 10, 'p': 0, 'x': 3.21, 'y': -3.21}}, {'track': {'f': 20, 'p': 0, 'x': 2.76, 'y': -3.28}}, {'track': {'f': 30, 'p': 0, 'x': 2.3, 'y': -3.34}}, {'track': {'f': 40, 'p': 0, 'x': 1.87, 'y': -3.4}], person_dict 68: [(-4.17, 6.8), (-4.15, 6.54), (-4.12, 6.31), (-4.09, 6.09), (-4.06, 5.89), (-4.03, 5.71), (-4.01, 5.54), (-3.97, 5.39), (-3.87, 5.27), (-3.73, 5.16), (-3.56, 5.05), (-3.39, 4.93), (-3.23, 4.81), (-3.07, 4.69), (-2.91, 4.57), (-2.75, 4.45), (-2.59, 4.33), (-2.43, 4.21), (-2.27, 4.09), (-2.11, 3.97), (-1.95, 3.85), (-1.79, 3.73), (-1.63, 3.61), (-1.47, 3.49), (-1.31, 3.37), (-1.15, 3.25), (-0.99, 3.13), (-0.83, 3.01), (-0.67, 2.89), (-0.51, 2.77), (-0.35, 2.65), (-0.19, 2.53), (0.05, 2.41), (0.21, 2.29), (0.37, 2.17), (0.53, 2.05), (0.69, 1.93), (0.85, 1.81), (1.01, 1.69), (1.17, 1.57), (1.33, 1.45), (1.49, 1.33), (1.65, 1.21), (1.81, 1.09), (1.97, 0.97), (2.13, 0.85), (2.29, 0.73), (2.45, 0.61), (2.61, 0.49), (2.77, 0.37), (2.93, 0.25), (3.09, 0.13), (3.25, 0.01), (3.41, -0.11), (3.57, -0.23), (3.73, -0.35), (3.89, -0.47), (4.05, -0.59), (4.21, -0.71), (4.37, -0.83), (4.53, -0.95), (4.69, -1.07), (4.85, -1.19), (5.01, -1.31), (5.17, -1.43), (5.33, -1.55), (5.49, -1.67), (5.65, -1.79), (5.81, -1.91), (5.97, -2.03), (6.13, -2.15), (6.29, -2.27), (6.45, -2.39), (6.61, -2.51), (6.77, -2.63), (6.93, -2.75), (7.09, -2.87), (7.25, -2.99), (7.41, -3.11), (7.57, -3.23), (7.73, -3.35), (7.89, -3.47), (8.05, -3.59), (8.21, -3.71), (8.37, -3.83), (8.53, -3.95), (8.69, -4.07), (8.85, -4.19), (9.01, -4.31), (9.17, -4.43), (9.33, -4.55), (9.49, -4.67), (9.65, -4.79), (9.81, -4.91), (9.97, -5.03), (10.13, -5.15), (10.29, -5.27), (10.45, -5.39), (10.61, -5.51), (10.77, -5.63), (10.93, -5.75), (11.09, -5.87), (11.25, -5.99), (11.41, -6.11), (11.57, -6.23), (11.73, -6.35), (11.89, -6.47), (12.05, -6.59), (12.21, -6.71), (12.37, -6.83), (12.53, -6.95), (12.69, -7.07), (12.85, -7.19), (13.01, -7.31), (13.17, -7.43), (13.33, -7.55), (13.49, -7.67), (13.65, -7.79), (13.81, -7.91), (13.97, -8.03), (14.13, -8.15), (14.29, -8.27), (14.45, -8.39), (14.61, -8.51), (14.77, -8.63), (14.93, -8.75), (15.09, -8.87), (15.25, -8.99), (15.41, -9.11), (15.57, -9.23), (15.73, -9.35), (15.89, -9.47), (16.05, -9.59), (16.21, -9.71), (16.37, -9.83), (16.53, -9.95), (16.69, -10.07), (16.85, -10.19), (17.01, -10.31), (17.17, -10.43), (17.33, -10.55), (17.49, -10.67), (17.65, -10.79), (17.81, -10.91), (17.97, -11.03), (18.13, -11.15), (18.29, -11.27), (18.45, -11.39), (18.61, -11.51), (18.77, -11.63), (18.93, -11.75), (19.09, -11.87), (19.25, -11.99), (19.41, -12.11), (19.57, -12.23), (19.73, -12.35), (19.89, -12.47), (20.05, -12.59), (20.21, -12.71), (20.37, -12.83), (20.53, -12.95), (20.69, -13.07), (20.85, -13.19), (20.1, -13.31), (20.26, -13.43), (20.42, -13.55), (20.58, -13.67), (20.74, -13.79), (20.9, -13.91), (21.06, -14.03), (21.22, -14.15), (21.38, -14.27), (21.54, -14.39), (21.7, -14.51), (21.86, -14.63), (22.02, -14.75), (22.18, -14.87), (22.34, -14.99), (22.5, -15.11), (22.66, -15.23), (22.82, -15.35), (22.98, -15.47), (23.14, -15.59), (23.3, -15.71), (23.46, -15.83), (23.62, -15.95), (23.78, -16.07), (23.94, -16.19), (24.1, -16.31), (24.26, -16.43), (24.42, -16.55), (24.58, -16.67), (24.74, -16.79), (24.9, -16.91), (25.06, -17.03), (25.22, -17.15), (25.38, -17.27), (25.54, -17.39), (25.7, -17.51), (25.86, -17.63), (26.02, -17.75), (26.18, -17.87), (26.34, -17.99), (26.5, -18.11), (26.66, -18.23), (26.82, -18.35), (26.98, -18.47), (27.14, -18.59), (27.3, -18.71), (27.46, -18.83), (27.62, -18.95), (27.78, -19.07), (27.94, -19.19), (28.1, -19.31), (28.26, -19.43), (28.42, -19.55), (28.58, -19.67), (28.74, -19.79), (28.9, -19.91), (29.06, -20.03), (29.22, -20.15), (29.38, -20.27), (29.54, -20.39), (29.7, -20.51), (29.86, -20.63), (29.1, -20.75), (29.26, -20.87), (29.42, -20.99), (29.58, -21.11), (29.74, -21.23), (29.9, -21.35), (30.06, -21.47), (30.22, -21.59), (30.38, -21.71), (30.54, -21.83), (30.7, -21.95), (30.86, -22.07), (30.1, -22.19), (30.26, -22.31), (30.42, -22.43), (30.58, -22.55), (30.74, -22.67), (30.9, -22.79), (31.06, -22.91), (31.22, -23.03), (31.38, -23.15), (31.54, -23.27), (31.7, -23.39), (31.86, -23.51), (31.1, -23.63), (31.26, -23.75), (31.42, -23.87), (31.58, -23.99), (31.74, -24.11), (31.9, -24.23), (32.06, -24.35), (32.22, -24.47), (32.38, -24.59), (32.54, -24.71), (32.7, -24.83), (32.86, -24.95), (33.1, -25.07), (33.26, -25.19), (33.42, -25.31), (33.58, -25.43), (33.74, -25.55), (33.9, -25.67), (34.06, -25.79), (34.22, -25.91), (34.38, -26.03), (34.54, -26.15), (34.7, -26.27), (34.86, -26.39), (35.1, -26.51), (35.26, -26.63), (35.42, -26.75), (35.58, -26.87), (35.74, -26.99), (35.9, -27.11), (36.06, -27.23), (36.22, -27.35), (36.38, -27.47), (36.54, -27.59), (36.7, -27.71), (36.86, -27.83), (37.1, -27.95), (37.26, -28.07), (37.42, -28.19), (37.58, -28.31), (37.74, -28.43), (37.9, -28.55), (38.06, -28.67), (38.22, -28.79), (38.38, -28.91), (38.54, -29.03), (38.7, -29.15), (38.86, -29.27), (39.1, -29.39), (39.26, -29.51), (39.42, -29.63), (39.58, -29.75), (39.74, -29.87), (39.9, -29.99), (40.06, -30.11), (40.22, -30.23), (40.38, -30.35), (40.54, -30.47), (40.7, -30.59), (40.86, -30.71), (41.1, -30.83), (41.26, -30.95), (41.42, -31.07), (41.58, -31.19), (41.74, -31.31), (41.9, -31.43), (42.06, -31.55), (42.22, -31.67), (42.38, -31.79), (42.54, -31.91), (42.7, -32.03), (42.86, -32.15), (43.1, -32.27), (43.26, -32.39), (43.42, -32.51), (43.58, -32.63), (43.74, -32.75), (43.9, -32.87), (44.06, -32.99), (44.22, -33.11), (44.38, -33.23), (44.54, -33.35), (44.7, -33.47), (44.86, -33.59), (45.1, -33.71), (45.26, -33.83), (45.42, -33.95), (45.58, -34.07), (45.74, -34.19), (45.9, -34.31), (46.06, -34.43), (46.22, -34.55), (46.38, -34.67), (46.54, -34.79), (46.7, -34.91), (46.86, -35.03), (47.1, -35.15), (47.26, -35.27), (47.42, -35.39), (47.58, -35.51), (47.74, -35.63), (47.9, -35.75), (48.06, -35.87), (48.22, -35.99), (48.38, -36.11), (48.54, -36.23), (48.7, -36.35), (48.86, -36.47), (49.1, -36.59), (49.26, -36.71), (49.42, -36.83), (49.58, -36.95), (49.74, -37.07), (49.9, -37.19), (50.06, -37.31), (50.22, -37.43), (50.38, -37.55), (50.54, -37.67), (50.7, -37.79), (50.86, -37.91), (51.1, -38.03), (51.26, -38.15), (51.42, -38.27), (51.58, -38.39), (51.74, -38.51), (51.9, -38.63), (52.06, -38.75), (52.22, -38.87), (52.38, -38.99), (52.54, -39.11), (52.7, -39.23), (52.86, -39.35), (53.1, -39.47), (53.26, -39.59), (53.42, -39.71), (53.58, -39.83), (53.74, -39.95), (53.9, -40.07), (54.06, -40.19), (54.22, -40.31), (54.38, -40.43), (54.54, -40.55), (54.7, -40.67), (54.86, -40.79), (55.1, -40.91), (55.26, -41.03), (55.42, -41.15), (55.58, -41.27), (55.74, -41.39), (55.9, -41.51), (56.06, -41.63), (56.22, -41.75), (56.38, -41.87), (56.54, -41.99), (56.7, -42.11), (56.86, -42.23), (57.1, -42.35), (57.26, -42.47), (57.42, -42.59), (57.58, -42.71), (57.74, -42.83), (57.9, -42.95), (58.06, -43.07), (58.22, -43.19), (58.38, -43.31), (58.54, -43.43), (58.7, -43.55), (58.86, -43.67), (59.1, -43.79), (59.26, -43.91), (59.42, -44.03), (59.58, -44.15), (59.74, -44.27), (59.9, -44.39), (60.06, -44.51), (60.22, -44.63), (60.38, -44.75), (60.54, -44.87), (60.7, -44.99), (60.86, -45.11), (61.1, -45.23), (61.26, -45.35), (61.42, -45.47), (61.58, -45.59), (61.74, -45.71), (61.9, -45.83), (62.06, -45.95), (62.22, -46.07), (62.38, -46.19), (62.54, -46.31), (62.7, -46.43), (62.86, -46.55), (63.1, -46.67), (63.26, -46.79), (63.42, -46.91), (63.58, -47.03), (63.74, -47.15), (63.9, -47.27), (64.06, -47.39), (64.22, -47.51), (64.38, -47.63), (64.54, -47.75), (64.7, -47.87), (64.86, -47.99), (65.1, -48.11), (65.26, -48.23), (65.42, -48.35), (65.58, -48.47), (65.74, -48.59), (65.9, -48.71), (66.06, -48.83), (66.22, -48.95), (66.38, -49.07), (66.54, -49.19), (66.7, -49.31), (66.86, -49.43), (67.1, -49.55), (67.26, -49.67), (67.42, -49.79), (67.58, -49.91), (67.74, -50.03), (67.9, -50.15), (68.06, -50.27), (68.22, -50.39), (68.38, -50.51), (68.54, -50.63), (68.7, -50.75), (68.86, -50.87), (69.1, -50.99), (69.26, -51.11), (69.42, -51.23), (69.58, -51.35), (69.74, -51.47), (69.9, -51.59), (70.06, -51.71), (70.22, -51.83), (70.38, -51.95), (70.54, -52.07), (70.7, -52.19), (70.86, -52.31), (71.1, -52.43), (71.26, -52.55), (71.42, -52.67), (71.58, -52.79), (71.74, -52.91), (71.9, -53.03), (72.06, -53.15), (72.22, -53.27), (72.38, -53.39), (72.54, -53.51), (72.7, -53.63), (72.86, -53.75), (73.1, -53.87), (73.26, -53.99), (73.42, -54.11), (73.58, -54.23), (73.74, -54.35), (73.9, -54.47), (74.06, -54.59), (74.22, -54.71), (74.38, -54.83), (74.54, -54.95), (74.7, -55.07), (74.86, -55.19), (75.1, -55.31), (75.26, -55.43), (75.42, -55.55), (75.58, -55.67), (75.74, -55.79), (75.9, -55.91), (76.06, -56.03), (76.22, -56.15), (76.38, -56.27), (76.54, -56.39), (76.7, -56.51), (76.86, -56.63), (77.1, -56.75), (77.26, -56.87), (77.42, -56.99), (77.58, -57.11), (77.74, -57.23), (77.9, -57.35), (78.06, -57.47), (78.22, -57.59), (78.38, -57.71), (78.54, -57.83), (78.7, -57.95), (78.86, -58.07), (79.1, -58.19), (79.26, -58.31), (79.42, -58.43), (79.58, -58.55), (79.74, -58.67), (79.9, -58.79), (80.06, -58.91), (80.22, -59.03), (80.38, -59.15), (80.54, -59.27), (80.7, -59.39), (80.86, -59.51), (81.1, -59.63), (81.26, -59.75), (81.42, -59.87), (81.58, -59.99), (81.74, -60.11), (81.9, -60.23), (82.06, -60.35), (82.22, -60.47), (82.38, -60.59), (82.54, -60.71), (82.7, -60.83), (82.86, -60.95), (83.1, -61.07), (83.26, -61.19), (83.42, -61.31), (83.58, -61.43), (83.74, -61.55), (83.9, -61.67), (84.06, -61.79), (84.22, -61.91), (84.38, -62.03), (84.54, -62.15), (84.7, -62.27), (84.86, -62.39), (85.1, -62.51), (85.26, -62.63), (85.42, -62.75), (85.58, -62.87), (85.74, -62.99), (85.9, -63.11), (86.06, -63.23), (86.22, -63.35), (86.38, -63.47), (86.54, -63.59), (86.7, -63.71), (86.86, -63.83), (87.1, -63.95), (87.26, -64.07), (87.42, -64.19), (87.58, -64.31), (87.74, -64.43), (87.9, -64.55), (88.06, -64.67), (88.22, -64.79), (88.38, -64.91), (88.54, -65.03), (88.7, -65.15), (88.86, -65.27), (89.1, -65.39), (89.26, -65.51), (89.42, -65.63), (89.58, -65.75), (89.74, -65.87), (89.9, -65.99), (90.06, -66.11), (90.22, -66.23), (90.38, -66.35), (90.54, -66.47), (90.7, -66.59), (90.86, -66.71), (91.1, -66.83), (91.26, -66.95), (91.42, -67.07), (91.58, -67.19), (91.74, -67.31), (91.9, -67.43), (92.06, -67.55), (92.22, -67.67), (92.38, -67.79), (92.54, -67.91), (92.7, -68.03), (92.86, -68.15), (93.1, -68.27), (93.26, -68.39), (93.42, -68.51), (93.58, -68.63), (93.74, -68.75), (93.9, -68.87), (94.06, -68.99), (94.22, -69.11), (94.38, -69.23), (94.54, -69.35), (94.7, -69.47), (94.86, -69.59), (95.1, -69.71), (95.26, -69.83), (95.42, -69.95), (95.58, -70.07), (95.74, -70.19), (95.9, -70.31), (96.06, -70.43), (96.22, -70.55), (96.38, -70.67), (96.54, -70.79), (96.7, -70.91), (96.86, -71.03), (97.1, -71.15), (97.26, -71.27), (97.42, -71.39), (97.58, -71.51), (97.74, -71.63), (97.9, -71.75), (98.06, -71.87), (98.22, -71.99), (98.38, -72.11), (98.54, -72.23), (98.7, -72.35), (98.86, -72.47), (99.1, -72.59), (99.26, -72.71), (99.42, -72.83), (99.58, -72.95), (99.74, -73.07), (99.9, -73.19), (100.06, -73.31), (100.22, -73.43), (100.38, -73.55), (100.54, -73.67), (100.7, -73.79), (100.86, -73.91), (101.1, -74.03), (101.26, -74.15), (101.42, -74.27), (101.58, -74.39), (101.74, -74.51), (101.9, -74.63), (102.06, -74.75), (102.22, -74.87), (102.38, -74.99), (102.54, -75.11), (102.7, -75.23), (102.86, -75.35), (103.1, -75.47), (103.26, -75.59), (103.42, -75.71), (103.58, -75.83), (103.74, -75.95), (103.9, -76.07), (104.06, -76.19), (104.22, -76.31), (104.38, -76.43), (104.54, -76.55), (104.7, -76.67), (104.86, -76.79), (105.1, -76.91), (105.26, -77.03), (105.42, -77.15), (105.58, -77.27), (105.74, -77.39), (105.9, -77.51), (106.06, -77.63), (106.22, -77.75), (106.38, -77.87), (106.54, -77.99), (106.7, -78.11), (106.86, -78.23), (107.1, -78.35), (107.26, -78.47), (107.42, -78.59), (107.58, -78.71), (107.74, -78.83), (107.9, -78.95), (108.06, -79.07), (108.22, -79.19), (108.38, -79.31), (108.54, -79.43), (108.7, -79.55), (108.86, -79.67), (109.1, -79.79), (109.26, -79.91), (109.42, -80.03), (109.58, -80.15), (109.74, -80.27), (109.9, -80.39), (110.06, -80.51), (110.22, -80.63), (110.38, -80.75), (110.54, -80.87), (110.7, -80.99), (110.86, -81.11), (111.1, -81.23), (111.26, -81.35), (111.42, -81.47), (111.58, -81.59), (111.74, -81.71), (111.9, -81.83), (112.06, -81.95), (112.22, -82.07), (112.38, -82.19), (112.54, -82.31), (112.7, -82.43), (112.86, -82.55), (113.1, -82.67), (113.26, -82.79), (113.42, -82.91), (113.58, -83.03), (113.74, -83.15), (113.9, -83.27), (114.06, -83.39), (114.22, -83.51), (114.38, -83.63), (114.54, -83.75), (114.7, -83.87), (114.86, -83.99), (115.1, -84.11), (115.26, -84.23), (115.42, -84.35), (115.58, -84.47), (115.74, -84.59), (115.9, -84.71), (116.06, -84.83), (116.22, -84.95), (116.38, -85.07), (116.54, -85.19), (116.7, -85.31), (116.86, -85.43), (117.1, -85.55), (117.26, -85.67), (117.42, -85.79), (117.58, -85.91), (117.74, -86.03), (117.9, -86.15), (118.06, -86.27), (118.22, -86.39), (118.38, -86.51), (118.54, -86.63), (118.7, -86.75), (118.86, -86.87), (119.1, -86.99), (119.26, -87.11), (119.42, -87.23), (119.58, -87.35), (119.74, -87.47), (119.9, -87.59), (120.06, -87.71), (120.22, -87.83), (120.38, -87.95), (120.54, -88.07), (120.7, -88.19), (120.86, -88.31), (121.1, -88.43), (121.26, -88.55), (121.42, -88.67), (121.58, -88.79), (121.74, -88.91), (121.9, -89.03), (122.06, -89.15), (122.22, -89.27), (122.38, -89.39), (122.54, -89.51), (122.7, -89.63), (122.86, -89.75), (123.1, -89.87), (123.26, -89.99), (123.42, -90.11), (123.58, -90.23), (123.74, -90.35), (123.9, -90.47), (124.06, -90.59), (124.22, -90.71), (124.38, -90.83), (124.54, -90.95), (124.7, -91.07), (124.86, -91.19), (125.1, -91.31), (125.26, -91.43), (125.42, -91.55), (125.58, -91.67), (125.74, -91.79), (125.9, -91.91), (126.06, -92.03), (126.22, -92.15), (126.38, -92.27), (126.54, -92.39), (126.7, -92.51), (126.86, -92.63), (127.1, -92.75), (127.26, -92.87), (127.42, -92.99), (127.58, -93.11), (127.74, -93.23), (127.9, -93.35), (128.06, -93.47), (128.22, -93.59), (128.38, -93.71), (128.54, -93.83), (128.7, -93.95), (128.86, -94.07), (129.1, -94.19), (129.26, -94.31), (129.42, -94.43), (129.58, -94.55), (129.74, -94.67), (129.9, -94.79), (130.06, -94.91), (130.22, -95.03), (130.38, -95.15), (130.54, -95.27), (130.7, -95.39), (130.86, -95.51), (131.1, -95.63), (131.26, -95.75), (131.42, -95.87), (131.58, -95.99), (131.74, -96.11), (131.9, -96.23), (132.06, -96.35), (132.22, -96.47), (132.38, -96.59), (132.54, -96.71), (132.7, -96.83), (132.86, -96.95), (133.1, -97.07), (133.26, -97.19), (133.42, -97.31), (133.58, -97.43), (133.74, -97.55), (133.9, -97.67), (134.06, -97.79), (134.22, -97.91), (134.38, -98.03), (134.54, -98.15), (134.7, -98.27), (134.86, -98.39), (135
```

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # Implement the kalman filter as a function using the numpy library, the matrix products can be done using the np.dot function
5 # The matrix inversion can be done using the np.linalg.inv
6
7 def kalman_filter(y, mu_tmin1, sigma_tmin1, A, C, Q, R):
8
9     # Predict the state at the next time step
10    mu_tmin1 = np.dot(A, mu_tmin1)      # (4,4)*(4,1) = (4,1)
11    sigma_tmin1 = np.dot(np.dot(A, sigma_tmin1), A.T) + Q # (4*4)
12    y_hat = np.dot(C, mu_tmin1) # (2*4)(4*1) = (2*1)
13
14    # Calculate the Kalman gain
15    S = R + np.dot(C, np.dot(sigma_tmin1, C.T)) # (2,2) + (2,4)*(4,4)*(4,2) = (2,2)
16    k_gain = np.dot(np.dot(sigma_tmin1, C.T), np.linalg.inv(S)) # (4,4)(4*2)(2*2) = (4*2)
17
18
19    # Update the estimate of the state
20    y = np.array([y]).T
21    r = np.subtract(y,y_hat) # (2*1)
22
23    mu_est = mu_tmin1 + np.dot (k_gain , r) # (4,1) + (4,2)(2,1) = (4,1)
24    I = np.identity(4)
25    sigma_est = np.dot((I - np.dot(k_gain, C)), sigma_tmin1) # ((4,4)- (4,4)(4,4))(4,4) = (4,4)
26
27    return mu_est, sigma_est
28
29
30
31 true_positions = person_dict_numpy[68] #choose a single trajectory by taking an element of the list person_dict_numpy, select a random index between 0 and 14
32
33 # Define parameter delta
34 Delta = 1 # constant velocity
35
36 A = np.array([[1, 0, Delta, 0],
37                 [0, 1, 0, Delta],
38                 [0, 0, 1, 0 ],
39                 [0, 0, 0, 1 ]])
34
35
36 # Define the measurement matrix
37 C = np.array([[1, 0, 0, 0], # x is measured directly
38                 [0, 1, 0, 0]]) # y is measured directly
39
40 # Set the standard deviation of the measurement noise
41
42 sigma_R = 0.5
43 sigma_Q = 0.003
44
45
46 # Define the process noise covariance matrix
47 Q = (sigma_Q**2)*np.eye(4)
48
49
50 # Define the measurement noise covariance matrix
51 R = (sigma_R**2)*np.eye(2)
52
53 # Set the initial state and covariance
54 mu_0 = [0.1,0.03,0.1,0.9]
55 mu_0 = np.array([mu_0])
56 mu_0 = mu_0.T
57 sigma_0 = ([[1, 0, 1, 0],
58                 [0, 1, 0, 1],
59                 [0, 0, 1, 0 ],
60                 [0, 0, 0, 1 ]])
51
52
53 # Iterate over the observed coordinates
54 y_est = []
55 y_est.append((mu_0,sigma_0))
56
57 t_series = []
58 x_kalman = []
59 y_kalman = []
60 x_observed = []
61 y_observed = []
62 x_true = []
63 y_true = []
64
65
66 for t in range(1,len(true_positions)):
67     # Get the observed coordinates at time t
68     # Note: the observed position is simulated by adding gaussian noise to the true_positions
69     y = true_positions[t] + np.random.normal(0,sigma_R,2)
70     #apply the kalman filter on the observed coordinates
71
72     y_est.append(kalman_filter(y, y_est[t-1][0], y_est[t-1][1], A, C, Q, R))
73     x_kalman.append(y_est[t][0].tolist())      # Creating x series
74     y_kalman.append(y_est[t][1].tolist())      # Creating y series
75
76     t_series.append(t)                      # Creating a time series for plotting
77
78     x_true.append(true_positions[t][0])
79     y_true.append(true_positions[t][1])
80
81     x_observed.append(y[0])
82     y_observed.append(y[1])
83
84
85 print()
86
87 print('t_series :',t_series)
88 print('\n x_kalman :',x_kalman)
89 print('\n y_kalman :',y_kalman)
90 print('x_true :',x_true)
91 print('y_true :',y_true)

```

```

104 print('\n y_true :', y_true)
105 print('\n x_observed :', x_observed)
106 print('\n y_observed :', y_observed)
107
108
109 #Use the matplotlib library to plot the true positions, observed positions y and the results of the kalman filtering
110 #You should obtain a plot which resembles the one in the figure
111 #Plot the trajectory of the x coordinate over time and the trajectory of the y coordinate over time into two separate plot
112

t_series : [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45]

x_kalman : [[-4.255485475408207], [-5.6585366593730235], [-5.281772407415517], [-4.973906366166316], [-5.039233049046499], [-4.734219723440784], [-4.076267638507], [-3.780997340821777], [-3.825637

y_kalman : [[6.960592743462526], [9.909720986350843], [10.75381601625452], [9.493094513759166], [8.23563598804174], [7.016455056104876], [6.645319751710337], [6.2370244878535965], [5.685297481031

x_true : [-4.15, -4.12, -4.09, -4.06, -4.03, -4.01, -3.97, -3.87, -3.73, -3.56, -3.39, -3.24, -3.14, -3.12, -3.16, -3.22, -3.3, -3.35, -3.37, -3.35, -3.28, -3.16, -2.99, -2.75, -2.47, -2.15, -1.81

y_true : [6.54, 6.31, 6.09, 5.89, 5.71, 5.54, 5.39, 5.27, 5.16, 5.05, 4.93, 4.8, 4.64, 4.46, 4.26, 4.03, 3.76, 3.46, 3.14, 2.82, 2.54, 2.26, 2.0, 1.73, 1.46, 1.19, 0.91, 0.63, 0.35, 0.1, -0.14, -0

x_observed : [-4.726774817824198, -4.278665979030187, -3.1218432666372156, -3.833635430463483, -4.613177373336612, -3.9971115367479637, -3.0612939494674256, -3.540978375919125, -4.24210179051258, -4.24210179051258

y_observed : [6.563140631107409, 7.011033198578276, 6.92107599555592, 5.52038121249522, 5.737102115256241, 5.354001020363288, 6.545660925692861, 6.144584680217733, 5.3976227925041345, 4.658056620

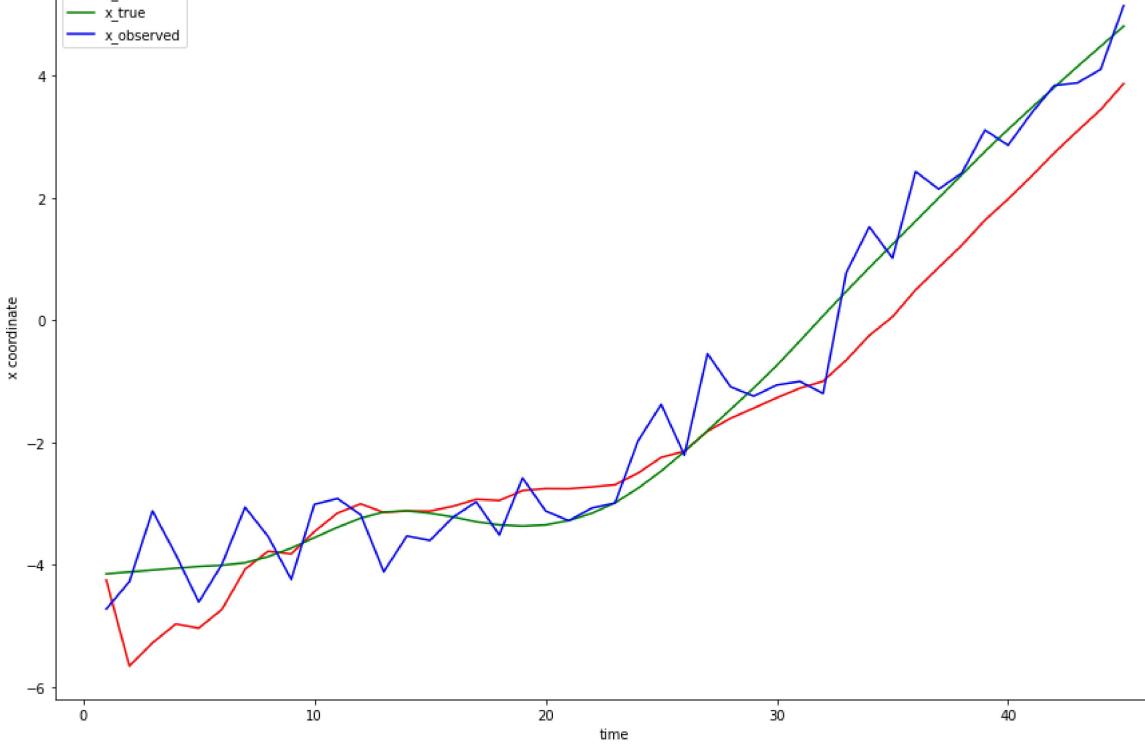
```

```

1 # Plotting x
2 f = plt.figure()
3 f.set_figwidth(15)
4 f.set_figheight(10)
5 plt.plot(t_series, x_kalman , color = 'red')
6 plt.plot(t_series, x_true , color = 'green')
7 plt.plot(t_series, x_observed , color = 'blue')
8 plt.legend(['x_kalman', 'x_true', 'x_observed'])
9 plt.title('x_kalman vs x_true vs x_observed')
10 plt.xlabel ('time')
11 plt.ylabel('x coordinate')
12 plt.show()
13
14
15
16
17

```

x_kalman vs x_true vs x_observed

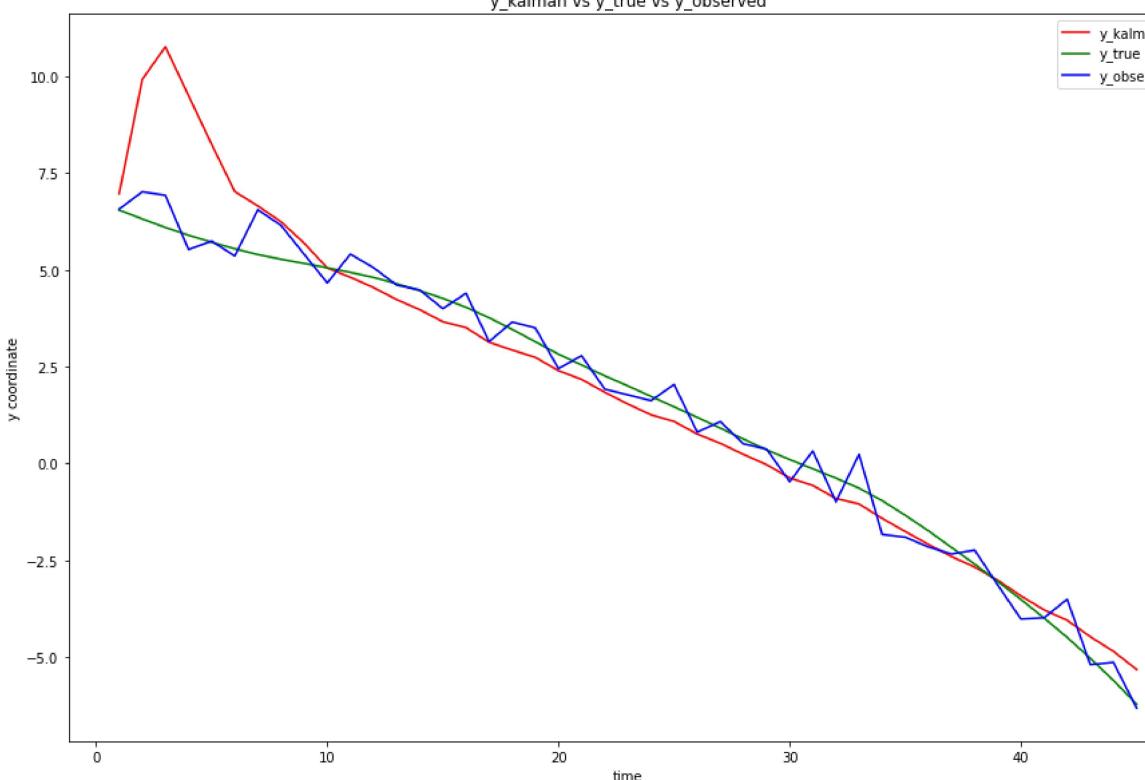


```

1 # Plotting y
2 f = plt.figure()
3 f.set_figwidth(15)
4 f.set_figheight(10)
5 plt.plot(t_series, y_kalman , color = 'red')
6 plt.plot(t_series, y_true , color = 'green')
7 plt.plot(t_series, y_observed , color = 'blue')
8 plt.legend(['y_kalman', 'y_true', 'y_observed'])
9 plt.title('y_kalman vs y_true vs y_observed')
10 plt.xlabel ('time')
11 plt.ylabel('y coordinate')
12 plt.show()

```

y_kalman vs y_true vs y_observed



Student's comments to COMPUTER LAB 4 - Kalman filter

- Kalman filter is used to find/predict the optimal estimate of a system when we can not directly measure or estimate it.
- It does it by combining the 1-measurement/observation(y) and 2-prediction(z) in two different phases : I.predict II.update

- We can make kalman filter follow the measurement or the model by:

1. Tuning measurement noise: The measurement noise is the uncertainty associated with the measurement, and it is used to determine the weight assigned to the measurement update. By decreasing the measurement noise, you increase the weight assigned to the measurement, and make the Kalman estimate follow the measured value more closely
2. Tuning process/model noise: The process noise is the uncertainty associated with the model, and it is used to determine the weight assigned to the model prediction. By decreasing the process noise, you increase the weight assigned to the model, and make the Kalman estimate trust the model more.

- We choose the delta = 1 since it is a constant velocity model

NOTE for figure:

- The Kalman filter may not have precise estimation at the beginning due to uncertainties in the initial conditions, model, and measurements, and **it may take some time for the estimate to converge to the true state of the system**

[Colab notebook](#) [GitHub notebook](#)

