

# Modelling and Analysis of Data

Exam no. 39

23. januar 2022

## Indhold

<b>1 Problem 1 - Maximum Likelihood Estimation</b>	<b>1</b>
1.1 . . . . .	1
<b>2 Problem 2 - Statistics</b>	<b>3</b>
<b>3 Problem 3 - Principal Component Analysis</b>	<b>3</b>
3.1 1) - Mean point of the data . . . . .	3
3.2 2) - Two eigenvalues . . . . .	4
3.3 3) - Two eigenvectors . . . . .	4
<b>4 Problem 4 - Regression</b>	<b>6</b>
4.1 a) . . . . .	6
4.2 b) . . . . .	6
4.3 c) . . . . .	7
<b>5 Problem 5 - Random Forests</b>	<b>8</b>
5.1 Finding Thresholds . . . . .	8
<b>6 Problem 6 - Classification &amp; Validation</b>	<b>11</b>
6.1 a) - Implement random forest training . . . . .	11
6.2 b) - Finding the optimal set of random forest classifier parameters . . .	11
6.3 c) - My results . . . . .	12
6.4 Code for a) & b) . . . . .	12
<b>7 Problem 7 - K-means Clustering &amp; Principal Component Analysis</b>	<b>15</b>
7.1 a) . . . . .	15
7.2 b) . . . . .	15
7.3 c) . . . . .	16
7.4 d) . . . . .	16
7.5 e) . . . . .	17
<b>8 Appendix A</b>	<b>18</b>
8.1 Q6_c . . . . .	18

## 1 Problem 1 - Maximum Likelihood Estimation

### 1.1

Assessing the following:

$$f(x) \begin{cases} \frac{1}{\sigma\sqrt{2\pi}} \cdot \frac{1}{x} \cdot \exp\left(-\frac{1}{2}\left(\frac{\log(x)-\mu}{\sigma}\right)^2\right) & \text{for } x \in \mathbb{R}_+ \\ 0 & \text{otherwise} \end{cases}$$

Since I can fix the parameter  $\sigma$  I can express the likelihood  $L$  of observing the given dataset  $x_1, \dots, x_N$  of  $N$  i.i.d samples from  $X$  as the joint distribution all the  $N$  i.i.d samples conditioning on  $\mu$  as:

$$L = p(x_1, \dots, x_N | \mu) = \prod_{n=1}^N p(x_n | \mu) = \prod_{n=1}^N f(x_n; \mu)$$

The PDF is now considered a function of both  $x_n$  and  $\mu$  I'm going to make take the logarithm of the expression, to make it more manageable, without changing the estimate. When trying to find the optimal  $\hat{\mu}$  that maximizes the likelihood  $L$ .

$$\begin{aligned} \log L &= \log\left(\prod_{n=1}^N f(x_n; \mu)\right) \\ &= \sum_{n=1}^N \log(f(x_n; \mu)) \\ &= \sum_{n=1}^N \left[ \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right) + \log\left(\frac{1}{x_n}\right) + \log\left(\exp\left(-\frac{1}{2}\left(\frac{\log(x)-\mu}{\sigma}\right)^2\right)\right) \right] \\ &= \sum_{n=1}^N \left[ \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right) + \log\left(\frac{1}{x_n}\right) - \frac{1}{2}\left(\frac{\log(x)-\mu}{\sigma}\right)^2 \right] \end{aligned}$$

To find the MLE  $\hat{\mu}$  for the parameter  $\mu$  I optimize  $L$  by solving  $\frac{\partial \log L}{\partial \mu} = 0$  When taking derivatives I'm only interested in terms that involve  $\mu$  and so I have the following expression.

$$\begin{aligned} \frac{\partial \log L}{\partial \mu} &= \frac{\partial}{\partial \mu} \left[ \sum_{n=1}^N -\frac{1}{2} \left( \frac{\log(x)-\mu}{\sigma} \right)^2 \right] \\ &= \sum_{n=1}^N -\left( \frac{\log(x)-\mu}{\sigma} \right) \cdot \left[ \frac{\log(x)-\mu}{\sigma} \cdot \frac{\partial}{\partial \mu} \right] \quad (\text{Chain rule}) \\ &= \sum_{n=1}^N -\left( \frac{\log(x)-\mu}{\sigma} \right) \cdot \left( -\frac{1}{\sigma} \right) \\ &= \frac{1}{\sigma^2} \sum_{n=1}^N \log(x_n) - \mu \end{aligned}$$

solving  $\frac{\partial \log L}{\partial \mu} = 0$  with respect to  $\mu$

$$\begin{aligned}
 0 &= \frac{1}{\sigma^2} \sum_{n=1}^N \log(x_n) - \mu = \sum_{n=1}^N \log(x_n) - \mu = \left[ \sum_{n=1}^N \log(x_n) \right] - N\mu \\
 &\Downarrow \\
 N\mu &= \sum_{n=1}^N \log(x_n) \\
 &\Downarrow \\
 \mu &= \frac{1}{N} \sum_{n=1}^N \log(x_n) = \frac{1}{N} \log(x_1 \cdot x_2 \cdot \dots \cdot x_N) = \log\left((x_1 \cdot x_2 \cdot \dots \cdot x_N)^{\frac{1}{N}}\right) \\
 &= \log(\sqrt[N]{x_1 \cdot x_2 \cdot \dots \cdot x_N})
 \end{aligned}$$

The above calculations is a necessary condition in order to find an optimal solution. Furthermore, to check if I'm dealing with a global maximum or a global minimum I have to take the second derivative of my function.

$$\frac{\partial^2 \log L}{\partial \mu^2} = -\frac{N}{\sigma^2}$$

Assessing the result from the second derivative, I can see that my function will always be negative. Thus by the second derivative test

$$\hat{\mu} = \log(\sqrt[N]{x_1 \cdot x_2 \cdot \dots \cdot x_N})$$

is a global maximum

## 2 Problem 2 - Statistics

## 3 Problem 3 - Principal Component Analysis

### 3.1 1) - Mean point of the data

Step one is to calculate the mean point of the data with respect to both x and y

$$\bar{x} = \frac{1}{8}(0.6 + 0.5 + 1.1 + (-0.5) + 0.8 + 0.2 + (-0.1) + 1.0) = 0.4500$$

$$\bar{y} = \frac{1}{8}(1.1 + 1.0 + 2.0 + 0.2 + (-0.1) + (-0.1) + (-1.5) + 2.5) = 0.6375$$

Step two is to normalize the dataset. This is done by taking each original point and subtract the mean.  $p - [\bar{x}, \bar{y}]$

Normalizing each point

$$p1 = (0.6, 1.1) - [0.4500, 0.6375] = [0.1500, 0.4625] \quad (1)$$

$$p2 = (0.5, 1.0) - [0.4500, 0.6375] = [0.0500, 0.3625] \quad (2)$$

$$p3 = (1.1, 2.0) - [0.4500, 0.6375] = [0.6500, 1.3625] \quad (3)$$

$$p4 = (-0.5, 0.2) - [0.4500, 0.6375] = [-0.9500, -0.4375] \quad (4)$$

$$p5 = (0.8, -0.1) - [0.4500, 0.6375] = [0.3500, -0.7375] \quad (5)$$

$$p6 = (0.2, -0.1) - [0.4500, 0.6375] = [-0.2500, -0.7375] \quad (6)$$

$$p7 = (-0.1, -1.5) - [0.4500, 0.6375] = [-0.5500, -2.1375] \quad (7)$$

$$p8 = (1.0, 2.5) - [0.4500, 0.6375] = [0.5500, 1.8625] \quad (8)$$

Step three is to calculate the covariances for XX, XY and YY using the normalized values for p1-p8

$$\begin{aligned} cov(xx) = \frac{1}{8} & \left( 0.15^2 + 0.05^2 + 0.65^2 + (-0.95)^2 \right. \\ & \left. + (-0.25)^2 + (-0.55)^2 + 0.55^2 \right) \approx 0.27 \end{aligned} \quad (9)$$

$$\begin{aligned} cov(yy) = \frac{1}{8} & \left( 0.4625^2 + 0.3625^2 + 1.3625^2 + (-0.4375)^2 + \right. \\ & \left. (-0.7375)^2 + (-0.7375)^2 + (-2.1375)^2 + 1.8625^2 \right) \approx 1.44 \end{aligned} \quad (10)$$

$$\begin{aligned} cov(xy) = \frac{1}{8} & \left( (0.15 \cdot 0.4625 + 0.05 \cdot 0.3625 + 0.65 \cdot 1.3625 + ((-0.95) \cdot (-0.4375)) \right. \\ & + 0.35 \cdot (-0.7375) + ((-0.25) \cdot (-0.7375)) + ((-0.55) \cdot (-2.1375)) \\ & \left. + 0.55 \cdot 1.8625 \right) \approx 0.44 \end{aligned} \quad (11)$$

having done the three calculations I can now form the covariance matrix inserting the calculated values into their respective place:

$$cov = \begin{pmatrix} xx & xy \\ xy & yy \end{pmatrix} =$$
$$cov = \begin{pmatrix} 0.27 & 0.44 \\ 0.44 & 1.44 \end{pmatrix}$$

### 3.2 2) - Two eigenvalues

I order to find the eigenvalues and the eigenvectors I need to solve **Characteristic equation** of the covariance matrix. The Characteristic equation is denoted as:

$$\det(A - \lambda I) = 0$$

To calculate the determinant you calculate the following:

$$\text{determinant} = \begin{pmatrix} a & b \\ c & d \end{pmatrix} = ad - bc$$

Solving the system equation that I get from using the Characteristics equation on the covariance matrix I get the following eigenvalues.

$$\begin{pmatrix} 0.27 - \lambda & 0.44 \\ 0.44 & 1.44 - \lambda \end{pmatrix} = 0$$

First I get the equation by calculating the determinant

$$\det(cov) \rightarrow ((0.27 - \lambda) \cdot (1.44 - \lambda)) - (0.44 \cdot 0.44)$$

Then setting the equation equal to zero and solving it will give me the two values for  $\lambda$  which is the eigenvalues.

$$((0.2675 - \lambda) \cdot (1.44 - \lambda)) - (0.44 \cdot 0.44) = 0$$

$$\text{eigenvalue}_1 = 1.59 \quad \text{eigenvalue}_2 = 0.12$$

### 3.3 3) - Two eigenvectors

I order to find the eigenvectors I use the definition of an eigenvector  $A\bar{x} = \lambda\bar{x}$  and solve this equation.

$$cov \cdot \bar{x} = \lambda \cdot \bar{x}$$

finding the eigenvectors, using the equation stated above, I have two equations systems.

$$0.27x_1 + 0.44x_2 = 1.59x_1$$

$$0.44x_1 + 1.44x_2 = 1.59x_2$$

and

$$0.27x_1 + 0.44x_2 = 0.12x_1$$

$$0.44x_1 + 1.44x_1 = 0.12x_2$$

solving these systems of equations using basic equation arithmetic

$$0.44x_2 = (1.59 - 0.27)x_1$$

$$x_2 = \frac{1.32}{0.44}x_1$$

$$x_2 = 3x_1$$

thus, if  $x_1 = 1$

then  $x_2 = 3$

Now that I have the values for  $x_1$  and  $x_2$  I can find the actual eigenvector, by finding the unit vector for the vector  $\bar{x} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$  to find the unit vector I'm using the formula  $\frac{1}{\|\bar{v}\|} \cdot \bar{v}$  which is equivalent to

$$\frac{1}{\sqrt{(1^2 + 3^2)}} \cdot \begin{bmatrix} 1 \\ 3 \end{bmatrix}$$

$$\frac{3}{\sqrt{(1^2 + 3^2)}} \cdot \begin{bmatrix} 1 \\ 3 \end{bmatrix}$$

$$x_1 = 0.316 \quad x_2 = 0.948$$

I've shown the calculations for the first system, with the first equation and I think it's fair to say that finding the second eigenvector is by doing exactly the same, with the only exception being that lambda now is the second eigenvalue. Thus my eigenvectors are:

$$Eigenvector_1 = \begin{bmatrix} 0.316 \\ 0.949 \end{bmatrix} \quad Eigenvector_2 = \begin{bmatrix} -0.949 \\ 0.316 \end{bmatrix}$$

note that signs of the eigenvector can be different. Depending on how you choose to solve your equation systems.

## 4 Problem 4 - Regression

### 4.1 a)

Using the euclidean distance I get an RSME value of

$$RMSE = 0.82430$$

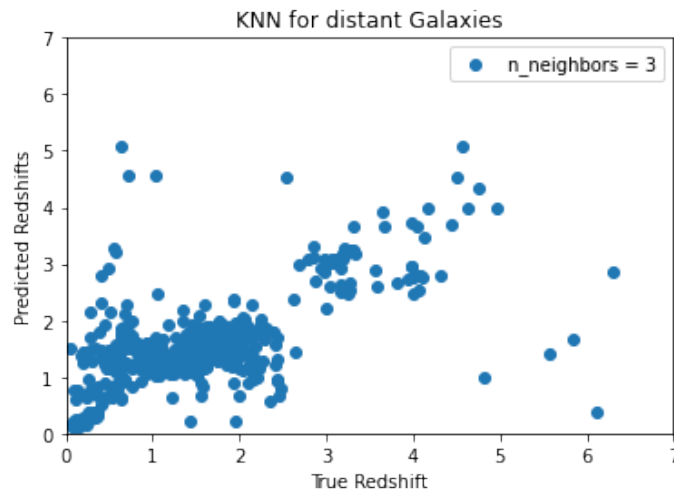


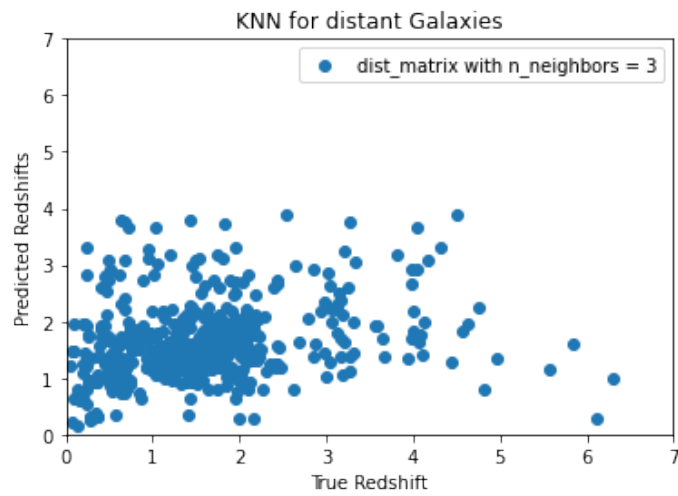
Figure 1: Plot of kNN with  $k = 3$

### 4.2 b)

Using the other distance calculation as given in the exam set I get the following RSME value

$$RMSE = 1.100519$$





Figur 2: Plot of kNN with  $k = 3$  and distance matrix

The distance matrix will lower/lessen the influence of all the features/attributes except for the last two, since they are set to 1.0.

#### 4.3 c)

## 5 Problem 5 - Random Forests

### 5.1 Finding Thresholds

There are four places that I should take into consideration as possible thresholds.

option	$x_i$	$x_{i+1}$
1	20	30
2	60	70
3	90	100
4	100	110

There are four positions that could be possible threshold, however, I'm going to ignore option 4, since it will be the same as option 1, except with reversed values. These four places have been found using the same approach as shown in our lecture, and lecture video (L10\_Classification\_Regression\_2\_Part3, starting at 17min 50sec.). This means, that I'm only going to look at positions where there is a "change" in labels for a given datapoint compared to its predecessor or successor. Since there is no reason to look at possible thresholds where the neighbors to a given feature have the same labels.

This process can be formalized:

$$x_i + \left(\frac{x_{i+1} - x_i}{2}\right) : y_{i+1} \neq y_i$$

As stated, option 4 is the same as option 1, so I will only perform the calculations for option 1-3.

$$\text{Option 1} = 20 + \frac{(30 - 20)}{2} = 25$$

$$\text{Option 2} = 60 + \frac{(70 - 60)}{2} = 65$$

$$\text{Option 3} = 90 + \frac{(100 - 90)}{2} = 95$$

Using my calculated values for possible thresholds I have three different decision trees to take into consideration.

Tree 1:

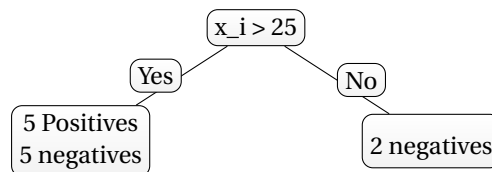
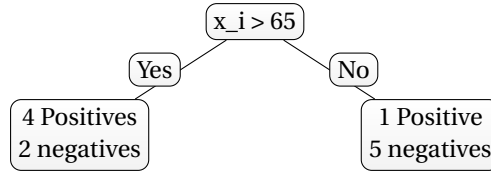


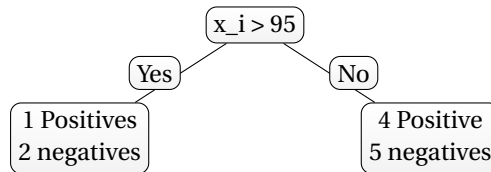
Figure 3: Decision tree with threshold = 25

Tree 2:



Figur 4: Decision three with threshold = 65

Tree 3:



Figur 5: Decision three with threshold = 95

From here, I will start by calculating the Entropy values for splits in my trees.  
Aswell as the information gain using the functions given in our lecture slides (L10\_Classification\_Regression\_2\_Part3 p. 8)

$$\text{Entropy } H(T) = -p_+ \cdot \log_2 \cdot p_+ - p_- \cdot \log_2 \cdot p_-$$

$$\text{Gain } (T, j) = H(T) - \sum_{i=0,1} \frac{|T_i|}{|T|} H(T_i)$$

Calculating the information gain for tree 1 where the threshold is = 25:

$$\begin{aligned} \text{Gain}(T, j) &= H(5 \text{ pos}, 7 \text{ neg}) - \frac{10}{12} H(5 \text{ pos}, 5 \text{ neg}) - \frac{2}{12} H(0 \text{ pos}, 2 \text{ neg}) \\ &= 0.979 - \frac{10}{12} 1 - \frac{2}{12} 0 \approx 0.145 \end{aligned}$$

Calculating the information gain for tree 2 where the threshold is = 65:

$$\begin{aligned} \text{Gain}(T, j) &= H(5 \text{ pos}, 7 \text{ neg}) - \frac{6}{12} H(4 \text{ pos}, 2 \text{ neg}) - \frac{6}{12} H(1 \text{ pos}, 5 \text{ neg}) \\ &= 0.979 - \frac{10}{12} 0.918 - \frac{2}{12} 0.650 \approx 0.195 \end{aligned}$$

Calculating the information gain for tree 3 where the threshold is = 95:

$$\begin{aligned} \text{Gain}(T, j) &= H(5 \text{ pos}, 7 \text{ neg}) - \frac{6}{12} H(4 \text{ pos}, 2 \text{ neg}) - \frac{6}{12} H(1 \text{ pos}, 5 \text{ neg}) \\ &= 0.979 - \frac{10}{12} 0.918 - \frac{2}{12} 0.991 \approx 0.0062 \end{aligned}$$

Since the goal is to find the tree which gives us the highest information gain. I can see that the second decision tree with a threshold of 65 is the best one.

Thus, the optimal threshold is 65.

## 6 Problem 6 - Classification & Validation

NOTE that I've included the full python code in the appendix

### 6.1 a) - Implement random forest training

I've implemented the Random Forest Training using the existing library from *Skikit-learn*, more specific the package: `sklearn.ensemble.RandomForestClassifier`. Below I've added the most notable part from my code in order to load the data. I'm reusing the "loaddata" function that's been handed out, though I've modified it slightly in line 12 in the code snippet.

```
1  # Modified the loaddata function to make the training features
   ↳ into 2d np.array, more specifically the 't value'
2  def loaddata(filename):
3      """Load the balloon data set from filename and return t, X
4          t - N-dim. vector of target (temperature) values
5          X - N-dim. vector containing the inputs (lift) x for
   ↳ each data point
6      """
7      # Load data set from CSV file
8      Xt = np.loadtxt(filename, delimiter=',')
9
10     # Split into data matrix and target vector
11     X = Xt[:,0] # Matrix
12     t = Xt[:,1:] # Vector
13
14     return t, X
15
16 # Loads the training data
17 X_train, t_train =
   ↳ loaddata("./accent-mfcc-data_shuffled_train.txt")
18
19 # Loads the validation data
20 X_validation, t_validation =
   ↳ loaddata("./accent-mfcc-data_shuffled_validation.txt")
```

To see how I've chosen to set up the Random Forest Classifier, training the algorithm using the it to predict - see the submitted code in section 6.4

### 6.2 b) - Finding the optimal set of random forest classifier parameters

I've chosen to use two extra packages in order to implement my Random Forest algorithm. It hasn't been specified that I'm not allowed to use other packages from the **Sklean** library so I assume that this is a possible solution.

I've made use of **ParameterGrid** from the *sklearn.model\_selection* and **accuracy\_score** from the *sklearn.metrics*.

I'm using ParameterGrid to generate a list of all the possible permutations of the specified parameters used for a given iteration to perform the optimized search. With each iteration the parameters and the resulting performance metrics are added to a list which can be sorted easily. I'm currently using the accuracy\_score as the first metric.

see the submitted code in section 6.4

### 6.3 c) - My results

Looking at the results, it's possible to see that the performance of the algorithm increases when tree depth, number of features that's taken into consideration and the complexity of the criterion is also increased. However it comes with the cost of a more computational complexity aswell. The second listing in section 6.4 shows the terminal output from the optimization loop. I've appended a sorted print of the permuted parameters in the appendix.

Assessing the result, I get the optimal results from the algorithm with the parameters:

```
{'criterion': 'entropy', 'max_tree_depth': 10, 'max_features': 'sqrt'}
```

### 6.4 Code for a) & b)

testtesttest

```
1  # Adding the parameters to test given in the exam question.
   ↳ Using these to find the optimal set of random forest
   ↳ classifier parameters.
2  random_f_parameters = {
3      'criterion'      : ['gini', 'entropy'],
4      'max_tree_depth' : [2,5,6,10,15],
5      'max_features'   : ['sqrt', 'log2']
6  }
7
8  # Empty array for the result metrics
9  res = np.empty((0,3))
10
11 # Looping through the chosen(given) parameters and setting up
   ↳ the Random forest classifier each time.
12 for params in list(ParameterGrid(random_f_parameters)):
13     clf = RandomForestClassifier(
14         criterion      = params['criterion'],
15         max_depth      = params['max_tree_depth'],
16         max_features   = params['max_features'])
17
```

```
18     # Training using the created classifier with the given
19     ↪ parameters.
20     clf.fit(X_train, t_train)
21
22     # number of correctly classified validation samples
23     t_prediction = clf.predict(X_validation)
24     acc_score = accuracy_score(t_validation, t_prediction)
25
26     # probability associated with classification
27     t_probability = clf.predict_proba(X_validation)
28     probability_score = np.mean([t_probability[int(t_val)]
29     ↪ for (t_probability, t_val)
30     ↪ in zip(t_probability,
31     ↪ t_validation)])
32
33     print("Accuracy score: %.2f"
34     ↪ %acc_score)
35     print("Average probability assigned to correct classes:
36     ↪ %.2f"
37     ↪ %probability_score)
38
39     # print the parameters if new ones are more optimal than
40     ↪ previously tried ones.
41     if len(res) > 0 and (acc_score > res[-1,1]
42     ↪ or (acc_score == res[-1,1]
43     ↪ and probability_score >
44     ↪ res[-1,2])):
45         print(params)
46
47     # accumulate results
48     res = np.append(res, np.array([[params, acc_score,
49     ↪ probability_score]]), axis=0)
50     # sort the results in ascending order
51     res = res[np.lexsort((res[:,1], res[:,2]))]
```

Accuracy score: 0.51  
Average probability assigned to correct classes: 0.36  
Accuracy score: 0.70  
Average probability assigned to correct classes: 0.48  
{'criterion': 'gini', 'max\_features': 'sqrt', 'max\_tree\_depth': 5}  
Accuracy score: 0.73  
Average probability assigned to correct classes: 0.51  
{'criterion': 'gini', 'max\_features': 'sqrt', 'max\_tree\_depth': 6}  
Accuracy score: 0.78  
Average probability assigned to correct classes: 0.56  
{'criterion': 'gini', 'max\_features': 'sqrt', 'max\_tree\_depth': 10}  
Accuracy score: 0.78  
Average probability assigned to correct classes: 0.57  
{'criterion': 'gini', 'max\_features': 'sqrt', 'max\_tree\_depth': 15}  
Accuracy score: 0.52  
Average probability assigned to correct classes: 0.36  
Accuracy score: 0.70  
Average probability assigned to correct classes: 0.49  
Accuracy score: 0.73  
Average probability assigned to correct classes: 0.51  
Accuracy score: 0.78  
Average probability assigned to correct classes: 0.56  
Accuracy score: 0.78  
Average probability assigned to correct classes: 0.55  
Accuracy score: 0.53  
Average probability assigned to correct classes: 0.37  
Accuracy score: 0.73  
Average probability assigned to correct classes: 0.51  
Accuracy score: 0.79  
Average probability assigned to correct classes: 0.54  
{'criterion': 'entropy', 'max\_features': 'sqrt', 'max\_tree\_depth': 6}  
Accuracy score: 0.82  
Average probability assigned to correct classes: 0.57  
{'criterion': 'entropy', 'max\_features': 'sqrt', 'max\_tree\_depth': 10}  
Accuracy score: 0.82  
Average probability assigned to correct classes: 0.57  
{'criterion': 'entropy', 'max\_features': 'sqrt', 'max\_tree\_depth': 15}  
Accuracy score: 0.56  
Average probability assigned to correct classes: 0.37  
Accuracy score: 0.71  
Average probability assigned to correct classes: 0.51  
Accuracy score: 0.75  
Average probability assigned to correct classes: 0.55  
Accuracy score: 0.84  
Average probability assigned to correct classes: 0.58  
{'criterion': 'entropy', 'max\_features': 'log2', 'max\_tree\_depth': 10}  
Accuracy score: 0.82  
Average probability assigned to correct classes: 0.57

Side 14 of 17

Figur 6: Caption of figure here



## 7 Problem 7 - K-means Clustering & Principal Component Analysis

### 7.1 a)

I've decided to make my own function to normalize my data. It normalizes the input value  $x$  by subtracting the mean and division by the standard deviation.

```
1  def normalizeData(x):  
2      '''  
3      Normalizes the input using mean and the standard deviation.  
4      Normalization is done on all the columns in the 2d  
5      ↪ numpy-array 'x'  
6      Returns the normalized 'x'  
7      '''  
8      return (x - np.mean(x, axis=0)) / np.std(x, axis=0)
```

### 7.2 b)

Below is a codesnippet of my implementation of the k-mean clustering. I've defined multiple helperfunctions, which can be found in the source code or in the appendix.

```
1  def k_mean_clustering(k, data, centroids):  
2      """  
3      K-means clustering  
4      Params:  
5      -----  
6      k : number of clusters  
7      data : (n_samples, n_features) (normalized) data matrix  
8      Returns:  
9      -----  
10     assignments, intra_cluster_dist : tuple  
11     """  
12     current_assignments = assign_datapoints_to_centroids(data,  
13     ↪ centroids)  
14     new_assignments = [] # initial empty value  
15     while not np.array_equal(current_assignments,  
16     ↪ new_assignments):
```

```
20     # repeat until assignments does not change
21     centroids = calculate_new_centroids(data,
22         ↪ current_assignments, centroids)
23     current_assignments = new_assignments
24     new_assignments = assign_datapoints_to_centroids(data,
25         ↪ centroids)
26
27     intra_cluster_dist = compute_sum_intra_cluster_dist(data,
28         ↪ current_assignments, centroids)
29
30     return current_assignments, intra_cluster_dist
```

The results of the k-means cluster algorithm with  $k = 3$ . The algorithm has been implemented with ***np.random.default\_rng*** which constructs a random generator so that a desired "seed\_value" can be reproduced.

### 7.3 c)

Here's a print of the values that I get from the k-means clustering algorithm, when using seed: 1

```
Calculations with seed value: 1
Smallest Intra-Cluster Distance: 268.5510
Number of samples in cluster 0: 67
Number of samples in cluster 1: 66
Number of samples in cluster 2: 67
```

### 7.4 d)

I've reused code that has been handed out in the assignments aswell as my own implementation of PCA from Assignment 5.

The code below shows how I've reused code in order to reduce the dimensions of data from the k-means clustering algorithm.

```
1     def __transformData(features, PCevecs):
2         """
3         Reused from A5
4         """
5         return np.dot(features, PCevecs[:, 0:2])
6
7     PCEvals, PCevecs = __PCA(normalized_data)
8
9     # Convert data to two dimemsions using PCA
10    features2D = __transformData(normalized_data, PCevecs)
11    centroids2D = __transformData(centroids, PCevecs)
```

I've also added a single line of code to the already provided "\_\_visualizeLabels()"function

```
1 plt.scatter(centroids[:, 0], centroids[:,1], c = 'black',  
    ↪ s=100)
```

### 7.5 e)

Figure 7 shows the plotted data of the k-means after PCA. There is a clear separation between the points around each centroid. My result look similar as the example show in our exam question, except that my clusters have been flipped vertically.

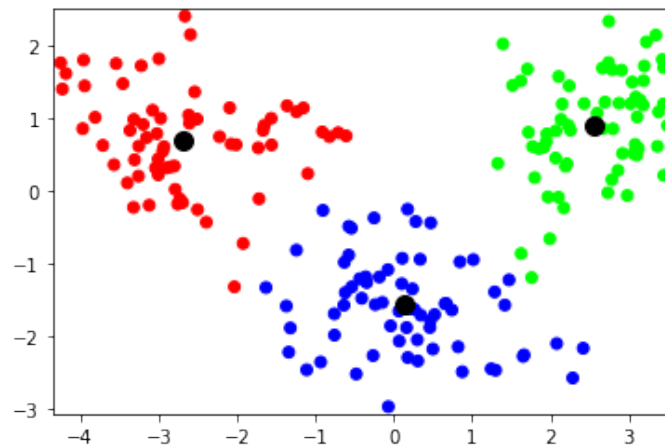


Figure 7: Plot of the K-means cluster after PCA

## 8 Appendix A

### 8.1 Q6\_c

This is a sorted print of the permuted parameters from the Random Forest algorithm. which which yields matrices that are increacinly better.

I've rearranged the terminal output slightly for better readability.

```
{'criterion': 'gini', 'max_features': 'log2', 'max_tree_depth': 2}
{'criterion': 'gini', 'max_features': 'sqrt', 'max_tree_depth': 2}
{'criterion': 'entropy', 'max_features': 'sqrt', 'max_tree_depth': 2}
{'criterion': 'entropy', 'max_features': 'log2', 'max_tree_depth': 2}
{'criterion': 'gini', 'max_features': 'sqrt', 'max_tree_depth': 5}
{'criterion': 'gini', 'max_features': 'log2', 'max_tree_depth': 5}
{'criterion': 'entropy', 'max_features': 'log2', 'max_tree_depth': 5}
{'criterion': 'gini', 'max_features': 'log2', 'max_tree_depth': 6}
{'criterion': 'entropy', 'max_features': 'sqrt', 'max_tree_depth': 5}
{'criterion': 'gini', 'max_features': 'sqrt', 'max_tree_depth': 6}
{'criterion': 'entropy', 'max_features': 'sqrt', 'max_tree_depth': 6}
{'criterion': 'entropy', 'max_features': 'log2', 'max_tree_depth': 6}
{'criterion': 'gini', 'max_features': 'log2', 'max_tree_depth': 15}
{'criterion': 'gini', 'max_features': 'sqrt', 'max_tree_depth': 10}
{'criterion': 'gini', 'max_features': 'log2', 'max_tree_depth': 10}
{'criterion': 'gini', 'max_features': 'sqrt', 'max_tree_depth': 15}
{'criterion': 'entropy', 'max_features': 'sqrt', 'max_tree_depth': 10}
{'criterion': 'entropy', 'max_features': 'sqrt', 'max_tree_depth': 15}
{'criterion': 'entropy', 'max_features': 'log2', 'max_tree_depth': 15}
{'criterion': 'entropy', 'max_features': 'log2', 'max_tree_depth': 10}
```