

Modelling and Analysis of Data

Exam no. 39

20. januar 2022

Indhold

1 Problem 1 - Maximum Likelihood Estimation	1
1.1	1
2 Problem 2 - Statistics	1
2.1	1
3 Problem 3 - Principal Component Analysis	1
3.1 1) - Mean point of the data	1
3.2 2) - Two eigenvalues	2
3.3 3) - Two eigenvectors	3
4 Problem 4 - Regression	3
4.1 a)	3
4.2 b)	4
4.3 c)	5
5 Problem 5	5
5.1	5
6 Problem 6 - Classification & Validation	5
6.1 a) - Implement random forest training	5
6.2 b) - Finding the optimal set of random forest classifier parameters . . .	6
6.3 c) - My results	6
6.4 Code for a) & b)	6
7 Problem 7 - K-means Clustering & Principal Component Analysis	9
7.1 a)	9
7.2 b)	9
8 Appendix A	10
8.1	10
9 Appendix B	11
9.1	11
10 Appendix Q6	12
11 Appendix X	15
11.1 Q6_c	15

1 Problem 1 - Maximum Likelihood Estimation

1.1

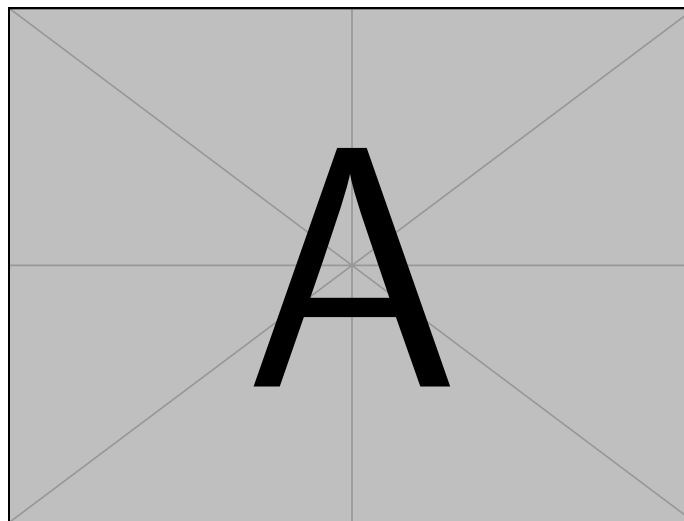
```
1 def minted_for_code_here
```

2 Problem 2 - Statistics

2.1

This is my implementation of KNN

```
1 def minted_for_code_here
```



Figur 1: Caption of figure here

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

$$\begin{aligned} cov(xx) = \frac{1}{8} \left(0.15^2 + 0.05^2 + 0.65^2 + (-0.95)^2 \right. \\ \left. + (-0.25)^2 + (-0.55)^2 + 0.55^2 \right) = 0.2675 \end{aligned} \quad (9)$$

$$\begin{aligned} cov(yy) = \frac{1}{8} \left(0.4625^2 + 0.3625^2 + 1.3625^2 + (-0.4375)^2 + \right. \\ \left. (-0.7375)^2 + (-0.7375)^2 + (-2.1375)^2 + 1.8625^2 \right) = 1.439843750 \end{aligned} \quad (10)$$

$$\begin{aligned} cov(xy) = \frac{1}{8} \left((0.15 \cdot 0.4625 + 0.05 \cdot 0.3625 + 0.65 \cdot 1.3625 + ((-0.95) \cdot (-0.4375)) \right. \\ \left. + 0.35 \cdot (-0.7375) + ((-0.25) \cdot (-0.7375)) + ((-0.55) \cdot (-2.1375)) \right. \\ \left. + 0.55 \cdot 1.8625 \right) = 0.439375 \end{aligned} \quad (11)$$

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

$$\begin{aligned} cov &= \begin{pmatrix} xx & xy \\ xy & yy \end{pmatrix} = \\ cov &= \begin{pmatrix} 0.2675 & 0.439375 \\ 0.439375 & 1.439843750 \end{pmatrix} \end{aligned}$$

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:

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

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

$$\begin{pmatrix} 0.2675 - \lambda & 0.439375 \\ 0.439375 & 1.439843750 - \lambda \end{pmatrix} = 0$$

First I get the equation by calculating the determinant

$$\det(\text{cov}) \rightarrow ((0.2675 - \lambda) \cdot (1.439843750 - \lambda)) - (0.439375 \cdot 0.439375)$$

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

$$((0.2675 - \lambda) \cdot (1.439843750 - \lambda)) - (0.439375 \cdot 0.439375) = 0$$

$$\text{eigenvalue}_1 = 0.1211093468 \quad \text{eigenvalue}_2 = 1.586234403$$

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 this equation.

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

I'm going to do the calculations in Maple since I'm dealing with odd numbers. I also going to reduce some of the decimals by rounding off just for readabilitys sake.

$$\text{eigenvector}_1 = [-0.9487, -0.3161] \quad \text{eigenvector}_2 = [0.3161, -0.9487]$$

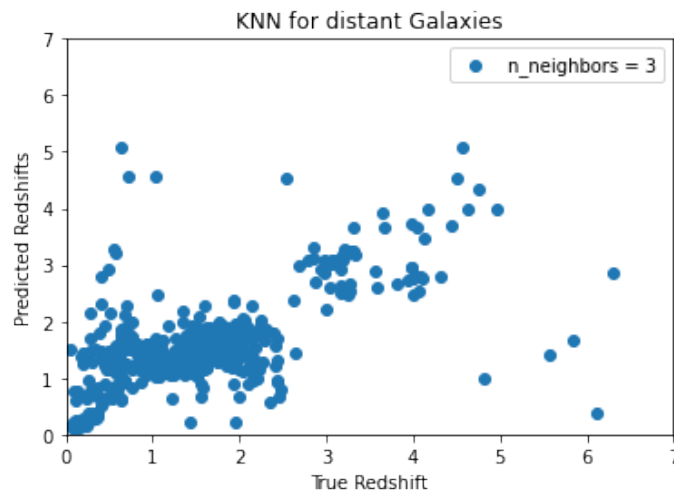
note that signs of the eigenvector can be different. e.g. eigenvector_1 could be $[0.9487, 0.3161]$ []

4 Problem 4 - Regression

4.1 a)

Using the euclidean distance I get an RSME value of

$$RMSE = 0.82430$$



Figur 2: Plot of kNN with $k = 3$

4.2 b)

Using the other distance calculation as given in the exam set I get the following RS-ME value

$$RMSE = 1.100519$$

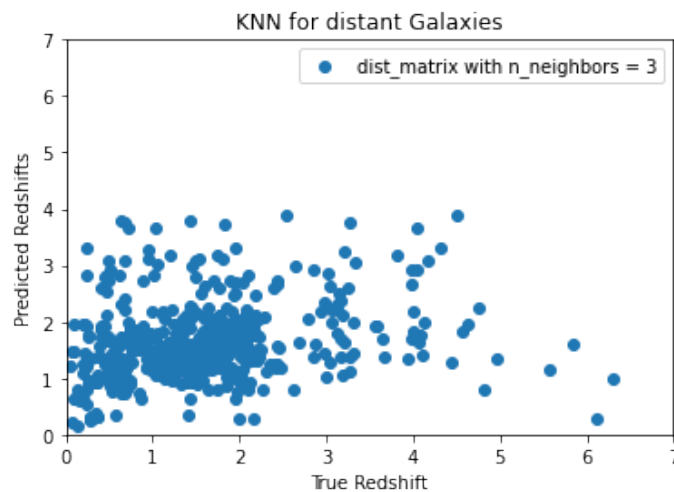


Figure 3: Plot of kNN with $k = 3$ and distance matrix

My assumption is, that 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

5.1 ...

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 *Sikit-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('D:\\Uddannelse\\Datalogi\\KU\\2_aar\\MAD\\MAD2021\\2022_Exam\\exam_data\\data
18
19 # Loads the validation data
20 X_validation, t_validation =
   ↳ loaddata('D:\\Uddannelse\\Datalogi\\KU\\2_aar\\MAD\\MAD2021\\2022_Exam\\exam_data\\data
```

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 as well. 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.

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: 0.48
Average probability assigned to correct classes: 0.36
Accuracy score: 0.66
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: 0.77
Average probability assigned to correct classes: 0.56
{'criterion': 'gini', 'max_features': 'sqrt', 'max_tree_depth': 10}
Accuracy score: 0.77
Average probability assigned to correct classes: 0.56
{'criterion': 'gini', 'max_features': 'sqrt', 'max_tree_depth': 15}
Accuracy score: 0.48
Average probability assigned to correct classes: 0.36
Accuracy score: 0.69
Average probability assigned to correct classes: 0.48
Accuracy score: 0.71
Average probability assigned to correct classes: 0.51
Accuracy score: 0.82
Average probability assigned to correct classes: 0.56
{'criterion': 'gini', 'max_features': 'log2', 'max_tree_depth': 10}
Accuracy score: 0.79
Average probability assigned to correct classes: 0.57
{'criterion': 'gini', 'max_features': 'log2', 'max_tree_depth': 15}
Accuracy score: 0.52
Average probability assigned to correct classes: 0.37
Accuracy score: 0.73
Average probability assigned to correct classes: 0.51
Accuracy score: 0.74
Average probability assigned to correct classes: 0.54
Accuracy score: 0.81
Average probability assigned to correct classes: 0.57
{'criterion': 'entropy', 'max_features': 'sqrt', 'max_tree_depth': 10}
Accuracy score: 0.81
Average probability assigned to correct classes: 0.58
{'criterion': 'entropy', 'max_features': 'sqrt', 'max_tree_depth': 15}
Accuracy score: 0.51
Average probability assigned to correct classes: 0.38
Accuracy score: 0.75
Average probability assigned to correct classes: 0.51
Accuracy score: 0.78
Average probability assigned to correct classes: 0.54
Accuracy score: 0.81
Average probability assigned to correct classes: 0.58
Accuracy score: 0.78
Average probability assigned to correct classes: 0.57

Side 8 of 9

Figur 4: 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  
→  numpy-array 'x'  
5  
6      Returns the normalized 'x'  
7      '''  
8      return (x - np.mean(x, axis=0)) / np.std(x, axis=0)
```

7.2 b)

8 Appendix A

8.1 ...

```
1 def minted_for_code_here
```

9 Appendix B

9.1 ...

```
1 def minted_for_code_here
```

10 Appendix Q6

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.metrics import accuracy_score
4 from sklearn.model_selection import ParameterGrid
5 from sklearn.ensemble import RandomForestClassifier
6
7 # Modified the loaddata function to make the training features
8   ↳ into
9 # 2d np.array, more specifically the 't value'
10 def loaddata(filename):
11     """Load the balloon data set from filename and return t, X
12         t - N-dim. vector of target (temperature) values
13         X - N-dim. vector containing the inputs (lift) x for
14         ↳ each data point
15     """
16     # Load data set from CSV file
17     Xt = np.loadtxt(filename, delimiter=',')
18
19     # Split into data matrix and target vector
20     X = Xt[:,0] # Matrix
21     t = Xt[:,1:] # Vector
22
23     return t, X
24
25 '''
26 QUESTION 6 - A
27 '''
28
29 # Loads the training data
30 X_train, t_train =
31   ↳ loaddata('D:\\Uddannelse\\Datalogi\\KU\\2_aar\\MAD\\MAD2021\\2022_Exam\\exam_data\\data
32
33 # Loads the validation data
34 X_validation, t_validation =
35   ↳ loaddata('D:\\Uddannelse\\Datalogi\\KU\\2_aar\\MAD\\MAD2021\\2022_Exam\\exam_data\\data
36
37 print("Shape of training targets: %s" %str(t_train.shape))
38 print("Shape of training features: %s" %str(X_train.shape))
39 print("Shape of validation targets: %s"
40   ↳ %str(t_validation.shape))
41 print("Shape of validation features: %s"
42   ↳ %str(X_validation.shape))
43 print()
```

```
38
39 '''
40 QUESTION 6 - B
41 '''
42 # Adding the parameters to test given in the exam question.
43 # Using these to find the optimal set of random forest
44   ↳ classifier parameters.
45 random_f_parameters = {
46     'criterion'      : ['gini', 'entropy'],
47     'max_tree_depth' : [2,5,6,10,15],
48     'max_features'   : ['sqrt', 'log2']
49 }
50
51 # Empty array for the result metrics
52 res = np.empty((0,3))
53
54 # Looping through the chosen(given) parameters and setting up
55   ↳ the
56 # Random forest classifier each time.
57 for params in list(ParameterGrid(random_f_parameters)):
58     clf = RandomForestClassifier(
59         criterion      = params['criterion'],
60         max_depth      = params['max_tree_depth'],
61         max_features    = params['max_features'])
62
63     # Training using the created classifier with the given
64     ↳ parameters.
65     clf.fit(X_train, t_train)
66
67     # number of correctly classified validation samples
68     t_prediction = clf.predict(X_validation)
69     acc_score = accuracy_score(t_validation, t_prediction)
70
71     # probability associated with classification
72     t_probability = clf.predict_proba(X_validation)
73     probability_score = np.mean([t_probability[int(t_val)]
74                                   for (t_probability, t_val)
75                                   in zip(t_probability, t_validation)])
76
77     print("Accuracy score: %.2f"
78           %acc_score)
79     print("Average probability assigned to correct classes:
80           ↳ %.2f"
81           %probability_score)
```

```
79     # print the parameters if new ones are more optimal
80     # than previously tried ones.
81     if len(res) > 0 and (acc_score > res[-1,1]
82                          or (acc_score == res[-1,1]
83                              and probability_score >
84                               res[-1,2])):
85
86         print(params)
87
88     # accumulate results
89     res = np.append(res, np.array([[params, acc_score,
90                                     probability_score]]), axis=0)
91     # sort the results in ascending order
92     res = res[np.lexsort((res[:,1], res[:,2]))]
```

```
91 for x in res:
92     print(x[0])
```


11 Appendix X

11.1 Q6_c

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

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

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