

TTK4205 Mønstergjennskekning
Prosjektoppgave 1

Anders Johannessen

Høst 2017

Introduksjon

Denen rapporten er skrevet som et svar på Prosjektoppgave 1 i faget Unik4590 / 9590 / TTK4205 Mønsterkjennning ved UiO / NTNU. Oppgaven går ut på å implementere tre forskjellige klassifisatorer (nærmeste nabo, minimum feilrate og minste kvadraters metode) og trene/teste disse klassifisatorene på tre ulike datasett. Nærmeste nabo klassifisatoren skal, ut ifra feilrate, brukes til å finne den beste egenskapskombinasjon for hver dimensjon som videre brukes i minimum feilrate og minste kvadraters metode klassifisatorene. De tre klassifisatorene blir så evaluert basert på feilrate for hver av de beste egenskapskombinasjonene funnet av nærmeste nabo klassifisatoren.

Gjennomførelse

All programmvare er skrevet i programmeringsspråket Python. Kode som er brukt er listet i appendix: Python scripts. Først ble dataene fra tekst filene importert i programmvaren slik at det kunne bearbeides av klassifisatorene. Dette ble gjort av funksjonen `readData` i `main.py` (Listing: 1). Funksjonen sorterer dataene etter linjenummer hvor oddetallslinjene legges i en trenings-data matrise og partallslinjene legges i en test-data matrise. Som i data filene så forteller første linje i trenings- og test-data matrisene hvilken klasse objektet tilhører. I de påfølgende kolonnene blir egenskapsverdiene til objektet lagt inn. I test-data matrisen er det lagt inn en ekstra kolonne som klassifisatorene bruker til å lagre klassen den har kommet frem til at objektet tilhører.

Etter at dataene har blitt lagret ble funksjonen `constructDimentionMatrix` kalt. Denne funksjonen tar inn en antall dimensjoner for data settet ($n_{kolonner} - 1$ i data filene) og returnerer en matrise med alle mulige kombinasjoner av egenskaper for en gitt dimensjon ($dimensjon = 1, \dots, n_{dimensjoner}$). For alle mulige egenskapkombinasjonene for alle dimensjoner kalles nærmeste nabo klassifisatoren.

2.1 Nærmeste nabo klassifikator

Nærmeste nabo klassifisatoren (Listing: 2) går gjennom alle egenskapskombinasjonene og beregner avstanden mellom et objekt i test-data matrisen med alle objektene trenings-data matrisen:

$$d = ||x_{test} * dimensjon - x_{tren} * dimensjon||$$

Den minste avstanden lagres og det siste elementet i test-data listen oppdateres til samme klasse som trenings-data objektet (det første kolonnen i trenings-data matrisen).

Etter at alle objektene i test-data matrisen er sammenlignet med alle objektene i trenings-data matrisen går funksjonen gjennom alle objektene i test-data matrisen og sammenligner egentlig klasse med tilegnet klasse.

$$feilrate = \frac{n_{feil}}{n_{objekter}}$$

Test-data matrisen, den beste kombinasjonen for hver dimensjon og tilhørende feilrate retuneres for videre bruk i programmet.

Den beste kombinasjonen for hver av dimensjonene brukes videre i minimum feilrate klassifisatoren og minste kvadraters metode klassifisatoren.

2.2 Minimum feilrate klassifisator

Minimum feilrate klassifisatoren (Listing: 3) begynner med å rekonstruere test-data og trenings-data matrisene slik at de kun inneholder de aktuelle egenskapene definert av de beste dimensjonene som ble funnet i nærmeste nabo klassifisatoren. Deretter går programmet gjennom alle nødvendige kalkulasjoner ($\hat{\mu}_i$, $\hat{\Sigma}_i$, W_i , w_i og w_{i0}) for å kunne regne ut diskriminantfunksjonene for hver av objektene i test-data matrisen:

$$g_i(\mathbf{x}) = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}^T \mathbf{x} + w_{i0} \quad i = 1, 2$$

Dersom $g_1(\mathbf{x}) \geq g_2(\mathbf{x})$ tilegnes objektet til klasse 1, og dersom $g_1(\mathbf{x}) < g_2(\mathbf{x})$ tilegnes objektet til klasse 2. Feilraten kalkuleres på samme måte som i nærmeste nabo klassifisatoren. Funksjonen returnerer test-data matrisen.

2.3 Minste kvadraters metode

Minste kvadraters metode klassifisatoren (Listing: 4) starter med å rekonstruere test-data og trenings-data matrisene på samme måte som i minimum feilrate klassifisatoren. Så ved hjelp av hjelpefunksjonene `calculateY` og `calculateB` kalkuleres $\mathbf{Y} = [\mathbf{y}_1^T \dots \mathbf{y}_n^T]^T$ hvor $\mathbf{y} = [1 \ x_1 \dots \ x_d]^T$ og \mathbf{b} som brukes i

$$\mathbf{a} = (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{b}$$

`leastSquareMethod` funksjonen finner nå diskriminantfunksjonen

$$g(\mathbf{x}) = \mathbf{a}^T \mathbf{y}$$

og dersom $g(\mathbf{x}) \geq 0$ tilegnes objektet klassen 1, og dersom $g(\mathbf{x}) < 0$ tilegnes objektet klassen 2. Feilraten kalkuleres på samme måte som i nærmeste nabo og minimum feilrate. Test-data matrisen retuneres ut av funksjonen.

Resultatet fra alle klassifisatorene for den beste todimensjonale egenskapskombinasjonen plottes ved hjelp av hjelpefunksjonen `plotResults`.

Resultat

Tabell 1-4 viser resultatet av nærmeste nabo klassifisatoren for dimensjonene 1, ..., 4 til hver av datasettene. Vi kan se at for de beste éndimensjonale egenskapskombinasjonene er:

- $[1, 0, 0, 0]$ med feilrate 0.24 for datasett 1
- $[1, 0, 0]$ med feilrate 0.18 for datasett 2
- $[0, 1, 0, 0]$ med feilrate 0.31 for data set 3.

Table 1: Resultat ved kjøring av nærmeste nabo klassifisatoren for 1-dimensjonale kombinasjoner

Datasett 1		Datasett 2		Datasett 3	
Dimensjon	Feilrate	Dimensjon	Feilrate	Dimensjon	Feilrate
$[1, 0, 0, 0]$	0.24	$[1, 0, 0]$	0.18	$[1, 0, 0, 0]$	0.33
$[0, 1, 0, 0]$	0.36			$[0, 1, 0, 0]$	0.31
$[0, 0, 1, 0]$	0.43	$[0, 1, 0]$	0.28	$[0, 0, 1, 0]$	0.34
$[0, 0, 0, 1]$	0.39	$[0, 0, 1]$	0.49	$[0, 0, 0, 1]$	0.40

De beste todimensjonale egenskapskombinasjonene er:

- $[1, 0, 0, 1]$ med feilrate 0.18 for datasett 1
- $[1, 1, 0]$ med feilrate 0.01 for datasett 2
- $[0, 1, 1, 0]$ med feilrate 0.10 for datasett 3

Table 2: Resultat ved kjøring av nærmeste nabo klassifisatoren for 2-dimensjonale kombinasjoner

Datasett 1		Datasett 2		Datasett 3	
Dimensjon	Feilrate	Dimensjon	Feilrate	Dimensjon	Feilrate
$[1, 1, 0, 0]$	0.18	$[1, 1, 0]$	0.01	$[1, 1, 0, 0]$	0.21
$[1, 0, 1, 0]$	0.19			$[1, 0, 1, 0]$	0.17
$[1, 0, 0, 1]$	0.17	$[1, 0, 1]$	0.19	$[1, 0, 0, 1]$	0.28
$[0, 1, 1, 0]$	0.32			$[0, 1, 1, 0]$	0.10
$[0, 1, 0, 1]$	0.23	$[1, 0, 1]$	0.29	$[0, 1, 0, 1]$	0.24
$[0, 0, 1, 1]$	0.30			$[0, 0, 1, 1]$	0.19

De beste tredimensjonale egenskapskombinasjonene er:

- $[1, 1, 0, 1]$ med feilrate 0.10 for datasett 1
- $[1, 1, 1]$ med feilrate 0.02 for datasett 2
- $[0, 1, 1, 1]$ med feilrate 0.07 for datasett 3

For firedimensjonale egenskapskombinasjoner er det kun datasett 1 og 3 som er med. Dette er fordi datasett 2 har kun tre egenskaper per objekt. Det er også kun én mulig egenskapskombinasjon for hver av datasettene:

Table 3: Resultat ved kjøring av nærmeste nabo klassifisatoren for 3-dimensjonale kombinasjoner

Datasett 1		Datasett 2		Datasett 3	
Dimensjon	Feilrate	Dimensjon	Feilrate	Dimensjon	Feilrate
[1, 1, 1, 0]	0.15	[1, 1, 1]	0.02	[1, 1, 1, 0]	0.10
[1, 1, 0, 1]	0.10			[1, 1, 0, 1]	0.20
[1, 0, 1, 1]	0.13			[1, 0, 1, 1]	0.15
[0, 1, 1, 1]	0.21			[0, 1, 1, 1]	0.07

Table 4: Resultat ved kjøring av nærmeste nabo klassifisatoren for 4-dimensjonale kombinasjoner

Datasett 1		Datasett 2		Datasett 3	
Dimensjon	Feilrate	Dimensjon	Feilrate	Dimensjon	Feilrate
[1, 1, 1, 1]	0.09			[1, 1, 1, 1]	0.10

- [1, 1, 1, 1] med feilrate 0.09 for datasett 1
- [1, 1, 1, 1] med feilrate 0.10 for datasett 3

Tabell 5 og Tabell 6 viser resultatet av kjøring av henholdsvis minimum feilrate klassifisatoren og minste kvadraters metode klassifisatoren for de beste egen-skapskombinasjonene fra nærmeste nabo klassifisatoren for hver av datasettene.

Table 5: Resultat ved kjøring av minimum feilrate klassifisatoren for beste dimensjon kombinasjoner

Datasett 1		Datasett 2		Datasett 3	
Dimensjon	Feilrate	Dimensjon	Feilrate	Dimensjon	Feilrate
[1, 0, 0, 0]	0.19	[1, 0, 0]	0.11	[0, 1, 0, 0]	0.23
[1, 0, 0, 1]	0.11	[1, 1, 0]	0.02	[0, 1, 1, 0]	0.20
[1, 1, 0, 1]	0.10			[0, 1, 1, 1]	0.13
[1, 1, 1, 1]	0.08	[1, 1, 1]	0.02	[1, 1, 1, 1]	0.07

Table 6: Resultat ved kjøring av minste kvadraters metode klassifisatoren for beste dimensjon kombinasjoner

Datasett 1		Datasett 2		Datasett 3	
Dimensjon	Feilrate	Dimensjon	Feilrate	Dimensjon	Feilrate
[1, 0, 0, 0]	0.15	[1, 0, 0]	0.11	[0, 1, 0, 0]	0.33
[1, 0, 0, 1]	0.11	[1, 1, 0]	0.13	[0, 1, 1, 0]	0.23
[1, 1, 0, 1]	0.11			[0, 1, 1, 1]	0.14
[1, 1, 1, 1]	0.09	[1, 1, 1]	0.13	[1, 1, 1, 1]	0.10

Den beste klassifikatoren for hver av egenskapskombinasjonen var for datasett 1:

- [1, 0, 0, 0] med feilrate 0.15 fra minste kvadraters metode
- [1, 0, 0, 1] med feilrate 0.11 fra minste kvadraters metode
- [1, 1, 0, 1] med feilrate 0.10 fra både nærmeste nabo og minimum feilrate
- [1, 1, 1, 1] med feilrate 0.08 fra minimum feilrate

For datasett 2:

- [1, 0, 0] med feilrate 0.11 fra minimum feilrate
- [1, 1, 0] med feilrate 0.01 fra nærmeste nabo
- [1, 1, 1] med feilrate 0.02 fra både nærmeste nabo og minimum feilrate

For datasett 3:

- [0, 1, 0, 0] med feilrate 0.23 fra minimum feilrate
- [0, 1, 1, 0] med feilrate 0.10 fra nærmeste nabo
- [0, 1, 1, 1] med feilrate 0.07 nærmeste nabo
- [1, 1, 1, 1] med feilrate 0.07 fra minimum feilrate

For hvert datasett ble den beste todimensjonale egenskapskombinasjonen plottet. Figur 1-3 viser henholdsvis resultatet av nærmeste nabo, minimum feilrate og minste kvadrater metode fra datasett 1, figur 4-6 fra datasett 2 og figur 7-9 fra datasett 3.

Avsluttende spørsmål

1. Den asymptotiske feilraten til nærmeste nabo regelen er gitt av

$$P^* \leq P \leq P^* \left(2 - \frac{c}{c-1} P^*\right)$$

hvor P^* er den optimale feilraten og c er antall klasser. I denne situasjonen er $c = 2$, noe som gir oss

$$P^* \leq P \leq P^* (2 - 2P^*)$$

Dette er da en relativ lav asymptotisk feilrate, noe som gjør nærmeste nabo klassifikatoren til en fornuftig klassifikator til å finne gunstige egenskapskombinasjoner.

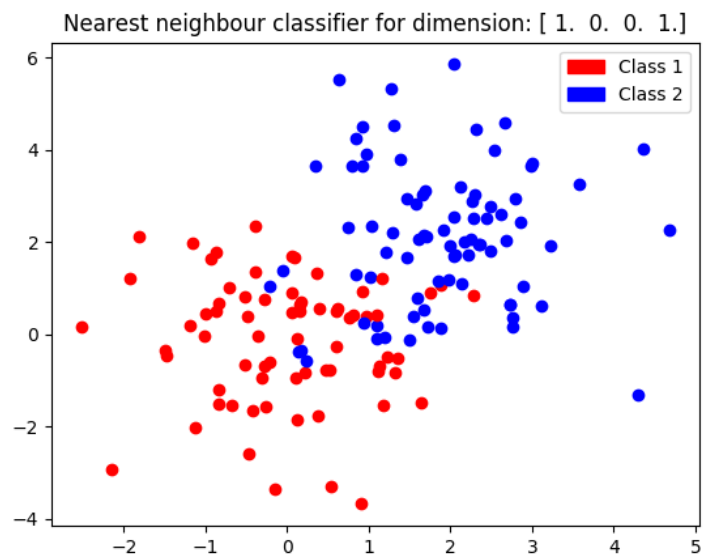


Figure 1: Nærmeste nabo klassifisering for dimensjonen: $[1, 0, 0, 1]$

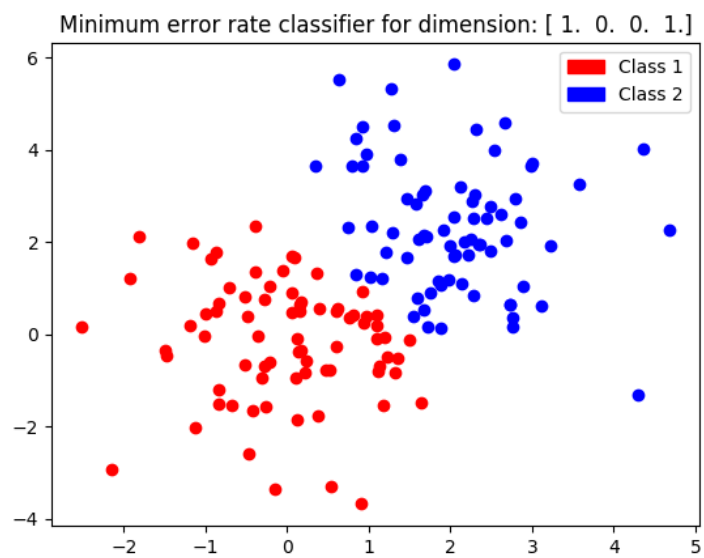


Figure 2: Minimum feilrate klassifisering for dimensjonen: $[1, 0, 0, 1]$

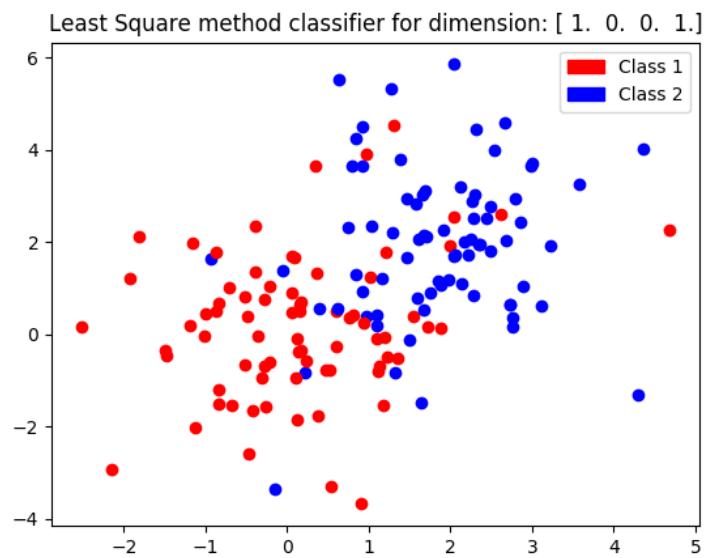


Figure 3: Minste kvadraters metode klassifisering for dimensjonen: $[1, 0, 0, 1]$

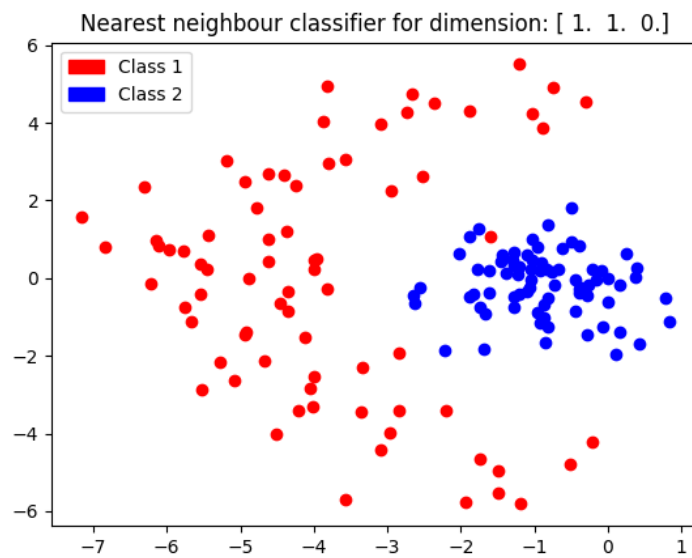


Figure 4: Nærmeste nabo klassifisering for dimensjonen: $[1, 1, 0]$

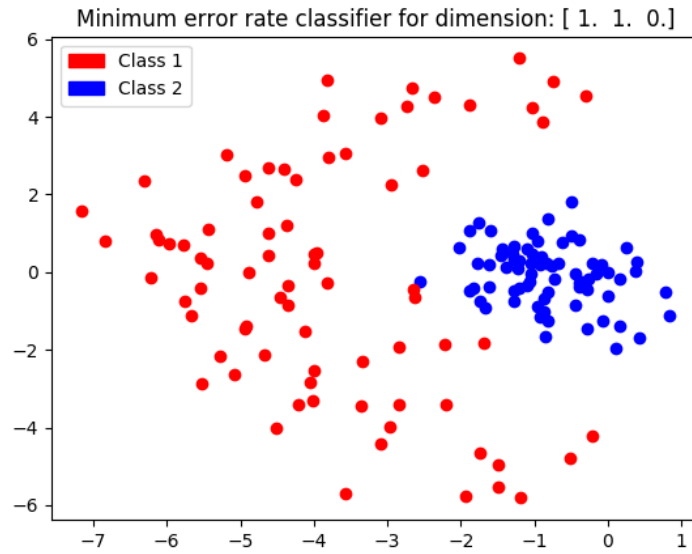


Figure 5: Minimum feilrate klassifisering for dimensjonen: $[1, 1, 0]$

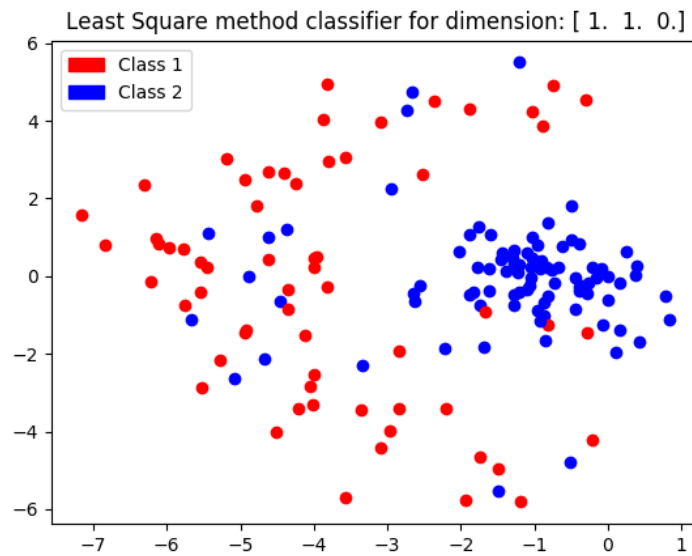


Figure 6: Minste kvadraters metode klassifisering for dimensjonen: $[1, 1, 0]$

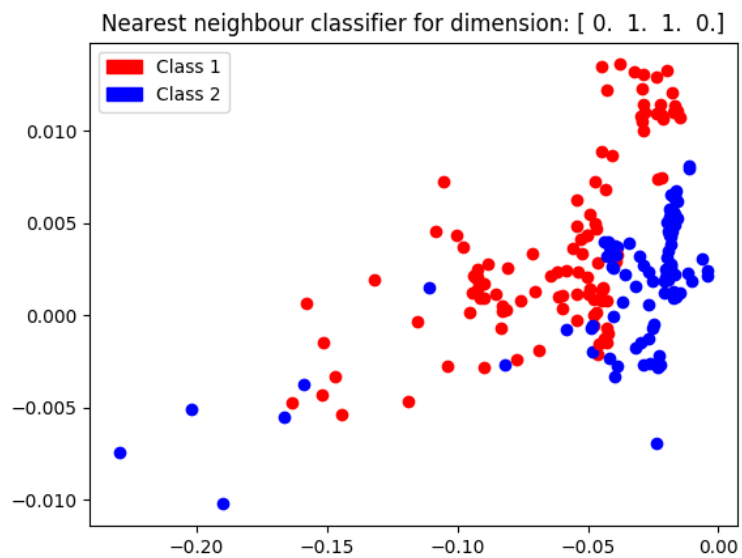


Figure 7: N rmeste nabo klassifisering for dimensjonen: [0, 1, 1, 0]

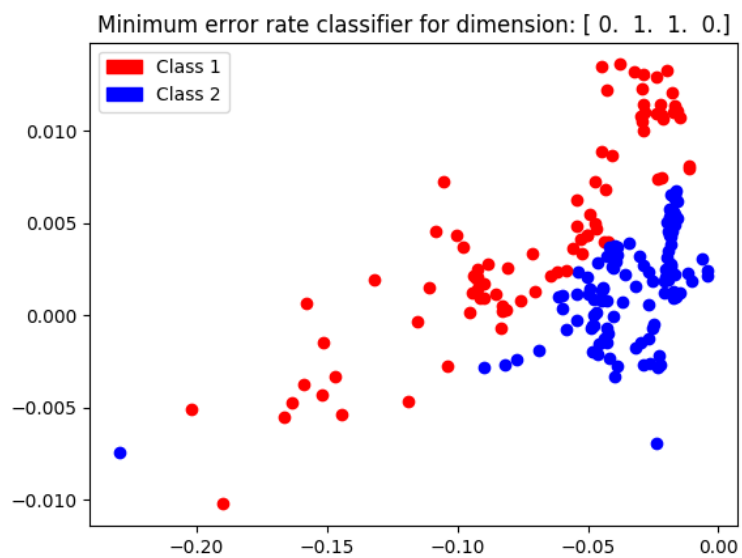


Figure 8: Minimum feilrate klassifisering for dimensjonen: [0, 1, 1, 0]

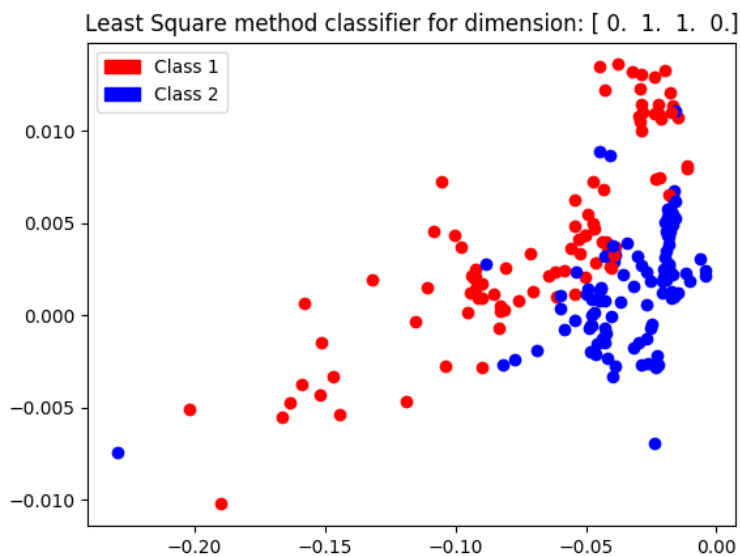


Figure 9: Minste kvadraters metode klassifisering for dimensjonen: $[0, 1, 1, 0]$

2. Nærmeste nabo klassifisatoren sammenligner alle test objektene med alle treningsobjektene. Når du i praktiske situasjoner ofte har større datasett krever denne klassifisatoren ekstremt mye regnekapasitet sammenlignet med f.eks minimum feilrate klassifisatoren som bruker en diskriminant-funksjon til å klassifisere et objekt. Et større datasett krever også større krav til lagringsplass.
3. Å bruke samme datasett til både trening og evaluering er en dårlig idé fordi da trenes klassifisatoren kun opp på en type data. Klassifisatoren vil da prestere bra på den test med samme type data, men presterer dårligere på data fra et annet set. Du vil altså få en lavere feilrate enn den faktiske feilraten ved evaluering på samme datasett.
4. At minste kvadraters metode klassifisatoren presterer dårligere enn nærmeste nabo og minste feilrate klassifisatorene er tydelig fremstilt i figur 6 sammenlignet med figur 4 og figur 5. Grunnen til dette kan være at datasett 2 ikke er lineært separabelt. Dette kan vi se fra figurene da det ikke er mulig å trekke en rett linje som skiller de to klassene.

Python scripts

A.1 Hoved script

Listing 1: main.py

```
1 import numpy as np
2 import itertools
3 import matplotlib.pyplot as plt
4 import matplotlib.patches as mpatches
5 from nearestNeighbour import nearestNeighbour
6 from minimumErrorRate import minimumErrorRate
7 from leastSquareMethod import leastSquareMethod
8
9
10
11 def readData(file_name):
12     file = open(file_name, 'r')
13     lines = file.readlines()
14
15     N_ROWS = len(lines)//2
16     N_COLUMNS = len(lines[0].split())
17
18     training_set = np.zeros([N_ROWS, N_COLUMNS])
19     test_set = np.zeros([N_ROWS, N_COLUMNS + 1])
20
21     for i in range(N_ROWS):
22         for j in range(N_COLUMNS):
23             training_set[i][j] = float(lines[i*2].split()[j])
24             test_set[i][j] = float(lines[i*2 + 1].split()[j])
25
26     return N_ROWS, N_COLUMNS, training_set, test_set
27
28
29
30 # Returns a matrix with all possible dimension combinations
31 def constructDimentionMatrix(dim, N_DIM):
32     for i in range(N_DIM):
33         combinations =
34             ↪ list(itertools.combinations(list(range(N_DIM)), dim +
35             ↪ 1))
36
37         dimension_matrix = np.zeros([len(combinations), N_DIM])
38
39         row = 0
40         for indexes in combinations:
```

```

39         for col in indexes:
40             dimension_matrix[row][col] = 1
41         row += 1
42
43     return dimension_matrix
44
45
46
47 def plotResults(best_error_rate_dimension, training_set,
48 ↪ test_set, N_ROWS, N_COLUMNS, data_file):
49     data_file,_ = data_file.split(".")
50     name = "results/" + data_file
51
52     dimension = [best_error_rate_dimension[1]]
53
54     result_NN,_,_ = nearestNeighbour(training_set, test_set,
55 ↪ N_ROWS, dimension)
56
57     x_test = np.zeros([N_ROWS, np.count_nonzero(dimension[0] ==
58 ↪ 1) + 2])
59     x_test[:, 0] = result_NN[:, 0]
60     x_test[:, -1] = result_NN[:, -1]
61     index = 1
62     for i in range(len(dimension[0])):
63         if dimension[0][i] == 1:
64             for j in range(N_ROWS):
65                 x_test[j][index] = result_NN[j][i + 1]
66             index += 1
67
68     result_MER = minimumErrorRate(training_set, test_set, N_ROWS,
69 ↪ dimension)
70     result_LSM = leastSquareMethod(training_set, test_set,
71 ↪ N_ROWS, N_COLUMNS, dimension)
72
73     plt.figure(1)
74     for row in x_test:
75         if row[0] == 1:
76             plt.plot(row[1], row[2], 'ro')
77         else:
78             plt.plot(row[1], row[2], 'bo')
79     plt.legend(handles=[mpatches.Patch(color='red', label='Class
80 ↪ 1'), mpatches.Patch(color='blue', label='Class 2')])
81     t1 = "Original plot before classifying for dimension: "
82     t2 = ''.join(str(i) for i in dimension)

```

```

79     plt.title(t1 + t2)
80     plt.savefig(name + " Original plot.png")
81
82     plt.figure(2)
83     for row in x_test:
84         if row[-1] == 1:
85             plt.plot(row[1],row[2], 'ro')
86         else:
87             plt.plot(row[1],row[2], 'bo')
88     plt.legend(handles=[mpatches.Patch(color='red', label='Class
89 ↪ 1'), mpatches.Patch(color='blue', label='Class 2')])
90     t1 = "Nearest neighbour classifier for dimension: "
91     t2 = ''.join(str(i) for i in dimension)
92     plt.title(t1 + t2)
93     plt.savefig(name + " NN plot.png")
94
95     plt.figure(3)
96     for row in result_MER:
97         if row[-1] == 1:
98             plt.plot(row[1],row[2], 'ro')
99         else:
100             plt.plot(row[1],row[2], 'bo')
101     plt.legend(handles=[mpatches.Patch(color='red', label='Class
102 ↪ 1'), mpatches.Patch(color='blue', label='Class 2')])
103     t1 = "Minimum error rate classifier for dimension: "
104     t2 = ''.join(str(i) for i in dimension)
105     plt.title(t1 + t2)
106     plt.savefig(name + " MER plot.png")
107
108     plt.figure(4)
109     for row in result_LSM:
110         if row[-1] == 1:
111             plt.plot(row[1],row[2], 'ro')
112         else:
113             plt.plot(row[1],row[2], 'bo')
114     plt.legend(handles=[mpatches.Patch(color='red', label='Class
115 ↪ 1'), mpatches.Patch(color='blue', label='Class 2')])
116     t1 = "Least Square method classifier for dimension: "
117     t2 = ''.join(str(i) for i in dimension)
118     plt.title(t1 + t2)
119     plt.savefig(name + " LSM plot.png")
120
121     plt.show()

```

```

122 if __name__ == "__main__":
123     data_file = "ds-1.txt"
124     #data_file = "ds-2.txt"
125     #data_file = "ds-3.txt"
126     N_ROWS, N_COLUMNS, training_set, test_set =
        ↪ readData(data_file)
127
128     best_error_rate_dimension = []
129     for current_dimension in range(N_COLUMNS - 1):
130         dimension_matrix =
            ↪ constructDimentionMatrix(current_dimension, N_COLUMNS
            ↪ - 1)
131         result_NN, min_error_rate, min_error_rate_dimension =
            ↪ nearestNeighbour(training_set, test_set, N_ROWS,
            ↪ dimension_matrix)
132
            ↪ best_error_rate_dimension.append(min_error_rate_dimension)
133         print("Best error rate: ", "%.2f" % min_error_rate, " for
            ↪ dimension: ", min_error_rate_dimension)
134
135         minimumErrorRate(training_set, test_set, N_ROWS,
            ↪ best_error_rate_dimension)
136         leastSquareMethod(training_set, test_set, N_ROWS, N_COLUMNS,
            ↪ best_error_rate_dimension)
137
138         plotResults(best_error_rate_dimension, training_set,
            ↪ test_set, N_ROWS, N_COLUMNS, data_file)

```

A.2 Nærmeste nabo klassifisator script

[h!t]

Listing 2: nearestNeighbour.py

```

1 import numpy as np
2
3
4
5 def nearestNeighbour(training_set, test_set, N_ROWS,
    ↪ dimension_matrix):
6     print("\n# ===== Running: Nearest neighbour classifier
    ↪ ===== #")
7     min_error_rate = float('inf')
8
9     for dimension in dimension_matrix:
10         error = 0

```

```

11
12     for test_row in range(N_ROWS):
13         min_dinstance = float('inf')
14
15         for training_row in range(N_ROWS):
16             distance =
17                 ↪ abs(np.linalg.norm(test_set[test_row][1:-1]*dimension
18                 ↪ - training_set[training_row][1:]*dimension))
19             if distance < min_dinstance:
20                 min_dinstance = distance
21                 test_set[test_row][-1] =
22                     ↪ training_set[training_row][0]
23
24         # Estimating error rate
25         if test_set[test_row][-1] != test_set[test_row][0]:
26             error += 1
27     error_rate = error/N_ROWS
28
29     if error_rate < min_error_rate:
30         min_error_rate = error_rate
31         min_error_rate_dimension = dimension
32
33     print("Minimum error rate: ", "%.2f" % error_rate, " for
34         ↪ dimension: ", dimension)
35
36     return test_set, min_error_rate, min_error_rate_dimension

```

A.3 Minimum feilrate klassifikator script

[h!t]

Listing 3: minimumErrorRate.py

```

1  import numpy as np
2
3
4
5  def minimumErrorRate(training_set, test_set, N_ROWS,
6  ↪ best_error_rate_dimension):
7      print("\n# ===== Running: Minimum error rate classifier
8      ↪ ===== #")
9
10     error_rate_storage = []
11
12     for dimension in best_error_rate_dimension:

```



```

11     # Restructuring data sets
12     x_training = np.zeros([N_ROWS, np.count_nonzero(dimension
13         ↪ == 1) + 1])
14     x_training[:, 0] = training_set[:, 0]
15     x_test = np.zeros([N_ROWS, np.count_nonzero(dimension ==
16         ↪ 1) + 2])
17     x_test[:, 0] = training_set[:, 0]
18     index = 1
19     for i in range(len(dimension)):
20         if dimension[i] == 1:
21             for j in range(N_ROWS):
22                 x_training[j][index] = training_set[j][i + 1]
23                 x_test[j][index] = test_set[j][i + 1]
24             index += 1
25     x_training = [x_training[x_training[:, 0] == 1],
26         ↪ x_training[x_training[:, 0] == 2]]
27
28     # Calculate mean
29     mu = [x_training[0][:, 1:].mean(axis=0), x_training[1][:,
30         ↪ 1:].mean(axis=0)]
31
32     # Calculate covariance
33     Sigma = [np.cov(x_training[0][:, 1:], rowvar=False),
34         ↪ np.cov(x_training[1][:, 1:], rowvar=False)]
35
36     #  $W_i = -1/2 * \Sigma_i^{-1}$       $i = 1, \dots, c$ 
37     W = [-np.linalg.inv(np.atleast_2d(Sigma[0]))/2,
38         ↪ -np.linalg.inv(np.atleast_2d(Sigma[1]))/2]
39
40     #  $w = \Sigma_i^{-1} * \mu_i$       $i = 1, \dots, c$ 
41     w = [np.matmul(np.linalg.inv(np.atleast_2d(Sigma[0])),
42         ↪ mu[0]),
43         ↪ np.matmul(np.linalg.inv(np.atleast_2d(Sigma[1])),
44         ↪ mu[1])]
45
46     #  $w_0 = -1/2 * \mu_i^T * \Sigma_i^{-1} * \mu_i - 1/2 * \ln(|\Sigma_i|) +$ 
47     ↪  $\ln(P(\omega_i))$       $i = 1, \dots, c$ 
48     w_0 = [-np.matmul(np.matmul(np.transpose(mu[0]),
49         ↪ np.linalg.inv(np.atleast_2d(Sigma[0]))), mu[0]) / 2
50         ↪ - np.log(np.linalg.det(np.atleast_2d(Sigma[0]))) /
51         ↪ 2 + np.log(0.5),
52         ↪ -np.matmul(np.matmul(np.transpose(mu[1]),
53         ↪ np.linalg.inv(np.atleast_2d(Sigma[1]))),
54         ↪ mu[1]) / 2

```

```

43         - np.log(np.linalg.det(np.atleast_2d(Sigma[1]))) /
44         ↪ 2 + np.log(0.5)]
45
46     for j in range(N_ROWS):
47         #  $g_i = x^T W_i x + w^T x + w_0_i \quad i = 1, \dots, c$ 
48         g_1 =
49         ↪ np.matmul(np.matmul(np.transpose(x_test[j][1:-1]),
50         ↪ W[0]), x_test[j][1:-1]) + \
51         ↪ np.matmul(np.transpose(w[0]), x_test[j][1:-1])
52         ↪ + w_0[0]
53         g_2 =
54         ↪ np.matmul(np.matmul(np.transpose(x_test[j][1:-1]),
55         ↪ W[1]), x_test[j][1:-1]) + \
56         ↪ np.matmul(np.transpose(w[1]), x_test[j][1:-1])
57         ↪ + w_0[1]
58
59         if (g_1 - g_2) >= 0:
60             x_test[j][-1] = 1
61         else:
62             x_test[j][-1] = 2
63
64     # Esimating error rate
65     error = 0
66     for i in range(N_ROWS):
67         if x_test[i][-1] != x_test[i][0]:
68             error += 1
69     error_rate = error / N_ROWS
70     error_rate_storage.append(error_rate)
71
72     print("Error rate: ", "%.2f" % error_rate, " for
73     ↪ dimension: ", dimension)
74
75     return x_test

```

A.4 Minste kvadraters metode script

[h!t]

Listing 4: leastSquareMethod.py

```

1  import numpy as np
2
3
4
5  def calculateY(x_training, N_ROWS):

```

```

6     Y = np.zeros([N_ROWS, len(x_training[0])])
7     for i in range(len(Y)):
8         Y[i][0] = 1
9         Y[i][1:] = x_training[i][1:]
10    return Y
11
12
13
14    def calculateB(x_training, N_ROWS):
15        b = np.zeros(N_ROWS)
16        for i in range(len(x_training)):
17            if x_training[i][0] == 1:
18                b[i] = 1
19            else:
20                b[i] = -1
21        return b
22
23
24
25    def leastSquareMethod(training_set, test_set, N_ROWS, N_COLUMNS,
26        ↪ best_error_rate_dimension):
27        print("\n# ===== Running: Least square method classifier
28        ↪ ===== # ")
29
30        error_rate_storage = []
31
32        for dimension in best_error_rate_dimension:
33            # Restructuring data sets
34            x_training = np.zeros([N_ROWS, np.count_nonzero(dimension
35                ↪ == 1) + 1])
36            x_training[:, 0] = training_set[:, 0]
37            x_test = np.zeros([N_ROWS, np.count_nonzero(dimension ==
38                ↪ 1) + 2])
39            x_test[:, 0] = training_set[:, 0]
40            index = 1
41            for i in range(len(dimension)):
42                if dimension[i] == 1:
43                    for j in range(N_ROWS):
44                        x_training[j][index] = training_set[j][i + 1]
45                        x_test[j][index] = test_set[j][i + 1]
46                    index += 1
47
48            # Create b
49            b = calculateB(x_training, N_ROWS)

```

```

48     # Create Y
49     Y = calculateY(x_training, N_ROWS)
50
51     # a = (Y^T*Y)^-1*Y*b^T
52     a =
53         ↪ np.matmul(np.matmul(np.linalg.inv(np.matmul(np.transpose(Y),
54         ↪ Y)), np.transpose(Y)), b)
55
56     for i in range(N_ROWS):
57         # g = a^T*y
58         g = np.matmul(np.transpose(a), Y[i])
59
60         if g >= 0:
61             x_test[i][-1] = 1
62         else:
63             x_test[i][-1] = 2
64
65     # Esimating error rate
66     error = 0
67     for i in range(N_ROWS):
68         if x_test[i][-1] != x_test[i][0]:
69             error += 1
70     error_rate = error/N_ROWS
71     error_rate_storage.append(error_rate)
72
73     print("Error rate: ", "%.2f" % error_rate, " for
74     ↪ dimension: ", dimension)
75
76     return x_test

```