

TEK5020 Prosjektoppgave 1 – Rapport

Fremgangsmåte

Som beskrevet i oppgaven, har vi brukt nærmeste-nabo klassifikatoren til å estimere feilraten for alle mulige kombinasjoner av egenskaper for et gitt antall dimensjoner, for hvert mulige antall av dimensjoner. For hvert antall av dimensjoner finner vi da den beste kombinasjonen av mulige egenskaper (den med minst estimert feilrate), og bruker de tre klassifikatorene på denne kombinasjonen av egenskaper for å finne den beste klassifikatoren for den egenskapskombinasjonen.

Resultater

Følgende feilrater ble oppnådd ved å utføre nærmeste nabo klassifikasjon på treningsdatasettet for ulike kombinasjoner av egenskaper. De beste (laveste) feilratene per datasett, per antall dimensjoner, er markert i grønt:

Egenskaper	Datasett 1	Datasett 2	Datasett 3
1	21.3%	18.0%	28.5%
2	36.0%	35.3%	37.0%
3	42.6%	48.7%	42.0%
4	28.7%	-	43.5%
(1, 2)	16.7%	1.3%	25.5%
(1, 3)	18.7%	18.6%	21.5%
(1, 4)	11.3%	-	26.0%
(2, 3)	39.3%	30.0%	15.5%
(3, 4)	22.0%	-	23.0%
(1, 2, 3)	12.7%	2.0%	13.0%
(1, 3, 4)	18.7%	-	13.5%
(2, 3, 4)	27.3%	-	11.5%
(1, 2, 3, 4)	13.3%	-	12.5%

Ved å trene de tre klassifikatorene på de beste egenskapene i treningssettet per antall dimensjoner, og deretter predikere på testsettet, fikk vi følgende feilrater. Beste feilrate per datasett, per antall dimensjoner, er markert i grønt.

Én dimensjon

Datasett	Nærmeste nabo	Lineærdiskriminant	Minste feilrate Gauss
1	28 . 7%	22 . 70%	15 . 3%
2	48 . 7%	51 . 3%	12 . 0%
3	43 . 5%	50 . 0%	35 . 0%

To dimensjoner

Datasett	Nærmeste nabo	Lineærdiskriminant	Minste feilrate Gauss
1	22 . 0%	23 . 3%	10 . 7%
2	30 . 0%	36 . 0%	2 . 7%
3	23 . 0%	50 . 00%	19 . 5%

Tre dimensjoner

Datasett	Nærmeste nabo	Lineærdiskriminant	Minste feilrate Gauss
1	27 . 3%	31 . 3%	14 . 7%
2	2 . 0%	11 . 3%	3 . 3%
3	11 . 5%	13 . 00%	11 . 5%

Fire dimensjoner

Datasett	Nærmeste nabo	Lineærdiskriminant	Minste feilrate Gauss
1	13 . 3%	10 . 0%	10 . 7%
2	-	-	-
3	12 . 5%	10 . 50%	6 . 5%

Avsluttende spørsmål

Hvorfor er det fornuftig å benytte nærmeste-nabo klassifikatoren til å finne gunstige egenskapskombinasjoner?

Nærmeste-nabo klassifikatoren gir er, som nevnt i oppgaveteksten, enkel å både forstå intuitivt og å programmere. Den gir et godt og lettforståelig mål på hvilke egenskapskombinasjoner som inneholder mest variabilitet.

Hvorfor kan det i en praktisk anvendelse være fornuftig å finne en lineær eller kvadratisk klassifikator til erstatning for nærmeste-nabo klassifikatoren?

I praksis er nærmeste-nabo klassifikatoren regnemessig kostbar, og tar mye tid å kalkulere. I tillegg er nærmeste-nabo klassifikatoren utsatt for overfitting; den er altså svært sensitiv for støy i datasettet. Lineære og kvadratiske klassifikatorer kan både lede til færre utregninger, og potensielt bidra til å redusere overfitting og gi en bedre balanse mellom bias og varians.

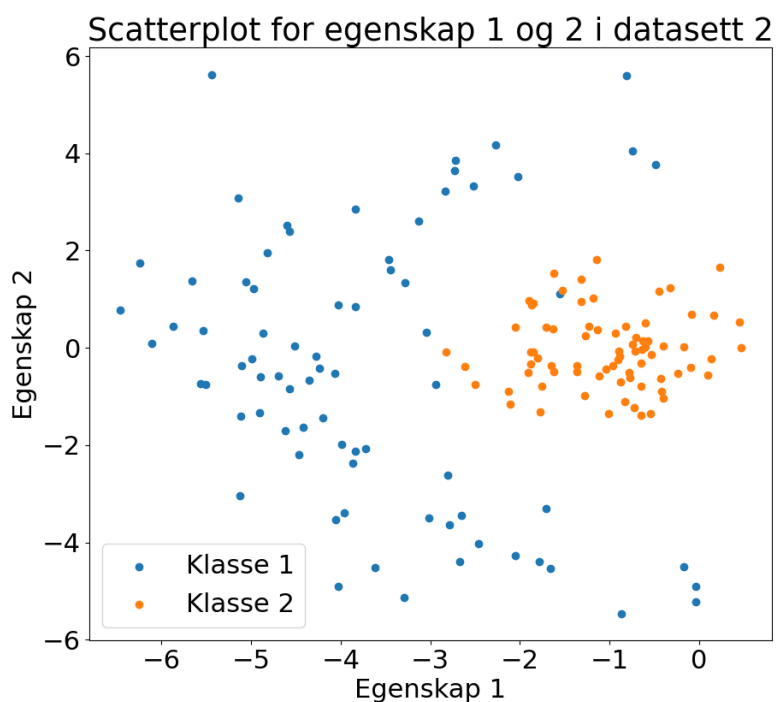
Hvorfor er det lite gunstig å bruke samme datasettet både til trening og evaluering av en klassifikator?

Å bruke samme datasett til trening og evaluering av en klassifikator gjør at når man evaluerer presisjonen av klassifikatoren på testsettet (som nå også er treningssettet), er ikke denne presisjonen representativ for hvilken presisjon klassifikatoren hadde fått på vilkårlige data samlet fra samme distribusjon som datasettet.

Dette er fordi klassifikatoren har sett, og trent på, dataen den blir evaluert på. Når vi evaluerer en klassifikator, ønsker vi å måle hvor presis den er på usette data fra samme distribusjon som treningsdataene. Klassifikatoren er laget for å optimalisere desisjongsgrensen for dataene den har sett, og når den trener på testsettet, blir den unaturlig god på å klassifisere akkurat testsettet, i forhold til hvordan den hadde vært på usette data.

Hvorfor gir en lineær klassifikator dårlige resultater for datasett 2?

Ved å plote egenskapskombinasjonen av datasett 2 som ga best resultat for nærmeste-nabo klassifikatoren, ser vi at datasettet i disse dimensjonene ikke er lineært separabelt (se figuren nedenfor). Det kreves altså en ikke-lineær desisjongsgrense for å optimalisere grensen, som lineære klassifikatorer ikke er i stand til å produsere.



Kode

1 utils.py

```
1 import numpy as np
2 import os
3 import cv2 as cv
4 import copy
5
6
7 def read_dataset(idx):
8     project_dir = os.path.dirname(os.path.dirname(__file__))
9     file_path = project_dir + f"/data/ds-{idx}.txt"
10    data_array = np.loadtxt(file_path)
11    targets, obs = data_array[:, 0].copy(), data_array[:, 1:].copy()
12    return targets, obs
13
14
15 def split_data(obs, targets):
16     train_obs, train_targets = obs[1::2], targets[1::2]
17     test_obs, test_targets = obs[0::2], targets[0::2]
18     return train_obs, test_obs, train_targets, test_targets
19
20
21 def least_params(train_obs, train_targets):
22     bias = np.ones((len(train_obs), 1))
23     ext_train_obs = np.concatenate((bias, train_obs), axis=1)
24
25     b = np.where(train_targets == 1, 1, -1)
26
27     params = np.linalg.inv(ext_train_obs.T @ ext_train_obs) @
28     ext_train_obs.T @ b
29     return params
30
31 def least_discriminant(params):
32     def discriminant(test_obs):
33         bias = np.ones((len(test_obs), 1))
34         ext_test_obs = np.concatenate((bias, test_obs), axis=1)
35         return np.where(ext_test_obs @ params > 0, 1, 2)
36
37     return discriminant
38
39
40 def create_dataset(pixels):
```

```

41 dataset = []
42 for i in range(len(pixels)):
43     pixels[i] = pixels[i].reshape(-1, 3)
44     pixels[i] = np.concatenate(
45         (np.ones((pixels[i].shape[0], 1)) * (i + 1), pixels[i])
46         , axis=1
47         )
48     dataset.extend(pixels[i])
49     return np.array(dataset)
50
51 def normalize_dataset(pixels):
52     r = pixels[:, 1]
53     g = pixels[:, 2]
54     b = pixels[:, 3]
55
56     t1 = r / np.sum(pixels[:, 1:], axis=1)
57     t2 = g / np.sum(pixels[:, 1:], axis=1)
58     t3 = b / np.sum(pixels[:, 1:], axis=1)
59
60     pixels[:, 1] = t1
61     pixels[:, 2] = t2
62     pixels[:, 3] = t3
63
64     return pixels
65
66
67 def estimate_pixels_apriori(pixels):
68     probs = []
69     for i in range(np.int64(np.max(pixels[:, 0], axis=0))):
70         prob = np.sum(pixels[:, 0] == (i + 1)) / pixels.shape[0]
71         probs.append(prob)
72     return np.array(probs)
73
74
75 def estimate_pixels_mean(pixels):
76     means = []
77     for i in range(np.int64(np.max(pixels[:, 0], axis=0))):
78         est_mean = pixels[pixels[:, 0] == (i + 1)].mean(axis=0)
79         means.append(est_mean)
80     return np.array(means)
81
82
83 def estimate_pixels_cov(pixels, pixel_class_means):
84     covs = []
85     for i in range(np.int64(np.max(pixels[:, 0], axis=0))):
86         N_class = np.sum(pixels[:, 0] == (i + 1))
87         class_dev = pixels[pixels[:, 0] == (i + 1), 1:] -
            pixel_class_means[i, 1:]
88         class_cov = (class_dev.T @ class_dev) / (N_class - 1)
89         covs.append(class_cov)
90
91     return np.array(covs)
92
93
94 def pixels_discriminants(pixel_means, pixel_covs, pixel_apriori):
95     discriminants = []

```

```

96     for i in range(len(pixel_means)):
97         discriminants.append(
98             class_discriminant(pixel_means[i], pixel_covs[i],
99                                 pixel_apriori[i], pixels=True)
100         )
101     return discriminants
102
103 def segment_image(image_path, discriminants):
104     img = cv.imread(image_path)
105     if img.shape[:2] > (800,800):
106         img = cv.resize(img, (600,600))
107     seg_img = np.zeros_like(img)
108     colors = np.array(
109         [
110             [255, 0, 0],
111             [0, 255, 0],
112             [0, 0, 255],
113             [255, 255, 0],
114             [255, 0, 255],
115             [0, 255, 255],
116             [128, 0, 128],
117             [255, 165, 0],
118             [0, 128, 0],
119             [128, 128, 128],
120         ]
121     )
122     for x in range(img.shape[0]):
123         for y in range(img.shape[1]):
124             c1 = np.argmax([disc(img[x, y]) for disc in
125                             discriminants])
126             seg_img[x,y] = colors[c1]
127     return seg_img
128
129 def measure_dist(obs_1, obs_2):
130     distance = np.linalg.norm(obs_1 - obs_2)
131     return distance
132
133
134 def nearest_neighbour(train_obs, train_targets, test_obs):
135     c_test_obs = np.zeros((len(test_obs), 1))
136
137     for i in range(len(test_obs)):
138         near_neigh = np.argmin(
139             [
140                 measure_dist(test_obs[i], train_obs[j])
141                 for j in range(len(train_obs))
142                 if i != j
143             ]
144         )
145         c_test_obs[i] = train_targets[near_neigh]
146
147     return c_test_obs.flatten()
148
149
150 def estimate_a_priori(train_targets):

```

```

151     class_one = np.sum(train_targets == 1)
152     prob_one = class_one / train_targets.shape[0]
153     prob_two = 1.0 - prob_one
154     return prob_one, prob_two
155
156
157 def estimate_class_mean(train_obs, train_targets):
158     class_one_mean = train_obs[train_targets == 1].mean(axis=0)
159     class_two_mean = train_obs[train_targets == 2].mean(axis=0)
160     return class_one_mean, class_two_mean
161
162
163 def estimate_class_cov(class_one_mean, class_two_mean, train_obs,
164     train_targets):
165     N_one = train_obs[train_targets == 1].shape[0]
166     N_two = train_obs.shape[0] - N_one
167     class_one_dev = train_obs[train_targets == 1] - class_one_mean
168     class_two_dev = train_obs[train_targets == 2] - class_two_mean
169     cov_one = (class_one_dev.T @ class_one_dev) / (N_one - 1)
170     cov_two = (class_two_dev.T @ class_two_dev) / (N_two - 1)
171
172     return cov_one, cov_two
173
174 def class_discriminant(class_mean, class_cov, a_priori_prob, pixels
175     =False):
176     W = -(1 / 2) * np.linalg.inv(class_cov)
177
178     w = np.linalg.inv(class_cov) @ class_mean
179
180     det_cov = np.log(np.linalg.det(class_cov))
181     det_cov = det_cov if det_cov > 1e-5 else 0
182
183     w_0 = (
184         -(1 / 2) * class_mean @ np.linalg.inv(class_cov) @
185         class_mean
186         - (1 / 2) * det_cov
187         + np.log(a_priori_prob)
188     )
189     if pixels:
190         return lambda test_obs: test_obs.T @ W @ test_obs +
191             test_obs @ w + w_0
192     else:
193         return lambda test_obs: np.sum(test_obs @ W * test_obs,
194             axis=1) + test_obs @ w + w_0
195
196
197 def minimum_error(train_obs, train_targets):
198     class_one_mean, class_two_mean = estimate_class_mean(train_obs,
199         train_targets)
200     cov_one, cov_two = estimate_class_cov(
201         class_one_mean, class_two_mean, train_obs, train_targets
202     )
203
204     a_priori_one, a_priori_two = estimate_a_priori(train_targets)
205
206     discriminant_one = _class_discriminant(class_one_mean, cov_one,

```

```

    a_priori_one)
discriminant_two = _class_discriminant(class_two_mean, cov_two,
    a_priori_two)
203
204     return gen_discriminant(discriminant_one, discriminant_two)
205
206
207 def gen_discriminant(c1_discr, c2_discr):
208     return lambda test_obs: np.where(c1_discr(test_obs) - c2_discr(
        test_obs) > 0, 1, 2)

```

2 oblig1.py

```

1 from snutils import *
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 dim_combinations_list = [
6     [(0,), (1,), (2,), (3,)],
7     [(0, 1), (0, 2), (0, 3), (1, 2), (2, 3)],
8     [(0, 1, 2), (0, 2, 3), (1, 2, 3)],
9     [(0, 1, 2, 3)],
10 ]
11
12 dim_combinations_for_dataset_2_list = [
13     [(0,), (1,), (2,)],
14     [(0, 1), (0, 2), (1, 2)],
15     [(0, 1, 2)],
16 ]
17
18 for dataset_idx in (1, 2, 3):
19     print(f"===== Dataset {dataset_idx} =====")
20     targets, obs = read_dataset(dataset_idx)
21     train_obs, test_obs, train_targets, test_targets = split_data(
        obs, targets)
22
23     for dim_combinations in (
24         dim_combinations_list
25         if dataset_idx != 2
26         else dim_combinations_for_dataset_2_list
27     ):
28         print(
29             f"===== Now testing for dimension {len(
30                 dim_combinations[0])} ====="
31         )
32         best_fail_rate = 1
33         best_dim = None
34         for dimensions in dim_combinations:
35             preds = nearest_neighbour(
36                 train_obs[:, dimensions], train_targets, train_obs
37                [:, dimensions]
38             )
39             fail_rate = (
40                 np.sum(np.where(preds != train_targets, 1, 0)) /

```



```

41         # print(f"{dimensions=} {fail_rate=}")
42
43         if fail_rate < best_fail_rate:
44             best_dim = dimensions
45             best_fail_rate = fail_rate
46     if dataset_idx == 2 and len(best_dim) == 2:
47         plt.scatter(
48             test_obs[:, best_dim[0]][test_targets == 1],
49             test_obs[:, best_dim[1]][test_targets == 1],
50         )
51         plt.scatter(
52             test_obs[:, best_dim[0]][test_targets == 2],
53             test_obs[:, best_dim[1]][test_targets == 2],
54         )
55         plt.show()
56
57     print(f"Lowest fail rate was {best_fail_rate:.3f}, for
58     features: {best_dim}")
59
60     # ----- here starts the actual grog way -----
61
62     print("\n\nNow testing all methods on test set:")
63     # Nearest neighbour
64     preds_test_nn = nearest_neighbour(
65         train_obs[:, best_dim], train_targets, test_obs[:,
66         best_dim]
67     )
68     fail_rate_test_nn = (
69         np.sum(np.where(preds != test_targets, 1, 0)) /
70         test_targets.shape[0]
71     )
72
73     # Linear discriminant
74     linear_discriminant = least_discriminant(
75         least_params(train_obs[:, best_dim], train_targets)
76     )
77
78     preds = linear_discriminant(test_obs[:, dimensions])
79     fail_rate_test_lindisc = (
80         np.sum(np.where(preds != test_targets, 1, 0)) /
81         test_targets.shape[0]
82     )
83
84     # Minimum error
85     minimum_error_discriminant = minimum_error(
86         train_obs[:, best_dim], train_targets
87     )
88
89     preds = minimum_error_discriminant(test_obs[:, best_dim])
90     fail_rate_test_minerror = (
91         np.sum(np.where(preds != test_targets, 1, 0)) /
92         test_targets.shape[0]
93     )
94
95     print(
96         f"Fail rates: NN: {fail_rate_test_nn:.3f} LINDISC: {
97         fail_rate_test_lindisc:.3f} MINERROR: {fail_rate_test_minerror
98         :.3f}"
99     )

```