

Homework 3

Shao-Ting Chiu (UIN:433002162)

10/26/22

Table of contents

Homework Description	1
Computational Enviromnent Setup	2
Third-party libraries	2
Version	2
Problem 5.1	3
Problem 5.2	4
(a)	4
(b)	4
Problem 5.6	4
(a)	4
(b)	5
Problem 5.10 (Python Assignment)	6
(a)	6
(b)	8
(c)	11
Appendix	12
Revised c05_kernel.py	12
References	14

Homework Description

- Course: ECEN649, Fall2022

Problems from the book:

5.1

5.2

5.6 (a,b)

5.10 (a,b,c)

Challenge (not graded):

5.4

5.6 (c,d)

- Deadline: Oct. 26th, 11:59 pm

Computational Environment Setup

Third-party libraries

```
1  %matplotlib inline
2  import sys
3  import matplotlib
4  import numpy as np
5  import scipy as sp
6  import pandas as pd
7  import sklearn as sk
8  import scipy.stats as st
9  import matplotlib.pyplot as plt
10 from scipy.stats import multivariate_normal as mvn
11 from sklearn.neighbors import KernelDensity as KD
12 from matplotlib.colors import ListedColormap
13 # Fix random state for reproducibility
14 np.random.seed(1978081)
15 # Matplotlib setting
16 plt.rcParams['text.usetex'] = True
17 matplotlib.rcParams['figure.dpi'] = 300
```

Version

```
1  print(sys.version)
2  print(matplotlib.__version__)
3  print(sp.__version__)
4  print(np.__version__)
5  print(pd.__version__)
```

```
6 print(sk.__version__)
```

3.8.14 (default, Sep 6 2022, 23:26:50)

[Clang 13.1.6 (clang-1316.0.21.2.5)]

3.3.1

1.5.2

1.19.1

1.1.1

1.1.2

Problem 5.1

Consider that an experimenter wants to use A 2-D cubic histogram classification rule, with square cells with side length h_n , and achieve consistency as the sample size n increases, for any possible distribution of the data. If the experimenter lets h_n decrease as $h_n = \frac{1}{\sqrt{n}}$, would they be guaranteed to achieve consistency and why? If not, how would they need to modify the rate of decrease of h_n to achieve consistency?

Use Braga-Neto (2020, Theorem 5.6).

Test of consistency

- $d = 2$
- $V_n = h_n^2 = \frac{1}{n}$
- $h_n \rightarrow 0, V_n \rightarrow 0$
- $nV_n = 1$ is not approaching to infinity as $n \rightarrow \infty$
- Thus, the consistency is not guaranteed.

Modification

- Let $h_n = \frac{1}{n^a}$
- $V_n = \frac{1}{n^{2p}}$, let $p > 0$
- $nV_n = \frac{1}{n^{2p-1}}$, $\lim_{n \rightarrow \infty} nV_n \rightarrow \infty$. $2p - 1 < 0$
 - $0 < p < \frac{1}{2}$
- The universal consistence of the cubic histogram rule is guaranteed.

Problem 5.2

Consider that an experimenter wants to use the kNN classification rule and achieve consistency as the sample size n increases. In each of the following alternatives, answer whether the experimenter is successful and why.

(a)

The experimenter does not know the distribution of (X, Y) and lets k increase as $k = \sqrt{n}$.

Use Braga-Neto (2020, Theorem 5.7)

- $k = \sqrt{n}$
- $\lim_{n \rightarrow \infty} k = \infty$
- $\lim_{n \rightarrow \infty} \frac{k}{n} = \lim_{n \rightarrow \infty} \frac{1}{\sqrt{n}} = 0$
- The kNN rule is universally consistent.

(b)

The experimenter does not know the distribution but knows that $\epsilon^* = 0$ and keeps k fixed, $k = 3$.

Because k is fixed and independent of n , the approach is not universally consistent. However, since $\epsilon^* = 0$, this approach is consistent.

Problem 5.6

Assume that the feature X in a classification problem is a real number in the interval $[0, 1]$. Assume that the classes are equally likely, with $p(x|Y = 0) = 2xI_{\{0 \leq x \leq 1\}}$ and $p(x|Y = 1) = 2(1 - x)I_{\{0 \leq x \leq 1\}}$.

(a)

Find the Bayes error ϵ^* .

Because the two classes are equally likely, $p(Y = 0) = p(Y = 1) = 0.5$.

$$\epsilon^* = E[\min(\eta(x), 1 - \eta(x))]$$

$$\eta(x) = p(Y = 1|x) \tag{1}$$

$$= \frac{p(x|Y = 1)p(Y = 1)}{p(x)} \tag{2}$$

$$= \frac{p(x|Y = 1)p(Y = 1)}{p(x|Y = 0)p(Y = 0) + p(x|Y = 1)p(Y = 1)} \tag{3}$$

$$= \frac{2(1-x) \cdot 0.5}{2x \cdot 0.5 + 2(1-x) \cdot 0.5} \tag{4}$$

$$= \frac{2(1-x)}{2x + 2 - 2x} \tag{5}$$

$$= 1 - x \tag{6}$$

$$1 - \eta(x) = x \tag{7}$$

$$p(x) = p(x|Y = 0)p(Y = 0) + p(x|Y = 1)p(Y = 1) = 2x \cdot 0.5 + 2(1-x) \cdot 0.5 = 1$$

$$\epsilon^* = E[\min(\eta(x), 1 - \eta(x))] \tag{8}$$

$$= E[\min(1 - x, x)] \tag{9}$$

$$= \int_0^1 \min(\eta(x), 1 - \eta(x))p(x)dx \tag{10}$$

$$= \int_0^1 \min(\eta(x), 1 - \eta(x))dx \tag{11}$$

$$= \int_0^{\frac{1}{2}} xdx + \int_{\frac{1}{2}}^1 (1-x)dx \tag{12}$$

$$= 0.25 \tag{13}$$

(b)

Find the asymptotic error rate ϵ_{NN} for the NN classification rule.

Use Cover-Hart Theorem (Braga-Neto 2020, Theorem 5.1).

$$\epsilon_{NN} = \lim_{n \rightarrow \infty} E[\epsilon_n] = E[2\eta(X)(1 - \eta(X))]$$

Use the result from [Problem 5.6\(a\)](#).

$$\eta(x) = 1 - x \quad (14)$$

$$1 - \eta(x) = x \quad (15)$$

$$\epsilon_{NN} = \lim_{n \rightarrow \infty} E[\epsilon_n] = E[2\eta(X)(1 - \eta(X))] \quad (16)$$

$$= E[2(1 - x)x] \quad (17)$$

$$= 2E[x - x^2] \quad (18)$$

$$= 2 \left(\int_0^1 xp(x)dx - \int_0^1 x^2p(x)dx \right) \quad (19)$$

$$= 2 \left(\int_0^1 xdx - \int_0^1 x^2dx \right) \quad (20)$$

$$= 2 \left(\left(\frac{1}{2}x^2 \right)_1^1 - \left(\frac{1}{3}x^3 \right)_0^1 \right) \quad (21)$$

$$= 2 \left(\frac{1}{2} - \frac{1}{3} \right) \quad (22)$$

$$= 2 \left(\frac{3-2}{6} \right) \quad (23)$$

$$= \frac{1}{3} \quad (24)$$

Problem 5.10 (Python Assignment)

(a)

Modify the code in `c05_kernel.py` (modified in [appendix](#)) to obtain plots for $h = 0.1, 0.3, 0.5, 1, 2, 5$ ¹ and $n = 50, 100, 250, 500$ per class. Plot the classifiers over the range $[-3, 9] \times [-3, 9]$ in order to visualize the entire data and reduce the marker size from 12 to 8 to facilitate visualization. Which classifiers are closest to the optimal classifier? How do you explain this in terms of underfitting/overfitting? See the coding hint in part (a) of Problem 5.8.

Since the optimal boundary is a straight line between two centroids. Those classifiers close to optimal decision tend to have large sample size.

¹In Problem 5.10, please replace “k=1,3,5,7,9,11” by “h=0.1,0.3,0.5,1,2,5” — [Ulisses on Slack](#)

The bandwidth (h) can have influences on the underfitting/overfitting. Small h may get overfitting; on the other hand, large h may get underfitting.

```

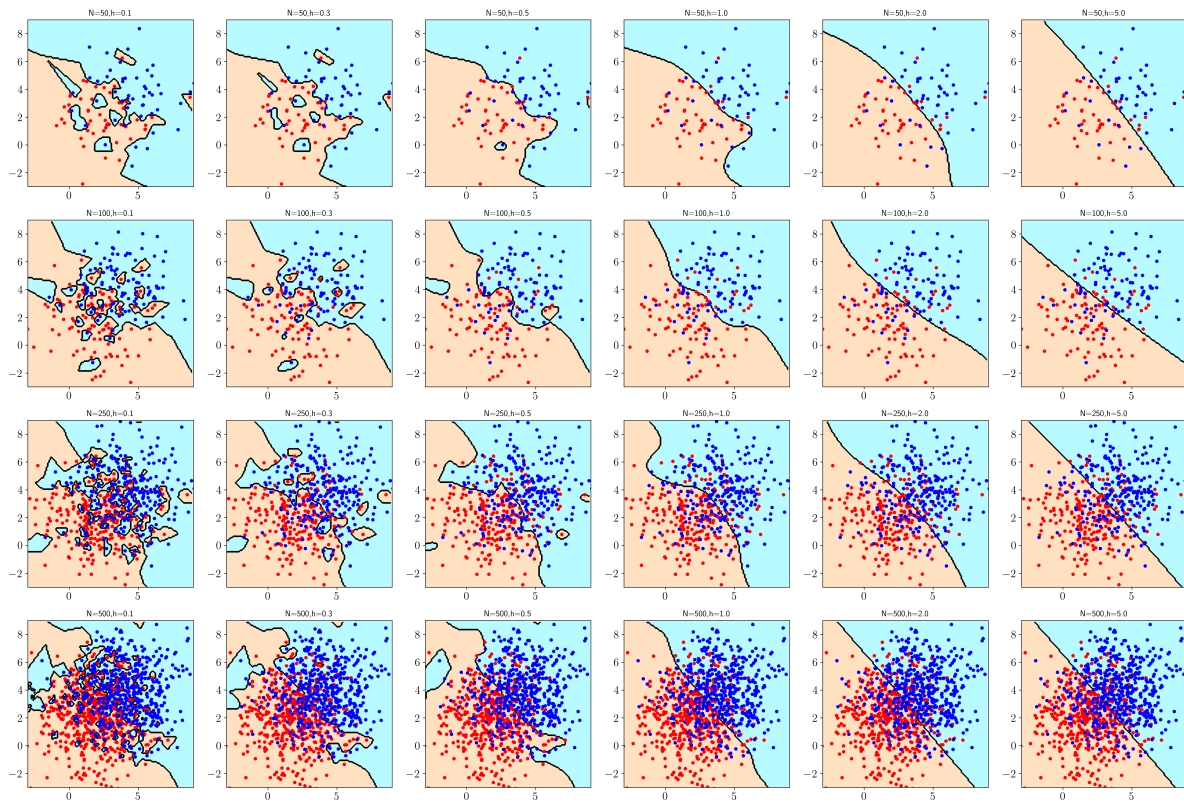
1  def plot_kd(ax, x0, y0, x1, y1, Z):
2      cmap_light = ListedColormap(["#FF0000", "#0000FF"])
3      plt.rc("xtick", labelsiz=16)
4      plt.rc("ytick", labelsiz=16)
5      ax.plot(x0, y0, ".r", markersiz=8) # class 0
6      ax.plot(x1, y1, ".b", markersiz=8) # class 1
7      ax.set_xlim([-3, 9])
8      ax.set_ylim([-3, 9])
9      ax.pcolormesh(xx, yy, Z, cmap=cmap_light, shading="nearest")
10     ax.contour(xx, yy, Z, colors="black", linewidths=0.5)
11     plt.close()
12     return ax
13
14
15     mm0 = np.array([2, 2])
16     mm1 = np.array([4, 4])
17     Sig0 = 4*np.identity(2)
18     Sig1 = 4*np.identity(2)
19     Ns = np.array([50, 100, 250, 500])
20     #Ns = [50]
21     hs = np.array([0.1, 0.3, 0.5, 1, 2, 5])
22     #hs = [0.1]
23     Xs = [[mvn.rvs(mm0, Sig0, n), mvn.rvs(mm1, Sig1, n)] for n in Ns]
24
25     clf0s = [[KD() for i in range(0, len(hs))] for j in range(0, len(Ns))]
26     clf1s = [[KD() for i in range(0, len(hs))] for j in range(0, len(Ns))]
27
28     # plotting
29     x_min, x_max = (-3, 9)
30     y_min, y_max = (-3, 9)
31     s = 0.1 #0.01 # mesh step size
32     xx, yy = np.meshgrid(np.arange(x_min, x_max, s), np.arange(y_min, y_max, s))
33     fig, axs = plt.subplots(len(Ns), len(hs), figsize=(30, 20), dpi=150)
34
35     for (i, X) in enumerate(Xs):
36         x0, y0 = np.split(X[0], 2, 1)
37         x1, y1 = np.split(X[1], 2, 1)
38         y = np.concatenate((np.zeros(Ns[i]), np.ones(Ns[i])))
39         for (j, h) in enumerate(hs):

```

```

40     clf0s[i][j] = KD(bandwidth=h)
41     clf0s[i][j].fit(X[0])
42     clf1s[i][j] = KD(bandwidth=h)
43     clf1s[i][j].fit(X[1])
44
45     Z0 = clf0s[i][j].score_samples(np.c_[xx.ravel(), yy.ravel()])
46     Z1 = clf1s[i][j].score_samples(np.c_[xx.ravel(), yy.ravel()])
47     Z = Z0 <= Z1
48     Z = Z.reshape(xx.shape)
49     plot_kd(axes[i][j], x0, y0, x1, y1, Z);
50     axes[i][j].set_title("N={},h={}".format(Ns[i],h))
51     plt.close()
52 fig.savefig("img/c05_kernel.png",bbox_inches="tight",facecolor="white");

```



(b)

Compute test set errors for each classifier in part (a), using the same procedure as in part (b) of Problem 5.8. Generate a table containing each classifier plot in

part (a) with its test set error rate. Which combinations of sample size and kernel bandwidth produce the top 5 smallest error rates?

Bayes error

$$\delta = \sqrt{(\mu_1 - \mu_0)^T \Sigma^{-1} (\mu_1 - \mu_0)} \quad (25)$$

$$\epsilon = \Phi\left(-\frac{\delta}{2}\right) \approx 0.23975$$

```

1 mu1 = np.matrix([[4],[4]])
2 mu0 = np.matrix([[2],[2]])
3 sig = np.matrix([[4,0],[0,4]])
4 delta = np.sqrt( (mu1-mu0).T @ np.linalg.inv(sig) @ (mu1-mu0))[0][0]
5 error = st.norm.cdf(-delta/2)
6
7 pd.DataFrame({"Bayes Error":[error]})

```

Bayes Error	
0	[[0.23975006109347669]]

Best Five

1. 0.269 (1,4) N=50, h=1
2. 0.271 (1,6) N=50, h=5
3. 0.273 (1,5) N=50, h=2
4. 0.274 (3,4) N=250, h=1
5. 0.276 (3,6) N=250, h=5

```

1 def measure_test_error(clf0, clf1, xxs, yys, ys):
2     Z0 = clf0.score_samples(np.c_[xxs, yys])
3     Z1 = clf1.score_samples(np.c_[xxs, yys])
4     Z = Z0<=Z1
5     return np.count_nonzero(Z != ys.astype(bool)) / len(ys)
6
7 nt = 500
8 X_test = [mvn.rvs(mm0, Sig0, nt), mvn.rvs(mm1, Sig1, nt)]
9 x0,y0 = np.split(X_test[0],2,1)
10 x1,y1 = np.split(X_test[1],2,1)
11 ys = np.concatenate((np.zeros(nt),np.ones(nt)))

```

```

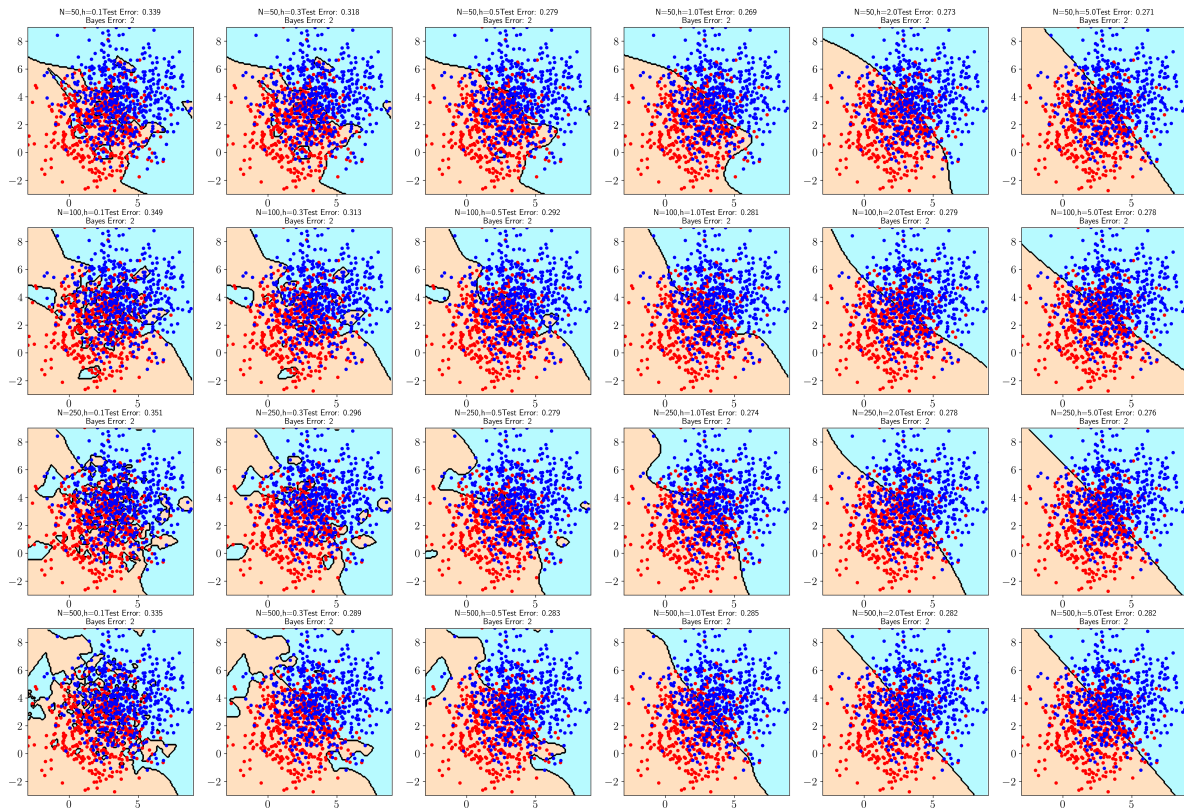
12  errs = np.zeros((len(Ns), len(hs)))
13
14  figt, axts = plt.subplots(len(Ns), len(hs), figsize=(30,20), dpi=150)
15
16  for i in range(0, len(Ns)):
17      xxs = np.concatenate((x0, x1))
18      yys = np.concatenate((y0, y1))
19      for (j, h) in enumerate(hs):
20          errs[i][j]= measure_test_error(clf0s[i][j], clf1s[i][j], xxs, yys, ys)
21
22          axts[i][j].set_title(axs[i][j].get_title() + "Test Error: {} \n Bayes Error: {}".fo
23          plt.close()
24
25          Z0 = clf0s[i][j].score_samples(np.c_[xx.ravel(), yy.ravel()])
26          Z1 = clf1s[i][j].score_samples(np.c_[xx.ravel(), yy.ravel()])
27          Z = Z0<=Z1
28          Z = Z.reshape(xx.shape)
29          plot_kd(axts[i][j], x0, y0, x1, y1, Z);
30
31  print(errs)
32
33  figt.savefig("img/c05_kernel_test.png",bbox_inches="tight",facecolor="white");

```

```

[[0.339 0.318 0.279 0.269 0.273 0.271]
 [0.349 0.313 0.292 0.281 0.279 0.278]
 [0.351 0.296 0.279 0.274 0.278 0.276]
 [0.335 0.289 0.283 0.285 0.282 0.282]]

```



(c)

Compute expected error rates for the Gaussian kernel classification rule in part (a), using the same procedure as in part (c) of Problem 5.8. Since error computation is faster here, a larger value $R = 200$ can be used, for better estimation of the expected error rates. Which kernel bandwidth should be used for each sample size?

```

1  errs = np.zeros((len(Ns), len(hs)))
2  R = 200
3  for jj in range(0, R):
4      nt = 500
5      X_test = [mvn.rvs(mm0, Sig0, nt), mvn.rvs(mm1, Sig1, nt)]
6      x0, y0 = np.split(X_test[0], 2, 1)
7      x1, y1 = np.split(X_test[1], 2, 1)
8      ys = np.concatenate((np.zeros(nt), np.ones(nt)))
9
10     xxs = np.concatenate((x0, x1))

```

```

11     yys = np.concatenate((y0, y1))
12     for i in range(0, len(Ns)):
13         for (j, h) in enumerate(hs):
14             errs[i][j] += measure_test_error(clf0s[i][j], clf1s[i][j], xxs, yys, ys)
15
16     errs = errs/R
17     best_hi = np.argmin(errs, axis=1)
18     print(errs)
19     print(best_hi)
20
21     pd.DataFrame({"Sample Size": Ns, "Best H": hs[best_hi ]})

```

```

[[0.309245 0.29467 0.26515 0.250625 0.24554 0.24833 ]
 [0.325025 0.2838 0.255735 0.2419 0.240805 0.240665]
 [0.33395 0.276555 0.25415 0.245395 0.23972 0.23971 ]
 [0.31032 0.25389 0.246315 0.240215 0.240195 0.239245]]
[4 5 5 5]

```

	Sample Size	Best H
0	50	2.0
1	100	5.0
2	250	5.0
3	500	5.0

Appendix

Revised c05_kernel.py²

```

1  """
2  Foundations of Pattern Recognition and Machine Learning
3  Chapter 5 Figure 5.5
4  Author: Ulisses Braga-Neto
5  Plot kernel classifiers
6  """
7
8  import numpy as np
9  import matplotlib.pyplot as plt
10 from scipy.stats import multivariate_normal as mvn

```

²All, be careful with script `c05_kernel.py`. If you replace the variable `b` to `h` in the script, it will shock with the grid step size parameter, which is also called `h`, and you will get erroneous results. — [Ulisses on Slack](#)

```

10 from sklearn.neighbors import KernelDensity as KD
11 from matplotlib.colors import ListedColormap
12 # Fix random state for reproducibility
13 np.random.seed(1978081)
14 mm0 = np.array([2,2])
15 mm1= np.array([4,4])
16 Sig0 = 4*np.identity(2)
17 Sig1 = 4*np.identity(2)
18 N = 50 # number of points in each class
19 X0 = mvn.rvs(mm0,Sig0,N)
20 x0,y0 = np.split(X0,2,1)
21 X1 = mvn.rvs(mm1,Sig1,N)
22 x1,y1 = np.split(X1,2,1)
23 X = np.concatenate((X0,X1),axis=0)
24 y = np.concatenate((np.zeros(N),np.ones(N)))
25 cmap_light = ListedColormap(["#FFEOCO", "#B7FAFF"])
26 s = .01 # mesh step size
27 x_min,x_max = (-0.5,6.5)
28 y_min,y_max = (-0.5,6.5)
29 for h in [0.1,0.3,0.5,1]:
30     clf0 = KD(bandwidth=h)
31     clf0.fit(X0)
32     clf1 = KD(bandwidth=h)
33     clf1.fit(X1)
34     xx,yy = np.meshgrid(np.arange(x_min,x_max,s),np.arange(y_min,y_max,s))
35     Z0 = clf0.score_samples(np.c_[xx.ravel(), yy.ravel()])
36     Z1 = clf1.score_samples(np.c_[xx.ravel(), yy.ravel()])
37     Z = Z0<=Z1
38     Z = Z.reshape(xx.shape)
39     fig,ax=plt.subplots(figsize=(8,8),dpi=150)
40     plt.rc("xtick",labelsize=16)
41     plt.rc("ytick",labelsize=16)
42     plt.plot(x0,y0,".r",markersize=16) # class 0
43     plt.plot(x1,y1,".b",markersize=16) # class 1
44     plt.xlim([-0.18,6.18])
45     plt.ylim([-0.18,6.18])
46     plt.pcolormesh(xx,yy,Z,cmap=cmap_light)
47     ax.contour(xx,yy,Z,colors="black",linewidths=0.5)
48     plt.show()
49     fig.savefig("c05_kernel"+str(int(10*h))+".png",bbox_inches="tight",facecolor="white")

```

References

Braga-Neto, Ulisses. 2020. *Fundamentals of Pattern Recognition and Machine Learning*. Springer.