

Statistical Machine Learning - Exercise 2



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Dominik Marino - 2468378, Pertamina Kunz - 2380210
8. August 2020

Inhaltsverzeichnis

1 Bayesian Decision Theory	2
1.1 Optimal Boundary	2
1.2 Decision Boundaries	2
1.3 Different Misclassification Costs	3
2 Density Estimation	4
2.1 Gaussian Maximization Likelihood Estimate	4
2.2 Prior Probabilities	5
2.3 Biased ML Estimate	5
2.4 Class Density	6
2.5 Posterior	7
2.6 Bayesian Estimation	11
3 Non-parametric Density Estimation	14
3.1 Histogram	14
3.2 Kernel Density Estimate	15
3.3 K-Nearest Neighbours	16
3.4 Comparision of the Non-Parametric Methods	17
4 Expectation Maximization	19
4.1 Gaussian Mixture Update Rules	19
4.2 EM	19

1 Bayesian Decision Theory

Consider data generated by a mixture of two Gaussian distributions,

$$\begin{aligned}C_1 : X &\sim \mathcal{N}(\mu_1, \sigma_1^2) \\C_2 : X &\sim \mathcal{N}(\mu_2, \sigma_2^2).\end{aligned}$$

1.1 Optimal Boundary

Tutor solution:

Bayesian Decision Theory is a framework for taking optimal decisions under uncertain conditions (see Bishop ch. 1.5). The goal is to minimize the cost of classification.

Let R_1 and R_2 denote regions and let C_1 and C_2 denote the respective class to which a point x is assigned. In the case of two classes, the missclassification error is:

$$\begin{aligned}p(\text{error}) &= p(x \in R_1, C_2) + p(x \in R_2, C_1) \\&= \int_{R_1} p(x, C_2) + \int_{R_2} p(x, C_1) \\&= \int_{R_1} p(x | C_2)p(C_2) + \int_{R_2} p(x | C_1)p(C_1)\end{aligned}$$

Bayesian decision theory is about finding the boundary (threshold) to decide which distribution a set of observations belongs to, based on Bayes' theorem on a-posteriori probability;

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k) p(C_k)}{p(\mathbf{x})},$$

where $p(\mathbf{x} | C_k)$ is likelihood, $p(C_k)$ is the class prior, and $p(\mathbf{x})$ is the normalization term. The goal is to minimize the misclassification rate. For two classes C_1 and C_2 , at the optimal decision boundary,

$$p(C_1 | \mathbf{x}) = p(C_2 | \mathbf{x}).$$

We decide for class C_1 over C_2 when $p(C_1 | \mathbf{x}) > p(C_2 | \mathbf{x})$, that is if class C_1 has the highest a-posteriori probability.

1.2 Decision Boundaries

At the decision boundary with $p(C_1) = p(C_2)$, $\sigma_1^2 = \sigma_2^2 = \sigma^2$ (and same misclassification costs):

$$\begin{aligned}p(C_1 | \mathbf{x}) &= p(C_2 | \mathbf{x}) \\ \frac{p(\mathbf{x} | C_1) p(C_1)}{p(\mathbf{x})} &= \frac{p(\mathbf{x} | C_2) p(C_2)}{p(\mathbf{x})} \\ \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (x - \mu_1)^2\right) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (x - \mu_2)^2\right) \\ \exp\left(-\frac{1}{2\sigma^2} (x - \mu_1)^2 + \frac{1}{2\sigma^2} (x - \mu_2)^2\right) &= 1 \\ \Leftrightarrow -(x - \mu_1)^2 + (x - \mu_2)^2 &= 0 \\ 2x\mu_1 - \mu_1^2 - 2x\mu_2 + \mu_2^2 &= 0 \\ \Leftrightarrow 2x(\mu_1 - \mu_2) &= \mu_1^2 - \mu_2^2 \\ \Leftrightarrow x &= \frac{\mu_1^2 - \mu_2^2}{2(\mu_1 - \mu_2)} = \frac{1}{2}(\mu_1 + \mu_2).\end{aligned}$$

1.3 Different Misclassification Costs

Loss Matrix:

$$\lambda = \begin{pmatrix} \lambda_{1,1} & \lambda_{1,2} \\ \lambda_{2,1} & \lambda_{2,2} \end{pmatrix} = \begin{pmatrix} \lambda(\alpha_1 | C_1) & \lambda(\alpha_1 | C_2) \\ \lambda(\alpha_2 | C_1) & \lambda(\alpha_2 | C_2) \end{pmatrix} = \begin{pmatrix} C_1 C_1 & C_1 C_2 \\ C_2 C_1 & C_2 C_2 \end{pmatrix} = \begin{pmatrix} 0 & 4 \\ 1 & 0 \end{pmatrix} \quad (1)$$

with:

$$R(\alpha_i | x) = \sum_j \lambda(\alpha_i | C_j) p(C_j | x)$$

where also:

$$\lambda(\alpha_i | C_j) = \lambda_{i,j}$$

Decide α_1 , if $R(\alpha_2 | x) > R(\alpha_1 | x)$

$$\begin{aligned} \lambda_{2,1} p(C_1 | x) + \lambda_{2,2} p(C_2 | x) &> \lambda_{1,1} p(C_1 | x) + \lambda_{1,2} p(C_2 | x) \\ \lambda_{2,1} p(C_1 | x) + \cancel{\lambda_{2,2} p(C_2 | x)} &> \cancel{\lambda_{1,1} p(C_1 | x)} + \lambda_{1,2} p(C_2 | x) \\ \lambda_{2,1} p(C_1 | x) &> \lambda_{1,2} p(C_2 | x) \\ p(C_1 | x) &> 4 * p(C_2 | x) \end{aligned}$$

With Bayes:

$$\begin{aligned} \frac{P(x | C_1) * P(C_1)}{P(x)} &> 4 * \frac{P(x | C_2) * P(C_2)}{P(x)} \\ \frac{P(x | C_1) * P(C_1)}{\cancel{P(x)}} &> 4 * \frac{P(x | C_2) * P(C_2)}{\cancel{P(x)}} \\ P(x | C_1) * \cancel{P(C_1)} &> 4 * P(x | C_2) * \cancel{P(C_2)} \\ P(x | C_1) &> 4 * P(x | C_2) \end{aligned}$$

With Gaussian:

$$p(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

Decision Boundary at $p(C_1) = p(C_2)$, $\sigma_1^2 = \sigma_2^2 = \sigma^2$, $\mu_1 > 0$, $\mu_1 = 2\mu_2$ (and same misclassification costs):

$$\begin{aligned} p(C_1 | x) &> 4 * p(C_2 | x) \\ \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu_1)^2}{2\sigma^2}\right) &> 4 * \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu_2)^2}{2\sigma^2}\right) \\ \cancel{\frac{1}{\sqrt{2\pi\sigma^2}}} \exp\left(-\frac{(x - \mu_1)^2}{2\sigma^2}\right) &> 4 * \cancel{\frac{1}{\sqrt{2\pi\sigma^2}}} \exp\left(-\frac{(x - \mu_2)^2}{2\sigma^2}\right) \\ \exp\left(-\frac{(x - \mu_2)^2}{2\sigma^2}\right) &> 4 * \exp\left(-\frac{(x - \mu_2)^2}{2\sigma^2}\right) \\ -\frac{(x - \mu_2)^2}{2\sigma^2} &> \ln(4) - \frac{(x - \mu_1)^2}{2\sigma^2} \\ -(x - \mu_2)^2 &> 2\sigma_2^2 \ln(4) - (x - \mu_1)^2 \\ -x^2 + 2x\mu_1 - \mu_1^2 &> 2\sigma_2^2 \ln(4) - (x^2 - 2x\mu_2 - \mu_2^2) \\ 2(\mu_1 - \mu_2)x &> 2\sigma_2^2 \ln(4) + \mu_1^2 - \mu_2^2 \\ x &> \frac{2\sigma_2^2 \ln(4) + \mu_1^2 - \mu_2^2}{2(\mu_1 - \mu_2)} \end{aligned}$$

with $\mu_1 = 2\mu_2$

$$x > \frac{2\sigma^2 \ln(4) + 3\mu_2^2}{2\mu_2}$$

When the cost of deciding C_1 when the true class is C_2 is higher than the other way round, make sense that the decision boundary is now shifted to the right closer to μ_1 (in this case by $\frac{\sigma^2 \ln(4)}{\mu_2}$), because now we are more conservative in choosing C_1 .

2 Density Estimation

2.1 Gaussian Maximization Likelihood Estimate

The p -variate Gaussian distribution, $\mathbf{x} \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, with $\mathbf{x}, \boldsymbol{\mu} \in \mathbb{R}^p$, $\boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$:

$$f(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{p}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right],$$

where $|\boldsymbol{\Sigma}|$ is the determinant of the covariance matrix

With a set of i.i.d data $\mathbf{x}_1, \dots, \mathbf{x}_N$ drawn from $\mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, the likelihood function is

$$\ell = p(\mathbf{x}_1, \dots, \mathbf{x}_N | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{Np}{2}} |\boldsymbol{\Sigma}|^{-\frac{N}{2}} \exp \left[-\frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right],$$

and the log likelihood

$$\mathcal{L} = \ln \ell = -\frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) + [\text{some constant}].$$

Note that the square mahalanobis distance $D^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = \text{Tr} \left[\boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})^T (\mathbf{x} - \boldsymbol{\mu}) \right]$, where $\text{Tr}[\cdot]$ is the trace function.

- Taking derivative w.r.t $\boldsymbol{\mu}$ and setting it to zero:

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\mu}} \left\{ -\frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \text{Tr} \left[\boldsymbol{\Sigma}^{-1} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T (\mathbf{x}_n - \boldsymbol{\mu}) \right] \right\} &\stackrel{!}{=} 0 \\ \cancel{\frac{1}{2}} \left\{ \cancel{2\boldsymbol{\Sigma}^{-1}} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}) \right\} &\stackrel{!}{=} 0 \\ \sum_{n=1}^N \mathbf{x}_n - N\boldsymbol{\mu} &\stackrel{!}{=} 0 \\ \Rightarrow \hat{\boldsymbol{\mu}} &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n. \end{aligned}$$

- Taking derivative w.r.t $\boldsymbol{\Sigma}^{-1}$ and setting it to zero:

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\Sigma}^{-1}} \left\{ -\frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \text{Tr} \left[\boldsymbol{\Sigma}^{-1} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T (\mathbf{x}_n - \boldsymbol{\mu}) \right] \right\} &\stackrel{!}{=} 0 \\ \frac{\partial}{\partial \boldsymbol{\Sigma}^{-1}} \left\{ N \ln |\boldsymbol{\Sigma}^{-1}| - \text{Tr} \left[\boldsymbol{\Sigma}^{-1} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T (\mathbf{x}_n - \boldsymbol{\mu}) \right] \right\} &\stackrel{!}{=} 0 \\ N\boldsymbol{\Sigma}^T - \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}) (\mathbf{x}_n - \boldsymbol{\mu})^T &\stackrel{!}{=} 0 \\ \Rightarrow \hat{\boldsymbol{\Sigma}} &= \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \hat{\boldsymbol{\mu}}) (\mathbf{x}_n - \hat{\boldsymbol{\mu}})^T. \end{aligned}$$

Using:

- $|\boldsymbol{\Sigma}|^{-1} = |\boldsymbol{\Sigma}^{-1}|$
- $\frac{\partial}{\partial \mathbf{A}} \text{Tr}[\mathbf{A}\mathbf{B}] = \mathbf{B}^T$
- $\frac{\partial}{\partial \mathbf{A}} \ln |\mathbf{A}| = (\mathbf{A}^{-1})^T$
- $\boldsymbol{\Sigma}^T = \boldsymbol{\Sigma}$.

2.2 Prior Probabilities

Probability Rule:

$$\sum P(X_n) = \sum_n P(X_n \in C_i) = 1, \forall i$$

Formula:

$$P(C_i) = \frac{|C_i|}{\sum_i |C_i|}$$

Prior Probability of densEst1	Prior Probability of densEst2
0.239	0.761

where C_i is the cardinality of class i

2.3 Biased ML Estimate

Tutor solution

The bias of an estimator $\hat{\theta}$ is the drift of the estimated value from the true value. It can be computed by:

$$\text{bias}(\hat{\theta}) = \mathbb{E}_{\theta}[\hat{\theta} - \theta]$$

where \mathbb{E}_{θ} denotes the expected value over the distribution $p_{\theta}(x)$.

Bias of an estimator is the difference between this estimator's expected value and the true value of the parameter being estimated.

	Definition of Biased Estimator	Definition of Unbiased Estimator
sample variance	$\hat{\sigma} = \frac{1}{N} \sum_{n=1}^N (x_n - \hat{\mu})^2$	$\hat{\sigma} = \frac{1}{N-1} \sum_{n=1}^N (x_n - \hat{\mu})^2$
sample mean	$\hat{\mu} = \frac{1}{N} \sum_{n=1}^N x_n$	$\hat{\mu} = \frac{1}{N} \sum_{n=1}^N x_n$

Multivariate Gaussian Distribution:

$$f(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{p}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right],$$

Parameters:

- $\boldsymbol{\mu}$: mean -> have to be calculated
- $\boldsymbol{\Sigma}$: Covariance -> have to be calculated
- $p = 2$: dimension -> given
- \mathbf{x} : datapoints -> given
- π : constant

Then

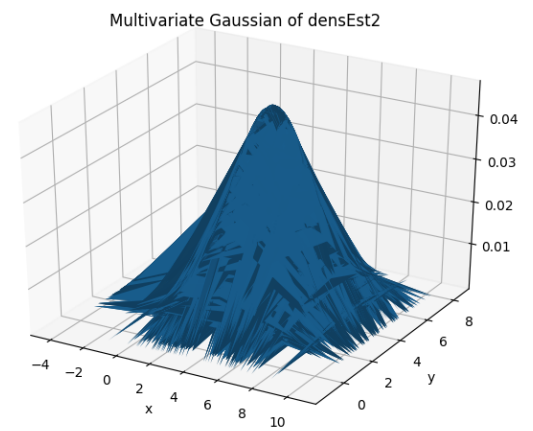
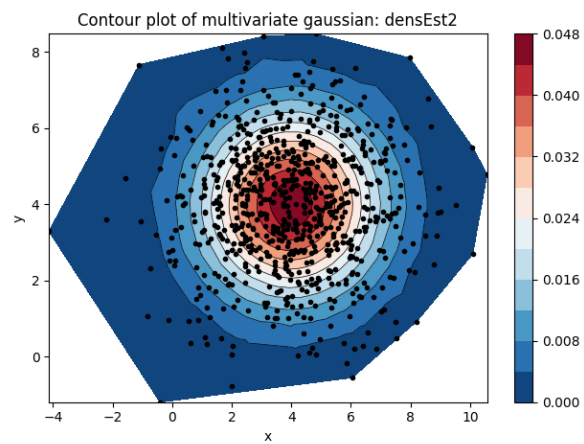
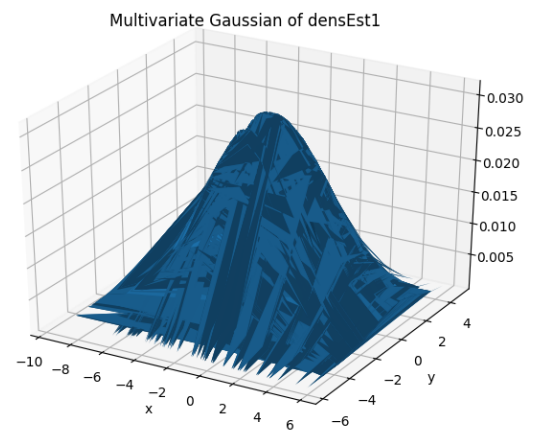
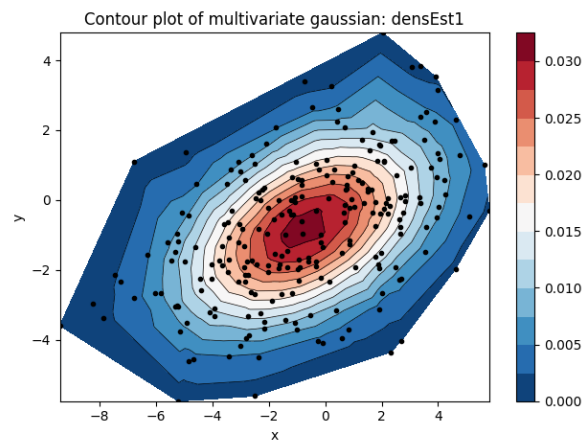
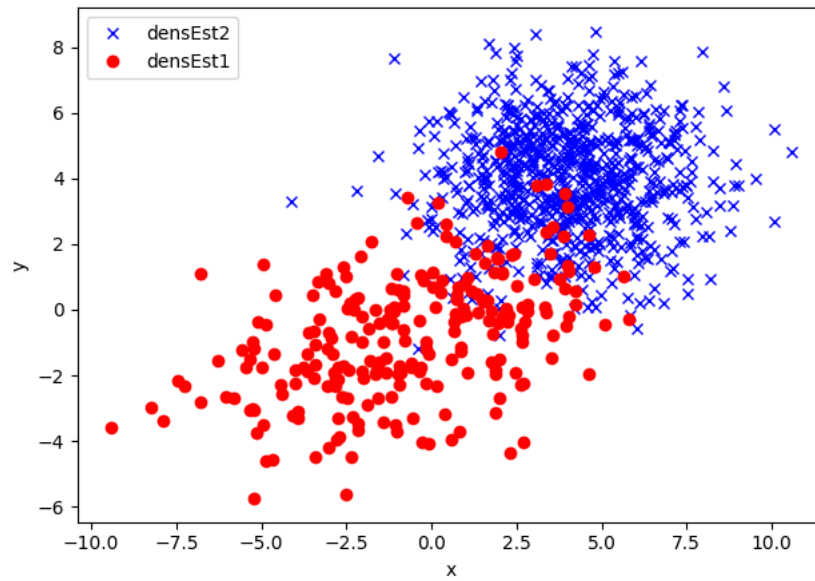
$$f(\mathbf{x}|C_i) = f(\mathbf{x}|\hat{\boldsymbol{\mu}}_i, \hat{\boldsymbol{\Sigma}}_i) = (2\pi)^{-\frac{p}{2}} |\hat{\boldsymbol{\Sigma}}_i|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \hat{\boldsymbol{\mu}}_i)^T \hat{\boldsymbol{\Sigma}}_i^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_i) \right]$$

	DensEst1	DensEst2
Mean	(-0.70681374, -0.81343083)	(3.98534252, 3.98438364)
Covariance Unbiased Estimator	$\begin{pmatrix} 9.05742302 & 2.6841014 \\ 2.6841014 & 3.61145033 \end{pmatrix}$	$\begin{pmatrix} 4.18087542 & 0.02761954 \\ 0.02761954 & 2.75658555 \end{pmatrix}$
Covariance biased estimator	$\begin{pmatrix} 9.01952586 & 2.67287085 \\ 2.67287085 & 3.59633965 \end{pmatrix}$	$\begin{pmatrix} 4.1753815 & 0.02758324 \\ 0.02758324 & 2.75296323 \end{pmatrix}$

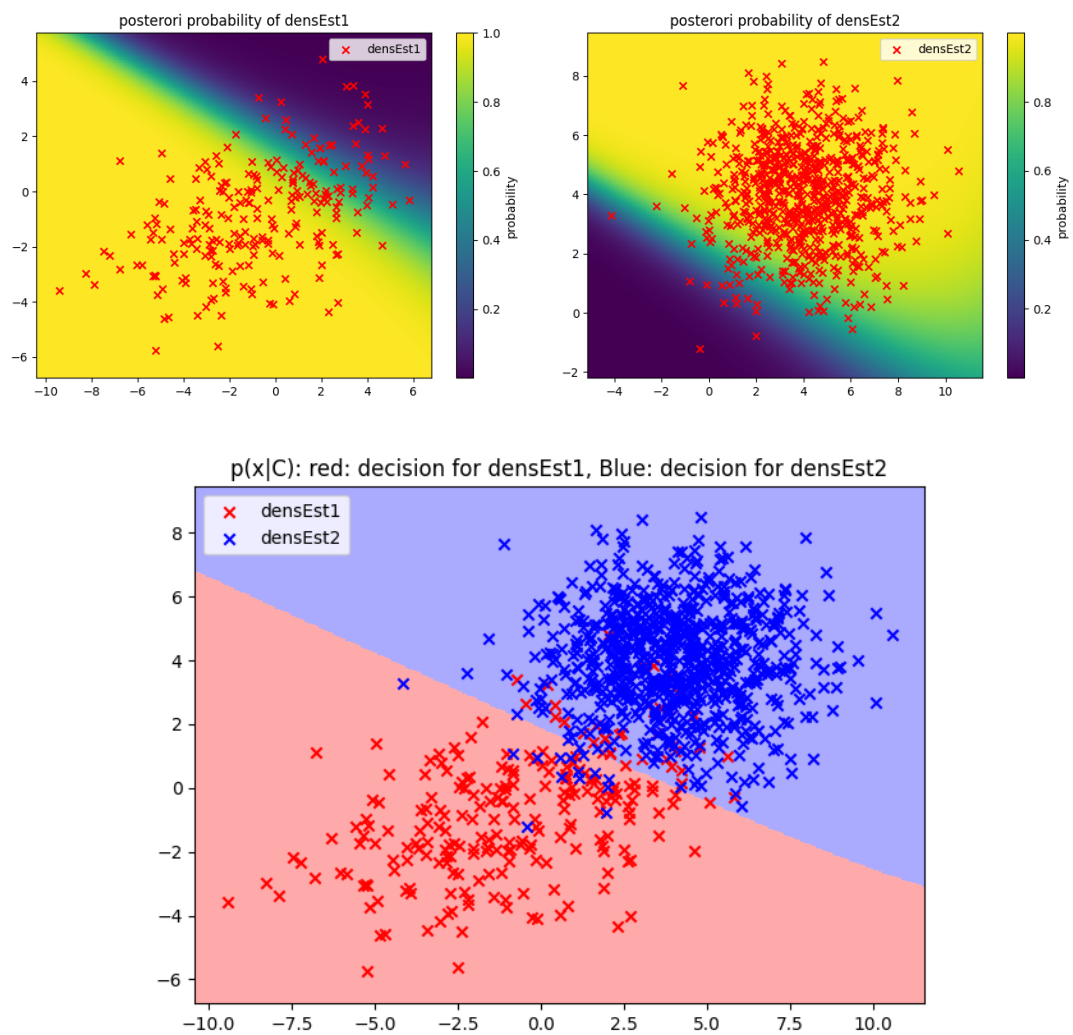
Tutor has different numbers

2.4 Class Density

Visualization of DensEst1 and DensEst2

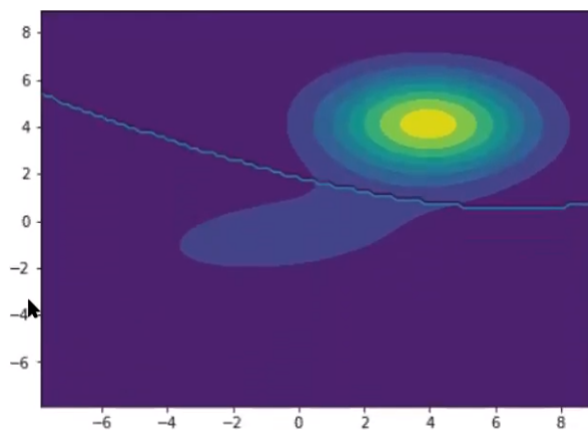
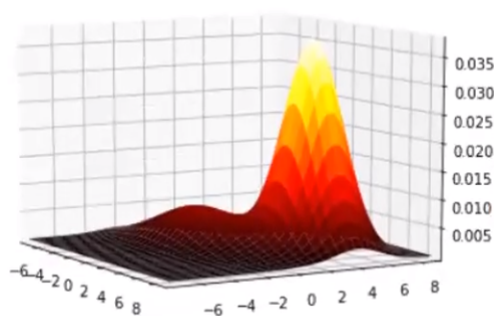


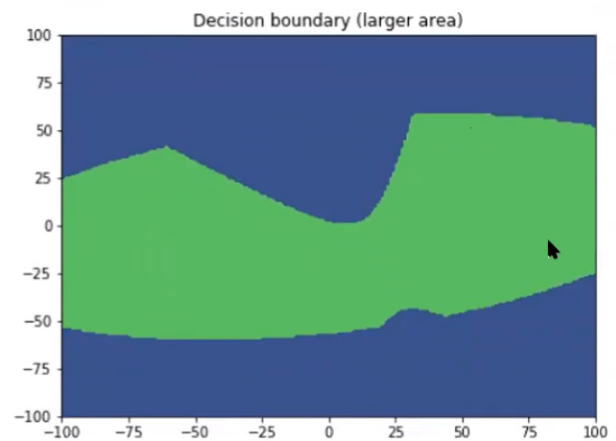
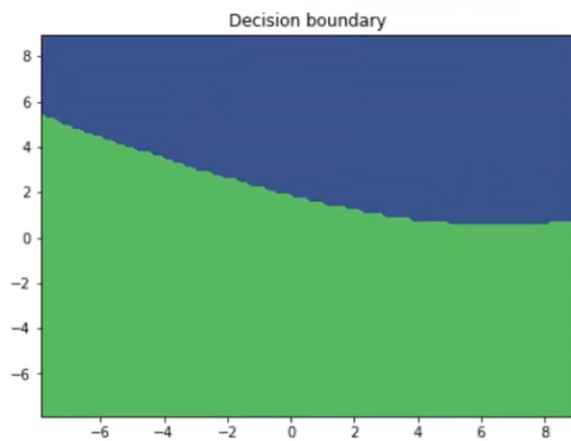
2.5 Posterior



Tutors solution plot:

Surface of the posterior of the mixture of two Gaussians





```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from matplotlib.colors import ListedColormap
4
5 """ 2c) Biased ML Estimate """
6 class Data:
7     mean = 0
8
9     variance_unbiased = 0
10    variance_biased = 0
11
12    cov_biased = 0
13    cov_unbiased = 0
14
15    def __init__(self, x, y, name):
16        self.x = x
17        self.y = y
18        self.name = name
19
20        self.calc_mean()
21        self.calc_sample_variance_biased_and_unbiased()
22        self.covariance_setup()
23
24    def calc_covariance_biased(self, x, y):
25        if (x == 0 and y == 0):
26            return self.variance_biased[x]
27        elif (x == 1 and y == 1):
28            return self.variance_biased[y]
29        else:
30            sum = 0
31            for i in range(self.x.shape[0]):
32                sum += (self.x[i] - self.mean[0])*(self.y[i] - self.mean[1])
33            return sum / self.x.shape[0]
34
35    def calc_covariance_unbiased(self, x, y):
36        if (x == 0 and y == 0):
37            return self.variance_unbiased[x]
38        elif (x == 1 and y == 1):
39            return self.variance_unbiased[y]
40        else:
41            sum = 0
42            for i in range(self.x.shape[0]):
43                sum += (self.x[i] - self.mean[0])*(self.y[i] - self.mean[1])
44            return sum / (self.x.shape[0] - 1)
45
46    def covariance_setup(self):
47        self.cov_biased = np.asarray([self.calc_covariance_biased(x, y) for x in range(2) for y in range(2)]).reshape(2, 2)
48        self.cov_unbiased = np.asarray([self.calc_covariance_unbiased(x, y) for x in range(2) for y in range(2)]).reshape(2, 2)
49
50    def calc_mean(self):
51        val_x = 0
52        val_y = 0
53        for i in range(self.x.shape[0]):
54            val_x += self.x[i]

```



```

56         val_y += self.y[i]
57         self.mean = np.asarray([val_x / len(self.x), val_y / len(self.y)])
58
59     def calc_sample_variance_biased_and_unbiased(self):
60         val_x = 0
61         val_y = 0
62         for i in range(self.x.shape[0]):
63             val_x += (self.x[i] - self.mean[0])**2
64             val_y += (self.y[i] - self.mean[1])**2
65
66         self.variance_unbiased = np.asarray([val_x / (len(self.x) - 1), val_y / (len(self.y) - 1)])
67         self.variance_biased = np.asarray([val_x / len(self.x), val_y / len(self.y)])
68     """ end 2c) """
69
70
71 def get_data(name :str) -> np.ndarray:
72     file = open(name, 'r')
73     x = []
74     y = []
75     for i in file.readlines():
76         line = i.split()
77         x.append(float(line[0]))
78         y.append(float(line[1]))
79     file.close()
80     return np.asarray(x), np.asarray(y)
81
82 def show_data(densEst1, densEst2):
83     plt.figure()
84     plt.plot(densEst2.x, densEst2.y, 'bo', marker='x', label='densEst2')
85     plt.plot(densEst1.x, densEst1.y, 'ro', marker='o', label='densEst1')
86     plt.title('Vizualization of DensEst1 and DensEst2')
87     plt.legend()
88     plt.xlabel('x')
89     plt.ylabel('y')
90     plt.show()
91
92 def calc_multi_gaussian(x, cov, mean):
93     factor = 1 / np.sqrt((2*np.pi)**2 * np.linalg.det(cov))
94     z = []
95     for pos in x:
96         multi = factor * np.exp(-0.5 * (pos - mean) @ np.linalg.inv(cov) @ (pos.T - mean.T))
97         z.append(multi)
98     return np.asarray(z)
99
100 def multi_gaussian(data):
101     z = calc_multi_gaussian(np.c_[data.x, data.y], data.cov_unbiased, data.mean)
102     fig, ax2 = plt.subplots()
103     ax2.tricontour(data.x, data.y, z, levels=14, linewidths=0.5, colors='k')
104     cntr2 = ax2.tricontourf(data.x, data.y, z, levels=14, cmap="RdBu_r")
105
106     fig.colorbar(cntr2, ax=ax2)
107     ax2.plot(data.x, data.y, 'ko', ms=3)
108     plt.xlabel('x')
109     plt.ylabel('y')
110     plt.title('Contour plot of multivariate gaussian: ' + data.name)
111     plt.subplots_adjust(hspace=0.5)
112     plt.show()
113
114 def prior_probability(densEst1, densEst2):
115     prior_dens1 = len(densEst1.x) / (len(densEst1.x) + len(densEst2.x))
116     prior_dens2 = len(densEst2.x) / (len(densEst1.x) + len(densEst2.x))
117     return np.asarray([prior_dens1, prior_dens2])
118
119 def calculate_posteriori(densEst1, densEst2, xx, yy):
120     def likelihood(xx, yy, cov, mean):
121         likeli = np.zeros(xx.shape)
122         for i in range(xx.shape[1]):
123             likeli[:, i] = calc_multi_gaussian(np.c_[xx[:, i], yy[:, i]], cov, mean)
124         return likeli
125
126     prior_dens1, prior_dens2 = prior_probability(densEst1, densEst2)
127
128     likelihood_dens1 = likelihood(xx, yy, densEst1.cov_unbiased, densEst1.mean)
129     likelihood_dens2 = likelihood(xx, yy, densEst2.cov_unbiased, densEst2.mean)
130     normalization = likelihood_dens1 * prior_dens1 + likelihood_dens2 * prior_dens2
131
132     posteriori_dens1 = likelihood_dens1 * prior_dens1 / normalization
133     posteriori_dens2 = likelihood_dens2 * prior_dens2 / normalization
134
135

```

```

136     return np.asarray(posteriori_dens1), np.asarray(posteriori_dens2)
137
138 def decision(Z1, Z2):
139     Z = np.zeros(Z1.shape)
140     for i in range(Z1.shape[0]):
141         for j in range(Z1.shape[1]):
142             if(Z1[i, j] >= Z2[i, j]):
143                 Z[i, j] = 0
144             else:
145                 Z[i, j] = 1
146     return Z
147
148
149 def multi_single_plot(densEst1, densEst2):
150     cmap_light = ListedColormap(["#FFAAAA", "#AAAAFF"])
151     X = np.c_[densEst1.x, densEst1.y]
152     x_min1, x_max1 = X[:, 0].min() - 1, X[:, 0].max() + 1
153     y_min1, y_max1 = X[:, 1].min() - 1, X[:, 1].max() + 1
154
155     X = np.c_[densEst2.x, densEst2.y]
156     x_min2, x_max2 = X[:, 0].min() - 1, X[:, 0].max() + 1
157     y_min2, y_max2 = X[:, 1].min() - 1, X[:, 1].max() + 1
158
159     xx, yy = np.meshgrid(np.arange(x_min1, x_max1, 0.05), np.arange(y_min1, y_max1, 0.05))
160     Z, _ = calculate_posteriori(densEst1, densEst2, xx, yy)
161     plt.figure()
162     plt.pcolormesh(xx, yy, Z, cmap='viridis')
163     plt.colorbar(label='probability')
164     plt.title('posteriori probability of ' + densEst1.name)
165     # Plot also the training points
166     plt.scatter(densEst1.x, densEst1.y, marker='x', color='red', label=densEst1.name)
167     plt.xlim(xx.min(), xx.max())
168     plt.ylim(yy.min(), yy.max())
169     plt.legend()
170     plt.show()
171
172
173     xx, yy = np.meshgrid(np.arange(x_min2, x_max2, 0.05), np.arange(y_min2, y_max2, 0.05))
174     _, Z = calculate_posteriori(densEst1, densEst2, xx, yy)
175     plt.figure()
176     plt.pcolormesh(xx, yy, Z, cmap='viridis')
177     plt.colorbar(label='probability')
178     plt.title('posteriori probability of ' + densEst2.name)
179     # Plot also the training points
180     plt.scatter(densEst2.x, densEst2.y, marker='x', color='red', label=densEst2.name)
181     plt.xlim(xx.min(), xx.max())
182     plt.ylim(yy.min(), yy.max())
183     plt.legend()
184     plt.show()
185
186     """Plot for decision Boundary """
187     x_min = min(x_min1, x_min2) #find min element (x_value) of densEst1 and densEst2
188     x_max = max(x_max1, x_max2) #find max element (x_value) " " " "
189     y_min = min(y_min1, y_max1)
190     y_max = max(y_max1, y_max2)
191
192     """ Create meshgrid of x and y axis"""
193     xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.05), np.arange(y_min, y_max, 0.05))
194     Z1, Z2 = calculate_posteriori(densEst1, densEst2, xx, yy)
195     Z = decision(Z1, Z2) #make a decision of the probability
196
197     """ create decision boundary with both dataset in one plot"""
198     plt.figure()
199     plt.pcolormesh(xx, yy, Z, cmap=cmap_light)
200     plt.title('p(x|C): red: decision for densEst1, Blue: decision for densEst2')
201     # Plot also the training points
202     plt.scatter(densEst1.x, densEst1.y, marker='x', color='red', label='densEst1')
203     plt.scatter(densEst2.x, densEst2.y, marker='x', color='blue', label='densEst2')
204     plt.xlim(xx.min(), xx.max())
205     plt.ylim(yy.min(), yy.max())
206     plt.legend()
207     plt.show()
208
209 def main():
210     #inicialization
211     x_densEst1, y_densEst1 = get_data('dataSets/densEst1.txt')
212     x_densEst2, y_densEst2 = get_data('dataSets/densEst2.txt')
213
214     densEst1 = Data(x_densEst1, y_densEst1, 'densEst1')
215     densEst2 = Data(x_densEst2, y_densEst2, 'densEst2')

```

```

216 """ 2d) Class Density """
217 show_data(densEst1, densEst2)
218 multi_gaussian(densEst1)
219 multi_gaussian(densEst2)
220
221 """ 2e) Posterior """
222 multi_single_plot(densEst1, densEst2)
223
224
225 if __name__ == '__main__':
226     main()

```

2.6 Bayesian Estimation

- State the generic case of Bayesian linear regression with data $\langle \vec{X}, \vec{Y} \rangle$ and parameter θ :
 - Separating two Classes form each other with a regression line. The aim is not not to find the single best value of the model parameters, but rather to determine the posterior distribution for the model parameters. Also, we want to maximize the Likelihood. This maximization of the likelihood can lead to excessively complex model and over-fitting. With Bayesian treatment of linear regression it avoid the over-fitting problem of the maximum likelihood and also determining the model complexity using the trainingdata alone.
- What do we assume about the data, the model and the parameter:
 - Data: $\vec{X}, \vec{Y} \in \mathbb{R}^{1 \times N}$ and they are independent
 - Model: Supervised Learning. The Label of the data are given.
 - Parameter θ : Is a normal distribution, with $\theta = (\mu, \sigma^2)$
- Formulate the posterior distribution for your model parameters given the data, i.e., $p(\theta | \vec{X}, \vec{Y})$, and derive its mean and covariance, assuming that the model of the output variable is a Gaussian distribution with a fixed variance.

Let $\theta = (\mu, \sigma^2)$ The posterior distribuiton can be written as:

$$p(\theta | X_1, \dots, X_n) = \frac{p(X_1, \dots, X_n | \theta)p(\theta)}{p(X_1, \dots, X_n)} = \frac{\mathcal{L}_n(\theta)p(\theta)}{p(X_1, \dots, X_n)} \propto \mathcal{L}_n(\theta)p(\theta)$$

where $\mathcal{L}_n(\theta) = \prod_{i=1}^n p(X_i | \theta)$ is the likelihood functon and

$$p(X_1, \dots, X_n) = \int p(X_1, \dots, X_n)p(\theta)d\theta = \int L_n(\theta)p(\theta)d\theta$$

is the normalizing constant (evidence).

We also know, if we estimate the distribution of the mean:

$$p(\mu | X) = \frac{p(X | \mu)p(\mu)}{p(X)}$$

with the prior:

$$p(\mu) \sim \mathcal{N}(\mu_0, \sigma_0^2) \Rightarrow p(\mu) = \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp\left(-\frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right)$$

$$p(X | \mu) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

With n independent observation $X = (X_1, X_2, \dots, X_n)$, such that:

$$p(\mu | x) \propto p(x | \mu)p(\mu) = p(\mu)p(x_1 | \mu)p(x_2 | \mu) \dots p(x_n | \mu)$$

and the Formula from above:

$$\begin{aligned}
p(\mu | x) &= \frac{p(x | \mu)p(\mu)}{\int p(x | \mu)p(\mu)d\mu} = \frac{p(x | \mu)p(\mu)}{p(x)} \propto p(x | \mu)p(\mu) \\
&= \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp\left(-\frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right) * \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \\
&= \frac{1}{(2\pi)^{\frac{n+1}{2}} \sqrt{\sigma_0^2\sigma^{2n}}} \exp\left(-\frac{\mu^2 + 2\mu\mu_0 - \mu_0^2}{2\sigma_0^2} - \sum_{i=1}^n \frac{x_i^2 - 2\mu x_i + \mu^2}{2\sigma^2}\right) \\
&= \text{const} * \exp\left(-\frac{\mu^2(\sigma^2 + n\sigma_0^2) + 2\mu(\mu_0\sigma^2 + \sigma_0^2x_1 + \dots + \sigma_0^2x_n) - (\mu_0^2\sigma^2 + \sigma_0^2x_1^2 + \dots + \sigma_0^2x_n^2)}{2\sigma_0^2\sigma^2}\right) \\
&\propto \exp\left(-\frac{-\mu^2 + 2\mu\frac{\mu_0\sigma^2 + \sum_{i=1}^n \sigma_0^2x_i}{\sigma^2 + n\sigma_0^2} - \left(\frac{\mu_0\sigma^2 + \sum_{i=1}^n \sigma_0^2x_i}{\sigma^2 + n\sigma_0^2}\right)^2}{2\frac{\sigma_0^2\sigma^2}{\sigma^2 + n\sigma_0^2}}\right) * \exp\left(-\frac{\mu_0^2\sigma^2 + \sum_{i=1}^n \sigma_0^2x_i^2}{2\sigma_0^2\sigma^2}\right) \\
&\propto \exp\left(-\frac{\left(\mu - \frac{\mu_0\sigma^2 + \sum_{i=1}^n \sigma_0^2x_i}{\sigma^2 + n\sigma_0^2}\right)^2}{2\frac{\sigma_0^2\sigma^2}{\sigma^2 + n\sigma_0^2}}\right)
\end{aligned}$$

Letting:

$$\begin{aligned}
\sigma_1 &= \frac{\sigma_0^2\sigma^2}{\sigma^2 + n\sigma_0^2} = \frac{1}{\sigma_0^{-2} + n\sigma^{-2}} \\
\mu_1 &= \frac{\mu_0^2\sigma^2 + \sum_{i=1}^n \sigma_0^2x_i^2}{\sigma^2 + n\sigma_0^2} = \frac{\mu_0\sigma_0^{-2} + \sum_{i=1}^n x_i\sigma^{-2}}{\sigma_0^{-2} + n\sigma^{-2}} = \sigma_1^2(\mu_0\sigma_0^{-2} + \sum_{i=1}^n x_i\sigma^{-2})
\end{aligned}$$

Now we can write it as:

$$\begin{aligned}
\sigma_1^2 &= \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}\right)^{-1} \\
\mu_1 &= \sigma_1^2\left(\frac{\mu_0}{\sigma_0^2} + \frac{\bar{x}n}{\sigma^2}\right)
\end{aligned}$$

We can alternatively write these formulas as:

$$\begin{aligned}
\mu_n &= \frac{N\sigma_0\bar{x} + \sigma^2\mu}{N\sigma_0^2 + \sigma^2} \\
\frac{1}{\sigma_N^2} &= \frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}
\end{aligned}$$

Multidimensional case:

$$p(\boldsymbol{\theta} | \mathbf{x}, y) \propto p(y | \mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta}).$$

Assuming multivariate Gaussian prior for $p(\boldsymbol{\theta})$,

$$p(\boldsymbol{\theta}) = p(\boldsymbol{\theta} | \boldsymbol{\mu}_\theta, \boldsymbol{\Sigma}_\theta) \sim \mathcal{N}_p(\boldsymbol{\theta} | \boldsymbol{\mu}_\theta, \boldsymbol{\Sigma}_\theta).$$

Assuming the output is Gaussian distribution with a fixed variance, $\hat{y} \sim \mathcal{N}(\hat{y} | \mu_y, \sigma_y^2) = \mathcal{N}(\hat{y} | \boldsymbol{\theta}^T \mathbf{x}, \sigma_y^2)$.

Given N samples $\{\mathbf{x}_n, y_n\}_{n=1}^N$ contained in output vector $\mathbf{y} \in \mathbb{R}^N$ and design matrix $\mathbf{X} \in \mathbb{R}^{N \times p}$, the likelihood is

$$p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) \sim \mathcal{N}_N(\mathbf{y} | \mathbf{X}\boldsymbol{\theta}, \sigma_y^2 \mathbf{I}).$$

The posterior is then a multiplication of two multivariate Gaussians,

$$\begin{aligned}
p(\boldsymbol{\theta} | \mathbf{x}, y) &\propto p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta}) \\
&\propto \mathcal{N}_N(\mathbf{y} | \mathbf{X}\boldsymbol{\theta}, \sigma_y^2 \mathbf{I}) \mathcal{N}_p(\boldsymbol{\theta} | \boldsymbol{\mu}_\theta, \boldsymbol{\Sigma}_\theta) \\
&\propto c \cdot \mathcal{N}(\cdot | \boldsymbol{\mu}_{\text{post}}, \boldsymbol{\Sigma}_{\text{post}}),
\end{aligned}$$

where the posterior mean

$$\mu_{\text{post}} = \frac{\left((\sigma_y^2 \mathbf{I})^{-1} \mathbf{X} \theta + \Sigma_{\theta}^{-1} \mu_{\theta} \right)}{\left((\sigma_y^2 \mathbf{I})^{-1} + \Sigma_{\theta}^{-1} \right)^{-1}},$$

the posterior covariance

$$\Sigma_{\text{post}} = \left((\sigma_y^2 \mathbf{I})^{-1} + \Sigma_{\theta}^{-1} \right)^{-1}.$$

with \mathbf{I} identity matrix in appropriate dimension¹.

- What do we do when we want to predict a new point?
 - We want the quality of the prediction. That means, we want to minimize the error of the calculation (i.e. RSS (Residual sum of squares), MSE (Mean squared error), etc.). And that means, we maximize the Likelihood, i.e. to approximate the probability of posterior. The Error minimization RSS corresponds to maximum likelihood if normal distribution is given (regression)
- Which are the advantages of being Bayesian?
 - By Bayesian regression, we introduce a prior to the parameter θ , $p(\theta)$, that may help against overfitting and instabilities
 - Bayesian helps us to use priors to help us with normalization.
 - If we have new information available, the previous posterior distribution can be used as a prior
 - it provides a wide range of models (hierarchical models and missing data problems)²

Tutor solution:

- Assumptions:
 - the data is identically independently distributed (i.i.d),
 - our model is linear in the parameters with some additional Gaussian noise, i.e. $y = f(\vec{x}, \vec{\theta}) + \epsilon$, where $\vec{\theta}$ are parameters we want to learn and $\epsilon \sim \mathcal{N}(0, \sigma_y^2)$
 - prior distribution over the parameters $p(\vec{\theta} | \vec{\mu}_{\vec{\theta}}, \vec{\Sigma}_{\vec{\theta}})$
- Given our second assumption, we can reformulate the model as Given our third assumption, we get the posterior distribution for $\vec{\theta}$ by using Bayes'theorem

$$\begin{aligned} p(\vec{\theta} | \vec{X}, \vec{Y}) &= \frac{p(\vec{Y} | \vec{X}, \vec{\theta}, \sigma_y^2) p(\vec{\theta} | \vec{\mu}_{\vec{\theta}}, \vec{\Sigma}_{\vec{\theta}})}{\vec{Y}} \\ &\Downarrow \\ p(\vec{\theta} | \vec{X}, \vec{Y}) &\propto p(\vec{Y} | \vec{X}, \vec{\theta}, \sigma_y^2) p(\vec{\theta} | \vec{\mu}_{\vec{\theta}}, \vec{\Sigma}_{\vec{\theta}}) = \mathcal{N}(\vec{\mu}, \vec{\Sigma}) \end{aligned}$$

- As everythin is Gaussian, we can compute the posterior in closed form (see Bishop, 2006, Eq. 1.68, 1.70, 2.115)(as stated in the assignment, we assume that σ_y^2 is constant)

$$\begin{aligned} \vec{\Sigma}^{-1} &= \vec{\Sigma}_{\vec{\theta}}^{-1} + \frac{1}{\sigma_y^2} \vec{X}^T \vec{X} \\ \vec{\mu} &= \vec{\Sigma}(\vec{\Sigma}_{\vec{\theta}}^{-1} \vec{\mu}_{\vec{\theta}} + \frac{1}{\sigma_y^2} \vec{X}^T \vec{X}) \end{aligned}$$

- When we want to predict a new point, we have to marginalize the model across the uncertainty in the parameters (we omit the dependencies on the means and covariances for simplicitys sake)

$$p(y | \vec{x}, \vec{Y}) = \int p(y | \vec{x}, \vec{\theta}) p(\vec{\theta} | \vec{X}, \vec{Y}) d\vec{\theta}$$

- The main advantage of this approach is that we do not have to fix the parameters of the model. instead of using one set of parameters, we average over all of the possible ones. Additionally, we are able to introduce a prior probability over the parameters. As a result we get a probability indicating the uncertainty of the model predictions. (more information in tutors exercise hour)

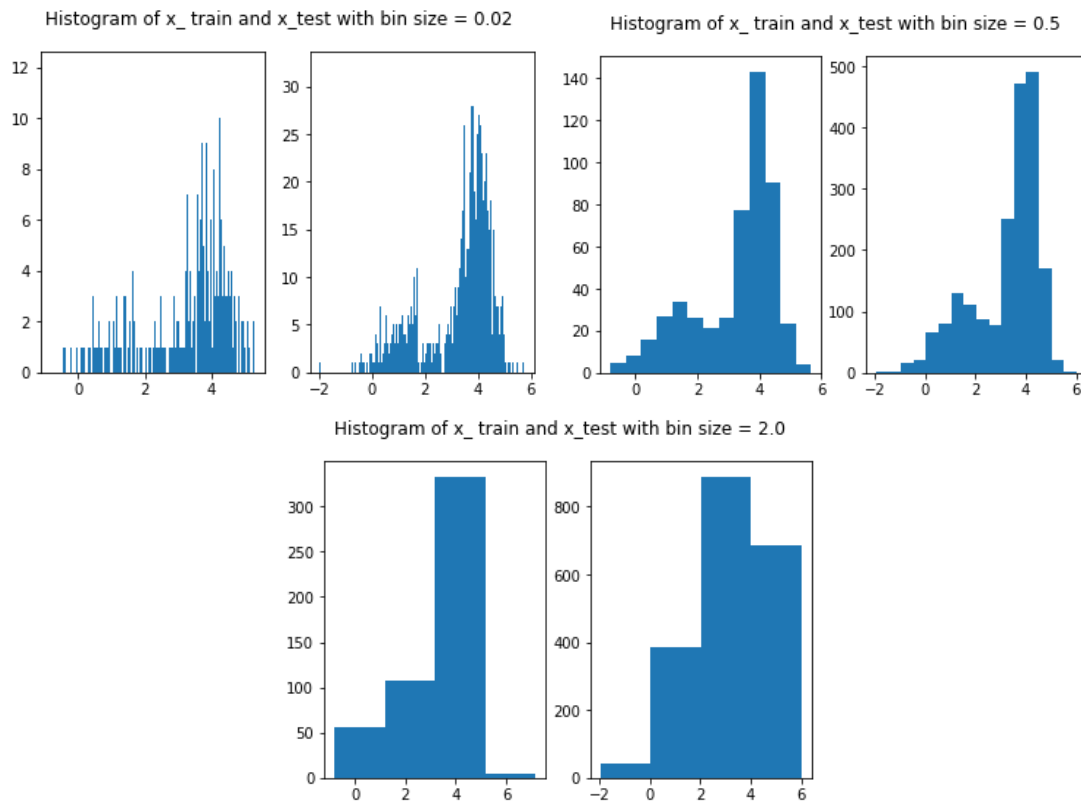
¹<http://compbio.fmph.uniba.sk/vyuka/ml/old/2008/handouts/matrix-cookbook.pdf>

²Source: https://www.cpp.edu/~djmorarty/wed/bayes_handout.pdf

3 Non-parametric Density Estimation

3.1 Histogram

The following figures show the histograms of the train and test data using different bin sizes (widths).



Intuitively, it is better to use the bin size of 0.5, because we can observe that using bin width that is too small (0.02) and seems to be over-fitting the data, the density is not smooth enough, while with bin width that is too big (2), the density is too smooth.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 def get_counts(x, binsize):
5     xmin = min(x)
6     xmax = max(x)
7     which_bin = np.floor((x-xmin)/binsize)
8     edges = np.arange(xmin, xmax + binsize, binsize)
9     Nbins = edges.shape[0]-1
10    counts = np.zeros(Nbins)
11    for n in range(Nbins):
12        counts[n] = np.count_nonzero(which_bin==n)
13    return (edges[1:]-binsize/2), counts
14
15 def plot_my_hist(x, y, binsize, title): # To plot two distributions side by side
16     fig, axes = plt.subplots(1,2)
17     fig.suptitle(title)
18
19     midpts, counts = get_counts(x, binsize)
20     axes[0].bar(x=midpts, height=counts, width=binsize)
21
22     midpts, counts = get_counts(y, binsize)
23     axes[1].bar(x=midpts, height=counts, width=binsize)
24
25     plt.show()
26
27 def main():
28     x_train = np.genfromtxt("nonParamTrain.txt")
29     x_test = np.genfromtxt("nonParamTest.txt")
30
```

```

31     for b in [0.02, 0.5, 2.0]:
32         plot_my_hist(x_train, x_test, b, "Histogram of x_train and x_test with bin size = "+ str(b))
33
34 if __name__ == '__main__':
35     main()

```

3.2 Kernel Density Estimate

Gaussian kernel in 1 dimension

$$k(u) = \frac{1}{\sqrt{2\pi}} \exp\{- (u)^2\},$$

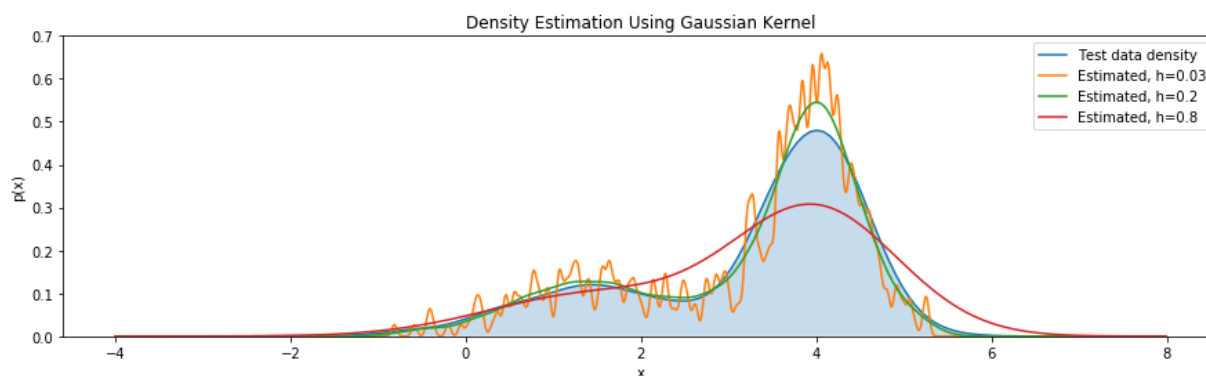
with

$$K(x; h) = \sum_{n=1}^N k\left(\frac{x - x^{(n)}}{h}\right)$$

where N is the number of total points and K is the number of points falling in the region R . Then, the density estimate (parameterized with h): is

$$p(x; h) \approx \frac{K(x; h)}{Nh} = \frac{1}{Nh} \sum_{n=1}^N k\left(\frac{x - x^{(n)}}{h}\right).$$

The estimated density was calculated using the training data with $h = 0.03, 0.2, 0.8$. The test data density was plotted using the kdeplot function from Seaborn. We can observe that for very small bandwidth ($h=0.03$), the curve is too grainy and the difference (error) w.r.t to the test data is relatively large. For $h=0.8$, the curve is overly smoothed and again the error w.r.t test data is large. In this case the best is $h=0.2$, where we can well distinguish the two modes (two peaks) and the closest one to the test data.



The likelihood function with a set of i.i.d data $x^{(1)}, \dots, x^{(M)}$ is

$$\mathcal{L} = \prod_{m=1}^M p(x^{(m)}; h) = \frac{1}{(Nh)^M} \prod_{m=1}^M K(x^{(m)}; h).$$

And the log likelihood

$$\ell = \ln \mathcal{L} = \sum_{m=1}^M \ln p(x^{(m)}; h) = -M \ln(Nh) + \sum_{m=1}^M \ln K(x^{(m)}; h).$$

Log likelihood:

- with $h=0.03$, training $\ell \approx -674.73$, test $\ell \approx -2812.11$
- with $h=0.2$, training $\ell \approx -717.02$, test $\ell \approx -2877.33$
- with $h=0.8$, training $\ell \approx -795.66$, test $\ell \approx -3192.49$

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 def H_gauss(u): #1-D Gaussian kernel
6     return np.exp(-(u)**2/2) / np.sqrt(2*np.pi)
7

```

```

8 def K_gauss(X, data, h):
9     M = X.shape[0]
10    N = data.shape[0]
11    K = np.zeros(M)
12    for m in range(M):
13        K_sum = 0
14        for n in range(N):
15            K_sum += H_gauss((X[m]-data[n])/h)
16        K[m] = K_sum
17    return K
18
19 def LL(p):
20     return np.sum(np.log(p))
21
22 def main():
23     x_train = np.genfromtxt("nonParamTrain.txt")
24     x_test = np.genfromtxt("nonParamTest.txt")
25     Ntrain = x_train.shape[0]
26     Ntest = x_test.shape[0]
27     X = np.arange(-4,8,0.01)
28
29     ll_train = np.zeros(3)
30     ll_test = np.zeros(3)
31     i = 0
32
33     plt.figure(figsize=(15,4))
34     sns.kdeplot(x_test, shade=True, label='Test data density')
35     for h in [0.03, 0.2, 0.8]:
36         p_train = K_gauss(x_train, x_train, h) / Ntrain / h
37         p_test = K_gauss(x_test, x_test, h) / Ntest / h
38         p_X_train = K_gauss(X, x_train, h) / Ntrain / h
39
40         # Log likelihood
41         ll_train[i] = LL(p_train)
42         print('Log likelihood of the training data with h='+str(h)+' is '+str(ll_train[i]))
43         ll_test[i] = LL(p_test)
44         print('Log likelihood of the test data with h='+str(h)+' is '+str(ll_test[i]))
45         i += 1
46
47         # Plotting
48         plt.plot(X, p_X_train, label='Estimated, h='+str(h))
49         plt.ylim((0.0,0.7))
50         plt.legend()
51         plt.title('Density Estimation Using Gaussian Kernel')
52         plt.ylabel('p(x)')
53         plt.xlabel('x')
54         plt.show()
55
56 if __name__ == '__main__':
57     main()

```

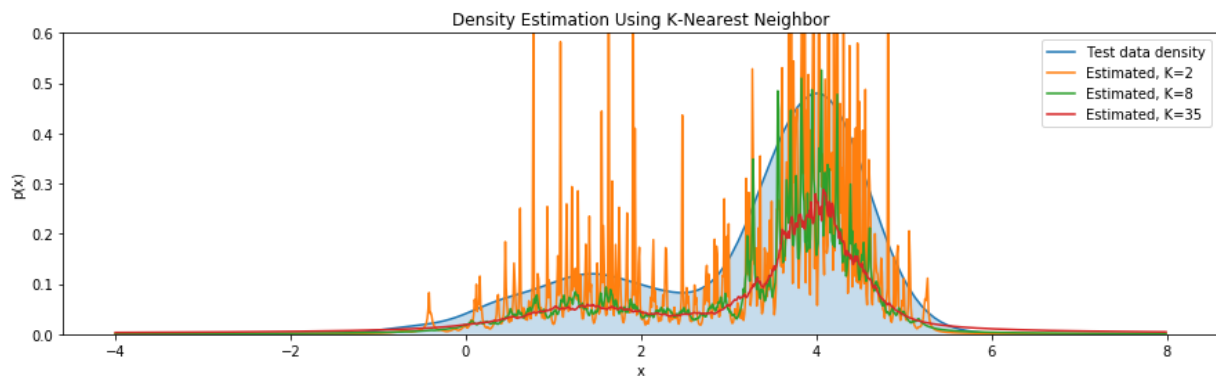
3.3 K-Nearest Neighbours

In the KNN method we grow the volume surrounding the estimation point x so that it encloses a total of K points. The density estimate then becomes

$$p(x; K) \approx \frac{K}{NV},$$

where for 1-D, $V = 2R_K$, where R_K is the distance to its K -th closest neighbor³.

³http://faculty.washington.edu/yenchic/18W_425/Lec7_knn-basis.pdf



The estimated density was calculated using the training data. The test data density was plotted using the `kdeplot` function from Seaborn. We can observe in the image that using $K = 2$, the estimate gives very bad result. With $K = 8$ is somewhat closer to the test density but still rather granular, not smooth enough. In this case, in our opinion, the best case is 35 where the result is smooth enough.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 def get_furthestKNNidx(x, data, K):
6     distances = np.abs(data-x)
7     return np.argsort(distances)[K-1] #index of furthest K nearest neighbor
8
9 def get_density(X, data, K):
10    NX = X.shape[0]
11    p_X = np.zeros(NX)
12    for n in range(NX):
13        furthest_knn_idx = get_furthestKNNidx(X[n], data, K) #index of furthest K nearest neighbor
14        R = np.abs(data[furthest_knn_idx] - X[n]) #max distance
15        p_X[n] = K/NX/2/R
16    return p_X
17
18 def main():
19    x_train = np.sort(np.genfromtxt("nonParamTrain.txt"))
20    x_test = np.sort(np.genfromtxt("nonParamTest.txt"))
21    X = np.arange(-4,8,0.01)
22
23    # Plotting
24    plt.figure(figsize=(15,4))
25    sns.kdeplot(x_test, shade=True, label='Test data density')
26    for K in [2,8,35]:
27        p_X_train = get_density(X, x_train, K)
28        plt.plot(X, p_X_train, label='Estimated, K='+str(K))
29    plt.ylim((0.0,0.6))
30    plt.legend()
31    plt.title('Density Estimation Using K-Nearest Neighbor')
32    plt.ylabel('p(x)')
33    plt.xlabel('x')
34    plt.show()
35
36 if __name__ == '__main__':
37     main()

```

3.4 Comparison of the Non-Parametric Methods

Estimate the log-likelihood of the testing data using the KDE estimators and the K-NN estimators. Why do we need to test them on a different data set? Compare the log-likelihoods of the estimators w.r.t both the training and testing sets in a table. Which estimator would you choose?

If we use only the training data for evaluating the performance of different parameters, then we are prone to overfit to the data. That means that we can explain the training data very well with our model, but we lose the generalization capability. As a result, when we evaluate our model on the testing data the performance degrades. Note that the testing data have been generated by the same process used for the training data!

The KDE estimator log-likelihoods on the testing data are $[-\infty, -2904.342, -3188.8334]$. The KNN estimator log-likelihoods are $[-2298.840, -2708.403, -2786.484]$. We observe that:

- We get the best performance using KNN estimator with $K=2$
- The log-likelihoods are considerably smaller than the ones on the training data
- However, from the plot we see that KNN with $K=2$ gives us a very noisy and overfitting density. The same for $K=8$
- With $K=35$ we get a good density, and its log-likelihood is higher than KDE with $\sigma = 0.2$, Therefore, KNN with $K=35$ would be a good choice

But why is the log-likelihood not truly a good indicator in this case?

- First, KNN does not give you a true distribution, i.e., the integral over the variable domain does not sum up to 1. For instance if we use $N=1$, we fit a distribution that has extremely high probability at each data point. As shown below, this density has even a higher log-likelihood than $N=2$ but is a terrible choice
- Second, the log-likelihood tells us how well we are fitting the data given some parameter. But it does not tell us how good the parameters are given the data (a.k.a posterior). Therefore, the likelihood is not always a good indicator for model selection. In general, MLE estimators are not always a good choice and sometimes it is better to have a prior over the parameters and use a MAP estimate.
 - The maximum likelihood estimate (MLE) of a parameter is the value of the parameter that maximizes the likelihood, where the likelihood is a function of the parameter and is actually equal to the probability of the data conditioning on the parameter.
 - Maximum a posteriori (MAP) estimation is the value of the parameter that maximizes the entire posterior distribution (which is calculated using the likelihood). A MAP estimate is the mode of the posterior distribution.
 - Note that there is no difference between the MLE and MAP estimate if the prior distribution we were assuming was a constant

```
plt.figure(figsize=(10,8))
p1_x = knn(x, 1, train_data)
p2_x = knn(x, 8, train_data)
p3_x = knn(x, 35, train_data)
line1, = plt.plot(x,p1_x,linewidth=2.0)
line2, = plt.plot(x,p2_x,linewidth=2.0)
line3, = plt.plot(x,p3_x,linewidth=2.0)
plt.legend([line1, line2, line3], ["K = 1", "K = 8", "K = 35"])
plt.xlabel('x')
plt.ylabel('p_x')

plt.axis([-4, 8, 0, 3])

p1_test = knn(test_data, 1, train_data)
log_like_test_1 = np.sum(np.log(p1_test))
p2_test = knn(test_data, 8, train_data)
log_like_test_2 = np.sum(np.log(p2_test))
p3_test = knn(test_data, 35, train_data)
log_like_test_3 = np.sum(np.log(p3_test))

print('\nlog likelihood on the test set')
print([log_like_test_1, log_like_test_2, log_like_test_3])
```

4 Expectation Maximization

4.1 Gaussian Mixture Update Rules

M -component Gaussian mixture model:

$$\begin{aligned} p(\mathbf{x}_n | \boldsymbol{\theta}) &= \sum_{j=1}^M p(\mathbf{x}_n | j) p(j) \\ &\equiv \sum_{j=1}^M \pi_j \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \end{aligned}$$

where the mixture parameters $\boldsymbol{\theta} = \{\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \pi_1, \dots, \boldsymbol{\mu}_M, \boldsymbol{\Sigma}_M, \pi_M\}$, $\sum_{j=1}^M \pi_j = 1$, and in this case $M = 4$, $\mathbf{x}, \boldsymbol{\mu} \in \mathbb{R}^2$, $\boldsymbol{\Sigma} \in \mathbb{R}^{2 \times 2}$. Given the observed data $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, we want to maximize the log-likelihood function of X w.r.t to the parameters:

$$\arg \max_{\boldsymbol{\theta}} \ell(X | \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \sum_{n=1}^N \log \sum_{j=1}^M \pi_j \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j).$$

E- (Expectation) step: Compute the posterior distribution $p(j | \mathbf{x}_n)$ for each mixture component j and for all data points $n = 1, \dots, N$. For time step i :

$$\hat{\alpha}_{nj}^{(i)} = [\hat{p}(j | \mathbf{x}_n)]^{(i)} = \frac{\hat{\pi}_j^{(i-1)} \mathcal{N}(\mathbf{x}_n; \hat{\boldsymbol{\mu}}_j^{(i-1)}, \hat{\boldsymbol{\Sigma}}_j^{(i-1)})}{\sum_{k=1}^M \hat{\pi}_k^{(i-1)} \mathcal{N}(\mathbf{x}_n; \hat{\boldsymbol{\mu}}_k^{(i-1)}, \hat{\boldsymbol{\Sigma}}_k^{(i-1)})}.$$

M- (Maximization) step: Estimate the parameters using the current values of $\hat{p}(j | \mathbf{x}_n)$. Let $N_j = \sum_{n=1}^N \hat{\alpha}_{nj}^{(i)}$, for time step i :

$$\begin{aligned} \hat{\boldsymbol{\mu}}_j^{(i)} &= \frac{1}{N_j} \sum_{n=1}^N \hat{\alpha}_{nj}^{(i)} \cdot \mathbf{x}_n \\ \hat{\boldsymbol{\Sigma}}_j^{(i)} &= \frac{1}{N_j} \sum_{n=1}^N \hat{\alpha}_{nj}^{(i)} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_j^{(i)}) (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_j^{(i)})^T \\ \hat{\pi}_j^{(i)} &= \frac{N_j}{N}. \end{aligned}$$

4.2 EM

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 def get_data(name :str) -> np.ndarray:
5     file = open(name, 'r')
6     x1 = []
7     x2 = []
8     for i in file.readlines():
9         line = i.split()
10        x1.append(float(line[0]))
11        x2.append(float(line[1]))
12    file.close()
13    return np.row_stack((x1,x2))
14
15 def get_cov(X):
16     return np.matmul(X,np.transpose(X))/X.shape[1]
17
18 def multivar_gauss(x,m,S):
19     p = m.shape[0]
20     return (2*np.pi)**(-p/2) / np.sqrt(np.linalg.det(S)) * np.exp(-0.5*np.matmul(np.transpose(x-m), np.matmul(np.linalg.
21        inv(S),(x-m))))
22
23 def get_mixpdf(x1,x2,t):
24     X=np.array([x1,x2])
25     M = int(len(t) / 3)
26     pdf = 0
27     for j in range(M):
28         m = t[j*3]
29         S = t[j*3+1]
```

```

29     pi = t[j*3+2]
30     pdf += pi * multvar_gauss(X,m,S)
31     return pdf
32
33 def plot_overlay(x,t,title):
34     nbins = 50
35     data_density = plt.hist2d(x[0,:],x[1:],bins=nbins)
36     x1_edges = data_density[1]
37     x2_edges = data_density[2]
38     X1,X2 = np.meshgrid(x1_edges, x2_edges)
39     mixpdf = np.zeros((nbins+1,nbins+1))
40     for i in range(nbins+1):
41         for j in range(nbins+1):
42             mixpdf[i,j] = get_mixpdf(X1[i,j],X2[i,j],t)
43     plt.contour(X1, X2, mixpdf, cmap="Wistia")
44     plt.title(title)
45     plt.show()
46
47 def main():
48     #read data
49     x = get_data('gmm.txt')
50     p = x.shape[0]
51     N = x.shape[1]
52
53     #initialization
54     M = 4
55     pi = 1/M
56     all_mu = [None]*M
57     all_Sigma = [None]*M
58     all_pi = [pi]*M
59     np.random.seed=1234
60     for j in range(M):
61         all_mu[j] = np.random.rand(2)
62         all_Sigma[j] = get_cov(x-all_mu[j][:,np.newaxis])
63
64     alpha = np.zeros((M,N))
65     steps = 30
66     theta = [None]*(steps) #tracking all the values for each step {mu_1, Sigma_1, pi_1, ...}
67     L = np.zeros(steps) #log likelihood
68
69     for i in range(steps):
70         t = [None] * (3*M)
71
72         #E-Step
73         for n in range(N):
74             alpha_denom = 0
75             beta = np.empty(M)
76             for j in range(M):
77                 bn = multvar_gauss(x[:,n],all_mu[j],all_Sigma[j])
78                 beta[j] = all_pi[j]*bn
79                 alpha_denom += beta[j]
80
81             alpha[:,n] = beta/alpha_denom
82
83         #M-step
84         Nj = np.sum(alpha, axis=1)
85         for j in range(M):
86             all_mu[j] = np.sum(alpha[j,:]*x, axis=1) / Nj[j]
87             s = np.zeros((p,p))
88             for n in range(N):
89                 X = x[:,n]-all_mu[j]
90                 X = X[:,np.newaxis]
91                 s = s + alpha[j,n]*np.matmul(X,np.transpose(X))
92             all_Sigma[j] = s / Nj[j]
93         all_pi = Nj / N
94
95         t[0:3*M:3] = all_mu
96         t[1:3*M:3] = all_Sigma
97         t[2:3*M:3] = all_pi
98         theta[i] = t
99
100     #log likelihood
101     ll = 0
102     for n in range(N):
103         l = 0
104         for j in range(M):
105             l += all_pi[j] * multvar_gauss(x[:,n],all_mu[j],all_Sigma[j])
106         ll += np.log(l)

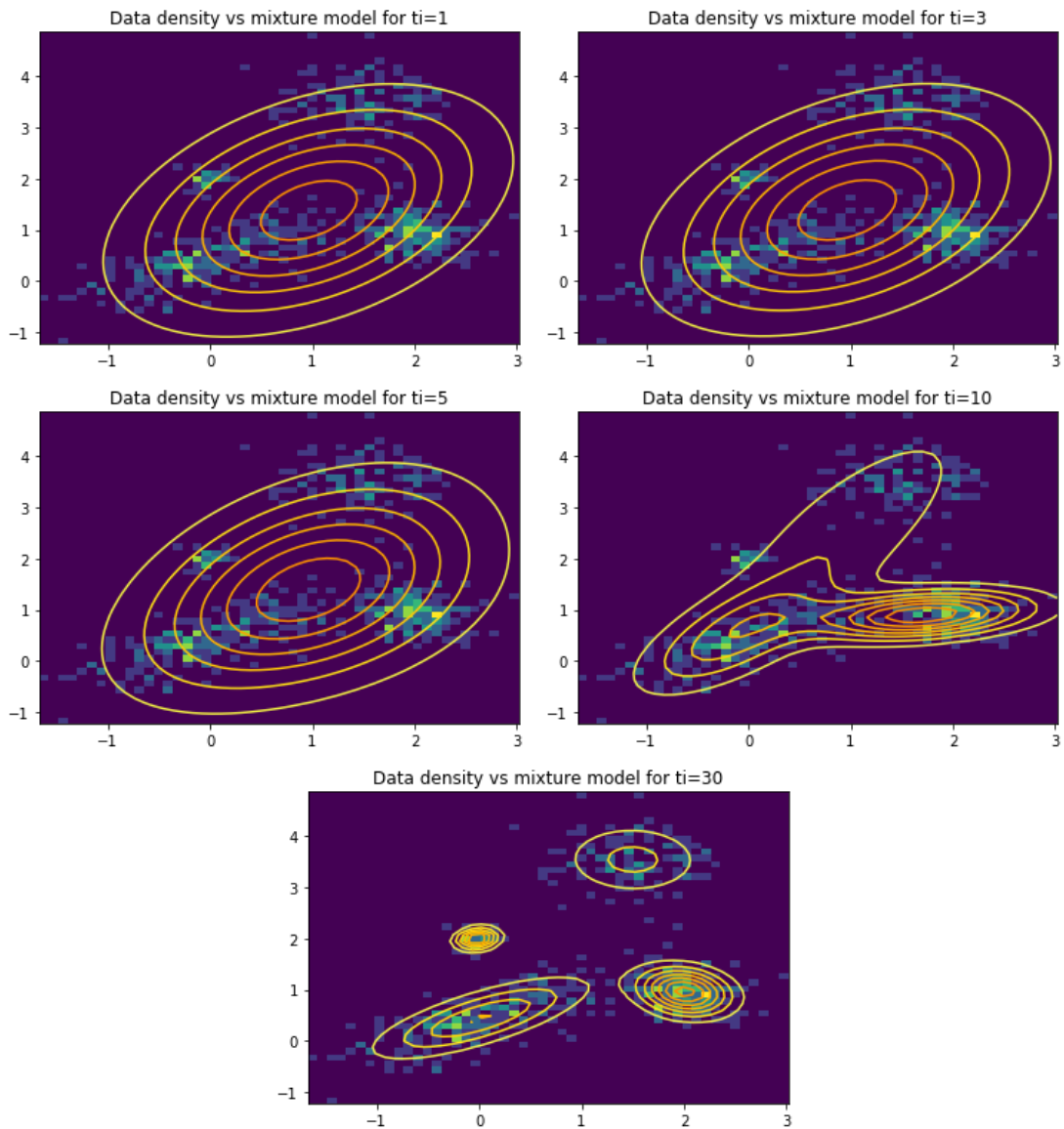
```

```

107     L[i] = ll
108
109     for i in [1,3,5,10,30]:
110         plot_overlay(x, theta[i-1], 'Data density vs mixture model for ti='+ str(i))
111
112     # plot log likelihood
113     plt.plot(L)
114     plt.ylabel("Log likelihood")
115     plt.xlabel("# iterations")
116     plt.title("Log likelihood updates")
117     plt.show()
118
119     # plot mixture scales
120     for j in range(M):
121         pi = [x[j*3+2] for x in theta]
122         plt.plot(np.arange(1,31), pi, label='pi'+ str(j+1))
123         plt.ylabel("pi")
124         plt.xlabel("# iterations")
125         plt.title("Updates of mixture scaling parameter")
126         plt.legend()
127         plt.show()
128
129 if __name__ == '__main__':
130     main()

```

Below we see the density of the observed data overlayed with the evolution of the mixture models at iterations $t_i = 1, 3, 5, 10, 30$.



The figure below shows the parameter updates for $t_i = 1, 3, 5, 10, 30$, where:

- indices 0,3,6,9 are $\mu_1, \mu_2, \mu_3, \mu_4$
- indices 1,4,7,10 are $\Sigma_1, \Sigma_2, \Sigma_3, \Sigma_4$
- indices 2,5,8,11 are $\pi_1, \pi_2, \pi_3, \pi_4$

After $t_i = 30$, $\pi_1, \pi_2, \pi_3, \pi_4$ are pretty stable at $\approx 0.31, 0.41, 0.06$, and 0.22 , respectively.

The figure displays five screenshots of parameter update windows for $t_i = 0, 2, 4, 9, 29$. Each window shows a table with 12 elements, indexed 0 to 11. The columns are Index, Typ, Größe, and Wert. The data represents parameter updates for four different components (μ, Σ, π) over time.

Index	Typ	Größe	Wert
0	float64	(2,)	[1.07750828 1.37224161]
1	float64	(2, 2)	[[0.9970463 0.42585596] [0.42585596 1.53723873]]
2	float64	1	0.263818898962714
3	float64	(2,)	[0.8328047 1.35941974]
4	float64	(2, 2)	[[1.09535564 0.55890766] [0.55890766 1.55728276]]
5	float64	1	0.26780671911158804
6	float64	(2,)	[0.94574193 1.38534719]
7	float64	(2, 2)	[[1.07842629 0.59356327] [0.59356327 1.70383808]]
8	float64	1	0.22609512801765128
9	float64	(2,)	[1.00313248 1.38201828]
10	float64	(2, 2)	[[1.04030247 0.51856691] [0.51856691 1.63560269]]
11	float64	1	0.2422792539744892

Index	Typ	Größe	Wert
0	float64	(2,)	[1.14888739 1.36394889]
1	float64	(2, 2)	[[0.98690898 0.34160726] [0.34160726 1.42389678]]
2	float64	1	0.2648358402995722
3	float64	(2,)	[0.77155899 1.28853069]
4	float64	(2, 2)	[[1.0981435 0.65614114] [0.65614114 1.62195117]]
5	float64	1	0.26820658652185114
6	float64	(2,)	[0.91905886 1.44072999]
7	float64	(2, 2)	[[1.04823075 0.59814205] [0.59814205 1.75515276]]
8	float64	1	0.22558604614850508
9	float64	(2,)	[1.01765507 1.4182123]
10	float64	(2, 2)	[[1.02832499 0.48538282] [0.48538282 1.62751828]]
11	float64	1	0.2413715270300716

Index	Typ	Größe	Wert
0	float64	(2,)	[1.26782116 1.27693154]
1	float64	(2, 2)	[[0.93479394 0.22161798] [0.22161798 1.16738919]]
2	float64	1	0.2678863979900569
3	float64	(2,)	[0.63517582 1.19505479]
4	float64	(2, 2)	[[1.09166105 0.72456739] [0.72456739 1.57552591]]
5	float64	1	0.26948807252030654
6	float64	(2,)	[0.90223136 1.58694835]
7	float64	(2, 2)	[[0.99788753 0.6700149] [0.6700149 1.96552881]]
8	float64	1	0.22472076551441142
9	float64	(2,)	[1.05340092 1.48544218]
10	float64	(2, 2)	[[0.98875459 0.46148113] [0.46148113 1.68574588]]
11	float64	1	0.2379047639752251

Index	Typ	Größe	Wert
0	float64	(2,)	[1.71115618 0.89754119]
1	float64	(2, 2)	[[0.45242787 0.05697317] [0.05697317 0.0886724]]
2	float64	1	0.38117663258626133
3	float64	(2,)	[-0.16207113 0.48172349]
4	float64	(2, 2)	[[0.35572482 0.28187744] [0.28187744 0.49544923]]
5	float64	1	0.26410627781455054
6	float64	(2,)	[0.93825421 2.60847134]
7	float64	(2, 2)	[[0.70413888 0.92942549] [0.92942549 1.87218618]]
8	float64	1	0.2169392531558372
9	float64	(2,)	[1.09720589 2.45984316]
10	float64	(2, 2)	[[0.66248457 0.59415525] [0.59415525 1.76881956]]
11	float64	1	0.13777783644335084

Index	Typ	Größe	Wert
0	float64	(2,)	[2.00779373 0.96110245]
1	float64	(2, 2)	[[0.08197349 -0.01569632] [-0.01569632 0.08947233]]
2	float64	1	0.31613479970550007
3	float64	(2,)	[0.01087482 0.4471223]
4	float64	(2, 2)	[[0.39324324 0.22017656] [0.22017656 0.22401418]]
5	float64	1	0.40653533555695787
6	float64	(2,)	[-0.01882301 1.99346402]
7	float64	(2, 2)	[[0.01774844 0.00299987] [0.00299987 0.019188]]
8	float64	1	0.05992940192274512
9	float64	(2,)	[1.50016772 3.53754339]
10	float64	(2, 2)	[[0.18275512 -0.00446836] [-0.00446836 0.18429649]]
11	float64	1	0.2174004628147969

The next figures show the evolution of the the mixture scaling parameters π_j and the log likelihood function.

