

SML: Exercise 2

Rinor Cakaj, Patrick Nowak



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Summer Term 2020
Sheet 1

Task 1.1: Density Estimation

We are given data C1 and C2, which we suppose to be generated by 2D-Gaussians with parameters μ_1, Σ_1 and μ_2, Σ_2 , respectively.

1.1a)

Assume we are given iid. datapoints $x_i, i = 1, \dots, n$ which are generated by a 2D-Gaussian. Following the max-likelihood principle, we maximize the log-likelihood function

$$l(\mu, \Sigma, x_1, \dots, x_n) = \ln\left(\prod_{i=1}^n p(x_i|\mu, \Sigma)\right) = \sum_{i=1}^n \ln(p(x_i|\mu, \Sigma))$$

for the Gaussian probability density

$$p(x|\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma (x - \mu)\right). \quad (1)$$

We receive

$$l(\mu, \Sigma) := l(\mu, \Sigma, x_1, \dots, x_n) = \sum_{i=1}^n \left(-\frac{k}{2} \ln(2\pi) - \frac{1}{2} \ln(|\Sigma|) - \frac{1}{2} (x_i - \mu)^T \Sigma (x_i - \mu) \right) \quad (2)$$

$$= -\frac{nk}{2} \ln(2\pi) - \frac{n}{2} \ln(|\Sigma|) - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma (x_i - \mu). \quad (3)$$

We compute the derivatives w.r.t. μ and Σ and set them equal to zero. This yields

$$\begin{aligned} \frac{d}{d\mu} l(\mu, \Sigma, x_1, \dots, x_n) &= \frac{d}{d\mu} \left(-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma (x_i - \mu) \right) \\ &= -\sum_{i=1}^n \frac{d}{d\mu} \frac{1}{2} (x_i - \mu)^T \Sigma (x_i - \mu). \end{aligned}$$

Using the matrix identity $\frac{d}{dw} \frac{w^T A w}{dw} = 2Aw$ which holds if w does not depend on A and if A is symmetric, we get (with $w = (x - \mu)$, $dw = -d\mu$)

$$\begin{aligned} 0 &\stackrel{!}{=} \frac{d}{d\mu} l(\mu, \Sigma, x_1, \dots, x_n) \\ 0 &\stackrel{!}{=} -\sum_{i=1}^n \Sigma^{-1} (x_i - \mu). \end{aligned}$$

Finally, we use that Σ^{-1} is positive definite, so we can leave it out here and get

$$0 \stackrel{!}{=} n\mu - \sum_{i=1}^n x_i ,$$

which is solved for the MLE-estimate

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i . \quad (4)$$

Secondly, we need to compute the derivative w.r.t Σ . To do that, we will need some results from mathematical classes. The following is used without prove:

- Cyclic permutations of a matrix product do not change the trace of it:

$$\text{tr} [ABC] = \text{tr} [CAB]$$

- The trace of a scalar is the scalar itself. In particular: the result of a quadratic form $x^T A x$ is a scalar, such that:

$$x^T A x = \text{tr} [x^T A x] = \text{tr} [x^T x A]$$

- $\frac{d}{dA} \text{tr} [AB] = B^T$
- $\frac{d}{dA} \ln |A| = A^{-T}$

As a first result of these assumptions, we can show, that

$$\frac{d}{dA} x^T A x = \frac{d}{dA} \text{tr} [x^T x A] = [x x^T]^T = x x^T .$$

We now got the tools to re-write the log-likelihood function in (3) to

$$\begin{aligned} l(\mu, \Sigma) &= -\frac{nk}{2} \ln(2\pi) - \frac{n}{2} \ln(|\Sigma|) - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma (x_i - \mu) \\ &= C + \frac{n}{2} \ln(|\Sigma^{-1}|) - \frac{1}{2} \sum_{i=1}^n \text{tr} [(x_i - \mu)(x_i - \mu)^T \Sigma^{-1}] \end{aligned}$$

for a constant C, and taking the derivative w.r.t Σ^{-1} yields

$$\frac{d}{d\Sigma^{-1}} l(\mu, \Sigma) = \frac{n}{2} \Sigma^T - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$$

and plugging in $\hat{\mu}$ as an estimation of μ and setting equal to zero finally gives us

$$\begin{aligned} 0 &\stackrel{!}{=} \frac{d}{d\Sigma^{-1}} l(\hat{\mu}, \Sigma) \\ 0 &\stackrel{!}{=} \frac{n}{2} \Sigma^T - \frac{1}{2} \sum_{i=1}^n (x_i - \hat{\mu})(x_i - \hat{\mu})^T \end{aligned}$$

which is solved for the (biased) MLE estimate

$$\tilde{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})(x_i - \hat{\mu})^T \quad (5)$$

1.1b)

We compute the prior probabilities of C1 and C2, using the following python code. We read the number of data points in each class and divide it by the sum of total data points in both classes.

```
import numpy as np

link1=" ../hw2/dataSets/densEst1.txt "
link2=" ../hw2/dataSets/densEst2.txt "

def get_lengths():
    l1=0;
    l2=0;
    for line in open(link1):
        l1=l1+1
    for line2 in open(link2):
        l2=l2+1
    return (l1,l2)

def get_priors(l1,l2):
    p_C1=l1/(l1+l2)
    p_C2=l2/(l1+l2)
    return(p_C1,p_C2)
```

Calling

```
lengths=get_lengths()
print(get_priors(lengths[0],lengths[1]))
```

we get the following results for the prior probabilities: $p(C1)=0.239$ and $p(C2)=0.761$.

1.1c)

Having a data set X and an estimator $\hat{\theta}$ on the true parameter θ , we define the bias of an estimator as the expected deviation from the true parameter. We get the formula

$$bias(\hat{\theta}) = \mathbb{E}_X [\hat{\theta}(X) - \theta]$$

We call an estimator unbiased iff $bias(\hat{\theta}) = 0$ and biased otherwise.

From the lecture we know that the MLE of the mean of a Gaussian is unbiased, but the MLE of the variance of a Gaussian is biased. In fact, an unbiased estimator on the variance would be the sample covariance matrix

$$\hat{\Sigma} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu})(x_i - \hat{\mu})^T. \quad (6)$$

To calculate the conditional distribution densities $p(x|C_i)$ we need to estimate the underlying parameters μ_i and Σ_i . For both classes we use the MLE-estimate of the mean, which is unbiased and calculated like in (4). We compute the biased MLE-estimate $\tilde{\Sigma}$ for the variance via (5) and the unbiased estimate $\hat{\Sigma}$ via (6). We wrote the following python code

```
def extract_data(t):
    a=np.empty((0,2),float)
    for line in t:
        v=np.array([[line.split()[0],line.split()[1]]],float)
        a=np.append(a,v,axis=0)
    return a
```

```

def get_mean_estimation(points):
    m=np.zeros(2)
    m[0]=sum(points[:,0])
    m[1]=sum(points[:,1])
    m=m/len(points)
    return m

def get_biased_var_estimation(mean_est, points):
    l=points-mean_est
    s=np.zeros((2,2))
    for line in l:
        s=s+np.outer(line, line)
    s=s/len(l)
    return s

def get_unbiased_var_estimation(mean_est, points):
    l=points-mean_est
    s=np.zeros((2,2))
    for line in l:
        s=s+np.outer(line, line)
    s=s/(len(l)-1)
    return s

def print_results(link):
    a=extract_data(open(link))
    mean=get_mean_estimation(a)
    print("mean=", mean)
    sigma=get_unbiased_var_estimation(mean, a)
    print("biased_Sigma=", get_biased_var_estimation(mean, a))
    print("sigma=", sigma)

```

which gives us the results for class C1

```

print_results(link1)

mean= [-0.70681374 -0.81343083]
biased_Sigma= [[9.01952586 2.67287085]
               [2.67287085 3.59633965]]
sigma= [[9.05742302 2.6841014 ]
        [2.6841014  3.61145033]]

```

and for class C2

```

print_results(link2)

mean= [3.98534252 3.98438364]
biased_Sigma= [[4.1753815  0.02758324]
               [0.02758324 2.75296323]]
sigma= [[4.18087542 0.02761954]
        [0.02761954 2.75658555]]

```

Transferred to the notation we introduced earlier, we get $\hat{\mu}_1 = \begin{pmatrix} -0.71 \\ -0.81 \end{pmatrix}$ with $\tilde{\Sigma}_1 = \begin{pmatrix} 9.02 & 2.67 \\ 2.67 & 3.6 \end{pmatrix}$ or with the unbiased $\hat{\Sigma}_1 = \begin{pmatrix} 9.06 & 2.68 \\ 2.68 & 3.61 \end{pmatrix}$ for C1.

And $\hat{\mu}_2 = \begin{pmatrix} 3.99 \\ 3.98 \end{pmatrix}$ with $\tilde{\Sigma}_2 = \begin{pmatrix} 4.18 & 0.03 \\ 0.03 & 2.75 \end{pmatrix}$ or with the unbiased $\hat{\Sigma}_2 = \begin{pmatrix} 4.18 & 0.03 \\ 0.03 & 2.76 \end{pmatrix}$ for C2.

We can now just plug in into formula (1) to get

for the biased estimate $p(x|C_i) = p(x|\hat{\mu}_i, \tilde{\Sigma}_i)$

and for the unbiased $p(x|C_i) = p(x|\hat{\mu}_i, \hat{\Sigma}_i)$, where the right hand sides are given like in (1).

1.1d)

Using the unbiased estimates $p(x|C_i) = p(x|\hat{\mu}_i, \hat{\Sigma}_i)$ from last task, for each class we plot the Gaussian in a single graph with the data points: Therefore we added 2 new functions to our code:

```
def gaussian(x,mu,detsigma , sigmainv):
    z=np.array(x-mu)
    enum=np.exp(-0.5*np.dot(np.dot(z , sigmainv),np.transpose(z)))
    denom=np.sqrt(np.power(2*np.pi , len(x))*detsigma)
    return enum/denom

def plot(link):
    a=extract_data(open(link))
    mean=get_mean_estimation(a)
    print("mean=",mean)
    sigma=get_unbiased_var_estimation(mean,a)
    print("biased_Sigma=",get_biased_var_estimation(mean,a))
    print("sigma=",sigma)
    detsigma=np.linalg.det(sigma)
    sigmainv=np.linalg.inv(sigma)
    x=np.linspace(min(a[:,0]),max(a[:,0]),300)
    y=np.linspace(min(a[:,1]),max(a[:,1]),300)
    X,Y=np.meshgrid(x,y)
    Z=np.empty_like(X)
    i=0
    while i<len(X):
        j=0
        while j<len(Y):
            xy=np.array([X[i , j],Y[i , j]])
            ergb=gaussian(xy,mean,detsigma , sigmainv)
            Z[i , j]=ergb
            j=j+1
        i=i+1
    plt.contourf(X,Y,Z,25)
    plt.colorbar()
    plt.scatter(a[:,0],a[:,1],alpha=1,c="white",s=0.8)
    return
```

Calling "plot(link1)" gives us Figure 1 and "plot(link2)" gives Figure 2.

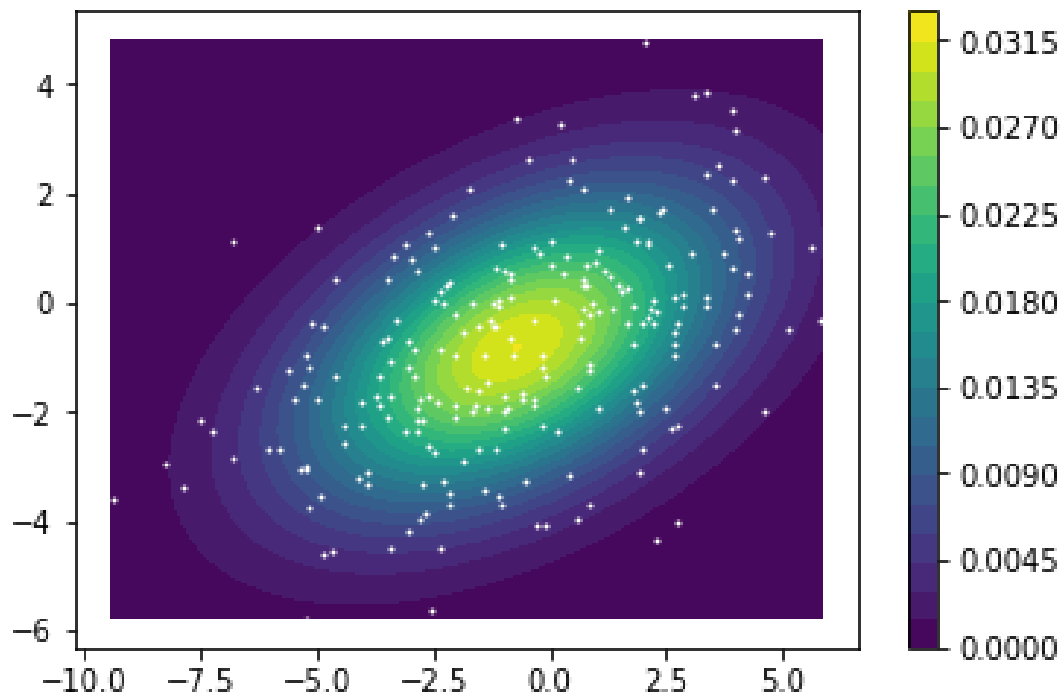


Figure 1: Gaussian ($\hat{\mu}_1, \hat{\Sigma}_1$) and data points C1 in white

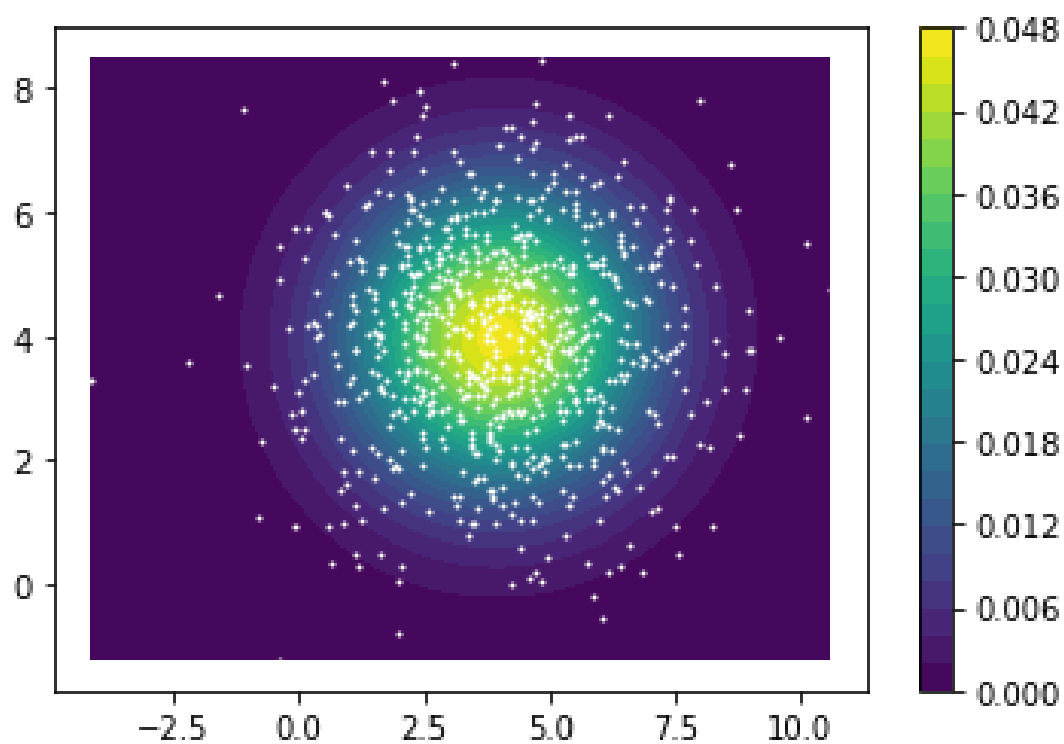


Figure 2: Gaussian ($\hat{\mu}_2, \hat{\Sigma}_2$) and data points C2 in white

1.1e)

We are interested in $p(C_i|x)$. From Bayes' rule we know that this equals

$$p(C_i|x) = \frac{p(x|C_i)p(C_i)}{\sum_j p(x|C_j)p(C_j)}$$

which helps us a lot because we calculated all terms occuring in the right hand side. We add two more functions to our code:

```
def postplot(link):
    g=extract_data(open(link3))
    a=extract_data(open(link))
    prior=len(a)/len(g)
    mean=get_mean_estimation(a)
    sigma=get_unbiased_var_estimation(mean,a)
    detsigma=np.linalg.det(sigma)
    sigmainv=np.linalg.inv(sigma)
    x=np.linspace(min(g[:,0]),max(g[:,0]),300)
    y=np.linspace(min(g[:,1]),max(g[:,1]),300)
    X,Y=np.meshgrid(x,y,sparse=False)
    Z=np.empty_like(X)
    i=0
    while i<len(X):
        j=0
        while j<len(Y):
            xy=np.array([X[i,j],Y[i,j]])
            ergb=prior*gaussian(xy,mean,detsigma,sigmainv)
            Z[i,j]=ergb
            j=j+1
        i=i+1
    return [X,Y,Z,prior]

def actual_plotting():
    X,Y,Z1,p1=postplot(link1)
    Z2,p2=postplot(link2)[2:4]
    S=np.add(Z1,Z2)
    Z1neu=np.divide(Z1,S)
    Z2neu=np.divide(Z2,S)
    plt.contour(X,Y,Z1,10)
    plt.contour(X,Y,Z2,10)
    plt.title("likelihood_x_prior")
    plt.figure()
    plt.contour(X,Y,Z1neu,10)
    plt.colorbar()
    plt.title("p(C1|x)")
    plt.figure()
    plt.contour(X,Y,Z2neu,10)
    plt.colorbar()
    plt.title("p(C2|x)")
    plt.figure()
    dec=np.greater(Z1neu,Z2neu)
    plt.scatter(X,Y,dec)
    plt.title("blue=decideC1 ,_white=decideC2")
    j=0
    mpoints=np.empty(shape=[0,2])
```

```

while j<len(X):
    i=0
    while i<len(Y):
        if not dec[i,j]:
            xy=np.array([[X[i,j],Y[i,j]]])
            mpoints=np.append(mpoints,xy,axis=0)
            break
        i=i+1
    j=j+1
plt.figure()
plt.scatter(mpoints[:,0],mpoints[:,1],s=0.5)
plt.title("decision_boundary")
plt.figure()
plt.contour(X,Y,Z1,10)
plt.contour(X,Y,Z2,10)
plt.scatter(mpoints[:,0],mpoints[:,1],s=0.5)
plt.title("likelihood_x_prior_and_decbound")

return

```

Where in link3 we stored to a .txt file containing all data from C1 and C2 combined. This can be generated by executing

```

def merge_txt(inp):
    with open('../hw2/dataSets/densEstCombined.txt', 'w') as outfile:
        for fname in inp:
            with open(fname) as infile:
                for line in infile:
                    outfile.write(line)

return

```

once, or just do it by hand. Calling actual_plotting() then gives us the following results:

