

Statistical Machine Learning 2018

Assignment 2

Deadline: 28th of October 2018

Christoph Schmidl

s4226887

c.schmidl@student.ru.nl

Mark Beijer

s4354834

mbeijer@science.ru.nl

November 2, 2018

Exercise 1 - weight 3

The financial services department of an insurance company receives numerous phone calls each day from people who want to make a claim against their policy. Most claims are genuine, however about 1 out of every 6 are thought to be fraudulent. To tackle this problem the company has installed a trial version of a software voice-analysis system that monitors each conversation and gives a numerical score z between 0 and 1, depending on allegedly suspicious vocal intonations of the customer. Unfortunately, nobody seems to know anymore how to interpret the score in this particular version of the system...

Tests revealed that the conditional probability density of z , given that a claim was valid ($c = 1$) or false ($c = 0$) are

$$\begin{aligned} p(z|c=0) &= \alpha_0(1-z^2), \\ p(z|c=1) &= \alpha_1 z(z+1). \end{aligned}$$

1.1

Compute the normalization constants α_0 and α_1 . How does the z score relate to the validity of the claim? What values for z would you expect when the claim is valid / false?

Answer:

We integrate over all possible values of z , which has to equal 1, for it is sure that give c some value of z will be chosen.

$$1 = \int_0^1 \alpha_0(1-z^2)dz \tag{1}$$

$$= \alpha_0 z - \frac{\alpha_0}{3} z^3 \Big|_0^1 \tag{2}$$

$$= \left(\alpha_0 - \frac{\alpha_0}{3} \right) = \frac{2}{3} \alpha_0 \tag{3}$$

$$\alpha_0 = \frac{3}{2} \tag{4}$$

And for the next constant:

$$1 = \int_0^1 \alpha_1 z(z+1) dz \quad (5)$$

$$= \left. \frac{\alpha_1}{3} z^3 + \frac{\alpha_1}{2} z^2 \right|_0^1 \quad (6)$$

$$= \frac{1}{3} \alpha_1 + \frac{1}{2} \alpha_1 = \frac{5}{6} \alpha_1 \quad (7)$$

$$\alpha_1 = \frac{6}{5} \quad (8)$$

From the plot at figure 1 we see that for higher values of z the probability $p(z|c=1)$ is higher. Therefore we think that higher values of z correspond with a higher chance that the claim is true/not fraudulent.

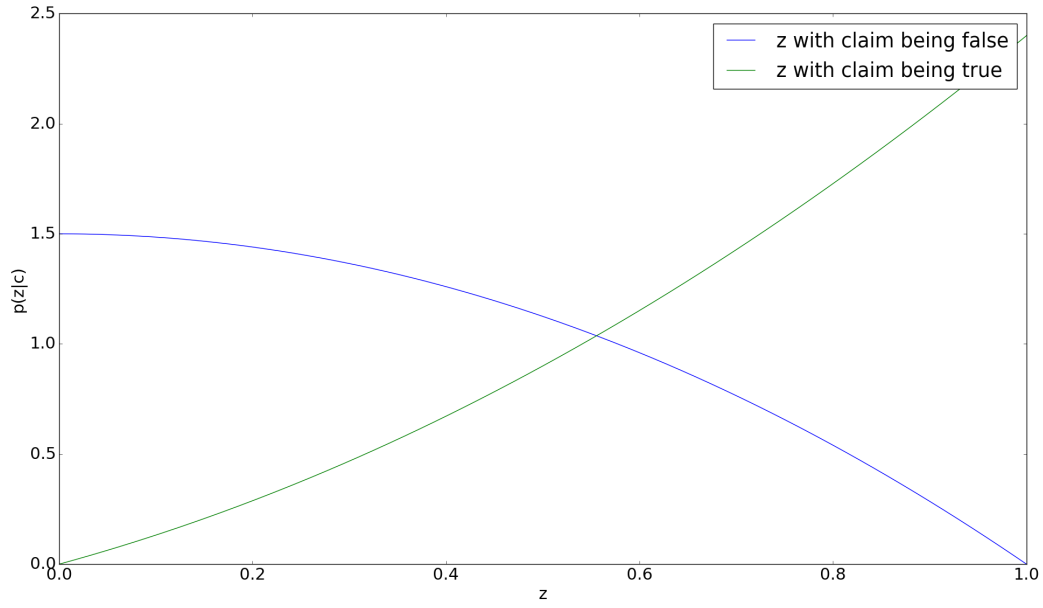


Figure 1: The probability $p(z|c)$ plotted against z , with $0 < z < 1$.

1.2

Use the sum and product rule to show that the probability distribution function $p(z)$ can be written as

$$p(z) = \frac{(3z+1)(z+1)}{4} \quad (9)$$

Answer:

For this we need to sum over the possible values of c and add the probability of c times the conditional probability over c .

$$p(z) = p(c=0)p(z|c=0) + p(c=1)p(z|c=1) \quad (10)$$

$$= \frac{1}{6} \frac{3}{2} (1 - z^2) + \frac{5}{6} \frac{6}{5} z(z+1) \quad (11)$$

$$= \frac{3}{12} (1 - z^2) + \frac{30}{30} z(z+1) \quad (12)$$

$$= \frac{1}{4} (1 - z^2) + z(z+1) \quad (13)$$

$$= \frac{1}{4} (1 - z^2) + \frac{4z(z+1)}{4} \quad (14)$$

$$= \frac{4z^2 + 4z - z^2 + 1}{4} \quad (15)$$

$$= \frac{3z^2 + 4z + 1}{4} \quad (16)$$

$$= \frac{(3z+1)(z+1)}{4} \quad (17)$$

1.3

Use Bayes' rule to compute the posterior probability distribution function $p(c|z)$. Plot these distributions in MATLAB as a function of z . How can these posterior probabilities help in making a decision regarding the validity of the claim?

Answer:

Bayes rule states that:

$$p(c|z) = \frac{p(z|c)p(c)}{p(z)} \quad (18)$$

So for both values of c we get:

$$p(c=0|z) = \frac{\frac{2}{3}(1-z^2)\frac{1}{6} \cdot 4}{(3z+1)(z+1)} \quad (19)$$

$$= \frac{8}{18} \frac{1-z^2}{(3z+1)(z+1)} \quad (20)$$

$$= \frac{4}{9} \frac{1-z^2}{(3z+1)(z+1)} \quad (21)$$

$$= \frac{4}{9} \frac{(1+z)(1-z)}{(3z+1)(z+1)} \quad (22)$$

$$= \frac{4}{9} \frac{1-z}{3z+1} \quad (23)$$

$$p(c=1|z) = \frac{\frac{6}{5}z(z+1)\frac{5}{6} \cdot 4}{(3z+1)(z+1)} \quad (24)$$

$$= \frac{4z(z+1)}{(3z+1)(z+1)} \quad (25)$$

$$= \frac{4z}{3z+1} \quad (26)$$

The plot of the distribution can be seen in figure 2.

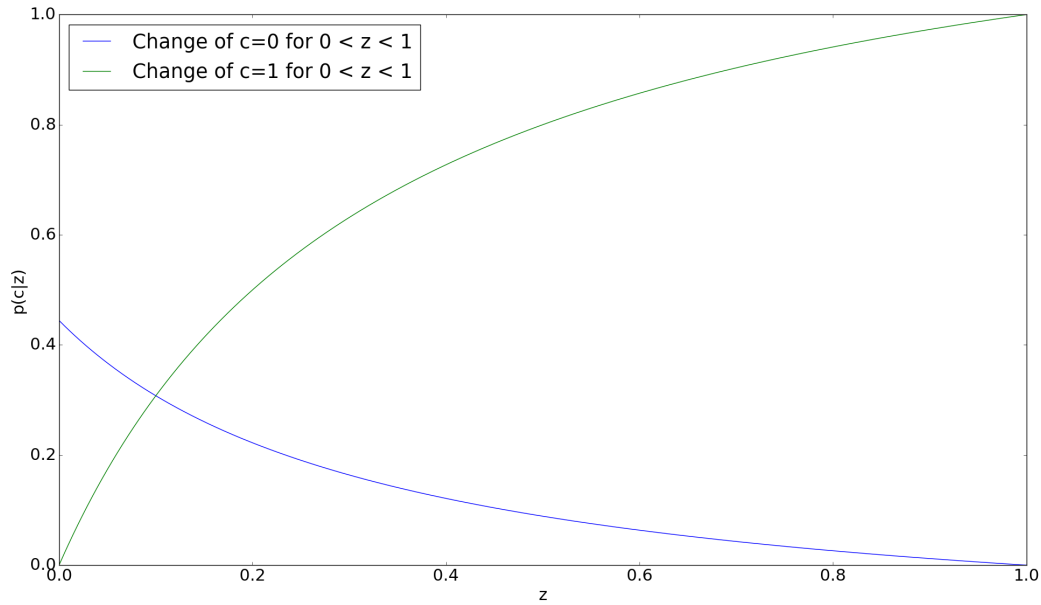


Figure 2: The probability $p(c|z)$ plotted against z , with $0 < z < 1$.

1.4

Compute the optimal decision boundary (based on our numerical score z) that minimizes the misclassification rate. For which z should we classify $c = 0$ (false) and for which z should we classify $c = 1$ (valid)? Explain your decision.

Answer:

We take the decision boundary at the point where the two functions, as plotted at Figure 2, intersect. This happens at the point where they equal:

$$\frac{4}{9} \frac{1-z}{3z+1} = \frac{4z}{3z+1} \quad (27)$$

$$\frac{4}{9}(1-z) = 4z \quad (28)$$

$$1-z = 9z \quad (29)$$

$$1 = 10z \quad (30)$$

$$z = 0.1 \quad (31)$$

We have also found the point 0.1 using our program.

1.5

Compute the misclassification rate, given the optimal decision boundary determined previously. Interpret the result you have obtained. Is the z score useful in determining the validity of the claim? Compare this with your prior guess from 1.

Answer:

We integrate the probabilities up to the decision point and from the decision point to 1. Then we compare the change of actually getting the right label vs us putting the decision boundary at 0.1 .

If z is less than 0.1 we classify it as $c=0$. The actual change is:

$$p(z < 0.1)p(c = 0|z < 0.1) = \int_0^{0.1} \frac{4}{9} \frac{1-z}{3z+1} dz \quad (32)$$

$$\frac{1-z}{3z+1} = \frac{4}{3(3z+1)} - \frac{1}{3} \quad (33)$$

$$p(z < 0.1)p(c = 0|z < 0.1) = \frac{4}{9} \int_0^{0.1} \left(\frac{4}{3} \frac{1}{3z+1} - \frac{1}{3} \right) dz \quad (34)$$

$$= \frac{4}{9} \left(\frac{4}{9} \ln(3z+1) - \frac{z}{3} \right) \Big|_0^{0.1} \quad (35)$$

$$\approx 0.0370 \quad (36)$$

Now for the other possibility:

$$p(z < 0.1)p(c = 1|z < 0.1) = \int_0^{0.1} \frac{4z}{3z+1} dz \quad (37)$$

$$\frac{4z}{3z+1} = \frac{4}{3} - \frac{4}{3} \frac{1}{3z+1} \quad (38)$$

$$= \frac{4}{3} \left(1 - \frac{1}{3z+1} \right) \quad (39)$$

$$p(z < 0.1)p(c = 1|z < 0.1) = \frac{4}{3} \int_0^{0.1} \left(1 - \frac{1}{3z+1} \right) dz \quad (40)$$

$$= \frac{4}{3} \left(z - \frac{1}{3} \ln(3z+1) \right) \Big|_0^{0.1} \quad (41)$$

$$\approx 0.0167 \quad (42)$$

For both $\ln(1) = 0$, so filling in $z = 0$ gives us 0.

Now the fraction of the two:

$$\text{Classification error} = \frac{p(c = 1|z < 0.1)}{p(c = 0|z < 0.1)} \approx 45.20\% \quad (43)$$

Now we do the same for $z > 0.1$. Here we already calculated the indefinite integrals, so we only have to fill in the boundaries:

$$p(z > 0.1)p(c = 0|z > 0.1) = \int_{0.1}^1 \frac{4}{9} \frac{1-z}{3z+1} dz \quad (44)$$

$$= \frac{4}{9} \left(\frac{4}{9} \ln(3z+1) - \frac{z}{3} \right) \Big|_{0.1}^1 \quad (45)$$

$$\approx 0.0887 \quad (46)$$

And for $c=1$:

$$p(z > 0.1)p(c = 1|z > 0.1) = \int_{0.1}^1 \frac{4z}{3z+1} dz \quad (47)$$

$$= \frac{4}{3} \left(\frac{1}{3} \ln(3z+1) - z \right) \Big|_{0.1}^1 \quad (48)$$

$$\approx 0.0467 \quad (49)$$

$$\text{Classification error} = \frac{p(c = 0|z > 0.1)}{p(c = 1|z > 0.1)} \approx 18.98\% \quad (50)$$

Now the total classification error is the sum of the two errors, weighted by the probability of z lying in these regions. So let's integrate z :

$$p(z) = \frac{(3z+1)(z+1)}{4} \quad (51)$$

$$= \frac{1}{4}(3z^2 + 4z + 1) \quad (52)$$

$$\int p(z)dz = \frac{1}{4}(z^3 + 2z^2 + z + C) \quad (53)$$

C is the integration constant independent of z .

Now we calculate the integrals:

$$p(z < 0.1) = \int_0^{0.1} \frac{(3z+1)(z+1)}{4} dz \approx 0.03025 \quad (54)$$

$$p(z > 0.1) = \int_{0.1}^1 \frac{(3z+1)(z+1)}{4} dz \approx 0.96975 \quad (55)$$

So the total error is:

$$\text{Misclassification error} = 0.03025 \cdot 45.20\% + 0.96975 \cdot 18.98\% \approx 19.77\% \quad (56)$$

Exercise 2 - weight 2

The government of the United Kingdom has decided to call a referendum regarding the country's European Union membership. The citizens of the UK will be asked the following question at the referendum: "Should the United Kingdom remain a member of the European Union or leave the European Union?". The European Commission (EC) is interested in the potential outcome of this referendum and has contracted a polling agency to study this issue.

Suppose that a person's vote follows a Bernoulli distribution with parameter θ and suppose the EC's prior distribution for θ , the proportion of British citizens that would be in favor of leaving the EU, is beta distribution with $\alpha = 90$ and $\beta = 110$.

2.1

Determine the mean and variance of the prior distribution. Plot the prior density function.

Answer:

The beta distribution is defined as:

$$\text{Beta}(\mu|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \mu^{\alpha-1} (1 - \mu)^{\beta-1} \quad (57)$$

*Not sure if this is needed, but we can do it some time regardless.*****

The mean is calculated as following:

$$\mathbb{E}[\mu] = \int_0^1 \text{Beta}(\mu|\alpha, \beta) \mu d\mu \quad (58)$$

$$= \int_0^1 \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \mu^{\alpha} (1 - \mu)^{\beta-1} d\mu \quad (59)$$

$$= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int_0^1 \mu^{\alpha} (1 - \mu)^{\beta-1} d\mu \quad (60)$$

***** The mean and variance are therefore:

$$\mathbb{E}[\mu] = \frac{\alpha}{\alpha + \beta} = 4.5 \cdot 10^{-1} \quad (61)$$

$$\text{var}[\mu] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)} \approx 1.23 \cdot 10^{-3} \quad (62)$$

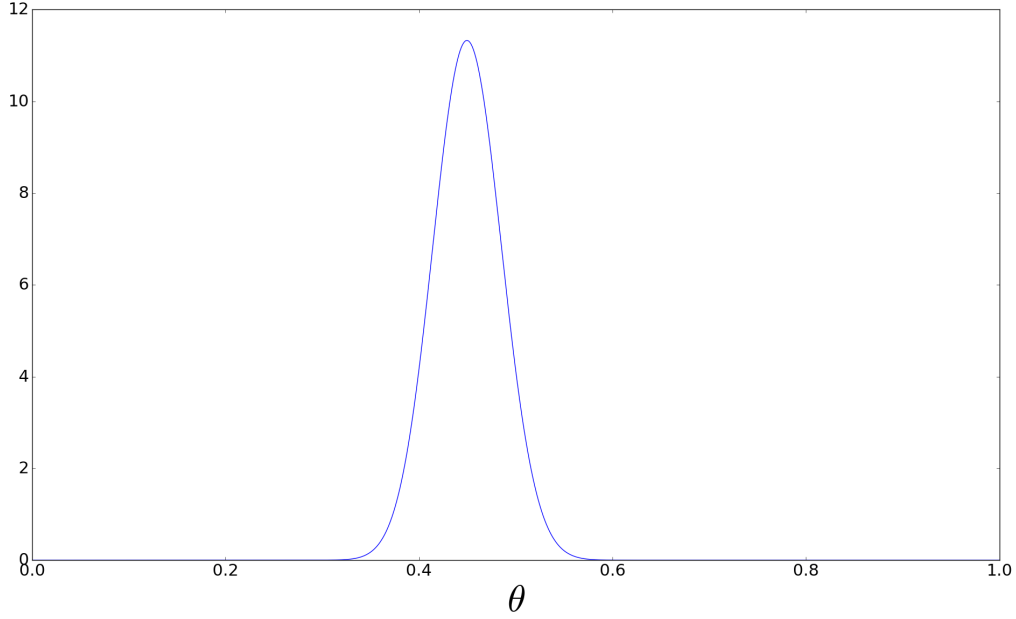


Figure 3: Plot of the beta distribution for $\alpha = 90$ and $\beta = 110$.

2.2

A random sample of 1000 British citizens is taken, and 60% of the people polled support leaving the European Union. What are the posterior mean and variance for θ ? Plot the posterior density function together with the prior density. Explain how the data from the sample changed the prior belief.

Answer:

The posterior probability is described as:

$$p(\theta|m, l, \alpha, \beta) = \frac{\Gamma(m + \alpha + l + \beta)}{\Gamma(m + \alpha)\Gamma(l + \beta)} \mu^{m+\alpha-1} (1 - \mu)^{l+\beta-1} \quad (63)$$

So the value of α is increased by m and the value of β is increased by l . The value of m is the amount of people in support, the value of l is how many people will vote in favour, while l is how many people will vote against. So the new mean is:

$$\mathbb{E}[\theta] = \frac{\alpha + m}{\alpha + m + \beta + l} = 5.75 \cdot 10^{-1} \quad (64)$$

$$\text{var}[\theta] = \frac{(\alpha + m)(\beta + l)}{(\alpha + \beta + l + m)^2(\alpha + \beta + m + l + 1)} = 2.03 \cdot 10^{-4} \quad (65)$$

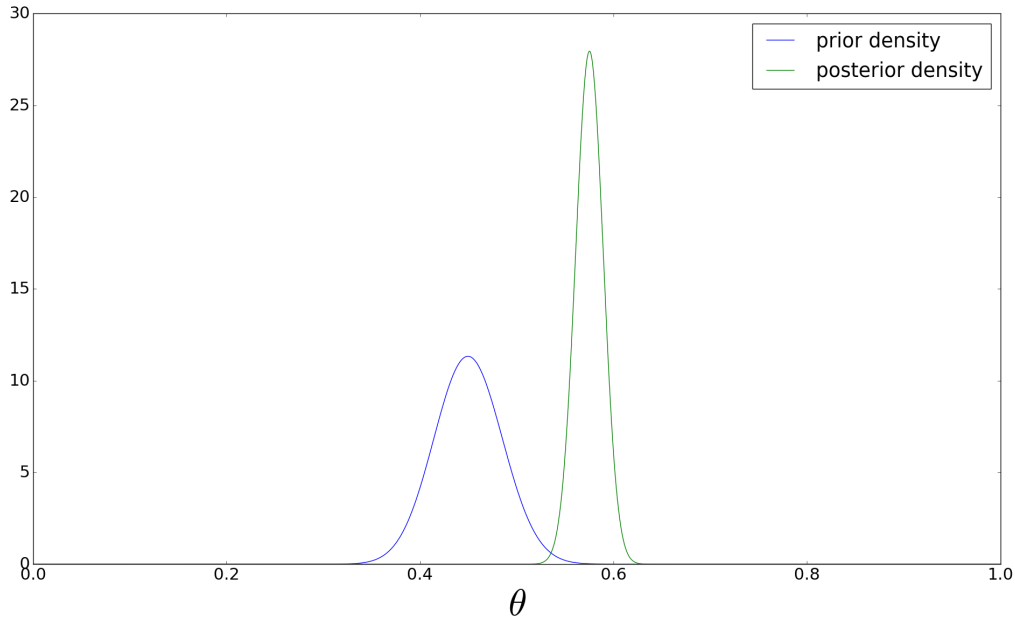


Figure 4: Plot of the prior and posterior probability distribution.

As you can see in the plot, the expected value (mean) of the plot is further to the right. For the new information showed more people supporting the brexit. You can also see that the variance has decreased, since we're more sure about our guess because we have more data.

2.3

Examine the effect of changing the prior hyperparameters (α, β) on the posterior by looking at several other hyperparamter configurations. Which values for α and β correspond to a non-informative prior? What is the interpretation of α and β for the beta prior? What does the choice of α and β in Question 1 tell you about the strength of the prior belief?

Answer:

I took different values of α and β and plotted them against each other:

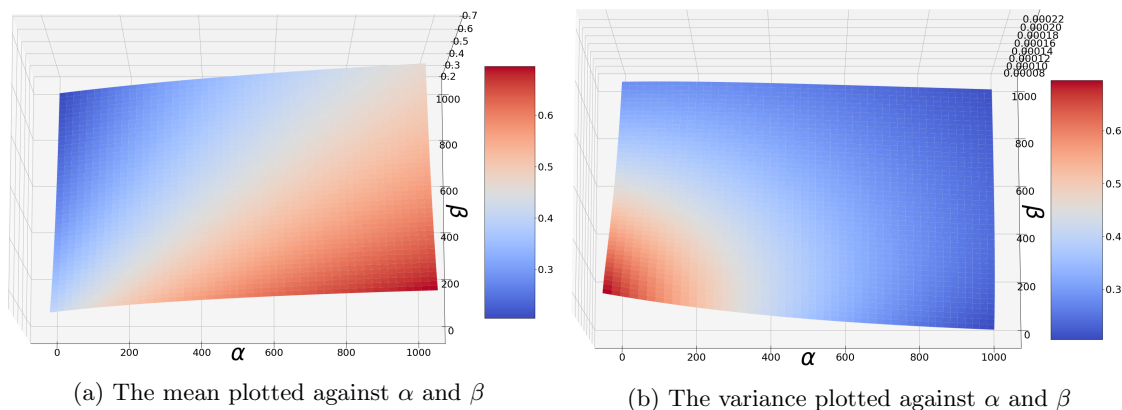


Figure 5: The mean and variance given different hyperparamters α and β .

At figure 5 you can see the mean is higher for higher *alpha* and lower *beta*. The variance increases when α or β is lowered. If you have a non-informed prior your guess is somewhere in the middle with a high variance. This corresponds to a low α and β , while $\alpha = \beta$ so the mean will be $\frac{\alpha}{2\alpha} = 0.5$. This makes sense, for if you gather more information you can update your hyperparameters, thus making them larger. So the hyperparameters in the prior tells us what we believe the distribution is likely to be.

The size of α and β tells us about the strength of our belief. If they are higher it is as if we had low hyperparameters to start with, but saw lots of evidence suggesting our current state of hyperparameters. They also give rise to a low variance as seen in figure 5.

2.4

Imagine you are now a reporter for the polling agency and you have been sent on field duty to gather more data. Your mission is to go out on the streets and randomly survey people on their thoughts regarding the upcoming referendum. Given all the available information you have acquired, what is the probability that the first person you talk to will vote 'Leave'?

Hint: Derive the predictive distribution for the next vote using the posterior distribution for θ computed in Question 2. For a reminder on predictive distribution, see subsection 1.2.6 in Bishop, in particular Equation (1.68).

Answer:

Equation (1.68) says:

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|x, w)p(w|\mathbf{x}, \mathbf{t})dw \quad (66)$$

Exercise 3 - Sequential learning (weight 5)

Part 1 - Obtaining the prior

Consider a four dimensional variable $[x_1, x_2, x_3, x_4]^T$, distributed according to a multivariate Gaussian with mean $\tilde{\mu} = [1, 0, 1, 2]^T$ and covariance matrix $\tilde{\Sigma}$ given as

$$\tilde{\Sigma} = \left(\begin{array}{cc|cc} 0.14 & -0.3 & 0.0 & 0.2 \\ -0.3 & 1.16 & 0.2 & -0.8 \\ \hline 0.0 & 0.2 & 1.0 & 1.0 \\ 0.2 & -0.8 & 1.0 & 2.0 \end{array} \right)$$

We are interested in the conditional distribution over $[x_1, x_2]^T$, given that $x_3 = x_4 = 0$. We know this conditional distribution will also take the form of a Gaussian:

$$p([x_1, x_2]^T | x_3 = x_4 = 0) = \mathcal{N}([x_1, x_2]^T | \mu_p, \Sigma_p) \quad (67)$$

for which the mean and covariance matrix are most easily expressed in terms of the (partitioned) precision matrix (see Bishop, 2.3.1).

3.1.1

Use the partitioned precision matrix $\tilde{\Lambda} = \tilde{\Sigma}^{-1}$ to give an explicit expression for the mean μ_p and covariance matrix Σ_p of this distribution and calculate their values. (This distribution will be taken as the prior information for the rest of this exercise, hence the subscript p). You may use the MATLAB command `inv` to calculate the matrix inverses.

Answer:

3.1.2

[MATLAB] - Create a function that can generate random number pairs, distributed according to the distribution in 67. Initialize your random generator and then draw a single pair

$$\mu_t = [\mu_{t1}, \mu_{t2}]^T \quad (68)$$

from this distribution. (These will be the 'true' means, hence the subscript t).

Hint: you can use the MATLAB function `mvnrnd` (which resides in the Statistics toolbox).

Answer:

3.1.3

[MATLAB] - Make a plot of the probability density of the distribution 67.

Hint: Use the MATLAB function `mvnpdf` (which resides in the Statistics toolbox) to calculate the probability density of a multivariate Gaussian random variable. The MATLAB functions `meshgrid` and `surf` may also prove useful.

Answer:

Part 2 - Generating the data

Here we assume we are dealing with a 2d-Gaussian data generating process

$$p(x) = \mathcal{N}(x|\mu, \Sigma) \quad (69)$$

For the mean μ , we will use the value μ_t drawn in ?? in order to generate the data. Subsequently, we will pretend that we do not know this "true" value μ_t of μ , and estimate μ from the data. For the covariance matrix Σ we will use the "true" value

$$\Sigma_t = \begin{pmatrix} 2.0 & 0.8 \\ 0.8 & 4.0 \end{pmatrix}$$

to generate the data.

3.2.1

[MATLAB] - Generate at least 1000 data pairs $\{x_i, y_i\}$, distributed according to ?? with $\mu = \mu_t$ and $\Sigma = \Sigma_t$ and save them to a file in plain-text format.

Answer:

3.2.2

From now on, we will assume (pretend) the 'true' mean μ_t is unknown and estimate μ from the data. Calculate the maximum likelihood estimate of μ_{ML} and Σ_{ML} for the data, and also an unbiased estimate of Σ (see Bishop, 2.3.4). Compare with the true values μ_t and Σ_t .

Answer:

Part 3 - Sequential learning algorithms

We will now estimate the mean μ from the generated data and the known variance Σ_t sequentially, i.e., by considering the data points one-by-one.

3.3.1

[MATLAB] - Write a procedure that processes the data points $\{x_n\}$ in the generated file one-by-one, and after each step computes an updated estimate of μ_{ML} , the maximum likelihood of the mean (using Bishop, eq. 2.126).

Answer:

3.3.2

Now we also use the prior information $p(\mu) = \mathcal{N}(\mu|\mu_p, \Sigma_p)$. From the prior, the generated data and the known variance Σ_t , we will estimate the mean μ .

Work out the details of sequential Bayesian inference (see eq. 2.144) for the mean μ . Apply Bayes' theorem in eq. 2.113 - 2.117 at each step $n = 1, \dots, N$ to compute the new posterior mean $\mu^{(n)}$ and covariance $\Sigma^{(n)}$ after a new point (x_n) has arrived from the old posterior step. The first step starts from the original prior ??.

Note: Do not confuse the posterior $\Sigma^{(n)}$ with the known Σ_t of the data generating process. For some more

hints, see appendix.

Answer:

3.3.3

[MATLAB] - Write a procedure that processes the data points $\{x_n\}$ in the generated file one-by-one, and after each step computes an updated estimate of μ_{MAP} - the maximum of the posterior distribution, using the results of the previous exercise.

Answer:

3.3.4

[MATLAB] - Plot both estimates (ML and MAP) in a single graph (1d or 2d) as a function of the number of data points observed. Indicate the true values $\{\mu_{t1}, \mu_{t2}\}$ as well. Evaluate your result.

Answer:

Hints

Below are some hints for **Exercise 3 - Part 3 - Question 2**.

Bayes' rule is also valid if earlier acquired information is taken into account. For example, if this is earlier seen data $D_{n-1} = \{x_1, \dots, x_{n-1}\}$. Bayes' rule conditioned on this earlier data is

$$P(\mu|x_n, D_{n-1}) \propto P(\mu|D_{n-1})P(x_n|\mu, D_{n-1})$$

Since $D_n = \{x_1, \dots, x_n\}$ this is written more conveniently as

$$P(\mu|D_n) \propto P(\mu|D_{n-1})P(x_n|\mu, D_{n-1})$$

If, given the model parameters μ , the probability distribution of x_n is independent of earlier data D_{n-1} , we can further reduce this to

$$P(\mu|D_n) \propto P(\mu|D_{n-1})P(x_n|\mu)$$

You should be able to see the relation with (2.144) and see in particular that the factor between brackets in (2.144) is to be identified with $P(\mu|D_{n-1})$.

Another important insight is that if $P(\mu|D_{n-1})$ and $P(x_n|\mu)$ are of the form (2.113) and (2.114), i.e., if $P(\mu|D_{n-1})$ is a Gaussian distribution over μ with a certain mean and covariance (you are free to give these any name, e.g. $\mu^{(n-1)}, \Sigma^{(n-1)}$) and if $P(x_n|\mu)$ is also Gaussian with a mean that is linear in μ , then you can use (2.116) and (2.117) to compute the posterior $P(\mu|D_n)$, which therefore is also Gaussian.

So it is your task to show this. To do this you have to figure out the mapping of the variables and parameters in the current exercise, i.e., what is the correspondence between $\mu, x_n, \Sigma_t, \mu^{(n-1)}, \Sigma^{(n-1)}$ etc. with x, μ, Λ, A, b, L . Don't forget that some quantities can also be zero or and other may be identity matrices.