

STATISTICAL MACHINE LEARNING

ASSIGNMENT 2

Inez Wijnands (s4149696) & Guido Zuidhof (s4160703)

Radboud University Nijmegen

04/11/2015

The entire code listing is in a separate file. The listings shown here are merely code snippets.

1 Sequential learning

1.1 Obtaining the prior

1.

$$\tilde{\Lambda}_{a,b} = \tilde{\Sigma}_{a,b}^{-1} \quad (1.1)$$

$$= \left(\begin{array}{cc|cc} 60 & 50 & -48 & 38 \\ 50 & 50 & -50 & 40 \\ \hline -48 & -50 & 52.4 & -41.4 \\ 38 & 40 & -41.4 & 33.4 \end{array} \right) \quad (1.2)$$

Using the precision matrix $\tilde{\Lambda}$ we can use equations 2.69, 2.73 and 2.75 from Bishop to obtain the mean and covariance of the conditional distribution $p([x_1, x_2]^T | x_3 = x_4 = 0)$.

$$\Sigma_p = \Lambda_{aa}^{-1} \quad (\text{Bishop 2.73})$$

$$\Lambda_{aa} = \begin{pmatrix} 60 & 50 \\ 50 & 50 \end{pmatrix} \quad (1.3)$$

$$\Sigma_p = \begin{pmatrix} 0.1 & -0.1 \\ -0.1 & 0.12 \end{pmatrix} \quad (1.4)$$

$$\mu_p = \mu_{a|b} = \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (\mathbf{x}_b - \mu_b) \quad (\text{Bishop 2.75})$$

We can fill in this equation, since $\tilde{\mu}$ and \mathbf{x}_b (the second partition of \mathbf{x}) are known.

$$\mu_p = \begin{pmatrix} 1 \\ 0 \end{pmatrix} - \begin{pmatrix} 60 & 50 \\ 50 & 50 \end{pmatrix}^{-1} \begin{pmatrix} -48 & 38 \\ -50 & 40 \end{pmatrix} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix} - \begin{pmatrix} 1 \\ 2 \end{pmatrix} \right) \quad (1.5)$$

$$= \begin{pmatrix} 1 \\ 0 \end{pmatrix} - \begin{pmatrix} 0.1 & -0.1 \\ -0.1 & 0.12 \end{pmatrix} \begin{pmatrix} -48 & 38 \\ -50 & 40 \end{pmatrix} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix} - \begin{pmatrix} 1 \\ 2 \end{pmatrix} \right) \quad (1.6)$$

$$= \begin{pmatrix} 0.8 \\ 0.8 \end{pmatrix} \quad (1.7)$$

2. Using the prior μ_p and Σ_p , we used the numpy-equivalent in Python for the MATLAB-function `mvnrnd` to obtain the μ_t we used for the remainder of this assignment:

```
np.random.multivariate_normal(mu_p, sigma_p, 1)
```

This resulted in:

$$\boldsymbol{\mu}_t = \begin{pmatrix} 0.28584241 \\ 1.42626702 \end{pmatrix} \quad (1.8)$$

3. The probability density is highest at the mean (as illustrated in Figure 1.1). The density decreases quickly as both x and y change, but less so when x XOR y change. In $\boldsymbol{\Sigma}_p$, the values for the x XOR y are lower, so this is consistent with our density plot.

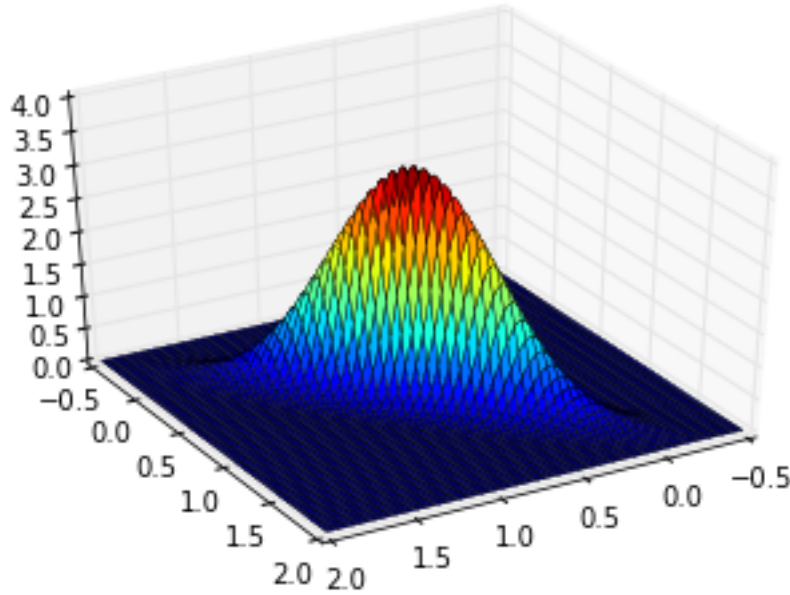


Figure 1.1: The probability density of the distribution.

1.2 Generating the data

1. We used the following function to generate our data:

```
np.random.multivariate_normal(mu_t, sigma_t, 1000)
```

- 2.

$$\boldsymbol{\mu}_{ML} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \quad (\text{Bishop 2.121})$$

$$= \left[\frac{1}{1000} \sum_{n=1}^{1000} \mathbf{x}_n, \frac{1}{1000} \sum_{n=1}^{1000} \mathbf{y}_n \right] \quad (1.9)$$

$$= \begin{pmatrix} 0.25383138 \\ 1.38260838 \end{pmatrix} \quad (1.10)$$

$$\boldsymbol{\Sigma}_{ML} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T \quad (\text{Bishop 1.22})$$

$$= \begin{pmatrix} \frac{1}{1000} \sum_{n=1}^{1000} (\mathbf{x}_n - \boldsymbol{\mu}_{(1)ML})(\mathbf{x}_n - \boldsymbol{\mu}_{(1)ML})^T & \frac{1}{1000} \sum_{n=1}^{1000} (\mathbf{x}_n - \boldsymbol{\mu}_{(1)ML})(\mathbf{y}_n - \boldsymbol{\mu}_{(2)ML})^T \\ \frac{1}{1000} \sum_{n=1}^{1000} (\mathbf{y}_n - \boldsymbol{\mu}_{(2)ML})(\mathbf{x}_n - \boldsymbol{\mu}_{(1)ML})^T & \frac{1}{1000} \sum_{n=1}^{1000} (\mathbf{y}_n - \boldsymbol{\mu}_{(2)ML})(\mathbf{y}_n - \boldsymbol{\mu}_{(2)ML})^T \end{pmatrix} \quad (1.11)$$

$$= \begin{pmatrix} 1.90513804 & 0.72479489 \\ 0.72479489 & 3.81690496 \end{pmatrix} \quad (1.12)$$

This is calculated using the code in Listing 1:

Listing 1: Python code to calculate μ_{ML} and Σ_{ML} .

```

1 mu_ml = sum(data)/len(data)
2
3 sse = [0,0]
4 for point in data:
5     point = np.matrix(point)
6     sse += (point-mu_ml).T*(point-mu_ml)
7 sigma_ml = sse/len(data)

```

The differences with the 'true' values are:

$$\mu_t - \mu_{ML} = \begin{pmatrix} 0.28584241 \\ 1.42626702 \end{pmatrix} - \begin{pmatrix} 0.25383138 \\ 1.38260838 \end{pmatrix} = \begin{pmatrix} 0.03201103 \\ 0.04365864 \end{pmatrix} \quad (1.13)$$

$$\Sigma_t - \Sigma_{ML} = \begin{pmatrix} 2.0 & 0.8 \\ 0.8 & 4.0 \end{pmatrix} - \begin{pmatrix} 1.90513804 & 0.72479489 \\ 0.72479489 & 3.81690496 \end{pmatrix} = \begin{pmatrix} 0.09486196 & 0.07520511 \\ 0.07520511 & 0.18309504 \end{pmatrix} \quad (1.14)$$

For the unbiased covariance:

$$\tilde{\Sigma} = \frac{1}{N-1} \sum_{n=1}^N (\mathbf{x}_n - \mu_{ML})(\mathbf{x}_n - \mu_{ML})^T \quad (\text{Bishop 2.125})$$

We added this line of code to the function of Listing 1:

```
sigma_ml_unbiased = sse * (1/(len(data)-1))
```

The difference between the unbiased covariance and the 'true' covariance is:

$$\Sigma_t - \tilde{\Sigma} = \begin{pmatrix} 2.0 & 0.8 \\ 0.8 & 4.0 \end{pmatrix} - \begin{pmatrix} 1.90704509 & 0.72552041 \\ 0.72552041 & 3.82072569 \end{pmatrix} = \begin{pmatrix} 0.09295491 & 0.07447959 \\ 0.07447959 & 0.17927431 \end{pmatrix} \quad (1.15)$$

The difference is slightly smaller than the difference between Σ_t and Σ_{ML} .

1.3 Sequential learning algorithms

1. See Listing 2 for our procedure to process all data points one-by-one to calculate an estimate of μ_{ML} .

Listing 2: Python code for function *sequential_learning_ml(data)*.

```

1 def sequential_learning_ml(data):
2     N = 0
3     mu_ml = 0
4     mus = []
5
6     for point in data:
7         N += 1
8         mu_ml = mu_ml + (1/N)*(point-mu_ml)
9         mus.append(mu_ml)
10
11     print "Sequential mu_ml:", mu_ml
12     return mus

```

This resulted in:

$$\mu_{ML} = \begin{pmatrix} 0.25383138 \\ 1.38260838 \end{pmatrix} \quad (1.16)$$

2.

$$p(\mathbf{x}|D_{n-1}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \quad (\text{Bishop 2.113})$$

where: $\mathbf{x} = \boldsymbol{\mu}, \boldsymbol{\mu} = \boldsymbol{\mu}_{(n-1)}, \boldsymbol{\Lambda}^{-1} = \boldsymbol{\Sigma}_{(n-1)}$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) \quad (\text{Bishop 2.114})$$

where $\mathbf{y} = \mathbf{x}_n, \mathbf{A} = \mathbf{I}, \mathbf{x} = \boldsymbol{\mu}, \mathbf{b} = 0, \mathbf{L}^{-1} = \boldsymbol{\Sigma}^t$

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^T \mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma}) \quad (\text{Bishop 2.116})$$

$$\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A})^{-1} \quad (\text{Bishop 2.117})$$

Matching the variables we get the following equations:

$$p(\boldsymbol{\mu}|\mathbf{x}_n) = \mathcal{N}(\boldsymbol{\mu}|\mathbf{S}\{\mathbf{I}^T \boldsymbol{\Sigma}_t^{-1}(\mathbf{x}_n - 0) + \boldsymbol{\Sigma}_{(n-1)}^{-1}\}, \mathbf{S}) \quad (1.17)$$

$$= \mathcal{N}(\boldsymbol{\mu}|\mathbf{S}\{\mathbf{I}^T \boldsymbol{\Sigma}_t^{-1} \mathbf{x}_n + \boldsymbol{\Sigma}_{(n-1)}^{-1}\}, \mathbf{S}) \quad (1.18)$$

$$= \mathcal{N}(\boldsymbol{\mu}|\mathbf{S}\{\boldsymbol{\Sigma}_t^{-1} \mathbf{x}_n + \boldsymbol{\Sigma}_{(n-1)}^{-1}\}, \mathbf{S}) \quad (1.19)$$

$$\mathbf{S} = (\boldsymbol{\Sigma}_{(n-1)}^{-1} + \mathbf{I}^T \boldsymbol{\Sigma}_t^{-1} \mathbf{I})^{-1} \quad (1.20)$$

$$= (\boldsymbol{\Sigma}_{(n-1)}^{-1} + \boldsymbol{\Sigma}_t^{-1})^{-1} \quad (1.21)$$

$\boldsymbol{\mu}_n$ is the mean of the distribution $p(\boldsymbol{\mu}|\mathbf{x}_n)$, so the functions we use for our sequential learning algorithm are:

$$\boldsymbol{\Sigma}_n = \mathbf{S} \quad (1.22)$$

$$\boldsymbol{\mu}_n = \boldsymbol{\Sigma}_n\{\boldsymbol{\Sigma}_t^{-1} \mathbf{x}_n + \boldsymbol{\Sigma}_{(n-1)}^{-1}\}, \boldsymbol{\Sigma}_n) \quad (1.23)$$

3. See Listing 3 for our procedure to make a MAP estimation of $\boldsymbol{\mu}$ by processing the data points one-by-one.

Listing 3: Python code for function *sequential_learning_map(data, mu_p, sigma_p, sigma_t)*.

```

1  def sequential_learning_map(data, mu_p, sigma_p, sigma_t):
2      sigma = sigma_p
3      mu = mu_p
4      mus = []
5
6      for point in data:
7          point = np.matrix(point).T
8          S = np.linalg.inv( np.linalg.inv(sigma) + np.linalg.inv(sigma_t))
9          mu = np.dot(S, np.dot( np.linalg.inv(sigma_t), point) + np.dot( np←
              .linalg.inv(sigma), mu))
10         sigma = S
11         mus.append(np.array(mu))
12
13     print "Sequential mu_map:", mu
14     return mus

```

This resulted in:

$$\boldsymbol{\mu}_{MAP} = \begin{pmatrix} 0.25941079 \\ 1.37796331 \end{pmatrix} \quad (1.24)$$

4. See Figure 1.2 for our results: We observed that the MAP estimate performs better for less data points. This seems logical enough, since the MAP estimate is a regularized version of the ML estimate. When more data points are observed, the difference is almost non-existent.

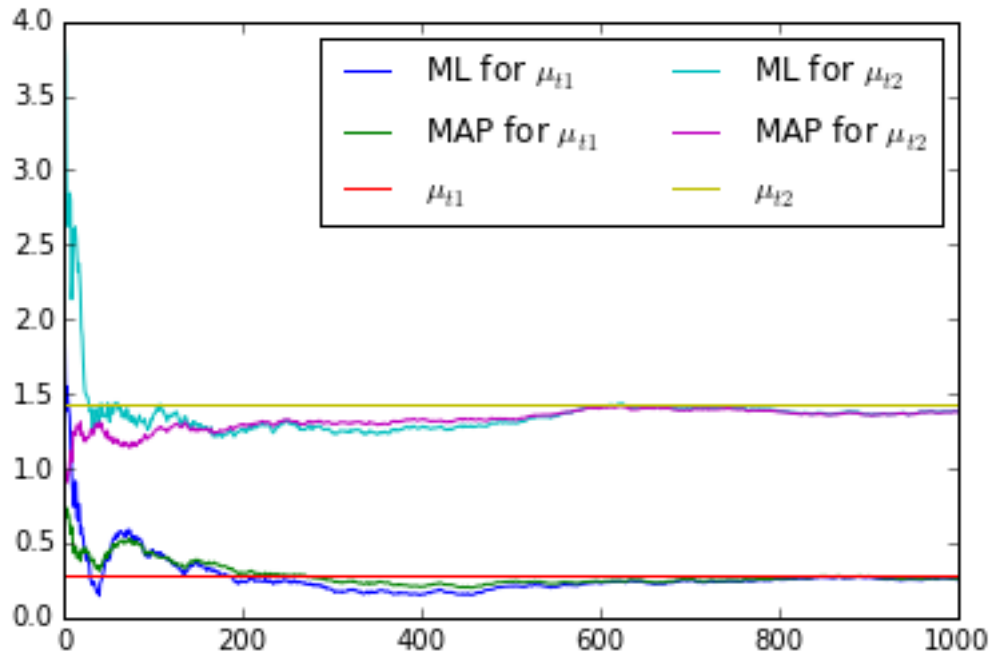


Figure 1.2: The ML and MAP estimates for μ are given, plotted against the amount of data points observed. The 'true' value μ_t is also indicated. This is done for both μ_{t1} and μ_{t2} .

2 The faulty lighthouse

2.1 Constructing the model

- 1.
2. First, using the given $\beta \tan(\theta_k) = x_k - \alpha$, we will calculate the derivation of θ .

$$\beta \tan(\theta_k) = x_k - \alpha \quad (\text{Assignment eq. 7})$$

$$\tan(\theta_k) = \frac{x_k - \alpha}{\beta} \quad (2.1)$$

$$\theta_k = \tan^{-1} \frac{x_k - \alpha}{\beta} \quad (2.2)$$

$$\left| \frac{d\theta}{dx} \right| = \frac{1}{1 + \left(\frac{x_k - \alpha}{\beta} \right)^2} \cdot \left| \frac{d \frac{x_k - \alpha}{\beta}}{dx} \right| \quad (2.3)$$

$$= \frac{1}{1 + \frac{(x_k - \alpha)^2}{\beta^2}} \cdot \frac{\beta}{\beta^2} \quad (2.4)$$

$$= \frac{\beta}{\beta^2 + \beta^2 \left(\frac{x_k - \alpha}{\beta} \right)^2} \quad (2.5)$$

$$= \frac{\beta}{\beta^2 + \beta^2 \frac{(x_k - \alpha)^2}{\beta^2}} \quad (2.6)$$

$$= \frac{\beta}{\beta^2 + (x_k - \alpha)^2} \quad (2.7)$$

Since the following equation holds:

$$p_x(x) = p_\theta(\theta_k) \left| \frac{d\theta}{dx} \right| \quad (\text{Bishop 1.27})$$

We need to multiply the derivation of θ with $p(\theta_k | \alpha, \beta)$ (Assignment eq. 6):

$$p(x_k | \alpha, \beta) = \frac{\beta}{\beta^2 + (x_k - \alpha)^2} \cdot \frac{1}{\pi} \quad (2.8)$$

$$= \frac{\beta}{\pi [\beta^2 + (x_k - \alpha)^2]} \quad (\text{Assignment eq. 8})$$

We have plotted the distribution for $\beta = 1$ and for α we chose the value 0.5, as illustrated in Figure 2.1:

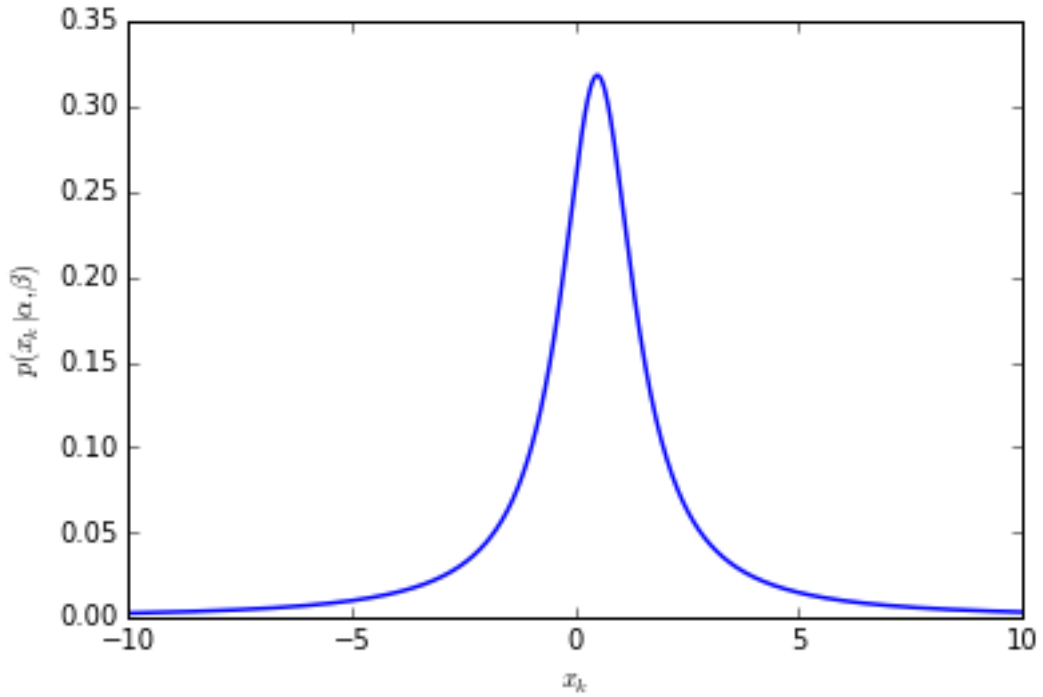


Figure 2.1: The probability distribution $p(x_k|\alpha, \beta)$ plotted against x_k , with $\alpha = 0.5$ and $\beta = 1$.

3.

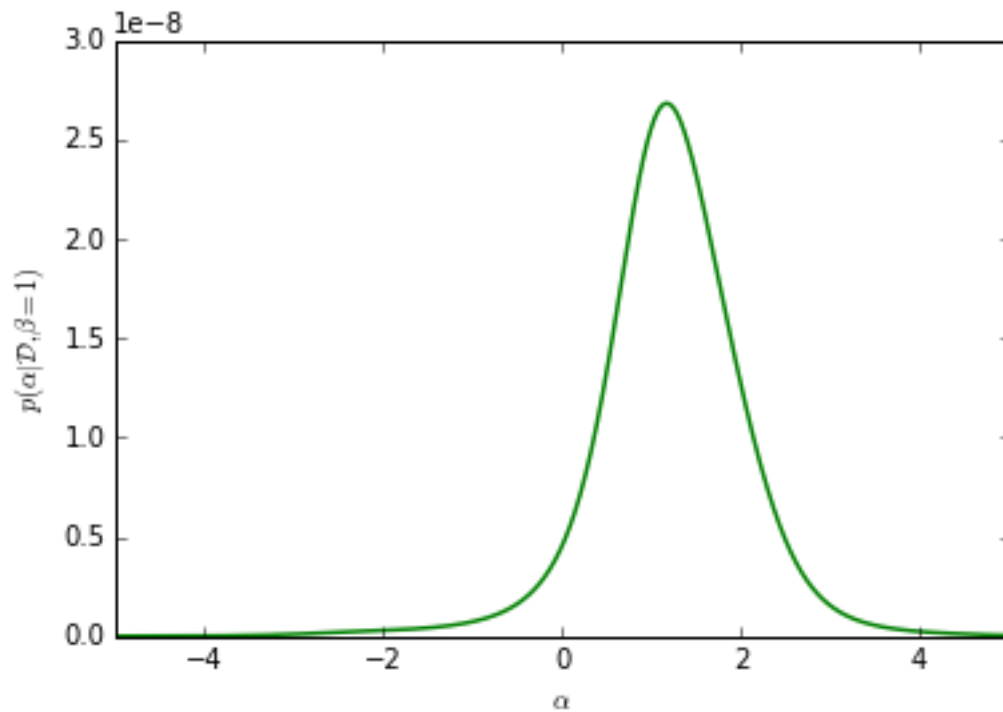


Figure 2.2: The probability density $p(\alpha|\mathcal{D}, \beta = 1)$ plotted against α .

4. The most likely estimate for $\hat{\alpha} = 1.17136$. The mean of $\alpha = -0.18333$. The difference is probably caused by the amount of data points we have, these are very limited. Outliers in the data will have a great effect on the mean.

2.2 Generate the lighthouse data

- 1.
- 2.

2.3 Find the lighthouse

- 1.
- 2.
- 3.