



University of Cyprus
MAI613 - Research Methodologies and Professional
Practices in AI
Lecturer: Stelios Timotheou

Exercises on Probability and Distributions

Exercises

Solve Exercises 6.1, 6.2a, 6.2b and 6.2c (compute only the mean not the mode and the median), 6.4, 6.5, 6.6, and 6.12 from the book:

- M.P. Deisenroth, A. A. Faisal, and C. S. Ong, “Mathematics for Machine Learning,” Cambridge University Press, 2020.

<https://mml-book.github.io/book/mml-book.pdf>

Note: These exercises are intended for self-assessment purposes. Their solutions will be uploaded in one week from now.

Probability and Distributions

Exercises

- 6.1 Consider the following bivariate distribution $p(x, y)$ of two discrete random variables X and Y .

Y	y_1	0.01	0.02	0.03	0.1	0.1
	y_2	0.05	0.1	0.05	0.07	0.2
	y_3	0.1	0.05	0.03	0.05	0.04
		x_1	x_2	x_3	x_4	x_5
		X				

Compute:

- The marginal distributions $p(x)$ and $p(y)$.
- The conditional distributions $p(x|Y = y_1)$ and $p(y|X = x_3)$.

The marginal and conditional distributions are given by

$$\begin{aligned}
 p(x) &= [0.16, 0.17, 0.11, 0.22, 0.34]^\top \\
 p(y) &= [0.26, 0.47, 0.27]^\top \\
 p(x|Y = y_1) &= [0.01, 0.02, 0.03, 0.1, 0.1]^\top \\
 p(y|X = x_3) &= [0.03, 0.05, 0.03]^\top.
 \end{aligned}$$

- 6.2 Consider a mixture of two Gaussian distributions (illustrated in Figure 6.4),

$$0.4\mathcal{N}\left(\begin{bmatrix} 10 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right) + 0.6\mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 8.4 & 2.0 \\ 2.0 & 1.7 \end{bmatrix}\right).$$

- Compute the marginal distributions for each dimension.
- Compute the mean, mode and median for each marginal distribution.
- Compute the mean and mode for the two-dimensional distribution.

Consider the mixture of two Gaussians,

$$p\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = 0.4\mathcal{N}\left(\begin{bmatrix} 10 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right) + 0.6\mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 8.4 & 2.0 \\ 2.0 & 1.7 \end{bmatrix}\right).$$

- a. Compute the marginal distribution for each dimension.

$$p(x_1) = \int 0.4 \left(\mathcal{N} \left(\begin{bmatrix} 10 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) + 0.6 \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 8.4 & 2.0 \\ 2.0 & 1.7 \end{bmatrix} \right) \right) dx_2 \quad (6.1)$$

$$= 0.4 \int \mathcal{N} \left(\begin{bmatrix} 10 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) dx_2 + 0.6 \int \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 8.4 & 2.0 \\ 2.0 & 1.7 \end{bmatrix} \right) dx_2 \quad (6.2)$$

$$= 0.4 \mathcal{N}(10, 1) + 0.6 \mathcal{N}(0, 8.4) .$$

From (6.1) to (6.2), we used the result that the integral of a sum is the sum of the integrals.

Similarly, we can get

$$p(x_2) = 0.4 \mathcal{N}(2, 1) + 0.6 \mathcal{N}(0, 1.7) .$$

- b. Compute the mean, mode and median for each marginal distribution.

Mean:

$$\begin{aligned} \mathbb{E}(x_1) &= \int x_1 p(x_1) dx_1 \\ &= \int x_1 (0.4 \mathcal{N}(x_1 | 10, 1) + 0.6 \mathcal{N}(x_1 | 0, 8.4)) dx_1 \quad (6.3) \end{aligned}$$

$$\begin{aligned} &= 0.4 \int x_1 \mathcal{N}(x_1 | 10, 1) dx_1 + 0.6 \int x_1 \mathcal{N}(x_1 | 0, 8.4) dx_1 \quad (6.4) \\ &= 0.4 \cdot 10 + 0.6 \cdot 0 = 4 . \end{aligned}$$

From step (6.3) to step (6.4), we use the fact that for $Y \sim \mathcal{N}(\mu, \sigma)$, where $\mathbb{E}[Y] = \mu$.

Similarly,

$$\mathbb{E}(x_2) = 0.4 \cdot 2 + 0.6 \cdot 0 = 0.8 .$$

Mode: In principle, we would need to solve for

$$\frac{dp(x_1)}{dx_1} = 0 \quad \text{and} \quad \frac{dp(x_2)}{dx_2} = 0 .$$

However, we can observe that the modes of each individual distribution are the peaks of the Gaussians, that is the Gaussian means for each dimension. Median:

$$\int_{-\infty}^a p(x_1) dx_1 = \frac{1}{2} .$$

- c. Compute the mean and mode for the 2 dimensional distribution.

Mean:

From (6.30), we know that

$$\mathbb{E} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} \mathbb{E}[x_1] \\ \mathbb{E}[x_2] \end{bmatrix} = \begin{bmatrix} 4 \\ 0.8 \end{bmatrix} .$$

Mode:

The two dimensional distribution has two peaks, and hence there are

two modes. In general, we would need to solve an optimization problem to find the maxima. However, in this particular case we can observe that the two modes correspond to the individual Gaussian means, that is

$$\begin{bmatrix} 10 \\ 2 \end{bmatrix} \text{ and } \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

- 6.4 There are two bags. The first bag contains four mangos and two apples; the second bag contains four mangos and four apples.

We also have a biased coin, which shows “heads” with probability 0.6 and “tails” with probability 0.4. If the coin shows “heads”, we pick a fruit at random from bag 1; otherwise we pick a fruit at random from bag 2.

Your friend flips the coin (you cannot see the result), picks a fruit at random from the corresponding bag, and presents you a mango.

What is the probability that the mango was picked from bag 2?

Hint: Use Bayes’ theorem.

We apply Bayes’ theorem and compute the posterior $p(b_2 | m)$ of picking a mango from bag 2.

$$p(b_2 | m) = \frac{p(m | b_2)p(b_2)}{p(m)}$$

where

$$p(m) = p(b_1)p(m|b_1) + p(b_2)p(m|b_2) = \frac{3}{5} \frac{2}{3} + \frac{2}{5} \frac{1}{2} = \frac{2}{5} + \frac{1}{5} = \frac{3}{5} \quad \text{Evidence}$$

$$p(b_2) = \frac{2}{5} \quad \text{Prior}$$

$$p(m|b_2) = \frac{1}{2} \quad \text{Likelihood}$$

Therefore,

$$p(b_2|m) = \frac{p(m|b_2)p(b_2)}{p(m)} = \frac{\frac{2}{5} \frac{1}{2}}{\frac{3}{5}} = \frac{1}{3}.$$

6.5 Consider the time-series model

$$\mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{w}, \quad \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q})$$

$$\mathbf{y}_t = \mathbf{C}\mathbf{x}_t + \mathbf{v}, \quad \mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{R}),$$

where \mathbf{w}, \mathbf{v} are i.i.d. Gaussian noise variables. Further, assume that $p(\mathbf{x}_0) = \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$.

- a. What is the form of $p(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_T)$? Justify your answer (you do not have to explicitly compute the joint distribution).
- b. Assume that $p(\mathbf{x}_t | \mathbf{y}_1, \dots, \mathbf{y}_t) = \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$.
 1. Compute $p(\mathbf{x}_{t+1} | \mathbf{y}_1, \dots, \mathbf{y}_t)$.
 2. Compute $p(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} | \mathbf{y}_1, \dots, \mathbf{y}_t)$.
 3. At time $t+1$, we observe the value $\mathbf{y}_{t+1} = \hat{\mathbf{y}}$. Compute the conditional distribution $p(\mathbf{x}_{t+1} | \mathbf{y}_1, \dots, \mathbf{y}_{t+1})$.
1. The joint distribution is Gaussian: $p(\mathbf{x}_0)$ is Gaussian, and \mathbf{x}_{t+1} is a linear/affine transformations of \mathbf{x}_t . Since affine transformations leave the Gaussianity of the random variable invariant, the joint distribution must be Gaussian.
2. 1. We use the results from linear transformation of Gaussian random variables (Section 6.5.3). We immediately obtain

$$p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{x}_{t+1} | \boldsymbol{\mu}_{t+1|t}, \boldsymbol{\Sigma}_{t+1|t})$$

$$\boldsymbol{\mu}_{t+1|t} := \mathbf{A}\boldsymbol{\mu}_t$$

$$\boldsymbol{\Sigma}_{t+1|t} := \mathbf{A}\boldsymbol{\Sigma}_t\mathbf{A}^\top + \mathbf{Q}.$$

2. The joint distribution $p(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} | \mathbf{y}_{1:t})$ is Gaussian (linear transformation of random variables). We compute every component of the Gaussian separately:

$$\mathbb{E}[\mathbf{y}_{t+1} | \mathbf{y}_{1:t}] = \mathbf{C}\boldsymbol{\mu}_{t+1|t} =: \boldsymbol{\mu}_{t+1|t}^y$$

$$\mathbb{V}[\mathbf{y}_{t+1} | \mathbf{y}_{1:t}] = \mathbf{C}\boldsymbol{\Sigma}_{t+1|t}\mathbf{C}^\top + \mathbf{R} =: \boldsymbol{\Sigma}_{t+1|t}^y$$

$$\text{Cov}[\mathbf{x}_{t+1}, \mathbf{y}_{t+1} | \mathbf{y}_{1:t}] = \text{Cov}[\mathbf{x}_{t+1}, \mathbf{C}\mathbf{x}_{t+1} | \mathbf{y}_{1:t}] = \boldsymbol{\Sigma}_{t+1|t}\mathbf{C}^\top =: \boldsymbol{\Sigma}_{xy}$$

Therefore,

$$p(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} | \mathbf{y}_{1:t}) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_{t+1|t} \\ \boldsymbol{\mu}_{t+1|t}^y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{t+1|t} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{t+1|t}^y \end{bmatrix}\right).$$

3. We obtain the desired distribution (again: Gaussian) by applying the rules for Gaussian conditioning of the joint distribution in B):

$$\begin{aligned}
 p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t+1}) &= \mathcal{N}(\boldsymbol{\mu}_{t+1|t+1}, \boldsymbol{\Sigma}_{t+1|t+1}) \\
 \boldsymbol{\mu}_{t+1|t+1} &= \boldsymbol{\mu}_{t+1|t} + \boldsymbol{\Sigma}_{xy}(\boldsymbol{\Sigma}_{t+1|t}^y)^{-1}(\hat{\mathbf{y}} - \boldsymbol{\mu}_{t+1|t}^y) \\
 \boldsymbol{\Sigma}_{t+1|t+1} &= \boldsymbol{\Sigma}_{t+1|t} - \boldsymbol{\Sigma}_{xy}(\boldsymbol{\Sigma}_{t+1|t}^y)^{-1}\boldsymbol{\Sigma}_{yx}
 \end{aligned}$$

- 6.6 Prove the relationship in (6.44), which relates the standard definition of the variance to the raw-score expression for the variance.

The formula in (6.43) can be converted to the so called raw-score formula for variance

$$\begin{aligned}
 \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 &= \frac{1}{N} \sum_{i=1}^N (x_i^2 - 2x_i\mu + \mu^2) \\
 &= \frac{1}{N} \sum_{i=1}^N x_i^2 - \frac{2}{N} \mu \sum_{i=1}^N x_i + \mu^2 \\
 &= \frac{1}{N} \sum_{i=1}^N x_i^2 - \mu^2 \\
 &= \frac{1}{N} \sum_{i=1}^N x_i^2 - \left(\frac{1}{N} \sum_{i=1}^N x_i \right)^2.
 \end{aligned}$$

6.11 Iterated Expectations.

Consider two random variables x, y with joint distribution $p(x, y)$. Show that

$$\mathbb{E}_X[x] = \mathbb{E}_Y[\mathbb{E}_X[x | y]] .$$

Here, $\mathbb{E}_X[x | y]$ denotes the expected value of x under the conditional distribution $p(x | y)$.

$$\begin{aligned} \mathbb{E}_X[x] &= \mathbb{E}_Y[\mathbb{E}_X[x | y]] \\ &\iff \int xp(x)dx = \iint xp(x | y)p(y)dxdy \\ &\iff \int xp(x)dx = \iint xp(x, y)dxdy \\ &\iff \int xp(x)dx = \int x \int p(x, y)dydx = \int xp(x)dx , \end{aligned}$$

which proves the claim.

6.12 Manipulation of Gaussian Random Variables.

Consider a Gaussian random variable $\mathbf{x} \sim \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x)$, where $\mathbf{x} \in \mathbb{R}^D$. Furthermore, we have

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b} + \mathbf{w} ,$$

where $\mathbf{y} \in \mathbb{R}^E$, $\mathbf{A} \in \mathbb{R}^{E \times D}$, $\mathbf{b} \in \mathbb{R}^E$, and $\mathbf{w} \sim \mathcal{N}(\mathbf{w} | \mathbf{0}, \mathbf{Q})$ is independent Gaussian noise. “Independent” implies that \mathbf{x} and \mathbf{w} are independent random variables and that \mathbf{Q} is diagonal.

a. Write down the likelihood $p(\mathbf{y} | \mathbf{x})$.

$$\begin{aligned} p(\mathbf{y} | \mathbf{x}) &= \mathcal{N}(\mathbf{y} | \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{Q}) \\ &= Z \exp \left(-\frac{1}{2}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b})^\top \mathbf{Q}^{-1}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}) \right) \\ \text{where } Z &= (2\pi)^{-E/2} |\mathbf{Q}|^{-1/2} \end{aligned}$$

- b. The distribution $p(\mathbf{y}) = \int p(\mathbf{y} | \mathbf{x})p(\mathbf{x})d\mathbf{x}$ is Gaussian. Compute the mean $\boldsymbol{\mu}_y$ and the covariance $\boldsymbol{\Sigma}_y$. Derive your result in detail.
 $p(\mathbf{y})$ is Gaussian distributed (affine transformation of the Gaussian random variable \mathbf{x}) with mean $\boldsymbol{\mu}_y$ and covariance matrix $\boldsymbol{\Sigma}_y$. We obtain

$$\begin{aligned}\boldsymbol{\mu}_y &= \mathbb{E}_Y[\mathbf{y}] = \mathbb{E}_{X,W}[\mathbf{A}\mathbf{x} + \mathbf{b} + \mathbf{w}] = \mathbf{A}\mathbb{E}_X[\mathbf{x}] + \mathbf{b} + \mathbb{E}_W[\mathbf{w}] \\ &= \mathbf{A}\boldsymbol{\mu}_x + \mathbf{b}, \quad \text{and} \\ \boldsymbol{\Sigma}_y &= \mathbb{E}_Y[\mathbf{y}\mathbf{y}^\top] - \mathbb{E}_Y[\mathbf{y}]\mathbb{E}_Y[\mathbf{y}]^\top \\ &= \mathbb{E}_{X,W}[(\mathbf{A}\mathbf{x} + \mathbf{b} + \mathbf{w})(\mathbf{A}\mathbf{x} + \mathbf{b} + \mathbf{w})^\top] - \boldsymbol{\mu}_y\boldsymbol{\mu}_y^\top \\ &= \mathbb{E}_{X,W}[\mathbf{A}\mathbf{x}\mathbf{x}^\top\mathbf{A}^\top + \mathbf{A}\mathbf{x}\mathbf{b}^\top + \mathbf{A}\mathbf{x}\mathbf{w}^\top + \mathbf{b}\mathbf{x}^\top\mathbf{A}^\top + \mathbf{b}\mathbf{b}^\top + \mathbf{b}\mathbf{w}^\top \\ &\quad + \mathbf{w}(\mathbf{A}\mathbf{x} + \mathbf{b})^\top + \mathbf{w}\mathbf{w}^\top] - \boldsymbol{\mu}_y\boldsymbol{\mu}_y^\top.\end{aligned}$$

We use the linearity of the expected value, move all constants out of the expected value, and exploit the independence of \mathbf{w} and \mathbf{x} :

$$\begin{aligned}\boldsymbol{\Sigma}_y &= \mathbf{A}\mathbb{E}_X[\mathbf{x}\mathbf{x}^\top]\mathbf{A}^\top + \mathbf{A}\mathbb{E}_X[\mathbf{x}]\mathbf{b}^\top + \mathbf{b}\mathbb{E}_X[\mathbf{x}^\top]\mathbf{A}^\top + \mathbf{b}\mathbf{b}^\top + \mathbb{E}_W[\mathbf{w}\mathbf{w}^\top] \\ &\quad - (\mathbf{A}\boldsymbol{\mu}_x + \mathbf{b})(\mathbf{A}\boldsymbol{\mu}_x + \mathbf{b})^\top,\end{aligned}$$

where we used our previous result for $\boldsymbol{\mu}_y$. Note that $\mathbb{E}_W[\mathbf{w}] = \mathbf{0}$. We continue as follows:

$$\begin{aligned}\boldsymbol{\Sigma}_y &= \mathbf{A}\mathbb{E}_X[\mathbf{x}\mathbf{x}^\top]\mathbf{A}^\top + \mathbf{A}\boldsymbol{\mu}_x\mathbf{b}^\top + \mathbf{b}\boldsymbol{\mu}_x^\top\mathbf{A}^\top + \mathbf{b}\mathbf{b}^\top + \mathbf{Q} \\ &\quad - \mathbf{A}\boldsymbol{\mu}_x\boldsymbol{\mu}_x^\top\mathbf{A}^\top - \mathbf{A}\boldsymbol{\mu}_x\mathbf{b}^\top - \mathbf{b}\boldsymbol{\mu}_x^\top\mathbf{A}^\top - \mathbf{b}\mathbf{b}^\top \\ &= \mathbf{A}\underbrace{(\mathbb{E}_X[\mathbf{x}\mathbf{x}^\top] - \boldsymbol{\mu}_x\boldsymbol{\mu}_x^\top)}_{=\boldsymbol{\Sigma}_x}\mathbf{A}^\top + \mathbf{Q} \\ &= \mathbf{A}\boldsymbol{\Sigma}_x\mathbf{A}^\top + \mathbf{Q}.\end{aligned}$$

Alternatively, we could have exploited

$$\begin{aligned}\mathbb{V}_y[\mathbf{y}] &= \mathbb{V}_{x,w}[\mathbf{A}\mathbf{x} + \mathbf{b} + \mathbf{w}] \stackrel{\text{i.i.d.}}{=} \mathbb{V}_x[\mathbf{A}\mathbf{x} + \mathbf{b}] + \mathbb{V}_w[\mathbf{w}] \\ &= \mathbf{A}\mathbb{V}_x\mathbf{A}^\top + \mathbf{Q} = \mathbf{A}\boldsymbol{\Sigma}_x\mathbf{A}^\top + \mathbf{Q}.\end{aligned}$$

- c. The random variable \mathbf{y} is being transformed according to the measurement mapping

$$\mathbf{z} = \mathbf{C}\mathbf{y} + \mathbf{v},$$

where $\mathbf{z} \in \mathbb{R}^F$, $\mathbf{C} \in \mathbb{R}^{F \times E}$, and $\mathbf{v} \sim \mathcal{N}(\mathbf{v} | \mathbf{0}, \mathbf{R})$ is independent Gaussian (measurement) noise.

- Write down $p(\mathbf{z} | \mathbf{y})$.

$$\begin{aligned}p(\mathbf{z} | \mathbf{y}) &= \mathcal{N}(\mathbf{z} | \mathbf{C}\mathbf{y}, \mathbf{R}) \\ &= (2\pi)^{-F/2} |\mathbf{R}|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{z} - \mathbf{C}\mathbf{y})^\top \mathbf{R}^{-1}(\mathbf{z} - \mathbf{C}\mathbf{y})\right).\end{aligned}$$

- Compute $p(\mathbf{z})$, i.e., the mean $\boldsymbol{\mu}_z$ and the covariance $\boldsymbol{\Sigma}_z$. Derive your result in detail.

Remark: Since \mathbf{y} is Gaussian and \mathbf{z} is a linear transformation of \mathbf{y} , $p(\mathbf{z})$ is Gaussian, too. Let's compute its moments (similar to the “time update”):

$$\begin{aligned}\boldsymbol{\mu}_z &= \mathbb{E}_Z[\mathbf{z}] = \mathbb{E}_{Y,V}[\mathbf{C}\mathbf{y} + \mathbf{v}] = \mathbf{C}\mathbb{E}_Y[\mathbf{y}] + \underbrace{\mathbb{E}_V[\mathbf{v}]}_{=0} \\ &= \mathbf{C}\boldsymbol{\mu}_y.\end{aligned}$$

For the covariance matrix, we compute

$$\begin{aligned}\boldsymbol{\Sigma}_z &= \mathbb{E}_Z[\mathbf{z}\mathbf{z}^\top] - \mathbb{E}_Z[\mathbf{z}]\mathbb{E}_Z[\mathbf{z}]^\top \\ &= \mathbb{E}_{Y,V}[(\mathbf{C}\mathbf{y} + \mathbf{v})(\mathbf{C}\mathbf{y} + \mathbf{v})^\top] - \boldsymbol{\mu}_z\boldsymbol{\mu}_z^\top \\ &= \mathbb{E}_{Y,V}[\mathbf{C}\mathbf{y}\mathbf{y}^\top\mathbf{C}^\top + \mathbf{C}\mathbf{y}\mathbf{v}^\top + \mathbf{v}(\mathbf{C}\mathbf{y})^\top + \mathbf{v}\mathbf{v}^\top] - \boldsymbol{\mu}_z\boldsymbol{\mu}_z^\top.\end{aligned}$$

We use the linearity of the expected value, move all constants out of the expected value, and exploit the independence of \mathbf{v} and \mathbf{y} :

$$\boldsymbol{\Sigma}_z = \mathbf{C}\mathbb{E}_Y[\mathbf{y}\mathbf{y}^\top]\mathbf{C}^\top + \mathbb{E}_V[\mathbf{v}\mathbf{v}^\top] - \mathbf{C}\boldsymbol{\mu}_y\boldsymbol{\mu}_y^\top\mathbf{C}^\top,$$

where we used our previous result for $\boldsymbol{\mu}_z$. Note that $\mathbb{E}_V[\mathbf{v}] = \mathbf{0}$. We continue as follows:

$$\begin{aligned}\boldsymbol{\Sigma}_y &= \mathbf{C} \underbrace{(\mathbb{E}_Y[\mathbf{y}\mathbf{y}^\top] - \boldsymbol{\mu}_y\boldsymbol{\mu}_y^\top)}_{=\boldsymbol{\Sigma}_x} \mathbf{C}^\top + \mathbf{R} \\ &= \mathbf{C}\boldsymbol{\Sigma}_x\mathbf{C}^\top + \mathbf{R}.\end{aligned}$$

- d. Now, a value $\hat{\mathbf{y}}$ is measured. Compute the posterior distribution $p(\mathbf{x} | \hat{\mathbf{y}})$. *Hint for solution:* This posterior is also Gaussian, i.e., we need to determine only its mean and covariance matrix. Start by explicitly computing the joint Gaussian $p(\mathbf{x}, \mathbf{y})$. This also requires us to compute the cross-covariances $\text{Cov}_{\mathbf{x},\mathbf{y}}[\mathbf{x}, \mathbf{y}]$ and $\text{Cov}_{\mathbf{y},\mathbf{x}}[\mathbf{y}, \mathbf{x}]$. Then apply the rules for Gaussian conditioning.

We derive the posterior distribution following the second hint since we do not have to worry about normalization constants:

Assume, we know the joint Gaussian distribution

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \middle| \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_x & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_y \end{bmatrix}\right), \quad (6.11)$$

where we defined $\boldsymbol{\Sigma}_{xy} := \text{Cov}_{\mathbf{x},\mathbf{y}}[\mathbf{x}, \mathbf{y}]$.

Now, we apply the rules for Gaussian conditioning to obtain

$$\begin{aligned}p(\mathbf{x} | \hat{\mathbf{y}}) &= \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_{x|y}, \boldsymbol{\Sigma}_{x|y}) \\ \boldsymbol{\mu}_{x|y} &= \boldsymbol{\mu}_x + \boldsymbol{\Sigma}_{xy}\boldsymbol{\Sigma}_y^{-1}(\hat{\mathbf{y}} - \boldsymbol{\mu}_y) \\ \boldsymbol{\Sigma}_{x|y} &= \boldsymbol{\Sigma}_x - \boldsymbol{\Sigma}_{xy}\boldsymbol{\Sigma}_y^{-1}\boldsymbol{\Sigma}_{yx}.\end{aligned}$$

Looking at (6.11), it remains to compute the cross-covariance term Σ_{xy} (the marginal distributions $p(\mathbf{x})$ and $p(\mathbf{y})$ are known and $\Sigma_{yx} = \Sigma_{xy}^\top$):

$$\Sigma_{xy} = \text{Cov}[\mathbf{x}, \mathbf{y}] = \mathbb{E}_{X,Y}[\mathbf{x}\mathbf{y}^\top] - \mathbb{E}_X[\mathbf{x}]\mathbb{E}_Y[\mathbf{y}]^\top = \mathbb{E}_{X,Y}[\mathbf{x}\mathbf{y}^\top] - \boldsymbol{\mu}_x\boldsymbol{\mu}_y^\top,$$

where $\boldsymbol{\mu}_x$ and $\boldsymbol{\mu}_y$ are known. Hence, it remains to compute

$$\begin{aligned}\mathbb{E}_{X,Y}[\mathbf{x}\mathbf{y}^\top] &= \mathbb{E}_X[\mathbf{x}(\mathbf{A}\mathbf{x} + \mathbf{b} + \mathbf{w})^\top] = \mathbb{E}_X[\mathbf{x}\mathbf{x}^\top]\mathbf{A}^\top + \mathbb{E}_X[\mathbf{x}]\mathbf{b}^\top \\ &= \mathbb{E}_X[\mathbf{x}\mathbf{x}^\top]\mathbf{A}^\top + \boldsymbol{\mu}_x\mathbf{b}^\top.\end{aligned}$$

With

$$\boldsymbol{\mu}_x\boldsymbol{\mu}_y^\top = \boldsymbol{\mu}_x\boldsymbol{\mu}_x^\top\mathbf{A}^\top + \boldsymbol{\mu}_x\mathbf{b}^\top$$

we obtain the desired cross-covariance

$$\begin{aligned}\Sigma_{xy} &= \mathbb{E}_X[\mathbf{x}\mathbf{x}^\top]\mathbf{A}^\top + \boldsymbol{\mu}_x\mathbf{b}^\top - \boldsymbol{\mu}_x\boldsymbol{\mu}_x^\top\mathbf{A}^\top - \boldsymbol{\mu}_x\mathbf{b}^\top \\ &= (\mathbb{E}_X[\mathbf{x}\mathbf{x}^\top] - \boldsymbol{\mu}_x\boldsymbol{\mu}_x^\top)\mathbf{A}^\top = \Sigma_x\mathbf{A}^\top.\end{aligned}$$

And finally

$$\begin{aligned}\boldsymbol{\mu}_{x|y} &= \boldsymbol{\mu}_x + \Sigma_x\mathbf{A}^\top(\mathbf{A}\Sigma_x\mathbf{A}^\top + \mathbf{Q})^{-1}(\hat{\mathbf{y}} - \mathbf{A}\boldsymbol{\mu}_x - \mathbf{b}), \\ \Sigma_{x|y} &= \Sigma_x - \Sigma_x\mathbf{A}^\top(\mathbf{A}\Sigma_x\mathbf{A}^\top + \mathbf{Q})^{-1}\mathbf{A}\Sigma_x.\end{aligned}$$