

Handout Probabilistic Graphical Models: Gaussian Computations

Volkan Cevher, Matthias Seeger

Abstract

This is a loosely organized set of notes complementing the lecture on *Gaussian Distributions*. It is meant to help you doing Gaussian computations in an economic way. A word of warning: there is no easy way into this domain, there is just practice.

1 Definitions

The density of a multivariate Gaussian variable $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ($\mathbf{x} \in \mathbb{R}^n$) is given as

$$P(\mathbf{x}) = |2\pi\boldsymbol{\Sigma}|^{-1/2} e^{-(1/2)(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}. \quad (1)$$

Here, $\boldsymbol{\mu} = E[\mathbf{x}]$ is the mean, $\boldsymbol{\Sigma} = \text{Cov}[\mathbf{x}]$ is the covariance matrix of $P(\mathbf{x})$. $|\cdot|$ denotes the determinant of a matrix, the product of its eigenvalues.

A Gaussian is completely defined by mean and covariance. In the maximum entropy sense, it is “nothing else but mean and covariance”, meaning that there is no additional structure. In *general* (not only for Gaussian variables), mean and covariance linearly transform in the following way. If $E[\mathbf{x}] = \boldsymbol{\mu}$, $\text{Cov}[\mathbf{x}] = \boldsymbol{\Sigma}$, and $\mathbf{y} = \mathbf{E}\mathbf{x} + \mathbf{b}$, then

$$E[\mathbf{y}] = \mathbf{E}\boldsymbol{\mu} + \mathbf{b}, \quad \text{Cov}[\mathbf{y}] = \mathbf{E}\boldsymbol{\Sigma}\mathbf{E}^T. \quad (2)$$

These properties hold simply because expectation is a linear transformation.

I will also frequently use the notation $A \doteq B$, which means that $A = CB$, where C is some constant that does not matter right now. If \doteq is used several times, the constants might be different.

2 Closure Properties. How to Determine a Gaussian Result

A family of distributions is *closed* under a set of operations on distributions or random variables if whenever you apply an operation to a family member, the outcome lies in the family as well. The Gaussian family is amazing in that sense.

- Gaussians are closed under linear (affine) transformations:

$$\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \mathbf{y} = \mathbf{E}\mathbf{x} + \mathbf{b} \quad \Rightarrow \quad \mathbf{y} \sim N(\mathbf{E}\boldsymbol{\mu} + \mathbf{b}, \mathbf{E}\boldsymbol{\Sigma}\mathbf{E}^T)$$

- Gaussians are closed under marginalization:

$$\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad I \subset \{1, \dots, n\} \quad \Rightarrow \quad \mathbf{x}_I = (x_i)_{i \in I} \text{ Gaussian}$$

In other words: the *sum* rule retains Gaussianity.

- Gaussians are closed under conditioning:

$$\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad I \subset \{1, \dots, n\}, \quad R = \{1, \dots, n\} \setminus I \quad \Rightarrow \quad P(\mathbf{x}_I | \mathbf{x}_R) \text{ Gaussian}$$

In other words: the *product* rule retains Gaussianity.

Beware the vocabulary. We say that a distribution is Gaussian, or that a random variable is Gaussian (the latter is sloppy for “is Gaussian distributed”).

Here is a very important recipe for *determining a Gaussian result*, which often works even in complicated models or situations. You combine what you know about closure properties with what you know about transformations of mean and covariance:

1. Establish that the desired result has to be a Gaussian.
2. Determine its mean and covariance. Since a Gaussian is defined by these, you are done. In particular, there is *no need to track multiplicative constants*: from mean, covariance, and (1), you figure them out at the end.

For the first part, there is a number of ways you can try:

- Use closure properties, possibly in combination.
- Write down the functional form of the result, then massage it into the *Gaussian form*: $Ce^{-q(\mathbf{x})}$, $q(\mathbf{x})$ a quadratic function. If you manage this, it’s a Gaussian.
- Use a sampling argument. If you can formulate the desired result as part of a joint sampling rule (say, a conditional distribution) where all input distributions are Gaussian and all transformation are affine, all intermediate distributions have to be Gaussian.

These points also help you with doing the second part (working out mean, covariance) in an economic fashion. A word of warning: there are books full of Gaussian formulae. You can learn them by heart and solve all these problems by long pages of algebra. But that’s not fun, it’s error-prone, and it’s simply very embarrassing in the long run. In the following, I will give examples for how to use these points, and then hope for really elegant solutions to exercise sheet problems.

2.1 Linear Transformations. Marginal Distributions

The first example is simple. Suppose $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, $\mathbf{y} = \mathbf{E}\mathbf{x} + \mathbf{b}$. What is $P(\mathbf{y})$? First step: $P(\mathbf{y})$ is Gaussian, because Gaussians are closed under affine transformations, and we start with a Gaussian.

At this point, you might be a bit sceptical. What if \mathbf{E} has a row that is zero? More generally, what if the rank of \mathbf{E} is less than the number of rows? The covariance of $P(\mathbf{y})$ will then be

singular, and (1) cannot be true (the determinant is zero). This is an example of a *degenerate Gaussian*. Strictly speaking, most closure properties are true only if we extend the Gaussian family by these members. Geometrically, a degenerate Gaussian is just a Gaussian that lives in a space that is too high-dimensional for it to have mass there. It is like a flat pancake in this space. Although these pancake Gaussians are important to understand certain models later on, we can avoid them by always starting with *proper Gaussians* (positive definite covariance), and requiring that \mathbf{E} in linear transforms has full rank (equal to number of rows).

Second step: use (2) to deduce mean and covariance of $P(\mathbf{y})$. Note that there is no need to compute the normalization constant of $P(\mathbf{y})$, because it is determined by mean and covariance.

What about the marginal distribution of \mathbf{x}_I , $I \subset \{1, \dots, n\}$? That is a special case, since $\mathbf{x}_I = \mathbf{I}_I \mathbf{x}$. A word on notation (more about this in Lecture 4). $\mathbf{x}_I = (x_i)_{i \in I}$, I is an ordered index set. $\mathbf{A}_{I,J} = (a_{ij})_{i \in I, j \in J}$, I for rows, J for columns. I also use the shorthand $\mathbf{A}_I := \mathbf{A}_{I,I}$. *Selection matrices* are obtained by applying indices to the identity matrix. Make yourself clear that $\mathbf{x}_I = \mathbf{I}_I \mathbf{x}$, or that $\mathbf{I}_{\cdot, I} \mathbf{x}_I$ is the vector of size n which is equal to \mathbf{x} at I , but zero elsewhere. Since marginalization is a special case of a linear transformation (for continuous variables), then $P(\mathbf{x}_I) = N(\boldsymbol{\mu}_I, \boldsymbol{\Sigma}_I)$.

3 Conditional Distribution. Natural Parameterization

In this section, we will see that conditional Gaussians (product rule) are a bit harder to determine than marginal Gaussians (sum rule). In a sense, this is an artefact of the parameterization we use (mean, covariance): in other *natural* parameters, conditioning is equally simple. In practice, we have one or the other, so we need to know how to convert between them painlessly. In this section, I will demonstrate the general ideas of the previous section, thereby obtaining some important concepts (Schur complement) and identities (partitioned inverse equations, determinant identity, Woodbury formula).

3.1 Towards Natural Parameterization

Given that $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, $I \subset \{1, \dots, n\}$, let $R = \{1, \dots, n\} \setminus I$. What is $P(\mathbf{x}_I | \mathbf{x}_R)$? Suppose we did not know that $P(\mathbf{x}_I | \mathbf{x}_R)$ is Gaussian. But we know that $P(\mathbf{x}) = P(\mathbf{x}_I, \mathbf{x}_R) = P(\mathbf{x}_I | \mathbf{x}_R) P(\mathbf{x}_R)$, so if we write down $P(\mathbf{x})$, the part depending on \mathbf{x}_I will be the conditional distribution. From (1), we see that we have to partition the inverse: $\mathbf{A} = \boldsymbol{\Sigma}^{-1}$. Then,

$$P(\mathbf{x}_I | \mathbf{x}_R) P(\mathbf{x}_R) \doteq e^{-(1/2)((\mathbf{x}_I - \boldsymbol{\mu}_I)^T \mathbf{A}_I (\mathbf{x}_I - \boldsymbol{\mu}_I)^T + 2(\mathbf{x}_R - \boldsymbol{\mu}_R)^T \mathbf{A}_{R,I} (\mathbf{x}_I - \boldsymbol{\mu}_I))}.$$

Note that I just write the quadratic and linear part w.r.t. \mathbf{x}_I . To determine what $P(\mathbf{x}_I | \mathbf{x}_R)$ is, I don't need to know the constant part: it'll be part of $P(\mathbf{x}_R)$. I *directly see* now that $P(\mathbf{x}_I | \mathbf{x}_R)$ must be Gaussian: it has the form $e^{-q(\mathbf{x}_I)}$, where q is a quadratic function. Moreover, I can read off its mean and covariance by matching the expression against (1). First, $\text{Cov}[\mathbf{x}_I | \mathbf{x}_R] = \mathbf{A}_I^{-1}$. Second,

$$(\mathbf{x}_R - \boldsymbol{\mu}_R)^T \mathbf{A}_{R,I} (\mathbf{x}_I - \boldsymbol{\mu}_I) = ((\mathbf{x}_R - \boldsymbol{\mu}_R)^T \mathbf{A}_{R,I} \mathbf{A}_I^{-1}) \mathbf{A}_I (\mathbf{x}_I - \boldsymbol{\mu}_I),$$

so that

$$E[\mathbf{x}_I|\mathbf{x}_R] = \boldsymbol{\mu}_I + \mathbf{A}_I^{-1}\mathbf{A}_{I,R}(\mathbf{x}_R - \boldsymbol{\mu}_R).$$

It is a good exercise to multiply out the quadratic for $P(\mathbf{x}_I|\mathbf{x}_R)$ from (1) and match terms against the expression for $P(\mathbf{x}_I|\mathbf{x}_R)P(\mathbf{x}_R)$ above, to retrieve these identities.

This means that while marginalization (sum rule) is simple when you know mean and covariance matrix, things are harder for the conditional distribution. However, they are *simple* again in a different parameterization. The one you know, $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$, is called *mean parameterization* (or *moment parameterization*). The novel one is obtained from this one as:

$$\mathbf{A} = \boldsymbol{\Sigma}^{-1}, \quad \mathbf{r} = \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}.$$

This is called *natural parameterization* (or *canonical parameterization*). The ultimate reason why Gaussian identities look somewhat messy, and unfortunately also why Gaussian information propagation is hard to do in a numerically stable way, is that we have to constantly convert between these parameterizations: inference involves both the sum and the product rule, and either is elementary in one of the parameterizations.

For a number of important reasons, among them numerical stability and processing time, it is *not* a good idea to constantly convert between parameterizations in practice. It's a useful way to think about the problem, but in practice you want to stay in one of the parameterizations and use direct identities. In the following, we will do that for the conditional distribution, given that we have mean parameters only. This way, you'll see some examples how these identities are obtained, illustrating ideas mentioned above already. We'll also derive some important general matrix identities this way, and see some "derivation saving" tricks in action.

3.2 Conditional Distribution by Sampling Argument

We derived mean and covariance of $P(\mathbf{x}_I|\mathbf{x}_R)$ above, but the result is in terms of natural parameters. What if I only have mean parameters? I could convert them, but that involves a matrix inversion. We'll see in Lecture 4 how *bad* matrix inversion is in practice. It's numerically unstable and expensive. You want to avoid it whenever possible. This is not something they tell you in your linear algebra lectures, but it is *Numerical Mathematics 101*.

To get to $P(\mathbf{x}_I|\mathbf{x}_R)$ directly, we can use a sampling argument. Let's first get rid of means by transforming $\mathbf{y} = \mathbf{x} - \boldsymbol{\mu}$, adding it back in later. We know that $P(\mathbf{y}) = P(\mathbf{y}_I|\mathbf{y}_R)P(\mathbf{y}_R)$, which tells us how to sample \mathbf{y} :

1. Draw $\mathbf{y}_R \sim N(\mathbf{0}, \boldsymbol{\Sigma}_R)$
2. Draw $\mathbf{y}_I \sim P(\mathbf{y}_I|\mathbf{y}_R) = N(?, ?)$

And we know the outcome, namely $\mathbf{y} \sim N(\mathbf{0}, \boldsymbol{\Sigma})$. We use an *ansatz*, which means that we guess a form for the solution. The ansatz is that $\mathbf{y}_I = \mathbf{u} + \mathbf{B}\mathbf{y}_R$, where $\mathbf{u} \sim N(\mathbf{0}, \mathbf{C})$ is independent of \mathbf{y}_R . But this is just a fancy way of saying that \mathbf{y} is a linear transform of \mathbf{y}_R and \mathbf{u} :

$$\begin{bmatrix} \mathbf{y}_I \\ \mathbf{y}_R \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{B} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{y}_R \end{bmatrix}. \quad (3)$$

We are on known ground now:

$$\text{Cov}[\mathbf{y}] = \begin{bmatrix} \mathbf{I} & \mathbf{B} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{C} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_R \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{B}^T & \mathbf{I} \end{bmatrix} = \begin{bmatrix} \mathbf{C} + \mathbf{B}\boldsymbol{\Sigma}_R\mathbf{B}^T & \mathbf{B}\boldsymbol{\Sigma}_R \\ \mathbf{B}^T\boldsymbol{\Sigma}_R & \boldsymbol{\Sigma}_R \end{bmatrix} \stackrel{!}{=} \begin{bmatrix} \boldsymbol{\Sigma}_I & \boldsymbol{\Sigma}_{I,R} \\ \boldsymbol{\Sigma}_{R,I} & \boldsymbol{\Sigma}_R \end{bmatrix} \quad (4)$$

Note that the matrix in the middle is block-diagonal, because \mathbf{u} and \mathbf{y}_R are independent. Now, just match terms to deduce \mathbf{B} and \mathbf{C} . First, $\mathbf{B} = \boldsymbol{\Sigma}_{I,R}\boldsymbol{\Sigma}_R^{-1}$. Second,

$$\mathbf{C} = \boldsymbol{\Sigma}_I - \mathbf{B}\boldsymbol{\Sigma}_R\mathbf{B}^T = \boldsymbol{\Sigma}_I - \boldsymbol{\Sigma}_{I,R}\boldsymbol{\Sigma}_R^{-1}\boldsymbol{\Sigma}_{R,I} =: \boldsymbol{\Sigma}/\boldsymbol{\Sigma}_R. \quad (5)$$

The form of \mathbf{C} in terms of blocks of $\boldsymbol{\Sigma}$ is the *Schur complement*, which we revisit below.

Now what is $P(\mathbf{x}_I|\mathbf{y}_R)$? Just a word on why the ansatz and the sampling argument worked out. The ansatz is a rule for sampling \mathbf{y} . Since we start with Gaussians and do linear transforms, \mathbf{y} must be Gaussian. Moreover, we managed to find \mathbf{B} and \mathbf{C} so that the “target parameters” $\boldsymbol{\mu}$, $\boldsymbol{\Sigma}$ are obtained. By the ansatz, $P(\mathbf{x}_I|\mathbf{y}_R)$ is Gaussian, and it is just one more step to its mean parameters. First, $\mathbb{E}[\mathbf{y}_I|\mathbf{y}_R] = \mathbb{E}[\mathbf{u} + \mathbf{B}\mathbf{y}_R|\mathbf{y}_R] = \mathbf{B}\mathbf{y}_R = \boldsymbol{\Sigma}_{I,R}\boldsymbol{\Sigma}_R^{-1}(\mathbf{x}_R - \boldsymbol{\mu}_R)$. Second,

$$\text{Cov}[\mathbf{x}_I|\mathbf{x}_R] = \text{Cov}[\mathbf{u}] = \mathbf{C} = \boldsymbol{\Sigma}/\boldsymbol{\Sigma}_R,$$

the Schur complement.

3.3 Schur Complement. Partitioned Inverse Equations

Recall that if $P(\mathbf{x}_I, \mathbf{x}_R) = N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then the Schur complement (5) is the covariance matrix of the conditional distribution $P(\mathbf{x}_I|\mathbf{x}_R)$. It exists and is invertible whenever $\boldsymbol{\Sigma}$ is invertible. The notation $\boldsymbol{\Sigma}/\boldsymbol{\Sigma}_R$ suggests something like “ $\boldsymbol{\Sigma} = (\boldsymbol{\Sigma}/\boldsymbol{\Sigma}_R)\boldsymbol{\Sigma}_R$ ”, which does not hold of course. It does however hold for the normalization constants:

$$|\boldsymbol{\Sigma}| = |\boldsymbol{\Sigma}/\boldsymbol{\Sigma}_R| |\boldsymbol{\Sigma}_R|. \quad (6)$$

Namely, recall the decomposition (4) of $\boldsymbol{\Sigma}$, and note that

$$\left| \begin{bmatrix} \mathbf{I} & \mathbf{B} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \right| = 1,$$

since the determinant of a triangular matrix is the product of the diagonal entries. (6) follows, since the right hand side is the determinant of the block-diagonal matrix in the middle in (4). Another way to check (6) is to recall that $P(\mathbf{x}_I, \mathbf{x}_R) = P(\mathbf{x}_I|\mathbf{x}_R)P(\mathbf{x}_R)$, then to match normalization constants.

Note that $\boldsymbol{\Sigma}/\boldsymbol{\Sigma}_R$ is defined over arbitrary square invertible matrices, they do not have to be symmetric. (6) holds in general as well. For the special case

$$\boldsymbol{\Sigma} = \begin{bmatrix} \mathbf{I} & -\mathbf{U} \\ \mathbf{V} & \mathbf{I} \end{bmatrix},$$

where the \mathbf{I} can be of different size, we have that

$$|\boldsymbol{\Sigma}/\boldsymbol{\Sigma}_R| = |\mathbf{I} + \mathbf{U}\mathbf{V}| = |\mathbf{I} + \mathbf{V}\mathbf{U}| = |\boldsymbol{\Sigma}/\boldsymbol{\Sigma}_I|,$$

a very useful determinant identity.

At this point, we have determined the conditional Gaussian in two different ways. We could match the results from Section 3.1 and Section 3.2 to obtain the partitioned inverse equations. But it is simpler to extend arguments from Section 3.2. Let's compute the inverse of the first equation in (4). Note that the matrices left and right have simple inverses:

$$\begin{bmatrix} \mathbf{I} & \mathbf{B} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{I} & -\mathbf{B} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}.$$

Therefore,

$$\Sigma^{-1} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\mathbf{B}^T & \mathbf{I} \end{bmatrix} \begin{bmatrix} (\Sigma/\Sigma_R)^{-1} & \mathbf{0} \\ \mathbf{0} & \Sigma_R^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{I} & -\mathbf{B} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \stackrel{!}{=} \begin{bmatrix} \mathbf{A}_I & \mathbf{A}_{I,R} \\ \mathbf{A}_{R,I} & \mathbf{A}_R \end{bmatrix},$$

where $\mathbf{B} = \Sigma_{I,R} \Sigma_R^{-1}$. Multiplying these out (which you should do as exercise) results in the *partitioned inverse equations*. They just do what the name suggests. For a matrix Σ , they are explicit formulae of blocks of the inverse in terms of corresponding blocks of Σ . Once more, they hold for any nonsingular Σ .

From the partitioned inverse equations, it is a simple step to the *Woodbury formula*. We simply apply the former with I and R interchanged, then equate terms. The significance of the Woodbury formula is that it allows you update the inverse of a matrix to which a low rank matrix is added. I strongly suggest that you derive the formula from the partitioned inverse equations as an exercise.

3.4 Another Sampling Argument

Here's another example for a sampling argument. We know that

$$N(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)N(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) = N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})C, \quad (7)$$

but what is C ? Recall that in most cases, we don't have to know: all that matters is $\boldsymbol{\mu}$, $\boldsymbol{\Sigma}$. But sometimes we have to know, examples will come up in later lectures. I invite you to do this very messy computation by hand, then to appreciate the following simple argument.

We can see (7) as a statement about three variables: \mathbf{x} , $\boldsymbol{\mu}_1$, $\boldsymbol{\mu}_2$. Formally, it does not matter whether we write $N(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ or $N(\boldsymbol{\mu}_1|\mathbf{x}, \boldsymbol{\Sigma}_1)$. Suppose we take the integral over \mathbf{x} on both sides. On the right side, we get C , since $\int N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x} = 1$. On the left side, flipping arguments in the first term, we have

$$\int N(\boldsymbol{\mu}_1|\mathbf{x}, \boldsymbol{\Sigma}_1)N(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) d\mathbf{x}.$$

But this is a sampling recipe for the distribution of $\boldsymbol{\mu}_1$ given $\boldsymbol{\mu}_2$:

1. Sample $\mathbf{x} \sim N(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$.
2. Sample $\boldsymbol{\mu}_1 \sim N(\boldsymbol{\mu}_1|\mathbf{x}, \boldsymbol{\Sigma}_1)$. Drop \mathbf{x} .

The combined sampling rule is $\boldsymbol{\mu}_1 = \mathbf{x} + \mathbf{v}$, $\mathbf{v} \sim N(\mathbf{0}, \boldsymbol{\Sigma}_1)$ independent of \mathbf{x} . But this means that $\boldsymbol{\mu}_1|\boldsymbol{\mu}_2 \sim N(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2)$, so that $C = N(\boldsymbol{\mu}_1|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2)$.

4 Linear-Gaussian Model

Recall the linear-Gaussian model:

$$\mathbf{y} = \mathbf{X}\mathbf{u} + \boldsymbol{\varepsilon}, \quad \mathbf{u} \sim N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0), \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \boldsymbol{\Psi}).$$

This is a fundamental latent variable model: \mathbf{u} is latent, \mathbf{y} is observed. For modelling data, what matters is the marginal distribution of \mathbf{y} , whose structure is determined by the dimensionality of \mathbf{u} , its prior, the mapping \mathbf{X} , and the noise covariance $\boldsymbol{\Psi}$. We can use the *tower formulae* to compute the marginal distribution. Here is the proof for the covariance. Let $\mathbf{v} = \mathbb{E}[\mathbf{y}|\mathbf{u}]$. Then,

$$\begin{aligned} \mathbb{E}[\text{Cov}[\mathbf{y}|\mathbf{u}]] &= \mathbb{E}[\mathbb{E}[\mathbf{y}\mathbf{y}^T|\mathbf{u}] - \mathbf{v}\mathbf{v}^T] \stackrel{*}{=} \mathbb{E}[\mathbf{y}\mathbf{y}^T] - \mathbb{E}[\mathbf{v}\mathbf{v}^T] = \text{Cov}[\mathbf{y}] + \mathbb{E}[\mathbf{y}]\mathbb{E}[\mathbf{y}]^T - \mathbb{E}[\mathbf{v}\mathbf{v}^T] \\ &\stackrel{*}{=} \text{Cov}[\mathbf{y}] + \mathbb{E}[\mathbf{v}]\mathbb{E}[\mathbf{v}]^T - \mathbb{E}[\mathbf{v}\mathbf{v}^T] = \text{Cov}[\mathbf{y}] - \text{Cov}[\mathbb{E}[\mathbf{y}|\mathbf{u}]]. \end{aligned}$$

Here, we used the tower formulae for expectation at each point “*” (make sure you understand each of these applications). For the linear-Gaussian model, $\mathbb{E}[\mathbf{y}] = \mathbb{E}[\mathbf{X}\mathbf{u}] = \mathbf{X}\boldsymbol{\mu}_0$, while

$$\text{Cov}[\mathbf{y}] = \text{Cov}[\mathbf{X}\mathbf{u}] + \mathbb{E}[\text{Cov}[\mathbf{y}|\mathbf{u}]] = \text{Cov}[\mathbf{X}\mathbf{u}] + \mathbb{E}[\text{Cov}[\boldsymbol{\varepsilon}]] = \mathbf{X}\boldsymbol{\Sigma}_0\mathbf{X}^T + \boldsymbol{\Psi}.$$

Another way is to obtain the joint distribution by using

$$\begin{bmatrix} \mathbf{u} \\ \mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{I} \\ \mathbf{X} \end{bmatrix} \mathbf{u} + \begin{bmatrix} \mathbf{0} \\ \boldsymbol{\varepsilon} \end{bmatrix}$$

together with what we know about linear transforms of Gaussians. Note that \mathbf{u} and $\boldsymbol{\varepsilon}$ are independent. For example, $\text{Cov}[\mathbf{u}, \mathbf{y}] = \boldsymbol{\Sigma}_0\mathbf{X}^T$.

Let us compute the posterior $P(\mathbf{u}|\mathbf{y})$ for this model, along the two different ways we derived above. We already know it must be Gaussian. Moreover,

$$\begin{aligned} \text{Cov}[\mathbf{u}|\mathbf{y}] &= \text{Cov}[(\mathbf{u}, \mathbf{y})] / \text{Cov}[\mathbf{y}] = \text{Cov}[\mathbf{u}] - \text{Cov}[\mathbf{u}, \mathbf{y}]\text{Cov}[\mathbf{y}]^{-1}\text{Cov}[\mathbf{u}, \mathbf{y}]^T \\ &= \boldsymbol{\Sigma}_0 - \boldsymbol{\Sigma}_0\mathbf{X}^T(\mathbf{X}\boldsymbol{\Sigma}_0\mathbf{X}^T + \boldsymbol{\Psi})^{-1}\mathbf{X}\boldsymbol{\Sigma}_0. \end{aligned}$$

Note that $\text{Cov}[(\mathbf{u}, \mathbf{y})]$ is the joint covariance, while $\text{Cov}[\mathbf{u}, \mathbf{y}]$ is the cross-covariance (a block within the former). And

$$\mathbb{E}[\mathbf{u}|\mathbf{y}] = \mathbb{E}[\mathbf{u}] + \text{Cov}[\mathbf{u}, \mathbf{y}]\text{Cov}[\mathbf{y}]^{-1}(\mathbf{y} - \mathbb{E}[\mathbf{y}]) = \boldsymbol{\mu}_0 + \boldsymbol{\Sigma}_0\mathbf{X}^T(\mathbf{X}\boldsymbol{\Sigma}_0\mathbf{X}^T + \boldsymbol{\Psi})^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\mu}_0).$$

We can also match terms in the joint distribution $P(\mathbf{y}, \mathbf{u}) = P(\mathbf{y}|\mathbf{u})P(\mathbf{u})$. As shown in the lecture, this leads to

$$\text{Cov}[\mathbf{u}|\mathbf{y}] = (\mathbf{X}^T\boldsymbol{\Psi}^{-1}\mathbf{X} + \boldsymbol{\Sigma}_0^{-1})^{-1}, \quad \mathbb{E}[\mathbf{u}|\mathbf{y}] = (\mathbf{X}^T\boldsymbol{\Psi}^{-1}\mathbf{X} + \boldsymbol{\Sigma}_0^{-1})^{-1}(\mathbf{X}^T\boldsymbol{\Psi}^{-1}\mathbf{y} + \boldsymbol{\Sigma}_0^{-1}\boldsymbol{\mu}_0).$$

These are different expressions for the same thing. In general, the former set is more memorable and often easier to compute. The latter is more direct to obtain (you just have to match terms), but contains more inverses. Still, in some special cases, it might be the preferred way for this computation.

An important take-home message for the linear-Gaussian model is as follows. Suppose that $\mathbf{u} \in \mathbb{R}^n$, $\mathbf{y} \in \mathbb{R}^m$. Then, you can always compute posterior quantities by doing expensive (superlinear) computations, such as inverses, in the *smaller* number only: $\min\{m, n\}$. The main vehicle to formally get from one set of expression to the other is the Woodbury formula.