

Handout Probabilistic Graphical Models: Gaussian Computations

1 Definitions

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

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

Here, $\boldsymbol{\mu} = E[\mathbf{x}]$, $\boldsymbol{\Sigma} = \text{Cov}[\mathbf{x}]$

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.

- 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:

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

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

3.2 Conditional Distribution by Sampling Argument

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 .

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

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

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

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

==>

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

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

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

where $\mathbf{B} = \boldsymbol{\Sigma}_{I,R}\boldsymbol{\Sigma}_R^{-1}$

4 Linear-Gaussian Model

the linear-Gaussian model:

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

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 “*”. For the linear-Gaussian model, $\mathbb{E}[\mathbf{y}] = \mathbb{E}[\mathbf{X}\mathbf{u}] = \mathbf{X}\boldsymbol{\mu}_0$, while

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

Another way is to obtain the joint distribution by using

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

together with what we know about linear transforms of Gaussians. Note that \mathbf{u} and ε 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,

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

==>

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

We can also match terms in the joint distribution $P(\mathbf{y}, \mathbf{u}) = P(\mathbf{y}|\mathbf{u})P(\mathbf{u})$. ==>

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

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.