**1.**

$$
\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \\ \mathbf{x}_c \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \\ \boldsymbol{\mu}_c \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} & \Sigma_{ac} \\ \Sigma_{ba} & \Sigma_{bb} & \Sigma_{bc} \\ \Sigma_{ca} & \Sigma_{cb} & \Sigma_{cc} \end{pmatrix}, \quad \Sigma^{-1} = \Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} & \Lambda_{ac} \\ \Lambda_{ba} & \Lambda_{bb} & \Lambda_{bc} \\ \Lambda_{ca} & \Lambda_{cb} & \Lambda_{cc} \end{pmatrix}.
$$

$$
\begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \\ \mathbf{x}_c \end{pmatrix} \sim \mathcal{N}\left( \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \\ \boldsymbol{\mu}_c \end{pmatrix}, \Sigma \right).
$$

Since $\mathbf{x}_c$ is marginalized, $\Lambda_{ac}, \Lambda_{bc}, \Lambda_{ca}, \Lambda_{cb}, \Lambda_{cc}$ do not affect the marginal distribution in $\mathbf{x}_a$ and $\mathbf{x}_b$. Thus,

$$
\begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \sim \mathcal{N}\left( \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}, \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \right).
$$

$$
\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left\{ -\tfrac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}.
$$

$$
-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})
$$

$$
= -\frac{1}{2} \begin{pmatrix} \mathbf{x}_a - \boldsymbol{\mu}_a \\ \mathbf{x}_b - \boldsymbol{\mu}_b \end{pmatrix}^\top \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix} \begin{pmatrix} \mathbf{x}_a - \boldsymbol{\mu}_a \\ \mathbf{x}_b - \boldsymbol{\mu}_b \end{pmatrix}
$$

$$
= -\tfrac{1}{2} (\mathbf{x}_a - \boldsymbol{\mu}_a)^\top \Lambda_{aa} (\mathbf{x}_a - \boldsymbol{\mu}_a) - \tfrac{1}{2} (\mathbf{x}_a - \boldsymbol{\mu}_a)^\top \Lambda_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b)
$$

$$
\quad - \tfrac{1}{2} (\mathbf{x}_b - \boldsymbol{\mu}_b)^\top \Lambda_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a) - \tfrac{1}{2} (\mathbf{x}_b - \boldsymbol{\mu}_b)^\top \Lambda_{bb} (\mathbf{x}_b - \boldsymbol{\mu}_b)
$$

$$
= -\tfrac{1}{2} \Big[ (\mathbf{x}_a - \boldsymbol{\mu}_a)^\top \Lambda_{aa} (\mathbf{x}_a - \boldsymbol{\mu}_a) + (\mathbf{x}_b - \boldsymbol{\mu}_b)^\top \Lambda_{bb} (\mathbf{x}_b - \boldsymbol{\mu}_b)
$$

$$
\quad + (\mathbf{x}_a - \boldsymbol{\mu}_a)^\top \Lambda_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) + (\mathbf{x}_b - \boldsymbol{\mu}_b)^\top \Lambda_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a) \Big]
$$

$$
= -\tfrac{1}{2} \Big[ (\mathbf{x}_a - \boldsymbol{\mu}_a)^\top \Lambda_{aa} (\mathbf{x}_a - \boldsymbol{\mu}_a) + (\mathbf{x}_b - \boldsymbol{\mu}_b)^\top \Lambda_{bb} (\mathbf{x}_b - \boldsymbol{\mu}_b)
$$

$$
\quad + 2 (\mathbf{x}_a - \boldsymbol{\mu}_a)^\top \Lambda_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \Big]
$$

$$
= -\tfrac{1}{2} (\mathbf{x}_a - \boldsymbol{\mu}_a)^\top \Lambda_{aa} (\mathbf{x}_a - \boldsymbol{\mu}_a) - \tfrac{1}{2} (\mathbf{x}_b - \boldsymbol{\mu}_b)^\top \Lambda_{bb} (\mathbf{x}_b - \boldsymbol{\mu}_b)
$$

$$
\quad - (\mathbf{x}_a - \boldsymbol{\mu}_a)^\top \Lambda_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b).
$$

Let

$$
\mathbf{a} = \mathbf{x}_a - \boldsymbol{\mu}_a \quad \text{and} \quad \mathbf{b} = \mathbf{x}_b - \boldsymbol{\mu}_b.
$$

$$-\frac{1}{2}\left(\mathbf{x}_a - \boldsymbol{\mu}_a\right)^\top \Lambda_{aa}\left(\mathbf{x}_a - \boldsymbol{\mu}_a\right) \;-\; \frac{1}{2}\left(\mathbf{x}_b - \boldsymbol{\mu}_b\right)^\top \Lambda_{bb}\left(\mathbf{x}_b - \boldsymbol{\mu}_b\right)$$

$$-\left(\mathbf{x}_a - \boldsymbol{\mu}_a\right)^\top \Lambda_{ab}\left(\mathbf{x}_b - \boldsymbol{\mu}_b\right)$$

$$= -\frac{1}{2}\,\mathbf{a}^\top \Lambda_{aa}\,\mathbf{a} \;-\; \frac{1}{2}\,\mathbf{b}^\top \Lambda_{bb}\,\mathbf{b} \;-\; \mathbf{a}^\top \Lambda_{ab}\,\mathbf{b}$$

$$= -\frac{1}{2}\,\mathbf{a}^\top \Lambda_{aa}\,\mathbf{a} \;-\; \mathbf{a}^\top \Lambda_{ab}\,\mathbf{b} \qquad \left(\tfrac{1}{2}\,\mathbf{b}^\top \Lambda_{bb}\,\mathbf{b}\text{ is independent of }\mathbf{a}\right)$$

$$= -\tfrac{1}{2}\left(\mathbf{x}_a - \boldsymbol{\mu}_a\right)^\top \Lambda_{aa}\left(\mathbf{x}_a - \boldsymbol{\mu}_a\right) \;-\; \left(\mathbf{x}_a - \boldsymbol{\mu}_a\right)^\top \Lambda_{ab}\,\mathbf{b}$$

$$= -\tfrac{1}{2}\,\mathbf{x}_a^\top \Lambda_{aa}\,\mathbf{x}_a \;+\; \mathbf{x}_a^\top \Lambda_{aa}\,\boldsymbol{\mu}_a \;-\; \tfrac{1}{2}\,\boldsymbol{\mu}_a^\top \Lambda_{aa}\,\boldsymbol{\mu}_a \;-\; \mathbf{x}_a^\top \Lambda_{ab}\,\mathbf{b} \;+\; \boldsymbol{\mu}_a^\top \Lambda_{ab}\,\mathbf{b}$$

$$= -\tfrac{1}{2}\,\mathbf{x}_a^\top \Lambda_{aa}\,\mathbf{x}_a \;+\; \mathbf{x}_a^\top \left[\Lambda_{aa}\,\boldsymbol{\mu}_a \;-\; \Lambda_{ab}\,\mathbf{b}\right] \;-\; \tfrac{1}{2}\,\boldsymbol{\mu}_a^\top \Lambda_{aa}\,\boldsymbol{\mu}_a \;+\; \boldsymbol{\mu}_a^\top \Lambda_{ab}\,\mathbf{b}. \longrightarrow (1)$$

$$-\tfrac{1}{2}\left(\mathbf{x}_a - \boldsymbol{\mu}\right)^\top \Lambda_{aa}\left(\mathbf{x}_a - \boldsymbol{\mu}\right) \;=\; -\tfrac{1}{2}\,\mathbf{x}_a^\top \Lambda_{aa}\,\mathbf{x}_a \;+\; \mathbf{x}_a^\top \Lambda_{aa}\,\boldsymbol{\mu} \;-\; \tfrac{1}{2}\,\boldsymbol{\mu}^\top \Lambda_{aa}\,\boldsymbol{\mu}. \longrightarrow (2)$$

From (1) and (2), comparing the term $\mathbf{x}_a^\top \Lambda_{aa}\,\boldsymbol{\mu}$ with $\mathbf{x}_a^\top\left[\Lambda_{aa}\,\boldsymbol{\mu}_a - \Lambda_{ab}\,\mathbf{b}\right]$ leads to:

$$\Lambda_{aa}\,\boldsymbol{\mu} \;=\; \Lambda_{aa}\,\boldsymbol{\mu}_a \;-\; \Lambda_{ab}\,\mathbf{b},$$

$$\boldsymbol{\mu} \;=\; \Lambda_{aa}^{-1}\left(\Lambda_{aa}\,\boldsymbol{\mu}_a - \Lambda_{ab}\,\mathbf{b}\right) \;=\; \boldsymbol{\mu}_a \;-\; \Lambda_{aa}^{-1}\Lambda_{ab}\,\mathbf{b}.$$

Substituting $\mathbf{b} \;=\; \mathbf{x}_b - \boldsymbol{\mu}_b$ gives the conditional mean:

$$\boldsymbol{\mu}_{a|b} \;=\; \boldsymbol{\mu}_a \;-\; \Lambda_{aa}^{-1}\Lambda_{ab}\left(\mathbf{x}_b - \boldsymbol{\mu}_b\right).$$

Focusing on the quadratic term in $\mathbf{x}_a$:

$$-\tfrac{1}{2}\left(\mathbf{x}_a - \boldsymbol{\mu}_a\right)^\top \Lambda_{aa}\left(\mathbf{x}_a - \boldsymbol{\mu}_a\right)$$

shows that $\Lambda_{aa}$ is the precision w.r.t $\mathbf{x}_a$, thus

$$\Sigma_{a|b}^{-1} \;=\; \Lambda_{aa} \quad \Longrightarrow \quad \Sigma_{a|b} \;=\; \Lambda_{aa}^{-1}.$$

Therefore, the conditional distribution is

$$p\left(\mathbf{x}_a \mid \mathbf{x}_b\right) \;=\; \mathcal{N}\left(\boldsymbol{\mu}_a - \Lambda_{aa}^{-1}\Lambda_{ab}\left(\mathbf{x}_b - \boldsymbol{\mu}_b\right),\ \Lambda_{aa}^{-1}\right).$$

## 2.

If $(A + BCD)^{-1} = A^{-1} - A^{-1}B\left(C^{-1} + DA^{-1}B\right)^{-1}DA^{-1}$, then by the definition of the inverse $AA^{-1}x = x$, for every vector $x$ it must be that $(A + BCD)\left[A^{-1} - A^{-1}B\left(C^{-1} + DA^{-1}B\right)^{-1}DA^{-1}\right]x \;=\; x$

Let $y = \left[A^{-1} - A^{-1}B\left(C^{-1} + DA^{-1}B\right)^{-1}DA^{-1}\right]x$

$$y = \left[A^{-1} - A^{-1}B\left(C^{-1} + DA^{-1}B\right)^{-1}DA^{-1}\right]x$$

$$= A^{-1}x \;-\; A^{-1}B\left(C^{-1} + DA^{-1}B\right)^{-1}DA^{-1}x$$

$$(A + BCD)\cdot y$$

$$= (A + BCD)\cdot \left[A^{-1}x \;-\; A^{-1}B\left(C^{-1} + DA^{-1}B\right)^{-1}DA^{-1}x\right]$$

$$= (A + BCD)\left[A^{-1}x\right] \;-\; (A + BCD)\left[A^{-1}B\left(C^{-1} + DA^{-1}B\right)^{-1}DA^{-1}x\right]$$

$$(A + BCD)\left[A^{-1}x\right]$$
$$= A\,A^{-1}x + BCD\left[A^{-1}x\right]$$
$$= x + BCDA^{-1}x$$

$$(A + BCD)\left[A^{-1}B\left(C^{-1} + D\,A^{-1}B\right)^{-1}D\,A^{-1}x\right]$$
$$= \left[A\,A^{-1}B + BCD\,A^{-1}B\right]\left(C^{-1} + D\,A^{-1}B\right)^{-1}D\,A^{-1}x$$
$$= \left[B + BCD\,A^{-1}B\right]\left(C^{-1} + D\,A^{-1}B\right)^{-1}D\,A^{-1}x$$
$$= B\left[1 + CD\,A^{-1}B\right]\left(C^{-1} + D\,A^{-1}B\right)^{-1}D\,A^{-1}x$$
$$= B\left[CC^{-1} + CD\,A^{-1}B\right]\left(C^{-1} + D\,A^{-1}B\right)^{-1}D\,A^{-1}x$$
$$= B\left[C(C^{-1} + D\,A^{-1}B)\right]\left(C^{-1} + D\,A^{-1}B\right)^{-1}D\,A^{-1}x$$
$$= BC\left[(C^{-1} + D\,A^{-1}B)\left(C^{-1} + D\,A^{-1}B\right)^{-1}\right]D\,A^{-1}x$$
$$= BCDA^{-1}x$$

$$(A + BCD)\cdot y$$
$$= \left[x + BCDA^{-1}x\right] - BCDA^{-1}x$$
$$= x$$

As for every vector $x$,

$$(A + BCD)\left[A^{-1} - A^{-1}B\left(C^{-1} + D\,A^{-1}B\right)^{-1}D\,A^{-1}\right]x = x,$$

and by definition

$$AA^{-1}x = x,$$

it follows that

$$A^{-1} - A^{-1}B\left(C^{-1} + D\,A^{-1}B\right)^{-1}D\,A^{-1} = (A + BCD)^{-1}$$

If $\mathbf{A}$ is an $n \times n$ diagonal matrix, then computing $\mathbf{A}^{-1}$ simply involves taking reciprocals of the diagonal entries. Even if $n$ is very large, only $n$ such operations are needed, eliminating the need for expensive factorizations or dense matrix inversions. As $\mathbf{A}^{-1}$ is known, its computational cost can be considered negligible.

Let $\mathbf{B}$ be an $n \times k$ matrix (with $n \gg k$) and $\mathbf{C}$ be an invertible $k \times k$ matrix. Instead of inverting the full $n \times n$ matrix $\mathbf{A} + \mathbf{BCD}$, the Woodbury identity can be applied which shifts the inversion problem to a smaller $k \times k$ matrix:

$$\mathbf{C}^{-1} + \mathbf{DA}^{-1}\mathbf{B}.$$

As $k$ is small compared to $n$—for example 3 or 10—this smaller system $3 \times 3$ or $10 \times 10$ is much cheaper to invert in comparison to the $n \times n$ system, especially if $n$ runs into the thousands or millions.

In terms of computational complexity, a naive inversion of $\mathbf{A} + \mathbf{BCD}$ costs $O(n^3)$ when $n \gg k$. With the Woodbury identity and a known $\mathbf{A}^{-1}$, the cost instead involves:

- Multiplications with $\mathbf{B}$ and $\mathbf{D}$, which is on the order of $O(nk^2)$.

- The inversion of the $k \times k$ matrix $\mathbf{C}$, which costs $O(k^3)$.

Because $k$ is smaller in comparison to $n$, this leads to a computational advantage and results in computing the right-hand side of r.h.s of the formula.

**3.**

$$\mathcal{N}(x \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$\mathcal{N}\left(x \mid \mu, \tfrac{1}{\lambda}\right) = \sqrt{\frac{\lambda}{2\pi}} \exp\left(-\frac{\lambda}{2}(x-\mu)^2\right)$$

$$p(\mathbf{X} \mid \mu, \lambda) = \prod_{n=1}^{N} \sqrt{\frac{\lambda}{2\pi}} \exp\left(-\frac{\lambda}{2}(x_n-\mu)^2\right)$$

$$= \left(\sqrt{\frac{\lambda}{2\pi}}\right)^N \exp\left(-\frac{\lambda}{2}\sum_{n=1}^{N}(x_n-\mu)^2\right)$$

$$= \left(\frac{\lambda}{2\pi}\right)^{\frac{N}{2}} \exp\left(-\frac{\lambda}{2}\sum_{n=1}^{N}(x_n-\mu)^2\right)$$

$$= \lambda^{\frac{N}{2}}(2\pi)^{-\frac{N}{2}} \exp\left(-\frac{\lambda}{2}\sum_{n=1}^{N}(x_n-\mu)^2\right).$$

$$p(\mathbf{X} \mid \mu, \lambda) \propto \lambda^{\frac{N}{2}} \exp\left(-\frac{\lambda}{2}\sum_{n=1}^{N}(x_n-\mu)^2\right).$$

$$p(\mu, \lambda) = \mathcal{N}\left(\mu \mid \mu_0, \ (\beta\lambda)^{-1}\right) \mathrm{Gam}(\lambda \mid a, \, b)$$

$$\mathcal{N}\left(\mu \mid \mu_0, \ (\beta\lambda)^{-1}\right) = \sqrt{\frac{\beta\lambda}{2\pi}} \exp\left(-\frac{\beta\lambda}{2}(\mu-\mu_0)^2\right)$$

$$\mathrm{Gam}(\lambda \mid a, \, b) = \frac{b^a}{\Gamma(a)} \lambda^{a-1} \exp(-b\lambda)$$

$$p(\mu, \lambda) = \sqrt{\frac{\beta\lambda}{2\pi}} \exp\left(-\frac{\beta\lambda}{2}(\mu-\mu_0)^2\right) \cdot \frac{b^a}{\Gamma(a)} \lambda^{a-1} \exp(-b\lambda)$$

$$= \sqrt{\frac{\beta}{2\pi}} \lambda^{\frac{1}{2}} \cdot \frac{b^a}{\Gamma(a)} \lambda^{a-1} \cdot \exp\left(-\frac{\beta\lambda}{2}(\mu-\mu_0)^2 - b\lambda\right)$$

$$= \sqrt{\frac{\beta}{2\pi}} \frac{b^a}{\Gamma(a)} \lambda^{\frac{1}{2}+(a-1)} \exp\left(-\frac{\beta\lambda}{2}(\mu-\mu_0)^2 - b\lambda\right)$$

$$= \sqrt{\frac{\beta}{2\pi}} \frac{b^a}{\Gamma(a)} \lambda^{a-\frac{1}{2}} \exp\left(-\frac{\beta\lambda}{2}(\mu-\mu_0)^2 - b\lambda\right)$$

$$= \sqrt{\frac{\beta}{2\pi}} \frac{b^a}{\Gamma(a)} \lambda^{a-\frac{1}{2}} \exp\left(-\frac{\beta\lambda}{2}(\mu^2 - 2\mu\mu_0 + \mu_0^2) - b\lambda\right)$$

$$= \sqrt{\frac{\beta}{2\pi}} \frac{b^a}{\Gamma(a)} \lambda^{a-\frac{1}{2}} \exp\left(-\frac{\beta\lambda}{2}\mu^2 + \beta\lambda\mu\mu_0 - \frac{\beta\lambda}{2}\mu_0^2 - b\lambda\right)$$

$$p(\mu, \lambda) = \sqrt{\frac{\beta}{2\pi}} \frac{b^a}{\Gamma(a)} \lambda^{a-\frac{1}{2}} \exp\left(-\frac{\beta\lambda\mu^2}{2} + \beta\mu_0\lambda\mu - \frac{\beta\lambda}{2}\mu_0^2 - b\lambda\right) \quad \longrightarrow (1)$$

Equation (2.153):

$$p(\mu, \lambda) \propto \left[\lambda^{\frac{1}{2}} \exp\left(-\frac{\lambda\mu^2}{2}\right)\right]^\beta \exp\left(c\lambda\mu - d\lambda\right)$$

$$= \left[\lambda^{\frac{1}{2}}\right]^\beta \left[\exp\left(-\frac{\lambda\mu^2}{2}\right)\right]^\beta \exp\left(c\lambda\mu - d\lambda\right)$$

$$= \lambda^{\frac{\beta}{2}} \exp\left(-\frac{\beta\lambda\mu^2}{2}\right) \exp\left(c\lambda\mu - d\lambda\right)$$

$$= \lambda^{\frac{\beta}{2}} \exp\left(-\frac{\beta\lambda\mu^2}{2} + c\lambda\mu - d\lambda\right) \quad \longrightarrow (2).$$

From (1) and (2):

From the power of $\lambda$ : $\quad a - \dfrac{1}{2} = \dfrac{\beta}{2} \quad \Longrightarrow \quad a = \dfrac{\beta}{2} + \dfrac{1}{2} = \dfrac{\beta+1}{2}.$

From the $\lambda\mu$ term: $\quad \beta\mu_0 = c \quad \Longrightarrow \quad \mu_0 = \dfrac{c}{\beta}.$

From the constant term in the exponent: $\quad -\dfrac{\beta\mu_0^2}{2} - b = -d \quad \Longrightarrow \quad d = b + \dfrac{\beta\mu_0^2}{2}.$

Substituting $\mu_0 = \frac{c}{\beta}$ gives $\frac{\beta\mu_0^2}{2} = \frac{\beta}{2}\left(\frac{c}{\beta}\right)^2 = \frac{c^2}{2\beta}$, so that $b = d - \frac{c^2}{2\beta}.$

$$\lambda^{\beta/2} = \lambda^{a-\frac{1}{2}} \quad \Longrightarrow \quad a - \frac{1}{2} = \frac{\beta}{2} \quad \Longrightarrow \quad a = \frac{1+\beta}{2}$$

The parameters of the distribution are:

$$\mu_0 = \frac{c}{\beta}, \quad a = \frac{\beta+1}{2}, \quad b = d - \frac{c^2}{2\beta}$$

$$p(\mu, \lambda \mid \mathbf{X}) \propto p(\mathbf{X} \mid \mu, \lambda) \cdot p(\mu, \lambda)$$

$$\propto \lambda^{\frac{N}{2}} \exp\left(-\frac{\lambda}{2}\sum_{n=1}^{N}(x_n - \mu)^2\right) \cdot \lambda^{\frac{\beta}{2}} \exp\left(-\frac{\beta\lambda\mu^2}{2} + c\lambda\mu - d\lambda\right)$$

$$\propto \lambda^{\frac{N+\beta}{2}} \exp\left[-\frac{\lambda}{2}\left(\sum_{n=1}^{N}(x_n - \mu)^2 + \beta\mu^2\right) + c\lambda\mu - d\lambda\right]$$

$$\propto \lambda^{\frac{N+\beta}{2}} \exp\left[-\frac{\lambda}{2}\left(\sum_{n=1}^{N}x_n^2 - 2\mu\sum_{n=1}^{N}x_n + N\mu^2 + \beta\mu^2\right) + c\lambda\mu - d\lambda\right]$$

$$\propto \lambda^{\frac{N+\beta}{2}} \exp\left[-\frac{\lambda}{2}\sum_{n=1}^{N}x_n^2 - d\lambda + \lambda\mu\left(\sum_{n=1}^{N}x_n + c\right) - \frac{\lambda}{2}(N+\beta)\mu^2\right]$$

$$\propto \lambda^{\frac{N+\beta}{2}} \exp\left[-\frac{\lambda(N+\beta)\mu^2}{2} + \left(c + \sum_{n=1}^{N}x_n\right)\lambda\mu - \left(d + \frac{1}{2}\sum_{n=1}^{N}x_n^2\right)\lambda\right]$$

$$d' = d + \frac{1}{2}\sum_{n=1}^{N}x_n^2$$

$$\mu_0' = \tfrac{c'}{\beta'} = \frac{c + \sum_{n=1}^{N} x_n}{N + \beta}$$

$$a' = \frac{\beta + N + 1}{2}$$

$$b' = d' - \frac{c'^2}{2\,\beta'}$$

By matching exponents and powers of $\lambda$ and $\mu$, the posterior distribution is the same functional form as the prior distribution but with updated parameters. Under the Normal–Gamma parameterization:

$$p(\mu, \lambda \mid \mathbf{X}) = \mathcal{N}\big(\mu \mid \mu_0', (\beta' \lambda)^{-1}\big)\, \Gamma(\lambda \mid a', b'),$$

where

$$\beta' = \beta + N, \quad c' = c + \sum_{n=1}^{N} x_n, \mu_0' = \tfrac{c'}{\beta'} = \frac{c + \sum_{n=1}^{N} x_n}{N + \beta}, \quad a' = \frac{\beta + N + 1}{2} \quad b' = d' - \frac{c'^2}{2\,\beta'}.$$

The posterior distribution is also a Gaussian–Gamma distribution of the same functional form as the prior, but with updated parameters, confirming that the posterior remains a Gaussian–Gamma distribution of the same functional form as the prior.

# 4.

## a.

Wishart as a conjugate prior to $\Lambda = \Sigma^{-1}$ for Gaussian distribution $\mathcal{N}(\mu, \Lambda^{-1})$

A Wishart distribution over the precision matrix $\Lambda \in \mathbb{R}^{D \times D}$ with parameters $(W, \nu)$ is given by:

$$\mathcal{W}(\Lambda \mid W, \nu) = B(W, \nu)\, |\Lambda|^{\frac{\nu - D - 1}{2}} \exp\!\Big[-\tfrac{1}{2}\operatorname{Tr}(W^{-1}\Lambda)\Big] \longrightarrow (1)$$

$$\mathcal{W}(\Lambda \mid W, \nu) \propto |\Lambda|^{\frac{\nu - D - 1}{2}} \exp\!\Big[-\tfrac{1}{2}\operatorname{Tr}(W^{-1}\Lambda)\Big]$$

$$p(X \mid \mu, \Lambda) \propto |\Lambda|^{\frac{N}{2}} \exp\!\Big[-\tfrac{1}{2}\sum_{n=1}^{N}(x_n - \mu)^T \Lambda (x_n - \mu)\Big]$$

Let $S = \tfrac{1}{N}\sum_{n=1}^{N}(x_n - \mu)\,(x_n - \mu)^T$. Then $\sum_{n=1}^{N}(x_n - \mu)^T \Lambda (x_n - \mu) = N\operatorname{Tr}(\Lambda\, S)$.

$$p(X \mid \mu, \Lambda) \propto |\Lambda|^{\frac{N}{2}} \exp\!\Big[-\tfrac{1}{2}\operatorname{Tr}(N\,S\,\Lambda)\Big] = |\Lambda|^{\frac{N}{2}} \exp\!\Big[-\tfrac{N}{2}\operatorname{Tr}(\Lambda S)\Big]$$

Posterior in $\Lambda$:

$$
\begin{aligned}
p(\Lambda \mid X, W, \nu) &\propto p(X \mid \mu, \Lambda)\, \mathcal{W}(\Lambda \mid W, \nu) \\
&= |\Lambda|^{\frac{N}{2}} \exp\!\Big[-\tfrac{1}{2}\operatorname{Tr}(N\,S\,\Lambda)\Big] \times |\Lambda|^{\frac{\nu - D - 1}{2}} \exp\!\Big[-\tfrac{1}{2}\operatorname{Tr}(W^{-1}\Lambda)\Big] \\
&= |\Lambda|^{\frac{N}{2}} \times |\Lambda|^{\frac{\nu - D - 1}{2}} \cdot \exp\!\Big[-\tfrac{1}{2}\operatorname{Tr}(N\,S\,\Lambda)\Big] \exp\!\Big[-\tfrac{1}{2}\operatorname{Tr}(W^{-1}\Lambda)\Big] \\
&= |\Lambda|^{\frac{N + \nu - D - 1}{2}} \cdot \exp\!\Big[-\tfrac{1}{2}\operatorname{Tr}\big((W^{-1} + N\,S)\,\Lambda\big)\Big] \longrightarrow (2)
\end{aligned}
$$

From (1) and (2),

$$\mathcal{W}\big(\Lambda \mid (W^{-1} + N\,S)^{-1}, N + \nu\big) = |\Lambda|^{\frac{N + \nu - D - 1}{2}} \exp\!\Big(-\tfrac{1}{2}\operatorname{Tr}\big((W^{-1} + N\,S)\,\Lambda\big)\Big)$$

$$\mathcal{W}(\Lambda \mid W, \nu) = B(W, \nu)\, |\Lambda|^{\frac{\nu - D - 1}{2}} \exp\!\Big[-\tfrac{1}{2}\operatorname{Tr}(W^{-1}\Lambda)\Big]$$
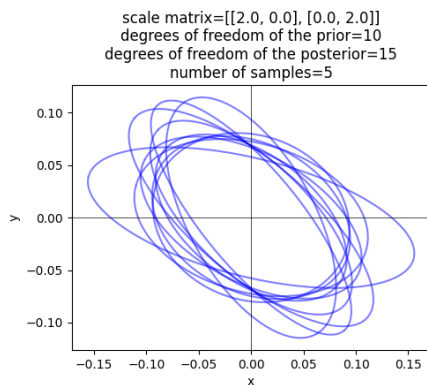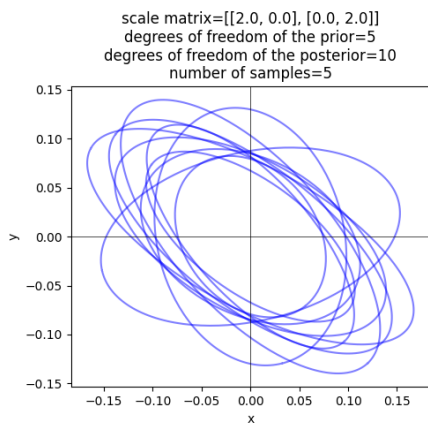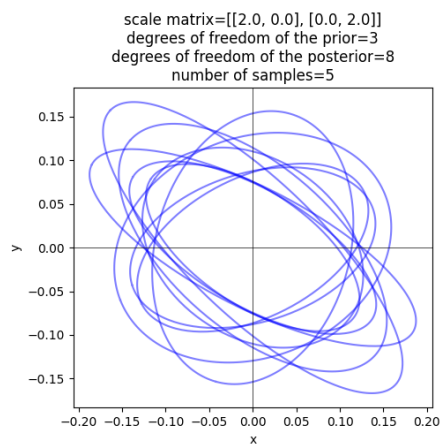
Therefore, the posterior is in the same form as a Wishart distribution with parameters:

$$\nu = N + \nu, \quad W = (W^{-1} + NS)^{-1}$$

As $\mathcal{W}(\boldsymbol{\Lambda} \mid W, \nu)$ (Wishart Prior) $\times\, p(X \mid \boldsymbol{\mu}, \boldsymbol{\Lambda})$ (Gaussian likelihood) $= \mathcal{W}(\boldsymbol{\Lambda} \mid W, \nu)$ (Wishart posterior), this closure property under posterior updates defines a conjugate prior for the Gaussian precision matrix.

**b.**

Varying Prior Degrees of Freedom



scale matrix=[[2.0, 0.0], [0.0, 2.0]]
degrees of freedom of the prior=3
degrees of freedom of the posterior=8
number of samples=5

scale matrix=[[2.0, 0.0], [0.0, 2.0]]
degrees of freedom of the prior=5
degrees of freedom of the posterior=10
number of samples=5

scale matrix=[[2.0, 0.0], [0.0, 2.0]]
degrees of freedom of the prior=10
degrees of freedom of the posterior=15
number of samples=5

Varying Number of Samples



scale matrix=[[1.0, 0.0], [0.0, 1.0]]
degrees of freedom of the prior=5
degrees of freedom of the posterior=7
number of samples=2

scale matrix=[[1.0, 0.0], [0.0, 1.0]]
degrees of freedom of the prior=5
degrees of freedom of the posterior=10
number of samples=5

scale matrix=[[1.0, 0.0], [0.0, 1.0]]
degrees of freedom of the prior=5
degrees of freedom of the posterior=25
number of samples=20

7

scale matrix=[[1.0, 0.0], [0.0, 1.0]]
degrees of freedom of the prior=5
degrees of freedom of the posterior=10
number of samples=5

scale matrix=[[2.0, 0.0], [0.0, 2.0]]
degrees of freedom of the prior=5
degrees of freedom of the posterior=10
number of samples=5

scale matrix=[[1.0, 0.0], [0.0, 2.0]]
degrees of freedom of the prior=5
degrees of freedom of the posterior=10
number of samples=5

```python
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import wishart

def generate_data(n, cov):
    return np.random.multivariate_normal(mean=[0, 0], cov=cov, size=n)

def sample_wishart(df, scale, num_samples=5):
    samples = []
    for _ in range(num_samples):
        W = wishart.rvs(df=df, scale=scale)
        samples.append(W)
    return samples

def plot_precision_ellipses(precisions, ax, title):
    for p in precisions:
        eigvals, eigvecs = np.linalg.eigh(p)
        if np.any(eigvals <= 0):
            continue

        angles = np.linspace(0, 2*np.pi, 200)
        circle = np.stack([np.cos(angles), np.sin(angles)], axis=1)
        scale_matrix = np.diag(1.0 / np.sqrt(eigvals))
        ellipse_y = circle @ scale_matrix
        ellipse_x = ellipse_y @ eigvecs.T
        ax.plot(ellipse_x[:, 0], ellipse_x[:, 1], 'b', alpha=0.5)

    ax.set_aspect('equal', 'box')
    ax.set_title(title)
    ax.set_xlabel('x')
    ax.set_ylabel('y')
    ax.axhline(0, color='black', linewidth=0.5)
    ax.axvline(0, color='black', linewidth=0.5)

def variation_1_distribution():
    n = 5
    scale_matrix = np.eye(2) * 2.0
```

```python
38        cov = np.array([[1.0, 0.5], [0.5, 1.0]])
39        X = generate_data(n, cov)
40        sum_T = X.T @ X
41        dfs = [3, 5, 10]
42
43        fig, axes = plt.subplots(1, 3, figsize=(15, 4))
44
45        for i, df_prior in enumerate(dfs):
46            df_post = df_prior + n
47            scale_post = scale_matrix + sum_T
48            posterior_samples = sample_wishart(df_post, scale_post, num_samples=10)
49
50            title = (
51                f"scale matrix={scale_matrix.tolist()}\n"
52                f"degrees of freedom of the prior={df_prior}\n"
53                f"degrees of freedom of the posterior={df_post}\n"
54                f"number of samples={n}"
55            )
56            plot_precision_ellipses(posterior_samples, axes[i], title)
57
58        fig.suptitle("Varying Prior Degrees of Freedom", fontsize=16)
59        plt.tight_layout()
60        plt.show()
61
62    def variation_2_distribution():
63        df_prior = 5
64        scale_matrix = np.eye(2)
65        cov = np.array([[1.0, 0.5],[0.5, 1.0]])
66        ns = [2, 5, 20]
67        fig, axes = plt.subplots(1, 3, figsize=(15, 4))
68
69        for i, n in enumerate(ns):
70            X = generate_data(n, cov)
71            sum_T = X.T @ X
72            df_post = df_prior + n
73            scale_post = scale_matrix + sum_T
74            posterior_samples = sample_wishart(df_post, scale_post, num_samples=10)
75
76            title = (
77                f"scale matrix={scale_matrix.tolist()}\n"
78                f"degrees of freedom of the prior={df_prior}\n"
79                f"degrees of freedom of the posterior={df_post}\n"
80                f"number of samples={n}"
81            )
82            plot_precision_ellipses(posterior_samples, axes[i], title)
83
84        fig.suptitle("Varying Number of Samples", fontsize=16)
85        plt.tight_layout()
86        plt.show()
87
88    def variation_3_distribution():
89        df_prior = 5
90        n = 5
91        scale_matrix_list = [np.eye(2), 2.0 * np.eye(2), np.diag([1.0, 2.0])]
92        cov = np.array([[1.0, 0.5],[0.5, 1.0]])
93        X = generate_data(n, cov)
94        sum_T = X.T @ X
95
96        fig, axes = plt.subplots(1, 3, figsize=(15, 4))
97
98        for i, scale_matrix in enumerate(scale_matrix_list):
99            df_post = df_prior + n
100            scale_post = scale_matrix + sum_T
101            posterior_samples = sample_wishart(df_post, scale_post, num_samples=10)
102
103            title = (
104                f"scale matrix={scale_matrix.tolist()}\n"
105                f"degrees of freedom of the prior={df_prior}\n"
106                f"degrees of freedom of the posterior={df_post}\n"
107                f"number of samples={n}"
108            )
```

```
109              plot_precision_ellipses(posterior_samples, axes[i], title)
110
111        fig.suptitle("Varying␣Scale␣Matrix", fontsize=16)
112        plt.tight_layout()
113        plt.show()
114
115  def main():
116        variation_1_distribution()
117        variation_2_distribution()
118        variation_3_distribution()
119
120  if __name__ == "__main__":
121        main()
```

## 5.

$$
\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{\frac{D}{2}} \, |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left\{ -\tfrac{1}{2} \left(\mathbf{x} - \boldsymbol{\mu}\right)^{\top} \boldsymbol{\Sigma}^{-1} \left(\mathbf{x} - \boldsymbol{\mu}\right) \right\}
$$

$$
= \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left\{ -\tfrac{1}{2} \left( \mathbf{x}^{\top} \boldsymbol{\Sigma}^{-1} \mathbf{x} - 2\,\boldsymbol{\mu}^{\top} \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}^{\top} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \right) \right\}
$$

$$
= \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left\{ -\tfrac{1}{2} \mathbf{x}^{\top} \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}^{\top} \boldsymbol{\Sigma}^{-1} \mathbf{x} - \tfrac{1}{2} \boldsymbol{\mu}^{\top} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \right\}
$$

$$
= \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left\{ -\tfrac{1}{2} \mathbf{x}^{\top} \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}^{\top} \boldsymbol{\Sigma}^{-1} \mathbf{x} \right\} \exp\left\{ -\tfrac{1}{2} \boldsymbol{\mu}^{\top} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \right\}
$$

$$
= (2\pi)^{-\frac{D}{2}} \cdot |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left\{ -\tfrac{1}{2} \boldsymbol{\mu}^{\top} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \right\} \exp\left\{ -\tfrac{1}{2} \mathbf{x}^{\top} \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}^{\top} \boldsymbol{\Sigma}^{-1} \mathbf{x} \right\}
$$

Exponential-family form:

$$
p(\mathbf{x} \mid \boldsymbol{\eta}) = h(\mathbf{x}) \cdot g(\boldsymbol{\eta}) \cdot \exp\!\left( \boldsymbol{\eta}^{\top} \mathbf{u}(\mathbf{x}) \right)
$$

$$
h(\mathbf{x}) = (2\pi)^{-\frac{D}{2}}
$$

$$
\mathbf{u}(\mathbf{x}) = \left( \mathbf{x}, \, \mathbf{x}\mathbf{x}^{\top} \right) \quad \text{(which generalizes the univariate } (x, x^2))
$$

$$
\boldsymbol{\eta} = \begin{pmatrix} \eta_1 \\ \eta_2 \end{pmatrix} = \begin{pmatrix} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \\ -\tfrac{1}{2} \boldsymbol{\Sigma}^{-1} \end{pmatrix}
$$

Solving for $\boldsymbol{\Sigma}$ and $\boldsymbol{\mu}$ in terms of $\boldsymbol{\eta}$.

$$
\eta_2 = -\frac{1}{2} \boldsymbol{\Sigma}^{-1} \quad \Longrightarrow \quad \boldsymbol{\Sigma}^{-1} = -2\, \eta_2 \quad \Longrightarrow \quad \boldsymbol{\Sigma} = \left(\boldsymbol{\Sigma}^{-1}\right)^{-1} = \left(-2\, \eta_2\right)^{-1} = -\tfrac{1}{2}\, \eta_2^{-1}.
$$

$$
\eta_1 = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \quad \Longrightarrow \quad \boldsymbol{\mu} = \boldsymbol{\Sigma}\, \eta_1 = \left( -\tfrac{1}{2}\, \eta_2^{-1} \right) \eta_1
$$

$$
|\boldsymbol{\Sigma}|^{-\frac{1}{2}} = \left| -2\, \eta_2 \right|^{\frac{D}{2}} \quad \text{(for } D \text{ dimensions)}
$$

$$
g(\boldsymbol{\eta}) = |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\!\left( -\tfrac{1}{2}\, \boldsymbol{\mu}^{\top} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \right)
$$

10

$$\boldsymbol{\mu}^\top \Sigma^{-1} \boldsymbol{\mu} \; = \; (\Sigma\,\boldsymbol{\eta}_1)^\top \Sigma^{-1} (\Sigma\,\boldsymbol{\eta}_1) \; = \; \boldsymbol{\eta}_1^\top \Sigma\,\boldsymbol{\eta}_1 \; = \; \boldsymbol{\eta}_1^\top \left( -\tfrac{1}{2}\,\boldsymbol{\eta}_2^{-1} \right) \boldsymbol{\eta}_1 \; = \; -\tfrac{1}{2}\,\boldsymbol{\eta}_1^\top\,\boldsymbol{\eta}_2^{-1}\,\boldsymbol{\eta}_1$$

Therefore,

$$
\begin{aligned}
g(\boldsymbol{\eta}) \; &= \; |\Sigma|^{-\frac{1}{2}} \, \exp\!\left( -\tfrac{1}{2}\,\boldsymbol{\mu}^\top \Sigma^{-1}\,\boldsymbol{\mu} \right) \\
&= \; |\Sigma|^{-\frac{1}{2}} \, \exp\!\left( -\tfrac{1}{2}\,\boldsymbol{\eta}_1^\top\,\boldsymbol{\eta}_2^{-1}\,\boldsymbol{\eta}_1 \right) \\
&= \; |-2\,\boldsymbol{\eta}_2|^{\frac{D}{2}} \, \exp\!\left( \tfrac{1}{4}\,\boldsymbol{\eta}_1^\top\,\boldsymbol{\eta}_2^{-1}\,\boldsymbol{\eta}_1 \right)
\end{aligned}
$$

Univariate $(D = 1)$ case:

$$g(\boldsymbol{\eta}) \; = \; |-2\,\eta_2|^{\frac{1}{2}} \, \exp\!\left( \tfrac{\eta_1^2}{4\,\eta_2} \right)$$

Multivariate $(D > 1)$ case:

$$g(\boldsymbol{\eta}) \; = \; |-2\,\boldsymbol{\eta}_2|^{\frac{D}{2}} \, \exp\!\left( \tfrac{1}{4}\,\boldsymbol{\eta}_1^\top\,\boldsymbol{\eta}_2^{-1}\,\boldsymbol{\eta}_1 \right)$$