



COLUMBIA UNIVERSITY  
IN THE CITY OF NEW YORK

STAT 4224/5224

*Bayesian Statistics*

Dobrin Marchev

# Recall: Multivariate Normal Distribution

Notation:  $\mathbf{X} \sim N_p(\boldsymbol{\mu}_X, \boldsymbol{\Sigma}_X)$ , where:

$$\mathbf{X} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix}, \boldsymbol{\mu}_X = E(\mathbf{X}) = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix}, \boldsymbol{\Sigma}_X = \sigma^2(\mathbf{X}) = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1p} & \sigma_{2p} & \cdots & \sigma_{pp} \end{bmatrix}$$

where  $\text{cov}(X_i, X_j) = \sigma_{ij}$ ,  $i \neq j$

Multivariate normal density function:

$$f(\mathbf{x}) = (2\pi)^{-p/2} \left| \boldsymbol{\Sigma}_X^{-1} \right|^{1/2} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_X)' \boldsymbol{\Sigma}_X^{-1} (\mathbf{x} - \boldsymbol{\mu}_X)}$$

# Multivariate Normal Model

Assume that we have multivariate observations

$$\mathbf{X}_1, \dots, \mathbf{X}_n \mid \boldsymbol{\theta}, \boldsymbol{\Sigma} \sim N_p(\boldsymbol{\theta}, \boldsymbol{\Sigma})$$

Then the likelihood is

$$\begin{aligned} f(\mathbf{x}_1, \dots, \mathbf{x}_n \mid \boldsymbol{\theta}, \boldsymbol{\Sigma}) \\ &= \prod_{i=1}^n (2\pi)^{-\frac{p}{2}} (\det \boldsymbol{\Sigma})^{-1/2} e^{-1/2 (\mathbf{x}_i - \boldsymbol{\theta})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\theta})} \\ &= (2\pi)^{-\frac{np}{2}} (\det \boldsymbol{\Sigma})^{-\frac{n}{2}} e^{-1/2 \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\theta})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\theta})} \\ &\propto e^{-1/2 (\sum_{i=1}^n \mathbf{x}_i' \boldsymbol{\Sigma}^{-1} \mathbf{x}_i + n \boldsymbol{\theta}' \boldsymbol{\Sigma}^{-1} \boldsymbol{\theta} - 2 \boldsymbol{\theta}' \boldsymbol{\Sigma}^{-1} \sum_{i=1}^n \mathbf{x}_i)} \\ &\propto e^{-\frac{1}{2} \sum_{i=1}^n \mathbf{x}_i' \boldsymbol{\Sigma}^{-1} \mathbf{x}_i - \frac{1}{2} \boldsymbol{\theta}' \mathbf{A}_1 \boldsymbol{\theta} + \boldsymbol{\theta}' \mathbf{b}_1} \end{aligned}$$

where  $\mathbf{A}_1 = n\boldsymbol{\Sigma}^{-1}$ ,  $\mathbf{b}_1 = n\boldsymbol{\Sigma}^{-1}\bar{\mathbf{x}}$

# Prior for the mean vector

Let

$$\boldsymbol{\theta} \sim N_p(\boldsymbol{\mu}_0, \boldsymbol{\Lambda}_0)$$

Then

$$\begin{aligned}\pi(\boldsymbol{\theta}) &= (2\pi)^{-\frac{p}{2}} (\det \boldsymbol{\Lambda}_0)^{-1/2} e^{-1/2(\boldsymbol{\theta} - \boldsymbol{\mu}_0)' \boldsymbol{\Lambda}_0^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}_0)} \\ &= (2\pi)^{-\frac{p}{2}} (\det \boldsymbol{\Lambda}_0)^{-1/2} e^{-1/2 \boldsymbol{\theta}' \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\theta} - \frac{1}{2} \boldsymbol{\mu}_0' \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0 + \boldsymbol{\theta}' \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0} \\ &\propto e^{-1/2 \boldsymbol{\theta}' \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\theta} + \boldsymbol{\theta}' \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0} \\ &= e^{-1/2 \boldsymbol{\theta}' \boldsymbol{A}_0 \boldsymbol{\theta} + \boldsymbol{\theta}' \boldsymbol{b}_0}\end{aligned}$$

where  $\boldsymbol{A}_0 = \boldsymbol{\Lambda}_0^{-1}$ ,  $\boldsymbol{b}_0 = \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0$ .

# Conditional Posterior of $\theta$

The conditional posterior of  $\theta | x_1, \dots, x_n, \Sigma$  is

$$\begin{aligned} f(\theta | x_1, \dots, x_n, \Sigma) &\propto e^{-\frac{1}{2}\theta' A_0 \theta + \theta' b_0} \times e^{-\frac{1}{2} \sum_{i=1}^n x_i' \Sigma^{-1} x_i - \frac{1}{2} \theta' A_1 \theta + \theta' b_1} \\ &\propto e^{-\frac{1}{2}\theta' A_n \theta + \theta' b_n} \end{aligned}$$

where

$$\begin{aligned} A_n &= A_0 + A_1 = \Lambda_0^{-1} + n\Sigma^{-1} \\ b_n &= b_0 + b_1 = \Lambda_0^{-1} \mu_0 + n\Sigma^{-1} \bar{x} \end{aligned}$$

The only distribution with such form of the density is the multivariate normal. Therefore,

$$\theta | x_1, \dots, x_n, \Sigma \sim N_p(\mu_n, \Lambda_n)$$

where

$$\begin{aligned} \mu_n &= A_n^{-1} b_n = (\Lambda_0^{-1} + n\Sigma^{-1})^{-1} (\Lambda_0^{-1} \mu_0 + n\Sigma^{-1} \bar{x}) \\ \Lambda_n &= A_n^{-1} = (\Lambda_0^{-1} + n\Sigma^{-1})^{-1} \end{aligned}$$

Notice the analogy with the univariate case!

# Posterior Predictive Distribution

If  $\Sigma$  is known, then it can be shown that

$$\mathbf{x}_{new}|\mathbf{x} \sim N_p(\boldsymbol{\mu}_n, \Sigma + \Lambda_n)$$

Proof:

$$f(\mathbf{x}_{new}|\mathbf{x}) = \int f(\mathbf{x}_{new}|\boldsymbol{\theta}, \Sigma) f(\boldsymbol{\theta}|\mathbf{x}, \Sigma) d\boldsymbol{\mu}$$

Convince yourself that only a multivariate normal density can be the answer, and then find the mean and variance with tricks we used before.

# Wishart Distribution

- It is a generalization to multidimensions of the Chi-Square distribution.
- The Wishart distribution is a sum of outer products of random vectors.
- It is a *random matrix* which is symmetric and positive definite.
- Let  $\mathbf{X}_1, \dots, \mathbf{X}_n \sim N_p(\mathbf{0}, \Sigma)$  be independent. Then the distribution of the  $p \times p$  random matrix  $\mathbf{M} = \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i'$  is said to have the Wishart distribution with  $\text{df} = n$ .
- It defines the distribution of the sample covariance matrix.
- Notation:  $\mathbf{M} \sim W_p(\Sigma, n)$

Obtaining samples in R:

To generate  $n$  random matrices, distributed according to the Wishart distribution with parameters  $\Sigma$  and  $\text{df}$ ,  $W_p(\Sigma, m)$ , where  $m = \text{df}$ .

Use Function: `rWishart(n, df, Σ)`

# Wishart Distribution Properties

- Let  $\mathbf{M} \sim W_p(\mathbf{\Sigma}, n)$
- $E(\mathbf{M}) = n\mathbf{\Sigma}$
- $\mathbf{M} \sim \mathbf{A}W_p(\mathbf{I}_p, n)\mathbf{A}'$ , where  $\mathbf{\Sigma} = \mathbf{A}\mathbf{A}'$  is the LU-decomposition
- Assume  $n > p$  and  $\mathbf{\Sigma}$  is invertible. Then the pdf of  $\mathbf{M}$  is

$$f(\mathbf{m}, n, \mathbf{\Sigma}) = \frac{|\mathbf{m}|^{\frac{n-p-1}{2}} e^{-\frac{1}{2} \text{tr}(\mathbf{m}\mathbf{\Sigma}^{-1})}}{2^{\frac{pn}{2}} \pi^{\frac{p(p-1)}{4}} |\mathbf{\Sigma}|^{\frac{n}{2}} \prod_{i=1}^p \Gamma\left(\frac{n+1-i}{2}\right)}$$

where the support is all symmetric positive definite matrices  $\mathbf{m}$ .



# Aside: Positive Definite Matrix

Any variance matrix must be positive definite, which is the analogy of the univariate variance  $\sigma^2 > 0$ .

**Definition:** A symmetric  $p \times p$  matrix  $\mathbf{M}$  is positive definite iff  
$$\mathbf{x}'\mathbf{M}\mathbf{x} > 0, \forall \mathbf{x} \in \mathbb{R}^p \setminus \mathbf{0}$$

**Definition:** If  $\mathbf{M}$  is a square  $p \times p$  matrix, then  $\mathbf{x} \neq \mathbf{0}$  is an eigenvector of  $\mathbf{M}$  if there  $\exists \lambda \in \mathbb{C}$  such that  $\mathbf{M}\mathbf{x} = \lambda\mathbf{x}$

**Theorem:** If  $\mathbf{M}$  is a symmetric matrix, then all its eigenvalues are real numbers.

**Theorem:** A symmetric matrix  $\mathbf{M}$  is positive definite iff all its eigenvalues are positive.

**Note:** The requirement that the covariance matrix  $\mathbf{\Sigma}$  is positive definite guarantees that all variances (on the diagonal) are positive and that all correlations are between -1 and 1.

# Connection with sample covariance

Let  $\mathbf{z}_1, \dots, \mathbf{z}_n \in \mathbb{R}^p$  be some observations from a multivariate distribution. The sum of squares (SS) *matrix* is defined as

$$SS = \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i' = \mathbf{Z}' \mathbf{Z}$$

where  $\mathbf{Z}$  is an  $n \times p$  matrix with  $i^{\text{th}}$  row  $\mathbf{z}_i'$ .

Notes:

- If  $n > p$  and the  $\mathbf{z}_i$ 's are *linearly independent*, then SS is a positive definite matrix.
- If we assume the population mean is 0, then  $SS/n$  is the MLE of the population covariance matrix.
- The Wishart distribution is constructed in this way and will be positive definite with probability 1 (the chance of linearly dependent observation vectors is 0).

# Inverse Wishart Distribution

- In the univariate case, for a conjugate prior, we had to use IG distribution on the population variance  $\sigma^2$ .
- In the multivariate case we need the *Inverse Wishart distribution* as a prior on  $\Sigma$ .
- We say that  $\mathbf{M} \sim W_p^{-1}(\Sigma, n)$ , if  $\mathbf{M}^{-1} \sim W_p(\Sigma^{-1}, n)$
- Property:  $E(\mathbf{M}) = \frac{\Sigma}{n-p-1}$
- Prior:  $\pi(\Sigma) \sim W_p^{-1}(\mathbf{S}_0, \nu_0)$  with pdf

$$f(\Sigma) = \left[ 2^{\frac{\nu_0 p}{2}} \pi^{\frac{\binom{p}{2}}{2}} (\det \mathbf{S}_0)^{-\frac{\nu_0}{2}} \prod_{j=1}^p \Gamma\left(\frac{\nu_0 + 1 - j}{2}\right) \right]^{-1} \times$$
$$(\det \Sigma)^{-\frac{\nu_0 + p + 1}{2}} e^{-\frac{\text{tr}(\mathbf{S}_0 \Sigma^{-1})}{2}}$$

# Aside: Trace of a matrix

**Definition:** Trace of a square  $p \times p$  matrix  $\mathbf{M}$  is defined as

$$\text{tr}(\mathbf{M}) = \sum_{j=1}^p m_{jj}$$

**Properties:**

- $\text{tr}(\mathbf{A} + \mathbf{B}) = \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B})$
- $\text{tr}(c\mathbf{A}) = c\text{tr}(\mathbf{A})$
- $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$ , where  $\mathbf{A}$  is  $m \times n$  and  $\mathbf{B}$  is  $n \times m$
- If  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{C}$  are symmetric, then
$$\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{CBA}) = \text{tr}(\mathbf{CAB}) = \dots$$
- Trace is equal to the sum of the eigenvalues.

# Back to Multivariate Normal Model

Assume that we have multivariate observations

$$\mathbf{X}_1, \dots, \mathbf{X}_n \mid \boldsymbol{\theta}, \boldsymbol{\Sigma} \sim N_p(\boldsymbol{\theta}, \boldsymbol{\Sigma})$$

Recall the likelihood is

$$\begin{aligned} & f(\mathbf{x}_1, \dots, \mathbf{x}_n \mid \boldsymbol{\theta}, \boldsymbol{\Sigma}) \\ &= (2\pi)^{-\frac{np}{2}} (\det \boldsymbol{\Sigma})^{-\frac{n}{2}} e^{-\frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\theta})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\theta})} \\ &\propto (\det \boldsymbol{\Sigma})^{-\frac{n}{2}} e^{-\frac{1}{2} \text{tr}(\mathbf{S}_\theta \boldsymbol{\Sigma}^{-1})} \end{aligned}$$

where we used

$$\begin{aligned} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\theta})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\theta}) &= \text{tr} \left[ \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\theta})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\theta}) \right] \\ &= \text{tr}(\mathbf{S}_\theta \boldsymbol{\Sigma}^{-1}) \end{aligned}$$

with  $\mathbf{S}_\theta = \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\theta})' (\mathbf{x}_i - \boldsymbol{\theta})$ .

# Conditional Posterior of $\Sigma$

The conditional posterior of  $\Sigma | \mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\theta}$  is

$$\begin{aligned} f(\Sigma | \mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\theta}) &\propto \pi(\Sigma) \times f(\mathbf{x}_1, \dots, \mathbf{x}_n | \boldsymbol{\theta}, \Sigma) \\ &\propto (\det \Sigma)^{-\frac{\nu_0 + p + 1}{2}} e^{-\frac{\text{tr}(\mathbf{S}_0 \Sigma^{-1})}{2}} \times (\det \Sigma)^{-\frac{n}{2}} e^{-\frac{1}{2} \text{tr}(\mathbf{S}_\theta \Sigma^{-1})} \\ &\quad (\det \Sigma)^{-\frac{\nu_0 + n + p + 1}{2}} e^{-\frac{1}{2} \text{tr}((\mathbf{S}_0 + \mathbf{S}_\theta) \Sigma^{-1})} \end{aligned}$$

Therefore,

$$\Sigma | \mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\theta} \sim W_p^{-1}(\mathbf{S}_0 + \mathbf{S}_\theta, \nu_0 + n)$$

and

$$E(\Sigma | \mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\theta}) = \frac{\mathbf{S}_0 + \mathbf{S}_\theta}{\nu_0 + n - p - 1}$$

More importantly, the fact that we know the distributions of each conditional posterior means we can develop a Gibbs sampler to obtain samples from the full joint posterior!

# Exercise 1

Consider the following model

$$\mathbf{X}_1, \dots, \mathbf{X}_n \mid \boldsymbol{\theta}, \boldsymbol{\Sigma} \sim N_p(\boldsymbol{\theta}, \boldsymbol{\Sigma})$$

$$\pi(\boldsymbol{\theta}, \boldsymbol{\Sigma}) \propto (\det \boldsymbol{\Sigma})^{-\frac{p+1}{2}}$$

This is known as the independence-Jeffreys prior (note that it is improper).

- a) Derive the conditional posteriors of  $\boldsymbol{\theta} \mid \mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\Sigma}$  and  $\boldsymbol{\Sigma} \mid \mathbf{x}_1, \dots, \mathbf{x}_n$
- b) Under what condition is the joint posterior proper?

Hint:

$$\begin{aligned} f(\boldsymbol{\theta}, \boldsymbol{\Sigma} \mid \mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\theta}) &\propto (\det \boldsymbol{\Sigma})^{-\frac{n+p+1}{2}} e^{-\frac{1}{2} \text{tr}(\mathbf{S}_{\boldsymbol{\theta}} \boldsymbol{\Sigma}^{-1})} \\ \mathbf{S}_{\boldsymbol{\theta}} &= \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\theta})(\mathbf{x}_i - \boldsymbol{\theta})' = \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})' + n(\boldsymbol{\theta} - \bar{\mathbf{x}})(\boldsymbol{\theta} - \bar{\mathbf{x}})' \\ \Rightarrow f(\boldsymbol{\theta}, \boldsymbol{\Sigma} \mid \mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\theta}) &\propto (\det \boldsymbol{\Sigma})^{-\frac{n+p+1}{2}} e^{-\frac{1}{2} \text{tr}(\mathbf{S}_{\bar{\mathbf{x}}} \boldsymbol{\Sigma}^{-1})} e^{-\frac{1}{2} \text{tr}(n(\boldsymbol{\theta} - \bar{\mathbf{x}})(\boldsymbol{\theta} - \bar{\mathbf{x}})' \boldsymbol{\Sigma}^{-1})} \end{aligned}$$

# Exercise 2

Consider the following model

$$\mathbf{X}_1, \dots, \mathbf{X}_n \mid \boldsymbol{\theta}, \boldsymbol{\Sigma} \sim N_p(\boldsymbol{\theta}, \boldsymbol{\Sigma})$$

$$\boldsymbol{\theta} \mid \boldsymbol{\Sigma} \sim N_p(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}/\kappa_0)$$

$$\boldsymbol{\Sigma} \sim W_p^{-1}(\Lambda_0, \nu_0)$$

- a) Derive the conditional posterior of  $\boldsymbol{\theta} \mid \mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\Sigma}$
- b) Derive the posterior of  $\boldsymbol{\Sigma} \mid \mathbf{x}_1, \dots, \mathbf{x}_n$

Answers:

- a) Multivariate normal
- b) Inverse Wishart



# Gibbs Sampler

We showed that:

$$\begin{aligned}\boldsymbol{\theta} | \mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\Sigma} &\sim N_p(\boldsymbol{\mu}_n, \boldsymbol{\Lambda}_n) \\ \boldsymbol{\Sigma} | \mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\theta} &\sim W_p^{-1}(\mathbf{S}_n, \nu_n)\end{aligned}$$

where

$$\begin{aligned}\boldsymbol{\mu}_n &= (\boldsymbol{\Lambda}_0^{-1} + n\boldsymbol{\Sigma}^{-1})^{-1}(\boldsymbol{\Lambda}_0^{-1}\boldsymbol{\mu}_0 + n\boldsymbol{\Sigma}^{-1}\bar{\mathbf{x}}) \\ \boldsymbol{\Lambda}_n &= (\boldsymbol{\Lambda}_0^{-1} + n\boldsymbol{\Sigma}^{-1})^{-1} \\ \mathbf{S}_n &= \mathbf{S}_0 + \mathbf{S}_\theta \\ \mathbf{S}_\theta &= \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\theta})'(\mathbf{x}_i - \boldsymbol{\theta}) \\ \nu_n &= \nu_0 + n\end{aligned}$$

Start with some starting value  $\boldsymbol{\Sigma}^{(0)}$  (say, the sample covariance matrix) and iterate between the above two conditional posteriors to obtain a sequence  $(\boldsymbol{\theta}^{(1)}, \boldsymbol{\Sigma}^{(1)}), \dots, (\boldsymbol{\theta}^{(S)}, \boldsymbol{\Sigma}^{(S)})$