MULTIVARIATE NORMAL MODEL

DR. OLANREWAJU MICHAEL AKANDE

Feb 19, 2020



ANNOUNCEMENTS

- Take Survey I
- Link: https://duke.qualtrics.com/jfe/form/SV_54rrMwDxp3hmagt
- Responses are anonymized.

OUTLINE

- Wrap up exercise from last class
- Multivariate normal/Gaussian model
 - Motivating example
 - Inference for mean
 - Inference for covariance



RECAP OF CONDITIONAL DISTRIBUTIONS

• Partition
$$oldsymbol{Y} = (Y_1, \dots, Y_p)^T$$
 as

$$oldsymbol{Y} = egin{pmatrix} oldsymbol{Y}_1 \ oldsymbol{Y}_2 \end{pmatrix} \sim \mathcal{N}_p \left[egin{pmatrix} oldsymbol{\mu}_1 \ oldsymbol{\mu}_2 \end{pmatrix}, egin{pmatrix} \Sigma_{11} & \Sigma_{12} \ \Sigma_{21} & \Sigma_{22} \end{pmatrix}
ight],$$

where

- $oldsymbol{Y}_1$ and $oldsymbol{\mu}_1$ are q imes 1,
- $oldsymbol{Y}_2$ and $oldsymbol{\mu}_2$ are (p-q) imes 1,
- Σ_{11} is q imes q, and
- Σ_{22} is (p-q) imes (p-q), with $\Sigma_{22} > 0$.
- Then,

$$oldsymbol{Y}_1|oldsymbol{Y}_2=oldsymbol{y}_2\sim\mathcal{N}_q\left(oldsymbol{\mu}_1+\Sigma_{12}\Sigma_{22}^{-1}(oldsymbol{y}_2-oldsymbol{\mu}_2),\Sigma_{11}-\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}
ight).$$



WORKING WITH NORMAL DISTRIBUTIONS

• Three real (univariate) random quantities x, y and z have a joint normal distribution given by p(x, y, z) = p(y|x)p(x|z)p(z).

Suppose

- $p(y|x) = \mathcal{N}(x, w)$ independently of z, for some known variance w;
- $p(x|z) = \mathcal{N}(\theta z, v)$ for some known parameter θ , and known variance v; and
- $p(z) = \mathcal{N}(m, M)$, with some known mean m, and known variance M.
- What is
 - p(x)? p(y)?
 - p(x|y)? p(z|x)?
- To be done on the board.



MULTIVARIATE DATA

- Survey data often yield multivariate data of varied types.
- Typical survey data: response vector $Y_i = (Y_{i1}, \ldots, Y_{ip})^T$ for each person i in a sample of survey respondents, $i = 1, \ldots, n$. For example, we could have
 - $y_{i1} = \mathsf{income}$
 - $y_{i2} =$ level of education
 - $y_{i3} =$ number of children
 - $y_{i4} = \mathsf{age}$
 - $y_{i5} = \mathsf{attitude}$
- Interest is then often on inferring the potential associations among these variables.
- See https://www.stat.washington.edu/people/pdhoff/public/coptalk.pdf







See https://www.stat.washington.edu/people/pdhoff/public/coptalk.pdf



CONDITIONAL MODELS

- Interest is often in conditional relationships between pairs of variables, accounting for heterogeneity in other variables of less interest.
- Consider the following models.
- GSS data:
 - Model 1

 $\text{INC}_i = \beta_0 + \beta_1 \frac{\text{CHILD}_i}{\beta_2 \text{DEG}_i} + \beta_3 \text{AGE}_i + \beta_4 \text{PCHILD}_i + \beta_5 \text{PINC}_i + \beta_6 \text{PDEG}_i + \epsilon_i$

p-value for β_1 here is 0.11: "little evidence" that $\beta_1 \neq 0$.

Model 2

 $\textbf{CHILD}_i \sim \textbf{Poisson} \left(\exp \left[\beta_0 + \beta_1 \textbf{INC}_i + \beta_2 \textbf{DEG}_i + \beta_3 \textbf{AGE}_i + \beta_4 \textbf{PCHILD}_i + \beta_5 \textbf{PINC}_i + \beta_6 \textbf{PDEG}_i \right] \right)$

p-value for β_1 here is 0.01: "strong evidence" that $\beta_1 \neq 0$.

- Not satisfactory; better to use multivariate models instead to do this jointly.
- See https://www.stat.washington.edu/people/pdhoff/public/coptalk.pdf
 STA 602L

MULTIVARIATE NORMAL DISTRIBUTION RECAP

• Recall that if $oldsymbol{Y} = (Y_1, \dots, Y_p)^T \sim \mathcal{N}_p(oldsymbol{ heta}, \Sigma)$, then

$$f(oldsymbol{y}) = (2\pi)^{-rac{p}{2}} |\Sigma|^{-rac{1}{2}} \exp\left\{-rac{1}{2}(oldsymbol{y}-oldsymbol{ heta})^T \Sigma^{-1}(oldsymbol{y}-oldsymbol{ heta})
ight\}.$$

- $oldsymbol{ heta}$ is the p imes 1 mean vector, that is, $oldsymbol{ heta} = (heta_1, \dots, heta_p)^T$.
- Σ is the $p \times p$ **positive definite** covariance matrix, that is, $\Sigma = \{\sigma_{jk}\}$, where σ_{jk} denotes the covariance between Y_j and Y_k .
- For each $j = 1, \ldots, p$, $Y_j \sim \mathcal{N}(\theta_j, \sigma_{jj})$.
- How to do posterior inference if this is our sampling model?



READING COMPREHENSION EXAMPLE

- Twenty-two children are given a reading comprehension test before and after receiving a particular instruction method.
 - Y_{i1} : pre-instructional score for student *i*.
 - Y_{i2}: post-instructional score for student *i*.
- Vector of observations for each student: $Y_i = (Y_{i1}, Y_{i2})^T$.
- Clearly, we should expect some correlation between Y_{i1} and Y_{i2} .



READING COMPREHENSION EXAMPLE

- Questions of interest:
 - Do students improve in reading comprehension on average?
 - If so, by how much?
 - Can we predict post-test score from pre-test score?
 - If there is a "significant" improvement, does that mean the instructional method is good?
 - If we have students with missing pre-test scores, can we predict the scores?
- We will come back to this example. First, let's specify priors and see what the implied (conditional) posteriors look like.



MULTIVARIATE NORMAL LIKELIHOOD

• For data $oldsymbol{Y}_i = (Y_{i1}, \dots, Y_{ip})^T \sim \mathcal{N}_p(oldsymbol{ heta}, \Sigma)$, the likelihood is

$$egin{aligned} L(m{Y};m{ heta},\Sigma) &= \prod_{i=1}^n (2\pi)^{-rac{p}{2}} |\Sigma|^{-rac{1}{2}} \exp\left\{-rac{1}{2}(m{y}_i-m{ heta})^T \Sigma^{-1}(m{y}_i-m{ heta})
ight\} \ &\propto |\Sigma|^{-rac{n}{2}} \exp\left\{-rac{1}{2}\sum_{i=1}^n (m{y}_i-m{ heta})^T \Sigma^{-1}(m{y}_i-m{ heta})
ight\}. \end{aligned}$$

• It will be super useful to be able to write the likelihood in two different formulations depending on whether we about the posterior of θ or Σ .



MULTIVARIATE NORMAL LIKELIHOOD

- For $\boldsymbol{\theta}$, it is convenient to write $L(\boldsymbol{Y}; \boldsymbol{\theta}, \boldsymbol{\Sigma})$ as

$$\begin{split} L(\boldsymbol{Y};\boldsymbol{\theta},\boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{-\frac{n}{2}} \exp\left\{-\frac{1}{2}\sum_{i=1}^{n}(\boldsymbol{y}_{i}-\boldsymbol{\theta})^{T}\boldsymbol{\Sigma}^{-1}(\boldsymbol{y}_{i}-\boldsymbol{\theta})\right\} \\ &\propto \exp\left\{-\frac{1}{2}\sum_{i=1}^{n}(\boldsymbol{y}_{i}^{T}-\boldsymbol{\theta}^{T})\boldsymbol{\Sigma}^{-1}(\boldsymbol{y}_{i}-\boldsymbol{\theta})\right\} \\ &= \exp\left\{-\frac{1}{2}\sum_{i=1}^{n}\left[\boldsymbol{y}_{i}^{T}\boldsymbol{\Sigma}^{-1}\boldsymbol{y}_{i}\underbrace{-\boldsymbol{y}_{i}^{T}\boldsymbol{\Sigma}^{-1}\boldsymbol{\theta}-\boldsymbol{\theta}^{T}\boldsymbol{\Sigma}^{-1}\boldsymbol{y}_{i}}_{\text{same term}}+\boldsymbol{\theta}^{T}\boldsymbol{\Sigma}^{-1}\boldsymbol{\theta}\right]\right\} \\ &\propto \exp\left\{-\frac{1}{2}\sum_{i=1}^{n}\left[\boldsymbol{\theta}^{T}\boldsymbol{\Sigma}^{-1}\boldsymbol{\theta}-2\boldsymbol{\theta}^{T}\boldsymbol{\Sigma}^{-1}\boldsymbol{y}_{i}\right]\right\} \\ &= \exp\left\{-\frac{1}{2}\sum_{i=1}^{n}\boldsymbol{\theta}^{T}\boldsymbol{\Sigma}^{-1}\boldsymbol{\theta}-\frac{1}{2}\sum_{i=1}^{n}(-2)\boldsymbol{\theta}^{T}\boldsymbol{\Sigma}^{-1}\boldsymbol{y}_{i}\right\} \\ &= \exp\left\{-\frac{1}{2}n\boldsymbol{\theta}^{T}\boldsymbol{\Sigma}^{-1}\boldsymbol{\theta}+\boldsymbol{\theta}^{T}\boldsymbol{\Sigma}^{-1}\sum_{i=1}^{n}\boldsymbol{y}_{i}\right\} \\ &= \exp\left\{-\frac{1}{2}\boldsymbol{\theta}^{T}(n\boldsymbol{\Sigma}^{-1})\boldsymbol{\theta}+\boldsymbol{\theta}^{T}(n\boldsymbol{\Sigma}^{-1}\boldsymbol{y})\right\}, \end{split}$$

where
$$ar{oldsymbol{y}}=({ar{y}}_1,\ldots,{ar{y}}_p)^T.$$



PRIOR FOR THE MEAN

- A convenient specification of the joint prior is $\pi(\theta, \Sigma) = \pi(\theta)\pi(\Sigma)$.
- As in the univariate case, a convenient conjugate prior distribution for θ is also normal (multivariate in this case).
- Assume that $\pi(\boldsymbol{\theta}) = \mathcal{N}_p(\boldsymbol{\mu}_0, \Lambda_0).$
- The pdf will be easier to work with if we write it as

$$\begin{aligned} \pi(\boldsymbol{\theta}) &= (2\pi)^{-\frac{p}{2}} |\Lambda_0|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\mu}_0)^T \Lambda_0^{-1}(\boldsymbol{\theta} - \boldsymbol{\mu}_0)\right\} \\ &\propto \exp\left\{-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\mu}_0)^T \Lambda_0^{-1}(\boldsymbol{\theta} - \boldsymbol{\mu}_0)\right\} \\ &= \exp\left\{-\frac{1}{2}\left[\boldsymbol{\theta}^T \Lambda_0^{-1} \boldsymbol{\theta} - \boldsymbol{\theta}^T \Lambda_0^{-1} \boldsymbol{\mu}_0 - \boldsymbol{\mu}_0^T \Lambda_0^{-1} \boldsymbol{\theta} + \boldsymbol{\mu}_0^T \Lambda_0^{-1} \boldsymbol{\mu}_0\right]\right\} \\ &\propto \exp\left\{-\frac{1}{2}\left[\boldsymbol{\theta}^T \Lambda_0^{-1} \boldsymbol{\theta} - 2\boldsymbol{\theta}^T \Lambda_0^{-1} \boldsymbol{\mu}_0\right]\right\} \\ &= \exp\left\{-\frac{1}{2}\boldsymbol{\theta}^T \Lambda_0^{-1} \boldsymbol{\theta} + \boldsymbol{\theta}^T \Lambda_0^{-1} \boldsymbol{\mu}_0\right\} \end{aligned}$$



PRIOR FOR THE MEAN

So we have

$$\pi(oldsymbol{ heta}) \propto \exp\left\{-rac{1}{2}oldsymbol{ heta}^T \Lambda_0^{-1}oldsymbol{ heta} + oldsymbol{ heta}^T \Lambda_0^{-1}oldsymbol{\mu}_0
ight\}.$$

- Key trick for combining with likelihood: When the normal density is written in this form, note the following details in the exponent.
 - In the first part, the inverse of the covariance matrix Λ_0^{-1} is "sandwiched" between θ^T and θ .
 - In the second part, the θ in the first part is replaced (sort of) with the mean μ_0 , with Λ_0^{-1} keeping its place.
- The two points above will help us identify updated means and updated covariance matrices relatively quickly.



CONDITIONAL POSTERIOR FOR THE MEAN

• Our conditional posterior (full conditional) $\boldsymbol{\theta}|\Sigma, \boldsymbol{Y}$, is then

$$\pi(\boldsymbol{\theta}|\boldsymbol{\Sigma}, \boldsymbol{Y}) \propto L(\boldsymbol{Y}; \boldsymbol{\theta}, \boldsymbol{\Sigma}) \cdot \pi(\boldsymbol{\theta})$$

$$\propto \underbrace{\exp\left\{-\frac{1}{2}\boldsymbol{\theta}^{T}(n\boldsymbol{\Sigma}^{-1})\boldsymbol{\theta} + \boldsymbol{\theta}^{T}(n\boldsymbol{\Sigma}^{-1}\bar{\boldsymbol{y}})\right\}}_{L(\boldsymbol{Y}; \boldsymbol{\theta}, \boldsymbol{\Sigma})} \cdot \underbrace{\exp\left\{-\frac{1}{2}\boldsymbol{\theta}^{T}\boldsymbol{\Lambda}_{0}^{-1}\boldsymbol{\theta} + \boldsymbol{\theta}^{T}\boldsymbol{\Lambda}_{0}^{-1}\boldsymbol{\mu}_{0}\right\}}_{\pi(\boldsymbol{\theta})}$$

$$= \exp\left\{\underbrace{-\frac{1}{2}\boldsymbol{\theta}^{T}(n\boldsymbol{\Sigma}^{-1})\boldsymbol{\theta} - \frac{1}{2}\boldsymbol{\theta}^{T}\boldsymbol{\Lambda}_{0}^{-1}\boldsymbol{\theta}}_{\text{First parts from } L(\boldsymbol{Y}; \boldsymbol{\theta}, \boldsymbol{\Sigma}) \text{ and } \pi(\boldsymbol{\theta})} + \underbrace{\boldsymbol{\theta}^{T}(n\boldsymbol{\Sigma}^{-1}\bar{\boldsymbol{y}}) + \boldsymbol{\theta}^{T}\boldsymbol{\Lambda}_{0}^{-1}\boldsymbol{\mu}_{0}}_{\text{Second parts from } L(\boldsymbol{Y}; \boldsymbol{\theta}, \boldsymbol{\Sigma}) \text{ and } \pi(\boldsymbol{\theta})}\right\}$$

$$= \exp\left\{-\frac{1}{2}\boldsymbol{\theta}^{T}\left[n\boldsymbol{\Sigma}^{-1} + \boldsymbol{\Lambda}_{0}^{-1}\right]\boldsymbol{\theta} + \boldsymbol{\theta}^{T}\left[n\boldsymbol{\Sigma}^{-1}\bar{\boldsymbol{y}} + \boldsymbol{\Lambda}_{0}^{-1}\boldsymbol{\mu}_{0}\right]\right\},$$

which is just another multivariate normal distribution.



CONDITIONAL POSTERIOR FOR THE MEAN

 To confirm the normal density and its parameters, compare to the prior kernel

$$\pi(oldsymbol{ heta}) \propto \exp\left\{-rac{1}{2}oldsymbol{ heta}^T \Lambda_0^{-1}oldsymbol{ heta} + oldsymbol{ heta}^T \Lambda_0^{-1}oldsymbol{\mu}_0
ight\}$$

and the posterior kernel we just derived, that is,

$$\pi(oldsymbol{ heta}|\Sigma,oldsymbol{Y}) \propto \exp\left\{-rac{1}{2}oldsymbol{ heta}^T \left[\Lambda_0^{-1} + n\Sigma^{-1}
ight]oldsymbol{ heta} + oldsymbol{ heta}^T \left[\Lambda_0^{-1}oldsymbol{\mu}_0 + n\Sigma^{-1}oldsymbol{ar{y}}
ight]
ight\}.$$

- Easy to see (relatively) that $oldsymbol{ heta}|\Sigma,oldsymbol{Y}\sim\mathcal{N}_p(oldsymbol{\mu}_n,\Lambda_n)$, with

$$\Lambda_n = \left[\Lambda_0^{-1} + n\Sigma^{-1}
ight]^{-1}$$

and

$$oldsymbol{\mu}_n = \Lambda_n \left[\Lambda_0^{-1} oldsymbol{\mu}_0 + n \Sigma^{-1} oldsymbol{ar{y}}
ight].$$



BAYESIAN INFERENCE

- As in the univariate case, we once again have that
 - Posterior precision is sum of prior precision and data precision:

 $\Lambda_n^{-1} = \Lambda_0^{-1} + n\Sigma^{-1}$

 Posterior expectation is weighted average of prior expectation and the sample mean:

$$oldsymbol{\mu}_n = \Lambda_n \left[\Lambda_0^{-1} oldsymbol{\mu}_0 + n \Sigma^{-1} oldsymbol{ar{y}}
ight]
onumber \ = \overbrace{\left[\Lambda_n \Lambda_0^{-1}
ight]}^{ ext{weight on prior mean}} oldsymbol{\mu}_0 + \overbrace{\left[\Lambda_n (n \Sigma^{-1})
ight]}^{ ext{weight on sample mean}} oldsymbol{ar{y}}
onumber \ ext{sample mean}$$

 Compare these to the results from the univariate case to gain more intuition.



WHAT ABOUT THE COVARIANCE MATRIX?

- In the univariate case with $y_i \sim \mathcal{N}(\mu, \sigma^2)$, the common choice for the prior is an inverse-gamma distribution for the variance σ^2 .
- As we have seen, we can rewrite as $y_i \sim \mathcal{N}(\mu, \tau^{-1})$, so that we have a gamma prior for the precision τ .
- In the multivariate normal case, we have a covariance matrix Σ instead of a scalar.
- Appealing to have a matrix-valued extension of the inverse-gamma (and gamma) that would be conjugate.



Positive definite and symmetric

- One complication is that the covariance matrix Σ must be **positive** definite and symmetric.
- "Positive definite" means that for all $x \in \mathcal{R}^p$, $x^T \Sigma x > 0$.
- Basically ensures that the diagonal elements of Σ (corresponding to the marginal variances) are positive.
- Also, ensures that the correlation coefficients for each pair of variables are between -1 and 1.
- Our prior for Σ should thus assign probability one to set of positive definite matrices.
- Analogous to the univariate case, the inverse-Wishart distribution is the corresponding conditionally conjugate prior for Σ (multivariate generalization of the inverse-gamma).
- The textbook covers the construction of Wishart and inverse-Wishart random variables. We will skip the actual development in class but will write code to sample random variates.



INVERSE-WISHART DISTRIBUTION

- A random variable $\Sigma \sim \mathrm{IW}_p(
u_0, m{S}_0)$, where Σ is positive definite and p imes p, has pdf

$$p(\Sigma) \propto |\Sigma|^{rac{-(
u_0+p+1)}{2}} \mathrm{exp} \left\{ -rac{1}{2} \mathrm{tr}(oldsymbol{S}_0 \Sigma^{-1})
ight\},$$

where

- $tr(\cdot)$ is the **trace function** (sum of diagonal elements),
- $u_0 > p-1$ is the "degrees of freedom", and
- S_0 is a $p \times p$ positive definite matrix.
- For this distribution, $\mathbb{E}[\Sigma] = rac{1}{
 u_0 p 1} oldsymbol{S}_0$, for $u_0 > p + 1$.
- Hence, S_0 is the scaled mean of the $IW_p(\nu_0, S_0)$.



WISHART DISTRIBUTION

- If we are very confidence in a prior guess Σ_0 , for Σ , then we might set
 - ν₀, the degrees of freedom to be very large, and

•
$$S_0 = (\nu_0 - p - 1)\Sigma_0.$$

In this case,
$$\mathbb{E}[\Sigma] = rac{1}{
u_0 - p - 1} S_0 = rac{1}{
u_0 - p - 1} (
u_0 - p - 1) \Sigma_0 = \Sigma_0$$
, and Σ is tightly (depending on the value of u_0) centered around Σ_0 .

- If we are not at all confident but we still have a prior guess Σ_0 , we might set

•
$$u_0 = p + 2$$
, so that the $\mathbb{E}[\Sigma] = rac{1}{
u_0 - p - 1} S_0$ is finite.
• $S_0 = \Sigma_0$

Here, $\mathbb{E}[\Sigma] = \Sigma_0$ as before, but Σ is only loosely centered around Σ_0 .



WISHART DISTRIBUTION

- Just as we had with the gamma and inverse-gamma relationship in the univariate case, we can also work in terms of the Wishart distribution (multivariate generalization of the gamma) instead.
- The Wishart distribution provides a conditionally-conjugate prior for the precision matrix Σ^{-1} in a multivariate normal model.
- Specifically, if $\Sigma \sim \mathrm{IW}_p(\nu_0, \boldsymbol{S}_0)$, then $\Phi = \Sigma^{-1} \sim \mathrm{W}_p(\nu_0, \boldsymbol{S}_0^{-1})$.
- A random variable $\Phi \sim \mathrm{W}_p(
 u_0, oldsymbol{S}_0^{-1})$, where Φ has dimension (p imes p), has pdf

$$f(\Phi) ~\propto ~ \left| \Phi
ight|^{rac{
u_0-p-1}{2}} ext{exp} \left\{ -rac{1}{2} ext{tr}(oldsymbol{S}_0 \Phi)
ight\}.$$

- Here, $\mathbb{E}[\Phi] = \nu_0 S_0$.
- Note that the textbook writes the inverse-Wishart as IW_p(\u03c6₀, S₀⁻¹). I prefer IW_p(\u03c6₀, S₀) instead. Feel free to use either notation but try not to get confused.



BACK TO INFERENCE ON COVARIANCE

- For inference on Σ, we need to rewrite the likelihood a bit to match the inverse-Wishart kernel.
- First a few results from matrix algebra:
 - 1. $\operatorname{tr}(\boldsymbol{A}) = \sum_{j=1}^p a_{jj}$, where a_{jj} is the *j*th diagonal element of a square $p \times p$ matrix \boldsymbol{A} .
 - 2. Cyclic property:

$$\operatorname{tr}(\boldsymbol{ABC}) = \operatorname{tr}(\boldsymbol{BCA}) = \operatorname{tr}(\boldsymbol{CAB}),$$

given that the product ABC is a square matrix.

3. If $oldsymbol{A}$ is a p imes p matrix, then for a p imes 1 vector $oldsymbol{x}$,

$$\boldsymbol{x}^T \boldsymbol{A} \boldsymbol{x} = \operatorname{tr}(\boldsymbol{x}^T \boldsymbol{A} \boldsymbol{x})$$

holds by (1), since $\boldsymbol{x}^T \boldsymbol{A} \boldsymbol{x}$ is a scalar.

4. $\operatorname{tr}(\boldsymbol{A} + \boldsymbol{B}) = \operatorname{tr}(\boldsymbol{A}) + \operatorname{tr}(\boldsymbol{B}).$



MULTIVARIATE NORMAL LIKELIHOOD AGAIN

It is thus convenient to rewrite L(Y; θ, Σ) as

$$\begin{split} \mathcal{L}(oldsymbol{Y};oldsymbol{ heta},\Sigma) \propto |\Sigma|^{-rac{n}{2}} \exp\left\{-rac{1}{2}\sum_{i=1}^{n}(oldsymbol{y}_{i}-oldsymbol{ heta})^{T}\Sigma^{-1}(oldsymbol{y}_{i}-oldsymbol{ heta})}
ight\} \ &= |\Sigma|^{-rac{n}{2}} \exp\left\{-rac{1}{2}\sum_{i=1}^{n} \operatorname{tr}\left[(oldsymbol{y}_{i}-oldsymbol{ heta})^{T}\Sigma^{-1}(oldsymbol{y}_{i}-oldsymbol{ heta})}
ight]
ight\} \ &= |\Sigma|^{-rac{n}{2}} \exp\left\{-rac{1}{2}\sum_{i=1}^{n} \operatorname{tr}\left[(oldsymbol{y}_{i}-oldsymbol{ heta})(oldsymbol{y}_{i}-oldsymbol{ heta})^{T}\Sigma^{-1}
ight]
ight\} \ &= |\Sigma|^{-rac{n}{2}} \exp\left\{-rac{1}{2}\operatorname{tr}\left[\sum_{i=1}^{n}(oldsymbol{y}_{i}-oldsymbol{ heta})(oldsymbol{y}_{i}-oldsymbol{ heta})^{T}\Sigma^{-1}
ight]
ight\}, \ &= |\Sigma|^{-rac{n}{2}} \exp\left\{-rac{1}{2}\operatorname{tr}\left[\sum_{i=1}^{n}(oldsymbol{y}_{i}-oldsymbol{ heta})(oldsymbol{y}_{i}-oldsymbol{ heta})^{T}\Sigma^{-1}
ight]
ight\}, \end{split}$$

where $m{S}_{ heta} = \sum_{i=1}^n (m{y}_i - m{ heta}) (m{y}_i - m{ heta})^T$ is the residual sum of squares matrix. 24 / 26



CONDITIONAL POSTERIOR FOR COVARIANCE

 Assuming π(Σ) = IW_p(ν₀, S₀), the conditional posterior (full conditional) Σ|θ, Y, is then

$$egin{aligned} \pi(\Sigma|m{ heta},m{Y}) &\propto L(m{Y};m{ heta},\Sigma)\cdot\pi(m{ heta}) \ &\propto |\Sigma|^{-rac{n}{2}}\exp\left\{-rac{1}{2} ext{tr}\left[m{S}_{ heta}\Sigma^{-1}
ight]
ight\}\cdot |\Sigma|^{rac{-(
u_0+p+1)}{2}}\exp\left\{-rac{1}{2} ext{tr}(m{S}_0\Sigma^{-1})
ight\} \ &\propto |\Sigma|^{rac{-(
u_0+p+n+1)}{2}}\exp\left\{-rac{1}{2} ext{tr}\left[m{S}_0\Sigma^{-1}+m{S}_{ heta}\Sigma^{-1}
ight]
ight\}, \ &\propto |\Sigma|^{rac{-(
u_0+n+p+1)}{2}}\exp\left\{-rac{1}{2} ext{tr}\left[m{S}_0\Sigma^{-1}+m{S}_{ heta}\Sigma^{-1}
ight]
ight\}, \end{aligned}$$

which is $IW_p(\nu_n, S_n)$, or using the notation in the book, $IW_p(\nu_n, S_n^{-1})$, with

- $u_n =
 u_0 + n$, and
- $\boldsymbol{S}_n = [\boldsymbol{S}_0 + \boldsymbol{S}_{ heta}]$



CONDITIONAL POSTERIOR FOR COVARIANCE

- We once again see that the "posterior sample size" or "posterior degrees of freedom" ν_n is the sum of the "prior degrees of freedom" ν₀ and the data sample size n.
- S_n can be thought of as the "posterior sum of squares", which is the sum of "prior sum of squares" plus "sample sum of squares".

• Recall that if
$$\Sigma \sim \mathrm{IW}_p(
u_0, oldsymbol{S}_0)$$
, then $\mathbb{E}[\Sigma] = rac{1}{
u_0 - p - 1} oldsymbol{S}_0.$

• \Rightarrow the conditional posterior expectation of the population covariance is

$$\begin{split} \mathbb{E}[\Sigma|\boldsymbol{\theta},\boldsymbol{Y}] &= \frac{1}{\nu_0 + n - p - 1} [\boldsymbol{S}_0 + \boldsymbol{S}_{\theta}] \\ &= \underbrace{\frac{\nu_0 - p - 1}{\nu_0 + n - p - 1}}_{\text{weight on prior expectation}} \left(\underbrace{\frac{1}{\nu_0 - p - 1} \boldsymbol{S}_0}_{\text{weight on sample estimate}} + \underbrace{\frac{n}{\nu_0 + n - p - 1}}_{\text{weight on sample estimate}} \underbrace{\left[\frac{1}{n} \boldsymbol{S}_{\theta}\right]}_{\text{weight on sample estimate}} \right), \end{split}$$

which is a weighted average of prior expectation and sample estimate.

