# IN-CLASS EXERCISE; INTRODUCTION TO MULTIVARIATE NORMAL DR. OLANREWAJU MICHAEL AKANDE

Feb 14, 2020



## ANNOUNCEMENTS

- No homework this week.
- Midterm in three weeks (might seem like a lot but it is NOT!).
- Spend time practicing how to manipulate the univariate and multivariate normal distributions.

## OUTLINE

- In-class exercise
- Multivariate normal distribution



#### **N-CLASS EXERCISE**

- Your friend agrees to conduct a poll for you, free of charge (lucky you!).
- You give the following instructions: "Please ask about 25 people whether they are in favor of more gun control, and report back to me the number who are in favor."
- After a few days your friend returns with the poll results: there were y = 20 in favor. "
- You then ask, "how many people did you ask?"Your friend says, "ummm, I dunno. You didn't ask me to record that. All I know is that it was about 25."
- What model can we use to do inference here?
- To be done on the board.



## PARTICIPATION EXERCISE

- You will work in groups of three. Work with the three students closest to you. Do the following:
  - 1. Using the full conditionals on the board, write a Gibbs sampler to sample from the joint posterior of N and  $\theta$ , using a starting value of N = 50 and  $\theta = 0.05$ . Set burn-in to 2000 and then proceed to generate 10000 draws.
  - 2. Look at the posterior densities for both parameters. Describe the distributions.
  - 3. Give the quantile-based equal-tailed posterior credible interval for  $\theta$ , rounded to two decimal places.
  - 4. What is the probability that exactly 20 people were polled? What can you takeaway from this?
  - 5. What is the probability that exactly 25 people were polled? What can you takeaway from this?



# MULTIVARIATE DATA

- So far we have only considered basic models with scalar/univariate outcomes, Y<sub>1</sub>,..., Y<sub>n</sub>.
- In practice however, outcomes of interest are actually often multivariate, e.g.,
  - Repeated measures of weight over time in a weight loss study
  - Measures of multiple disease markers
  - Tumor counts at different locations along the intestine
- Longitudinal data is just a special case of multivariate data.
- Interest then is often on how multiple outcomes are correlated, and on how that correlation may change across outcomes or time points.



## MULTIVARIATE NORMAL DISTRIBUTION

- The most common model for multivariate outcomes is the multivariate normal distribution.
- Next week, we will do actual inference with the multivariate normal distribution.
- We will explore the common choices for prior distributions and then derive the corresponding posterior distributions.
- Today, we'll start slow and simply explore some properties of the multivariate normal distribution.
- Let  $\mathbf{Y} = (Y_1, \dots, Y_p)^T$ , where p represents the dimension of the multivariate outcome variable for a single unit of observation.
- For multiple observations,  $Y_i = (Y_{i1}, \ldots, Y_{ip})^T$  for  $i = 1, \ldots, n$ .



## MULTIVARIATE NORMAL DISTRIBUTION

•  $m{Y}$  follows a multivariate normal distribution, that is,  $m{Y} \sim \mathcal{N}_p(m{\mu}, \Sigma)$ , if

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

where  $|\Sigma|$  denotes the determinant of A.

- $\mu$  is the  $p \times 1$  mean vector, that is,  $\mu = \mathbb{E}[\mathbf{Y}] = \{\mathbb{E}[Y_1], \dots, \mathbb{E}[Y_p]\} = (\mu_1, \dots, \mu_p)^T.$
- $\Sigma$  is the  $p \times p$  positive definite and symmetric covariance matrix, that is,  $\Sigma = \{\sigma_{jk}\}$ , where  $\sigma_{jk}$  denotes the covariance between  $Y_j$  and  $Y_k$ .
- Note that Y<sub>1</sub>,..., Y<sub>p</sub> may be linearly dependent depending on the structure of Σ, which characterizes association between them.

• For each 
$$j = 1, \ldots, p$$
,  $Y_j \sim \mathcal{N}(\mu_j, \sigma_{jj})$ .



#### **BIVARIATE NORMAL DISTRIBUTION**

In the bivariate case, we have

$$oldsymbol{Y} = egin{pmatrix} Y_1 \ Y_2 \end{pmatrix} \sim \mathcal{N}_2 \left[ \mu = egin{pmatrix} \mu_1 \ \mu_2 \end{pmatrix}, \Sigma = egin{pmatrix} \sigma_{11} = \sigma_1^2 & \sigma_{12} \ \sigma_{21} & \sigma_{22} = \sigma_2^2 \end{pmatrix} 
ight],$$

and  $\sigma_{12} = \sigma_{21} = \mathbb{C}\mathrm{ov}[Y_1,Y_2].$ 

- The correlation between  $Y_1$  and  $Y_2$  is defined as

$$ho_{1,2} = rac{\mathbb{C}\mathrm{ov}[Y_1,Y_2]}{\sqrt{\mathbb{V}\mathrm{ar}[Y_1]}\sqrt{\mathbb{V}\mathrm{ar}[Y_2]}} = rac{\sigma_{12}}{\sigma_1\sigma_2}.$$

- $-1 \le \rho_{1,2} \le 1.$
- Correlation coefficient is free of the measurement units.



#### BACK TO THE MULTIVARIATE NORMAL

- There are many special properties of the multivariate normal as we will see as we continue to work with the distribution.
- First, dependence between any  $Y_j$  and  $Y_k$  does not depend on the other p-2 variables.
- Second, while generally, independence implies zero covariance, for the normal family, the converse is also true. That is, independence implies zero covariance.
- Thus, the covariance Σ carries a lot of information about marginal relationships, especially marginal independence.
- If  $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_p) \sim \mathcal{N}_p(\boldsymbol{0}, \boldsymbol{I}_p)$ , that is,  $\epsilon_1, \dots, \epsilon_p \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ , then

 $oldsymbol{Y} = oldsymbol{\mu} + Aoldsymbol{\epsilon} \Rightarrow oldsymbol{Y} \sim \mathcal{N}_p(oldsymbol{\mu}, \Sigma)$ 

holds for any matrix square root A of  $\Sigma$ , that is,  $AA^T = \Sigma$  (see Cholesky decomposition).



#### 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, and  $oldsymbol{Y}_2$  and  $oldsymbol{\mu}_2$  are (p-q) imes 1;
- $\Sigma_{11}$  is  $q \times q$ , and  $\Sigma_{22}$  is  $(p-q) \times (p-q)$ , with  $\Sigma_{22} > 0$ .
- Then, it turns out that

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

- That is, the conditional distribution of  $Y_1$  given  $Y_2$  is also normal!
- Marginal distributions are once again normal, that is,

 $oldsymbol{Y}_1 \sim \mathcal{N}_q\left(oldsymbol{\mu}_1, \Sigma_{11}
ight); \hspace{1em} oldsymbol{Y}_2 \sim \mathcal{N}_{p-q}\left(oldsymbol{\mu}_2, \Sigma_{22}
ight).$ 



#### CONDITIONAL DISTRIBUTIONS

In the bivariate normal case with

$$oldsymbol{Y} = egin{pmatrix} Y_1 \ Y_2 \end{pmatrix} \sim \mathcal{N}_2 \left[ \mu = egin{pmatrix} \mu_1 \ \mu_2 \end{pmatrix}, \Sigma = egin{pmatrix} \sigma_{11} = \sigma_1^2 & \sigma_{12} \ \sigma_{21} & \sigma_{22} = \sigma_2^2 \end{pmatrix} 
ight],$$

we have

$$Y_1|Y_2=y_2\sim \mathcal{N}\left(\mu_1+rac{\sigma_{12}}{\sigma_2}(y_2-\mu_2),\sigma_1-rac{\sigma_{12}^2}{\sigma_2}
ight).$$

which can also be written as

$$Y_1|Y_2=y_2\sim \mathcal{N}\left(\mu_1+rac{\sigma_1}{\sigma_2}
ho(y_2-\mu_2),(1-
ho^2)\sigma_1^2
ight).$$



## 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  $\theta$ , 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.

