

# STA 360/602L: MODULE 5.1

## HIERARCHICAL NORMAL MODELS WITH CONSTANT VARIANCE: TWO GROUPS

DR. OLANREWAJU MICHAEL AKANDE

# MOTIVATION

- Sometimes, we may have a natural grouping in our data, for example
  - students within schools,
  - patients within hospitals,
  - voters within counties or states,
  - biology data, where animals are followed within natural populations organized geographically and, in some cases, socially.
- For such grouped data, we may want to do inference across all the groups, for example, comparison of the group means.
- Ideally, we should do so in a way that takes advantage of the relationship between observations in the same group, but we should also look to borrow information across groups when possible.
- **Hierarchical modeling** provides a principled way to do so.

# BAYES ESTIMATORS AND BIAS

- Recall the normal model:

$$y_i | \mu, \sigma^2 \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2).$$

- The MLE for the population mean  $\mu$  is just the sample mean  $\bar{y}$ .
- $\bar{y}$  is unbiased for  $\mu$ . That is, for any data  $y_i \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$ ,  $\mathbb{E}[\bar{y}] = \mu$ .
- However, recall that in the conjugate normal model with known variance for example, the posterior expectation is a **weighted average** of the prior mean and the sample mean.
- That is, the posterior mean is actually biased.

# SHRINKAGE

- Usually through the weighting of the sample data and prior, Bayes procedures have the tendency to pull the estimate of  $\mu$  toward the prior mean.
- Of course, the magnitude of the pull depends on the sample size.
- This "pulling" phenomenon is referred to as **shrinkage**.
- Why would we ever want to do this? Why not just stick with the MLE?
- Well, in part, because shrinkage estimators are often "more accurate" in prediction problems – i.e. they tend to do a better job of predicting a future outcome or of recovering the actual parameter values. Remember variance-bias trade off!
- The fact that a biased estimator would do a better job in many prediction problems can be proven rigorously, and is referred to as **Stein's paradox**.

# MODERN RELEVANCE

- Stein's result implies, in particular, that the sample mean is an *inadmissible* estimator of the mean of a multivariate normal distribution in more than two dimensions – i.e. there are other estimators that will come closer to the true value in expectation.
- In fact, these are Bayes point estimators (the posterior expectation of the parameter  $\mu$ ).
- Most of what we do now in high-dimensional statistics is develop biased estimators that perform better than unbiased ones.
- Examples: lasso regression, ridge regression, various kinds of hierarchical Bayesian models, etc.
- So, here we will get a very basic introduction to **Bayesian hierarchical models**, which provide a formal and coherent framework for constructing shrinkage estimators.

# WHY HIERARCHICAL MODELS?

- **Bayesian hierarchical models** is a sort of catch-all phrase for a large class of models that have several levels of conditional distributions making up the prior.
- Like simpler one-level priors, they also accomplish shrinkage. However, they are much more flexible.
- Why use them? Several reasons:
  - We may want to exploit more complex dependence structures.
  - We may have many parameters relative to the amount of data that we have, and want to borrow information in estimating them.
  - We may want to shrink toward something other than a simple prior mean/hyper-parameter.

# COMPARING TWO GROUPS

- Suppose we want to do inference on mean body mass index (BMI) for two groups (male or female).
- BMI is known to often follow a normal distribution, so let's assume the same here.
- We should expect some relationship between the mean BMI for the two groups.
- We may also think the shape of the two distributions would be relatively the same (at least as a simplifying assumption for now).
- Thus, a reasonable model might be

$$\begin{aligned} y_{i,\text{male}} &\overset{iid}{\sim} \mathcal{N}(\theta_m, \sigma^2); \quad i = 1, \dots, n_m; \\ y_{i,\text{female}} &\overset{iid}{\sim} \mathcal{N}(\theta_f, \sigma^2); \quad i = 1, \dots, n_f. \end{aligned}$$

but with some relationship between  $\theta_m$  and  $\theta_f$ .

# BAYESIAN INFERENCE

- One parameterization that can reflect some relationship between  $\theta_m$  and  $\theta_f$  is

$$\begin{aligned} y_{i,\text{male}} &\stackrel{iid}{\sim} \mathcal{N}(\mu + \delta, \sigma^2); \quad i = 1, \dots, n_m; \\ y_{i,\text{female}} &\stackrel{iid}{\sim} \mathcal{N}(\mu - \delta, \sigma^2); \quad i = 1, \dots, n_f. \end{aligned}$$

where

- $\theta_m = \mu + \delta$  and  $\theta_f = \mu - \delta$ ,
- $\mu = \frac{\theta_m + \theta_f}{2}$  is the average of the population means, and
- $2\delta = \theta_m - \theta_f$  is the difference in population means.



# BAYESIAN INFERENCE

- Convenient prior:
  - $\pi(\mu, \delta, \sigma^2) = \pi(\mu) \cdot \pi(\delta) \cdot \pi(\sigma^2)$ , where
    - $\pi(\mu) = \mathcal{N}(\mu_0, \gamma_0^2)$ ,
    - $\pi(\delta) = \mathcal{N}(\delta_0, \tau_0^2)$ , and
    - $\pi(\sigma^2) = \mathcal{IG}(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2})$ .
- We will set the hyper-parameters as:
  - $\mu_0 = 15, \gamma_0 = 5$ ,
  - $\delta_0 = 0, \tau_0 = 3$ ,
  - $\nu_0 = 1, \sigma_0 = 5$ .
- Do these values seem reasonable?

# BAYESIAN INFERENCE

- Note that we can rewrite

$$\begin{aligned} y_{i,\text{male}} &\stackrel{iid}{\sim} \mathcal{N}(\mu + \delta, \sigma^2); \quad i = 1, \dots, n_m; \\ y_{i,\text{female}} &\stackrel{iid}{\sim} \mathcal{N}(\mu - \delta, \sigma^2); \quad i = 1, \dots, n_f \end{aligned}$$

as

$$\begin{aligned} (y_{i,\text{male}} - \delta) &\stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2); \quad i = 1, \dots, n_m; \\ (y_{i,\text{female}} + \delta) &\stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2); \quad i = 1, \dots, n_f \end{aligned}$$

or

$$\begin{aligned} (y_{i,\text{male}} - \mu) &\stackrel{iid}{\sim} \mathcal{N}(\delta, \sigma^2); \quad i = 1, \dots, n_m; \\ (-1)(y_{i,\text{female}} - \mu) &\stackrel{iid}{\sim} \mathcal{N}(\delta, \sigma^2); \quad i = 1, \dots, n_f. \end{aligned}$$

as needed, so we can leverage past results for the full conditionals.

# FULL CONDITIONALS

- For the full conditionals we will derive here, we will take advantage of previous results from the regular univariate normal model.
- Recall that if we assume

$$y_i \sim \mathcal{N}(\mu, \sigma^2), \quad i = 1, \dots, n,$$

and set our priors to be

$$\begin{aligned}\pi(\mu) &= \mathcal{N}(\mu_0, \gamma_0^2) . \\ \pi(\sigma^2) &= \mathcal{IG}\left(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2}\right),\end{aligned}$$

then we have

$$\pi(\mu, \sigma^2 | Y) \propto \left\{ \prod_{i=1}^n p(y_i | \mu, \sigma^2) \right\} \cdot \pi(\mu) \cdot \pi(\sigma^2)$$

# FULL CONDITIONALS

- We have

$$\pi(\mu|\sigma^2, Y) = \mathcal{N}(\mu_n, \gamma_n^2).$$

where

$$\gamma_n^2 = \frac{1}{\frac{n}{\sigma^2} + \frac{1}{\gamma_0^2}}; \quad \mu_n = \gamma_n^2 \left[ \frac{n}{\sigma^2} \bar{y} + \frac{1}{\gamma_0^2} \mu_0 \right],$$

- and

$$\pi(\sigma^2|\mu, Y) = \mathcal{IG}\left(\frac{\nu_n}{2}, \frac{\nu_n \sigma_n^2}{2}\right),$$

where

$$\nu_n = \nu_0 + n; \quad \sigma_n^2 = \frac{1}{\nu_n} \left[ \nu_0 \sigma_0^2 + \sum_{i=1}^n (y_i - \mu)^2 \right].$$

# FULL CONDITIONALS

- With  $\pi(\mu) = \mathcal{N}(\mu_0, \gamma_0^2)$ , and

$$\begin{aligned}(y_{i,male} - \delta) &\stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2); \quad i = 1, \dots, n_m; \\ (y_{i,female} + \delta) &\stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2); \quad i = 1, \dots, n_f,\end{aligned}$$

we have

$$\mu|Y, \delta, \sigma^2 \sim \mathcal{N}(\mu_n, \gamma_n^2), \quad \text{where}$$

$$\begin{aligned}\gamma_n^2 &= \frac{1}{\frac{1}{\gamma_0^2} + \frac{n_m + n_f}{\sigma^2}} \\ \mu_n &= \gamma_n^2 \left[ \frac{\mu_0}{\gamma_0^2} + \frac{\sum_{i=1}^{n_m} (y_{i,male} - \delta) + \sum_{i=1}^{n_f} (y_{i,female} + \delta)}{\sigma^2} \right].\end{aligned}$$

# FULL CONDITIONALS

- With  $\pi(\delta) = \mathcal{N}(\delta_0, \tau_0^2)$ , and

$$\begin{aligned}(y_{i,\text{male}} - \mu) &\stackrel{iid}{\sim} \mathcal{N}(\delta, \sigma^2); \quad i = 1, \dots, n_m; \\ (-1)(y_{i,\text{female}} - \mu) &\stackrel{iid}{\sim} \mathcal{N}(\delta, \sigma^2); \quad i = 1, \dots, n_f,\end{aligned}$$

we have

$$\delta | Y, \mu, \sigma^2 \sim \mathcal{N}(\delta_n, \tau_n^2), \quad \text{where}$$

$$\tau_n^2 = \frac{1}{\frac{1}{\tau_0^2} + \frac{n_m + n_f}{\sigma^2}}$$

$$\delta_n = \tau_n^2 \left[ \frac{\delta_0}{\tau_0^2} + \frac{\sum_{i=1}^{n_m} (y_{i,\text{male}} - \mu) + (-1) \sum_{i=1}^{n_f} (y_{i,\text{female}} - \mu)}{\sigma^2} \right].$$

# FULL CONDITIONALS

- With  $\pi(\sigma^2) = \mathcal{IG}(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2})$ , and

$$y_{i,male} \stackrel{iid}{\sim} \mathcal{N}(\mu + \delta, \sigma^2); \quad i = 1, \dots, n_m;$$
$$y_{i,female} \stackrel{iid}{\sim} \mathcal{N}(\mu - \delta, \sigma^2); \quad i = 1, \dots, n_f$$

we have

$$\sigma^2 | Y, \mu, \delta \sim \mathcal{IG}(\frac{\nu_n}{2}, \frac{\nu_n \sigma_n^2}{2}), \quad \text{where}$$

$$\nu_n = \nu_0 + n_m + n_f$$

$$\sigma_n^2 = \frac{1}{\nu_n} \left[ \nu_0 \sigma_0^2 + \sum_{i=1}^{n_m} (y_{i,male} - [\mu + \delta])^2 + \sum_{i=1}^{n_f} (y_{i,female} - [\mu - \delta])^2 \right].$$

- We will use write a Gibbs sampler for this model and fit the model to real data in the next module.

# WHAT'S NEXT?

MOVE ON TO THE READINGS FOR THE NEXT MODULE!