

Consider the dataset in "mathtest.csu", which shows the scores on a standardized math test from a sample of 10th grade students at 100 US urban high schools. Let θ_i be the underlying mean test score for school i , and let y_{ij} be the score for the j th student at school i . Notice the extreme school-level averages \bar{y}_i (both high and low) tend to be at schools where fewer students were sampled.

1.) Explain briefly why this would be.

The sample variance is $\frac{\sigma^2}{n}$ where σ^2 is the true population variance. Therefore, it is easy to see that when n is small, the variance is larger. The variance as a function of sample size is a funnel plot.

The mean value theorem states that the sample mean would converge to the population mean for large n . That is, as the sample size grows, the probability that the sample mean is very different from the true population mean goes to 0. Now, we do not know what the true value of θ_i are for the schools that had the fewest samples but it is not unreasonable to assume that the variance of θ_i is not too large. Therefore, supposing that the θ_i are not too far from their overall mean μ , the \bar{y}_i that come from schools with large sample sizes would also be relatively close in value to their θ_i and μ by the mvt. However, when the sample size is small, the mvt no longer holds and \bar{y}_i can deviate greatly from θ_i , and by extension, $\bar{y}_{j|i}$.

2.) Fit this normal hierarchical model to the math test data via Gibbs sampling:

$$\begin{aligned} y_{ij} | \theta_i &\sim N(\theta_i, \sigma^2) \\ \theta_i &\sim N(\mu, \tau^2 \sigma^2) \end{aligned} \quad E(y) = E(E(y | \theta_i)) = \mu$$

Decide upon sensible priors for the unknown model parameters (μ, σ^2, τ^2) .

In order to implement gibbs sampling, we need to derive the full conditionals $p(\theta_i | \theta_{-i}, \text{data})$. First I will state the pieces of the joint posterior. To make the derivations easier, I will switch to working with precisions. Let $\lambda = \frac{1}{\sigma^2}$, $\nu = \frac{1}{\tau^2}$.

$$\begin{aligned} p(y_{ij} | \theta_i, \lambda) &\sim N(\theta_i, \frac{1}{\lambda}) \\ p(\theta_i | \mu, \nu, \lambda) &\sim N(\mu, \frac{1}{\lambda \nu}) \end{aligned}$$

I will assign conjugate priors to μ, λ , and ν . Let

$$\begin{aligned} p(\mu | 0, \gamma) &\sim N(0, \frac{1}{\gamma}) \\ p(\lambda | a, b) &\sim \text{Ga}(a, b) \\ p(\nu | c, d) &\sim \text{Ga}(c, d) \end{aligned}$$

We can now construct the full posterior

$$p(\theta, \lambda, \mu, \nu | y_{ij}) \propto \prod_{i=1}^I \prod_{j=1}^{n_i} [p(y_{ij} | \theta_i, \lambda)] \prod_{i=1}^I [p(\theta_i | \mu, \nu, \lambda)] p(\mu) p(\nu) p(\lambda)$$

$$\begin{aligned}
 p(\theta, \lambda, \mu, \nu | y_{ij}) &\propto \prod_{i=1}^n \prod_{j=1}^m [p(y_{ij} | \theta_i, \lambda)] \prod_{i=1}^n [p(\theta_i | \mu, \nu, \lambda)] p(\mu) p(\nu) p(\lambda) \\
 &\propto \prod_{i=1}^n \prod_{j=1}^m \lambda^{\frac{1}{2}} \exp(-\frac{1}{2} (y_{ij} - \theta_i)^2) \prod_{i=1}^n \left[(\nu \lambda)^{\frac{1}{2}} \exp(-\frac{\nu \lambda}{2} (\theta_i - \mu)^2) \right] \exp(-\frac{\gamma}{2} (\mu - \mu_0)^2) \lambda^{a-1} \exp(-b\lambda) \nu^{c-1} \exp(-d\nu) \\
 &= \lambda^{\frac{\sum_{i,j} m_i}{2}} \exp(-\frac{1}{2} \sum_{i,j} (y_{ij} - \theta_i)^2) (\nu \lambda)^{\frac{n}{2}} \exp(-\frac{\nu \lambda}{2} \sum_{i=1}^n (\theta_i - \mu)^2) \exp(-\frac{\gamma}{2} (\mu - \mu_0)^2) \lambda^{(a-1)} \exp(-b\lambda) \nu^{(c-1)} \exp(-d\nu)
 \end{aligned}$$

Then the full conditionals are of the form

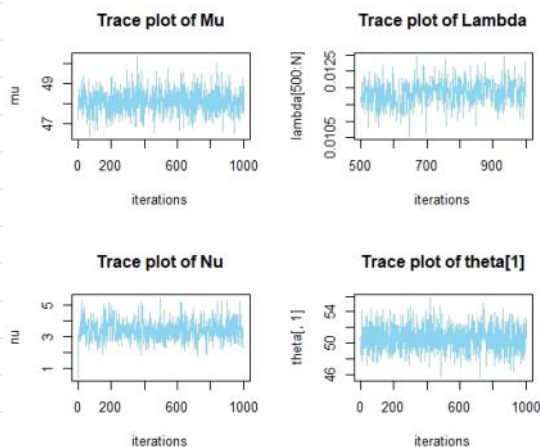
$$\begin{aligned}
 p(\theta_i | y_{ij}, \lambda, \mu, \nu) &\propto \exp(-\frac{1}{2} \sum_{j=1}^m (y_{ij} - \theta_i)^2) \exp(-\frac{\nu \lambda}{2} (\theta_i - \mu)^2) \\
 &= \exp(-\frac{1}{2} \sum_{j=1}^m y_{ij}^2 - 2\theta_i \sum_{j=1}^m y_{ij} + \theta_i^2) \exp(-\frac{\nu \lambda}{2} (\theta_i^2 - 2\theta_i \mu + \mu^2)) \\
 &\propto \exp(-\frac{1}{2} (m\theta_i^2 - 2\theta_i (\sum_{j=1}^m y_{ij})) - \frac{\nu \lambda}{2} (\theta_i^2 - 2\theta_i \mu)) \\
 &= \exp(-\frac{1}{2} [(m\lambda + \nu \lambda) \theta_i^2 - 2\theta_i (\lambda \sum_{j=1}^m y_{ij} + \nu \lambda \mu)]) \\
 &\propto \exp(-\frac{m\lambda + \nu \lambda}{2} (\theta_i - \frac{\lambda \sum_{j=1}^m y_{ij} + \nu \lambda \mu}{m\lambda + \nu \lambda})^2) \\
 &\sim N(\frac{\lambda \sum_{j=1}^m y_{ij} + \nu \lambda \mu}{m\lambda + \nu \lambda}, \frac{1}{m\lambda + \nu \lambda})
 \end{aligned}$$

$$\begin{aligned}
 p(\lambda | y_{ij}, \mu, \theta_i, \nu) &\propto \lambda^{\frac{\sum_{i,j} m_i}{2} + a - 1} \exp(-\lambda (\frac{1}{2} \sum_{i,j} (y_{ij} - \theta_i)^2 + \frac{\nu}{2} \sum_{i=1}^n (\theta_i - \mu)^2 + b)) \\
 &\sim \text{Ga}(\frac{\sum_{i,j} m_i}{2} + a, \frac{1}{2} \sum_{i,j} (y_{ij} - \theta_i)^2 + \frac{\nu}{2} \sum_{i=1}^n (\theta_i - \mu)^2 + b)
 \end{aligned}$$

$$\begin{aligned}
 p(\nu | \theta_i, \mu, c, d) &\propto \nu^{\frac{n}{2} + c - 1} \exp(-\nu (\frac{1}{2} \sum_{i=1}^n (\theta_i - \mu)^2 + d)) \\
 &\sim \text{Ga}(\frac{n}{2} + c, \frac{1}{2} \sum_{i=1}^n (\theta_i - \mu)^2 + d)
 \end{aligned}$$

$$\begin{aligned}
 p(\mu | \theta_i, \gamma) &\propto \exp(-\frac{\nu \lambda}{2} \sum_{i=1}^n (\theta_i - \mu)^2) \exp(-\frac{\gamma}{2} (\mu - \mu_0)^2) \\
 &\propto \exp(-\frac{\nu \lambda}{2} (n\mu^2 - 2\mu \sum_{i=1}^n \theta_i) - \frac{\gamma}{2} (\mu^2 - 2\mu \mu_0 + \mu_0^2)) \\
 &\propto \exp(-\frac{\nu \lambda n + \gamma}{2} (\mu - \frac{\nu \lambda \sum_{i=1}^n \theta_i + \gamma \mu_0}{\nu \lambda n + \gamma})^2) \\
 &\sim N(\frac{\nu \lambda \sum_{i=1}^n \theta_i + \gamma \mu_0}{\nu \lambda n + \gamma}, \frac{1}{\nu \lambda n + \gamma})
 \end{aligned}$$

For the hyperparameters, I will let $\mu_0 = \text{average}(\bar{y}_i)$. I will set γ to something small so that the variance is large. This will make the prior of μ more uninformative. I will let $a=b=c=d=1$.



The posterior mean of $\mu = 48.119$

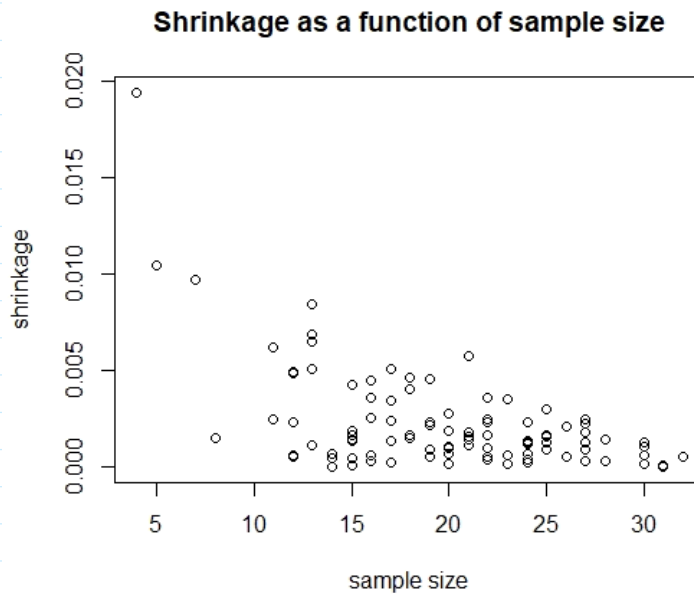
$$\lambda = 0.012$$

$$\nu = 3.410$$

3-) Define the shrinkage coefficient K_i as

$$K_i = \frac{\bar{y}_i - \hat{\theta}_i}{y_i}$$

which tells you how much the posterior mean shrinks the observed sample mean. Plot this shrinkage coefficient (in absolute value) for each school as a function of that school's sample size, and comment.



The shrinkage values are higher for schools with the smaller samples. This is in line with our expectations. With fewer samples, the $\sum_{j=1}^n y_{ij}$ would be smaller and so the posterior mean of θ_i would be pulled towards the prior mean μ . This is the effect we want for this dataset. Since having fewer samples makes the sampling distribution more variable, we want to shrink these means towards the overall mean to remove some of the unwarranted variance.