

MATH TEST DATASET

dataset summary:

- $P = 100$ highschools
- $n_i = \#$ student for $i = 1, \dots, P$
 $\hookrightarrow n_i \geq 400 \quad \forall i$

- $y_{ij} = \text{score for } j\text{th student in school } i$

$$y_{ij} \sim N(\theta_i, \sigma^2) \quad \forall i = 1, \dots, P, \quad \forall j = 1, \dots, n_i$$

model:

$$\theta_i \sim N(\mu, \tau^2 \sigma^2) \quad \forall i = 1, \dots, P$$

A) Find MLE for $\theta_i, \forall i = 1, \dots, P$

$$\mathcal{L}(\theta, y) = \prod_{i=1}^P \prod_{j=1}^{n_i} p(y_{ij} | \theta_i)$$

$$\propto \prod_{i=1}^P \prod_{j=1}^{n_i} \exp \left\{ -\frac{1}{2\sigma^2} (y_{ij} - \theta_i)^2 \right\}$$

$$= \prod_{i=1}^P \exp \left\{ -\frac{1}{2\sigma^2} \sum_{j=1}^{n_i} (y_{ij} - \theta_i)^2 \right\}$$

$$= \prod_{i=1}^P \exp \left\{ -\frac{1}{2\sigma^2} \sum_{j=1}^{n_i} (y_{ij}^2 - 2y_{ij}\theta_i + \theta_i^2) \right\}$$

$$\propto \prod_{i=1}^P \exp \left\{ -\frac{1}{2\sigma^2} \sum_{j=1}^{n_i} (\theta_i^2 - 2y_{ij}\theta_i) \right\}$$

so that

$$\log \mathcal{L}(\theta; y) \propto \sum_{i=1}^p \frac{1}{2\sigma^2} \sum_{j=1}^{n_i} (\theta_i^2 - 2y_{ij}\theta_i)$$

so that

$$\frac{\partial \mathcal{L}(\theta; y)}{\partial \theta_i} \propto \sum_{j=1}^{n_i} (\theta_i - y_{ij}) = n_i \theta_i - \sum_{j=1}^{n_i} y_{ij}$$

so that

$$\frac{\partial \mathcal{L}(\theta; y)}{\partial \theta_i} = 0 \rightarrow \theta_i = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij} = \bar{y}_i, \quad \forall i \in \{1, \dots, p\}$$

B) As suggested, we can notice that we observe extreme score (average) in schools in which n_i is higher. This is because smaller datasets will be more affected by outliers in their estimate.

(see notebook)

c) Consider the following two level hierarchical model:

- $y_{ij} | \theta_i, \delta^2 \sim N(\theta_i, \delta^2)$
- $\theta_i | \tau^2, \delta^2 \sim N(\mu, \tau^2 \delta^2)$

with priors:

- $p(\mu) \propto C$
- $p(\delta^2) \sim \text{Jeffrey's} \quad \longrightarrow \quad p(\delta^2) = \frac{1}{\delta^2}$
- $\tau^2 \sim \text{In-Gamma}(\frac{1}{2}, \frac{1}{2})$

In order to fit a Gibbs sampler for the posterior, we first need to derive the full conditionals:

- 1) $p(\theta_i | y_i, \delta^2, \mu, \tau^2)$
- 2) $p(\delta^2 | \theta_i, y_i, \mu, \tau^2)$
- 3) $p(\mu | \theta_i, \delta^2, y_i, \tau^2)$
- 4) $p(\tau^2 | \theta_i, \delta^2, y_i, \mu)$

$$1) p(\theta_i | -) \propto \prod_{j=1}^{n_i} p(y_{ij} | \theta_i, \sigma^2) \cdot p(\theta_i)$$

$$\propto \exp\left\{-\frac{1}{2\sigma^2} \sum_{j=1}^{n_i} (y_{ij} - \theta_i)^2\right\} \cdot \exp\left\{-\frac{1}{2\tau^2\sigma^2} (\theta_i - \mu)^2\right\}$$

$$\propto \exp\left\{-\frac{1}{2\sigma^2} \sum_{j=1}^{n_i} (\theta_i^2 - 2y_{ij}\theta_i)\right\} \cdot \exp\left\{-\frac{1}{2\tau^2\sigma^2} (\theta_i^2 - 2\theta_i\mu)\right\}$$

$$= \exp\left\{-\frac{1}{2\sigma^2} (n_i\theta_i^2 - 2\bar{y}_i \cdot n_i\theta_i)\right\} \cdot \exp\left\{-\frac{1}{2\tau^2\sigma^2} (\theta_i^2 - 2\theta_i\mu)\right\}$$

$$= \exp\left\{-\frac{n_i}{2\sigma^2\tau^2} \left(\theta_i^2\tau^2 - 2\bar{y}_i\tau^2\theta_i + \frac{\theta_i^2}{n_i} - \frac{2\theta_i\mu}{n_i}\right)\right\}$$

$$= \exp\left\{-\frac{n_i}{2\sigma^2\tau^2} \left(\theta_i^2(\tau^2 + n_i^{-1}) - 2\theta_i \left(\bar{y}_i \cdot \tau^2 + \mu \cdot n_i^{-1}\right)\right)\right\}$$

$$= \exp\left\{-\frac{n_i(\tau^2 + n_i^{-1})}{2\sigma^2\tau^2} \left(\theta_i^2 - 2\theta_i \left(\frac{\bar{y}_i \cdot \tau^2 + \mu \cdot n_i^{-1}}{\tau^2 + n_i^{-1}}\right)\right)\right\}$$

completing the square:

$$\theta_i | - \sim N\left(\frac{\bar{y}_i \cdot \tau^2 + \mu \cdot n_i^{-1}}{\tau^2 + n_i^{-1}}; \frac{\sigma^2\tau^2}{(n_i\tau^2 + 1)}\right)$$

⇒ Note: we can update all θ_i together using

$$\theta | - \sim \text{MVN}(m, V)$$

with:
 • m = vectors of means as in
 • V = diagonal matrix with elements as in

(1)

$$2) p(\sigma^2 | -) \propto \left(\prod_{i=1}^P \left(\prod_{j=1}^{n_i} p(y_{ij} | \theta_i, \sigma^2) \right) p(\theta_i | \sigma^2, \tau^2, \mu) \right) p(\sigma^2)$$

$$\propto \left(\prod_{i=1}^P \left(\prod_{j=1}^{n_i} \frac{1}{\sigma} \exp \left\{ -\frac{1}{2\sigma^2} (y_{ij} - \theta_i)^2 \right\} \right) \cdot \frac{1}{\sigma} \exp \left\{ -\frac{1}{2\sigma^2 \tau^2} (\theta_i - \mu)^2 \right\} \right) \frac{1}{\sigma^2} \quad \left. \begin{array}{l} \text{matrix} \\ \text{form} \\ \text{for } \theta \end{array} \right\}$$

$$= (\sigma^2)^{-\frac{N+P}{2}-1} \cdot \exp \left\{ -\frac{1}{2\sigma^2} \left(\sum_{i=1}^P \sum_{j=1}^{n_i} (y_{ij} - \theta_i)^2 + (\theta - \mu)' (\tau^2 I) (\theta - \mu) \right) \right\}$$

so that:

$$\sigma^2 | - \sim \text{IG} \left(\frac{\sum_{i=1}^P n_i + P}{2}, \frac{\sum_{i=1}^P \sum_{j=1}^{n_i} (y_{ij} - \theta_i)^2 + (\theta - \mu)' (\tau^2 I) (\theta - \mu)}{2} \right) \quad (2)$$

$$3) p(\mu | -) \propto \left(\prod_{i=1}^P p(\theta_i | \mu, \sigma^2, \tau^2) \right) \cdot p(\mu)$$

$$\propto \exp \left\{ -\frac{1}{2\sigma^2 \tau^2} \sum_{i=1}^P (\theta_i - \mu)^2 \right\}$$

$$\propto \exp \left\{ -\frac{1}{2\sigma^2 \tau^2} (P\mu^2 - 2P\bar{\theta} \cdot \mu) \right\}$$

$$= \exp \left\{ -\frac{P}{2\sigma^2 \tau^2} (\mu^2 - 2\bar{\theta} \mu) \right\}$$

completing the square:

$$\mu | - \sim N \left(\bar{\theta}, \frac{\sigma^2 \tau^2}{P} \right) \quad (3)$$

$$4) P(\gamma^2 | \rightarrow) \propto \left(\prod_{i=1}^P p(\theta_i | \mu, \sigma^2, \gamma^2) \right) p(\gamma^2)$$

$$\propto \left(\frac{1}{\gamma^2} \right)^P \cdot \exp \left[-\frac{1}{2\sigma^2 \gamma^2} \sum_{i=1}^P (\theta_i - \mu)^2 \right] \cdot (\gamma^2)^{-\frac{1}{2}-1} \cdot \exp \left\{ -\frac{1}{2} \gamma^2 \right\}$$

$$= (\gamma^2)^{\frac{-P-1}{2}-1} \cdot \exp \left\{ -\frac{1}{\gamma^2} \left(\frac{(\theta - \mu)' (\sigma^2 I)^{-1} (\theta - \mu) + 1}{2} \right) \right\}$$

using matrix notation for θ

so that

$$\gamma^2 | \rightarrow \sim \text{Inv-Ge} \left(\frac{P+1}{2}, \frac{(\theta - \mu)' (\sigma^2 I)^{-1} (\theta - \mu) + 1}{2} \right) \quad (4)$$

Here, a Gibbs sampler will iteratively update the parameters using

(1), (2), (3), (4)

D) We have, as derived in c):

$$\theta_i | - \sim N \left(\frac{\bar{y}_i \cdot \tau^2 + \mu \cdot n_i^{-1}}{\tau^2 + n_i^{-1}} ; \frac{\sigma^2 \tau^2}{(n_i \tau^2 + 1)} \right)$$

so that

$$E(\theta_i | -) = \frac{\bar{y}_i \cdot \tau^2 + \mu \cdot n_i^{-1}}{\tau^2 + n_i^{-1}}$$

which can be expressed as a convex combination of prior mean and data mean. Indeed:

$$= \underbrace{\bar{y}_i \frac{\tau^2}{\tau^2 + n_i^{-1}}}_{1 - \kappa_i} + \mu \underbrace{\frac{n_i^{-1}}{\tau^2 + n_i^{-1}}}_{\kappa_i} \rightarrow \text{SHRINKAGE COEFFICIENT}$$

(see notebook for the plot)

One can notice how, when the group size (sample size) increases, the shrinkage effect is more evident.

⇒ Note that you can also write the model as:

$$y_{ij} = \mu + \delta_i + e_{ij}$$

\swarrow
 $N(0, \tau^2 \sigma^2)$

\searrow
 $N(0, \sigma^2)$

so that:

CONDITIONALLY ON μ :

• $\text{cov}(y_{ij}, y_{ik})$ for $j \neq k$ (same school but different individual)

$$\begin{aligned}
 &= E[(\mu + \delta_i + e_{ij} - E(y_{ij}))(\mu + \delta_i + e_{ik} - E(y_{ik}))] \\
 &= E[(\mu + \delta_i + e_{ij} - \mu)(\mu + \delta_i + e_{ik} - \mu)] \\
 &= E[(\delta_i^2 + e_{ij}\delta_i + e_{ik}\delta_i + e_{ij} \cdot e_{ik})] \\
 &= \underbrace{E(\delta_i^2)}_{\tau^2 \sigma^2 = \text{VAR}(\delta_i)} + 2 \underbrace{E(e_{ij}\delta_i)}_{= \text{cov}(e_{ij}, \delta_i) = 0} + \underbrace{E(e_{ij} \cdot e_{ik})}_{= 0 \text{ since independent}} = \boxed{\tau^2 \sigma^2}
 \end{aligned}$$

• $\text{cov}(y_{ij}, y_{i'k})$ for $i \neq i'$

$$\begin{aligned}
 &= E[(\delta_i + e_{ij})(\delta_{i'} + e_{i'k})] \\
 &= E[(\delta_i \delta_{i'} + e_{ij} \delta_{i'} + e_{i'k} \delta_i + e_{ij} \cdot e_{i'k})] \\
 &= \underbrace{E(\delta_i \delta_{i'})}_{= \text{cov}(\delta_i, \delta_{i'}) = 0} + \underbrace{E(e_{ij} \delta_{i'})}_{= 0 \text{ (indep.)}} + \underbrace{E(e_{i'k} \delta_i)}_{= 0} + \underbrace{E(e_{ij} \cdot e_{i'k})}_{\text{indep by hp.} = 0} = \boxed{\emptyset}
 \end{aligned}$$

If we consider assuming that observations coming from different schools were completely independent, I would say that I am not convinced by the assumption. Indeed, think for example about the fact that all the observations come from the same country, the US, so some sort of correlation between different schools should be expected.

In practice, we are not making this assumption. Indeed, we are assuming independence once conditioning on the ground mean μ , which covered account for the common effect.

Hence, the assumption seems realistic.

¶) We now focus on the assumption that σ^2 is the same for all schools.

This assumption doesn't allow for schools to differentiate in terms of variety of students' input, even when the number of students is considerably different.

We could for example have an "elite" school in which all students were really well prepared for the math exam and also have a school in which people in some classes focused more on math and scored higher while people in other classes focused more on classical studies and scored lower.

The model should be such that different school variances are allowed

A more formal fitting strategy for this hypothesis is presented in the code.

The strategy proposed in class was to perform a POSTERIOR PREDICTIVE CHECK.

In other words, we could simulate a dataset obtained using the posterior parameters and see if such a dataset "resembles" the original one.

For implementation, see *mathscores notebook*
looking at the p-value obtained doing an F test for the standard deviations, the null hypothesis is not rejected and the assumption made by the model seems to be confirmed.

BLOOD PRESSURE DATASET

20 individuals



some are given a treatment and some are not.



each with repeated observations about blood pressure

A) Want to check if the treatment is effective in reducing blood pressure.

⇒ Naive solution: perform a t-test on the difference of the means

(treating all the observations from same patient as independent)

To perform this test is not a good strategy because it assumes that the sample are independent.

This is not the case: values are strongly correlated since they are longitudinal data coming from few individuals

(see notebook for the test)

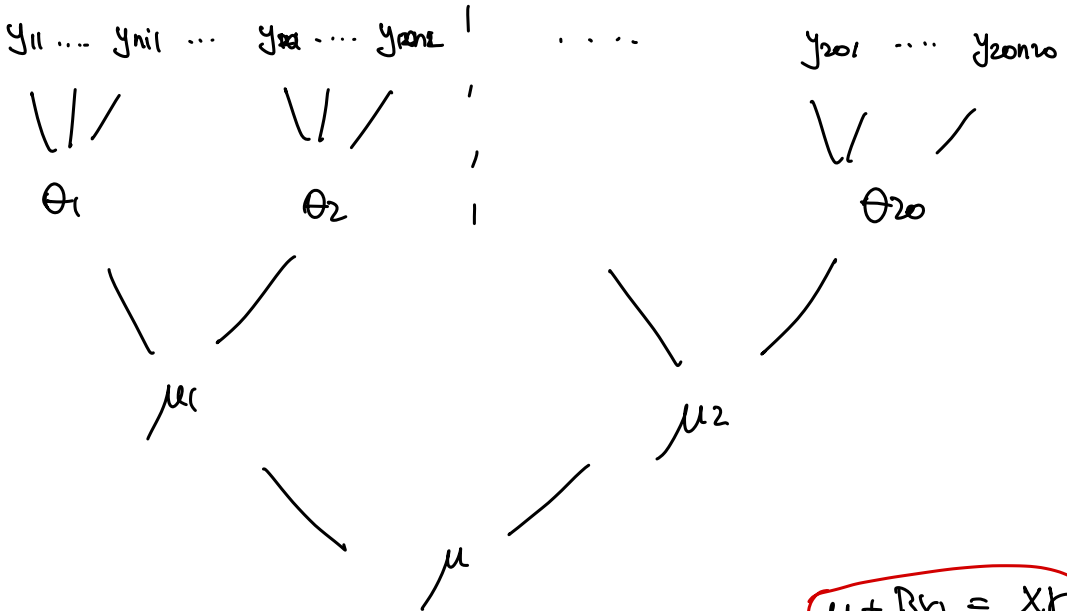
B) The standard error of the test is considerably higher because now we are considering fewer data points (by collapsing data from the same individual to only be represented by the mean of that individual). (see notebook)

Still, this test is not ideal because of **heteroskedasticity**.

20 different subjects \rightarrow groups: θ_i

$$y_{ij} | \theta_i, \sigma^2 \sim N(\theta_i, \sigma^2)$$

$$\theta_i | \gamma^2, \sigma^2, \beta \sim N(\mu + \beta x_i, \gamma^2 \sigma^2)$$



$$\mu + \beta x_i = x_i \gamma$$

- riors:
- $p(\mu, \beta) \propto C \rightarrow p(\gamma) \propto C$
 - $p(\sigma^2) \propto \frac{1}{\sigma^2}$
 - $p(\gamma^2) \propto \text{Inv-G}(\frac{1}{2}, \frac{1}{2})$

So that the derivations for the poster will be similar:

$$1) p(\theta_i | -) \propto \left(\prod_{j=1}^{n_i} p(y_{ij} | \theta_i, \delta^2) \right) p(\theta_i | \gamma, \delta^2 \tau^2)$$

$$\propto \left(\prod_{j=1}^{n_i} \exp \left\{ -\frac{1}{2\delta^2} (y_{ij} - \theta_i)^2 \right\} \right) \exp \left\{ -\frac{1}{2\delta^2 \tau^2} (\theta_i - \cancel{x_i' \gamma})^2 \right\}$$

will be the same if poster

$$= \exp \left\{ -\frac{1}{2\delta^2} \left(\sum_{j=1}^{n_i} (y_{ij} - \theta_i)^2 + \tau^2 (\theta_i - x_i' \gamma)^2 \right) \right\}$$

and following the same steps followed in A) for the previous dataset, but with now $\mu = x_i' \gamma$:

$$\Rightarrow \theta_i | - \sim N \left(\frac{\bar{y}_i \cdot \tau^2 + \cancel{x_i' \gamma} \cdot n_i}{\tau^2 + n_i} ; \frac{\delta^2 \tau^2}{(n_i \tau^2 + 1)} \right)$$

$$2) p(\delta^2 | -) \propto \prod_{i=1}^P \left(\prod_{j=1}^{n_i} p(y_{ij} | \theta_i, \delta^2) \right) \cdot p(\theta_i | \delta^2, \tau^2, \gamma) \cdot p(\delta^2)$$

$$\left(\prod_{i=1}^P \left(\prod_{j=1}^{n_i} \frac{1}{\delta} \cdot \exp \left\{ -\frac{1}{2\delta^2} (y_{ij} - \theta_i)^2 \right\} \right) \cdot \frac{1}{\delta} \cdot \exp \left\{ -\frac{1}{2\tau^2} (\theta_i - x_i' \gamma)^2 \right\} \right) \frac{1}{\delta^2}$$

$$= \frac{P+N}{(\delta^2)^2} \cdot \exp \left\{ -\frac{1}{2\delta^2} \left(\sum_{i=1}^P \sum_{j=1}^{n_i} (y_{ij} - \theta_i)^2 + (\theta - \cancel{x_f})' \tau^2 I (\theta - \cancel{x_f}) \right) \right\}$$

vector of 20

$$\Rightarrow \delta^2 | - \sim \text{Inv-Gol} \left(\frac{N+P}{2} ; \frac{\sum_{i=1}^P \sum_{j=1}^{n_i} (y_{ij} - \theta_i)^2 + (\theta - x_f)' \tau^2 I (\theta - x_f)}{2} \right)$$

3) $\gamma \rightarrow$ now bidimensional. still we are giving flat prior

$$p(\gamma|-) \propto \prod_{i=1}^P p(\theta_i | \gamma, \sigma^2, \tau^2) \cdot p(\gamma)$$

θ_i
condition
indep

$$\propto \exp \left\{ -\frac{1}{2} (\theta - X\gamma)' (\sigma^2 I + \tau^2 I) (\theta - X\gamma) \right\}$$

$$\propto \exp \left\{ -\frac{1}{2} (\theta - X\gamma)' (1) (\theta - X\gamma) \right\}$$

$$= \exp \left\{ -\frac{1}{2} [\theta' \theta - \gamma' X' \theta + \gamma' X' \theta - \gamma' X' \theta + \gamma' X' \theta - \gamma' X' \theta] \right\}$$

$$\propto \exp \left\{ -\frac{1}{2} [-2\gamma' X' \theta + \gamma' X' \theta] \right\}$$

so that completing the square:

$$\gamma|- \sim N(m, V)$$

$$\cdot V = (X' \lambda X)^{-1} = (X' (\sigma^2 \tau^2 I)^{-1} X)^{-1} = \sigma^2 \tau^2 (X' X)^{-1}$$

$$\cdot m = \cancel{\sigma^2 \tau^2} (X' X)^{-1} \cdot X' \frac{1}{\cancel{\sigma^2 \tau^2}} I \theta = (X' X)^{-1} X' \theta$$

$$4) p(\tau^2|-) \propto \left(\prod_{i=1}^P p(\theta_i | \gamma, \sigma^2, \tau^2) \right) p(\tau^2)$$

$$\propto \left(\frac{1}{\tau^2} \right)^P \cdot \exp \left\{ -\frac{1}{2} \frac{(\theta - X\gamma)' (\sigma^2 I)^{-1} (\theta - X\gamma)}{\tau^2} \right\} \cdot (\tau^2)^{-\frac{1}{2} - 1} \cdot \exp \left\{ -\frac{1}{2} \tau^2 \right\}$$

$$= (\tau^2)^{-P - \frac{1}{2} - 1} \cdot \exp \left\{ -\frac{1}{\tau^2} \left(\frac{(\theta - X\gamma)' (\sigma^2 I)^{-1} (\theta - X\gamma) + 1}{2} \right) \right\}$$

$$\Rightarrow \tau^2|- \propto \text{In-Ga} \left(P + \frac{1}{2}, \frac{(\theta - X\gamma)' (\sigma^2 I)^{-1} (\theta - X\gamma) + 1}{2} \right)$$

\Rightarrow Now looks more like the posterior you would have under a linear model $y = X\beta$

b) (see notebook)

looking at autocorrelation plot, the hypothesis of independence (and conditional on θ_i), seems to be realistic