

# Hierarchical Modelling

Arthur White

# Overview

- ▶ We have seen how to fit Bayesian models to data using MCMC
- ▶ This is a powerful computational tool that lets us fit a lot of different models in a Bayesian framework
- ▶ We are going to look at a simple hierarchical model that is a natural fit in this framework

# Statistical models – review

- ▶ So far we have looked at
  - ▶ Binomial models
  - ▶ Normal models
- ▶ We have seen how to fit these models to data in a Bayesian framework
  - ▶ Assign priors to parameters
  - ▶ Fit using MCMC if needed
  - ▶ Interpret and posterior distribution for parameters as needed

## Normal distribution - generative model

$$\mu \sim \mathcal{N}(\mu_0, 1/\tau_0)$$

$$\tau \sim \mathcal{G}(a, b)$$

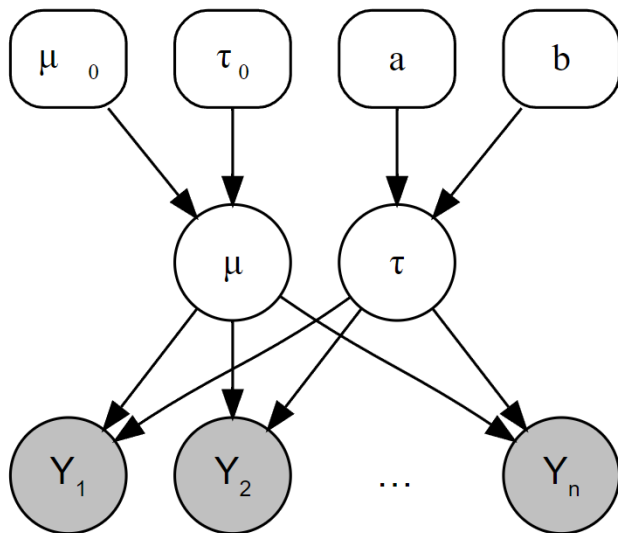
$$i = 1, \dots, n : y_i \sim \mathcal{N}(\mu, 1/\tau)$$

- Conditional posterior distributions for  $\mu$  and  $\tau$ :

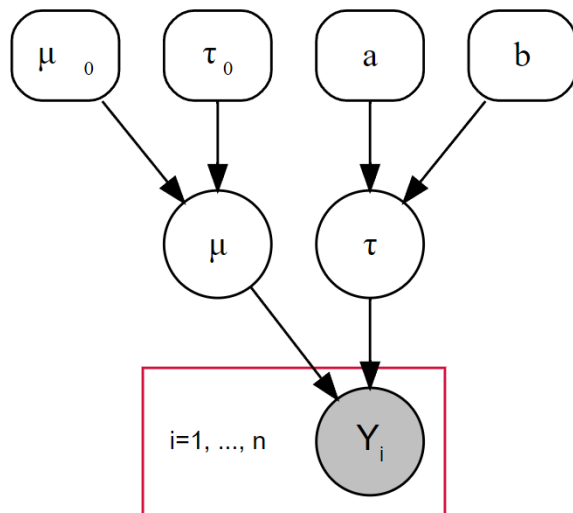
$$\mu|y, \tau \sim \mathcal{N}(\mu_n, 1/\tau_n)$$

$$\tau|y, \mu \sim \mathcal{G}(a_n, b_n)$$

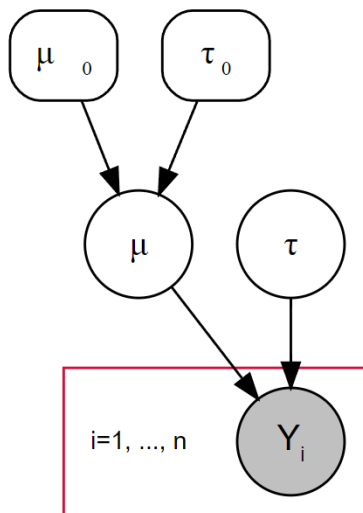
## Normal distribution - graph representation



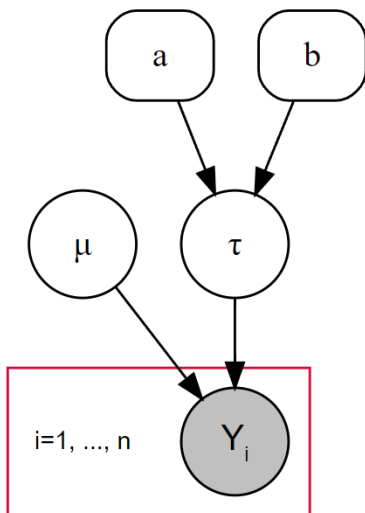
## Normal distribution - plate notation



## Normal distribution - conditional on $\tau$



## Normal distribution - conditional on $\mu$





## Comparing multiple means

- ▶ Now suppose we observe data samples from  $K$  multiple groups:

$$y_1 = y_{1,1}, \dots, y_{n_1,1}$$

$$\vdots$$

$$y_K = y_{1,K}, \dots, y_{n_K,K}.$$

- ▶ We assume that each sample is normally distributed:

$$y_{i,k} \sim \mathcal{N}(\theta_k, 1/\tau_w),$$

for each observation  $i$  in each group  $k$ .

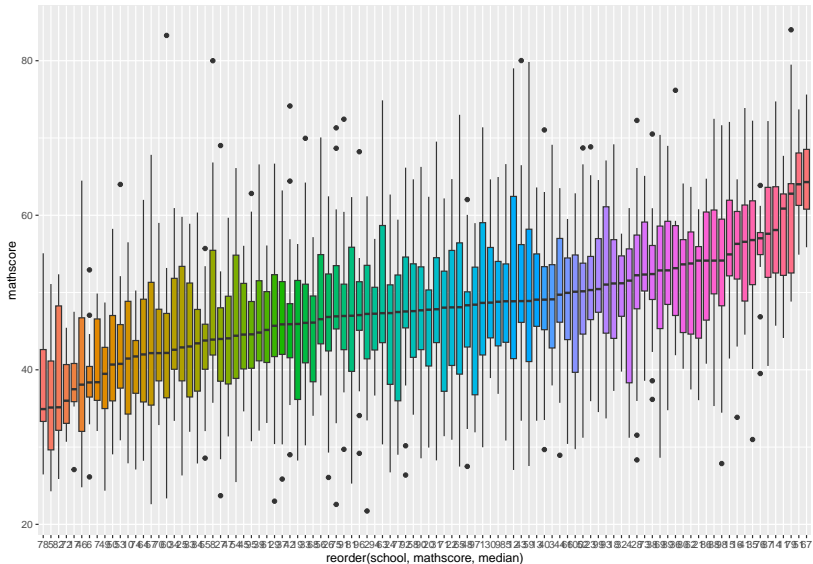
- ▶ We assume a common precision (variance) parameter  $\tau$  across all groups.

## Example - schools data

##	school	mathscore
## 1157	58	50.94
## 889	46	35.77
## 447	24	46.94
## 1649	86	64.73
## 962	50	41.55
## 1231	61	55.28

- ▶ We assess exam scores of  $n = 1993$  students from 100 different schools
- ▶ Same context as earlier example but many more schools

# Comparing multiple means - data



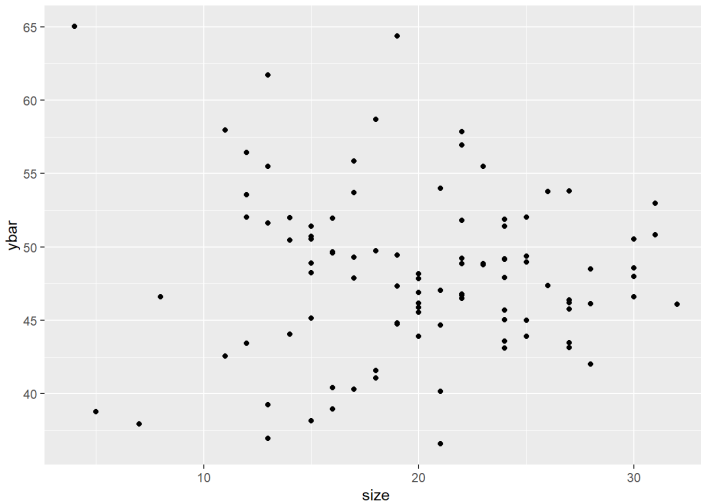
## Comparing multiple groups

- ▶ We want to compare the performance of the schools, i.e., at the level of population means

$$\theta_1, \dots, \theta_K.$$

- ▶ We could of course directly compare the sample means  $\bar{y}_1, \dots, \bar{y}_K$ .
- ▶ But we know there is variability in the way that these summary statistics have been collected:
  - ▶ A student's performance in an exam will vary from day to day;
  - ▶ Some students are weaker/stronger than others, so different samples may obtain a somewhat different scores.
- ▶ Our comparison needs to take this into account and distinguish between systematic and random variation.

# Schools data



# Modelling approach

- ▶ A classic approach to take here would be to perform a hypothesis test
  - ▶ Assume a null model where  $d = 0$  across all groups;
  - ▶ Take the variability of the data into account
  - ▶ If the scores are different enough that statistical significance is achieved  $\Rightarrow$  conclude that a difference exists.
  - ▶ Otherwise, if e.g.,  $p > 0.05$ , fail to reject  $H_0$  and conclude that no real differences can be detected in the data.

# Modelling approach

- ▶ This approach is well known and called ANOVA (analysis of variance)
- ▶ But we can criticise this approach.
- ▶ Suppose that  $p = 0.055$ . Is this result really so different from  $p = 0.045$ ?
- ▶ The initial ANOVA test also only gives an overall result. To compare all groups requires adjustments to avoid catastrophic Type 1 error.
- ▶ It is more flexible to explicitly model these means in a Bayesian framework.

## Comparing multiple means - model

- ▶ Let  $\mu$  be a **population mean**;
- ▶ Let  $\tau_b$  be precision **between** groups;
- ▶ Let  $\theta_k$  denote **mean** of group  $k$ ;
- ▶ Let  $\tau_w$  be precision **within** groups; this is common to all groups.



# Comparing multiple means - model

- ▶ Model the data as follows:

$$\theta_k = \mu + d_k, k = 1, \dots, K$$

- ▶  $d_k \sim \mathcal{N}(0, 1/\tau_b)$  for all  $k$ .

- ▶ Then:

$$y_{i,k} = \mu + d_k + \epsilon_{i,k}, i = 1, \dots, n_k$$

- ▶  $\epsilon_{i,k} \sim \mathcal{N}(0, 1/\tau_w)$  for all  $i, k$ , i.e., noise.

## Comparing multiple means - model

- ▶ Another way to write this is that we assume that:
- ▶ For  $k = 1, \dots, K$ :

$$\theta_k \sim \mathcal{N}(\mu, 1/\tau_b);$$

- ▶ For  $i = 1, \dots, n_k$ :

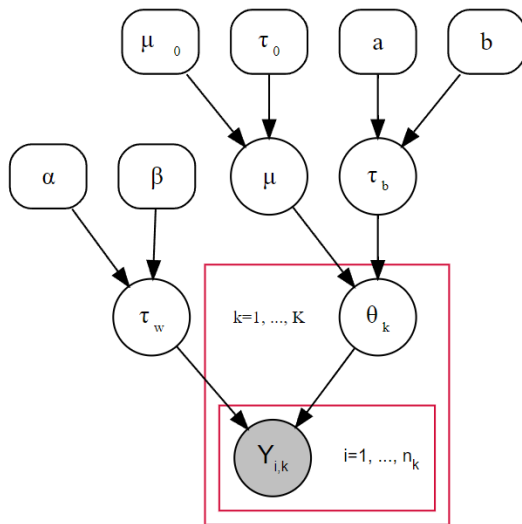
$$y_{i,k} \sim \mathcal{N}(\theta_k, 1/\tau_w),$$

- ▶ These representations are equivalent
- ▶ But I prefer the latter representation
  - ▶ makes dependency between variables more explicit
  - ▶ more parsimonious

## Choosing priors

- ▶ Let's choose the following priors for this model
  - ▶  $\mu \sim \mathcal{N}(\mu_0, 1/\tau_0)$ ;
  - ▶  $\tau_b \sim \mathcal{G}(a, b)$ ;
  - ▶ Let  $\tau_w \sim \mathcal{G}(\alpha, \beta)$ ;
- ▶ These choices should be unsurprising at this point
  - ▶ see notes on Normal model for details
- ▶ We don't need to specify a prior for group means  $\theta_1, \dots, \theta_K$ .
- ▶  $\mu$  and  $\tau_b$  act as the “prior” in this case – see graph.

## Comparing multiple means – graph



# Inference

- ▶ In practice, we can fit this model using Stan, or other specialist languages.
- ▶ However it is still instructive to review how to perform inference in this case
- ▶ So we will outline the key steps

## Inference - build the posterior

- ▶ Assume for now that  $\mu, \tau_b, \tau_w$  are known.
- ▶ We have a joint distribution for  $y$  and  $\theta$  of the form:

$$\begin{aligned} p(y, \theta | \mu, \tau_b, \tau_w) &= p(y | \theta, \tau_w) p(\theta | \mu, \tau_b) \\ &= \prod_{k=1}^K \prod_{i=1}^{n_k} p(y_{i,k} | \theta_k, \tau_w) p(\theta_k | \mu, \tau_b). \end{aligned}$$

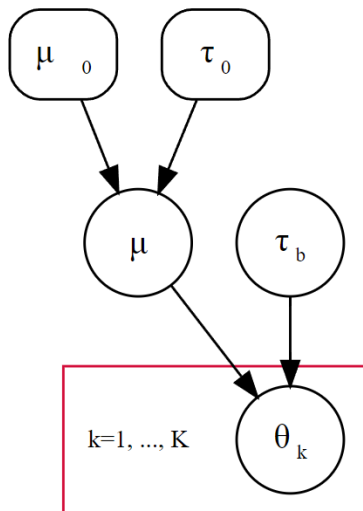
## Building the posterior

- We can then construct the full posterior by adding in prior terms for  $\mu, \tau_b$ , and  $\tau_w$  :

$$\begin{aligned} & p(\theta, \mu, \tau_b, \tau_w | y, \mu_0, \tau_0, a, b, \alpha, \beta) \\ & \propto p(y, \theta | \mu, \tau_b, \tau_w) p(\mu | \mu_0, \tau_0) p(\tau_b | a, b) p(\tau_w | \alpha, \beta) \\ & \propto \left\{ \prod_{k=1}^K \prod_{i=1}^{n_k} p(y_{i,k} | \theta_k, \tau_w) \right\} \times \prod_{k=1}^K p(\theta_k | \mu, \tau_b) \\ & \quad \times p(\mu | \mu_0, \tau_0) p(\tau_b | a, b) p(\tau_w | \alpha, \beta). \end{aligned}$$

- Note the relationship between the structure of this posterior and the graph on the previous slide.

## Conditional distributions – $\mu$





## Conditional distributions – $\mu$

- ▶ The conditional distribution  $p(\mu|\theta, \tau_b)$  for the population mean  $\mu$  is:

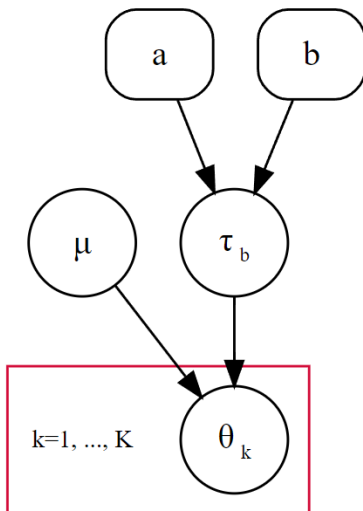
$$p(\mu|\theta, \tau_b, \mu_0, \tau_0) \propto \prod_{k=1}^K p(\theta_k|\mu, \tau_b)p(\mu|\mu_0, \tau_0),$$

- ▶ where  $p(\theta_k|\mu, \tau_b)$  is a normal distribution, for all  $k$ , and  $p(\mu|\mu_0, \tau_0)$  is also normal.
- ▶ Hence  $\mu|\theta, \tau_b, \mu_0, \tau_0 \sim \mathcal{N}(\mu_K, 1/\tau_K)$ , with

$$\mu_K = \frac{K\tau_b\bar{\theta} + \tau_0\mu_0}{K\tau_b + \tau_0};$$

$$\tau_K = K\tau_b + \tau_0.$$

## Conditional distributions – $\tau_b$



## Conditional distributions – $\tau_b$

- ▶ The conditional distribution for the between group precision  $\tau_b$  is:

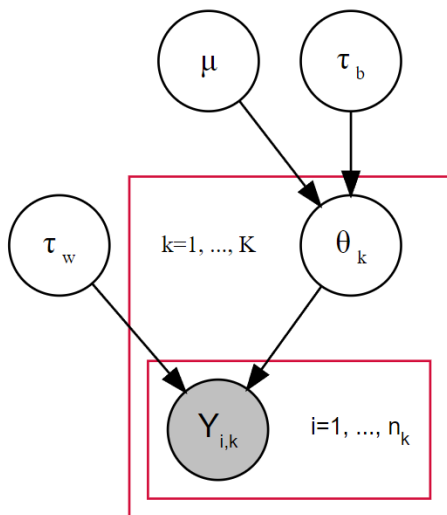
$$p(\tau_b|\theta, \mu, a, b) \propto \prod_{k=1}^K p(\theta_k|\mu, \tau_b)p(\tau_b|a, b),$$

- ▶ where  $p(\theta_k|\mu, \tau_b)$  is a normal distribution, for all  $k$ , and  $p(\tau_b|a, b)$  is Gamma.
- ▶ Then  $\tau_b|\theta, \mu, a, b \sim \mathcal{G}(a_K, b_K)$ , with

$$a_K = K/2 + a;$$

$$b_K = \frac{1}{2} \left\{ \sum_{k=1}^K (\theta_k - \mu)^2 \right\} + b.$$

## Conditional distributions – $\theta_k$



## Conditional distributions – $\theta_k$

- ▶ The conditional distribution for each  $\theta_k$  is:

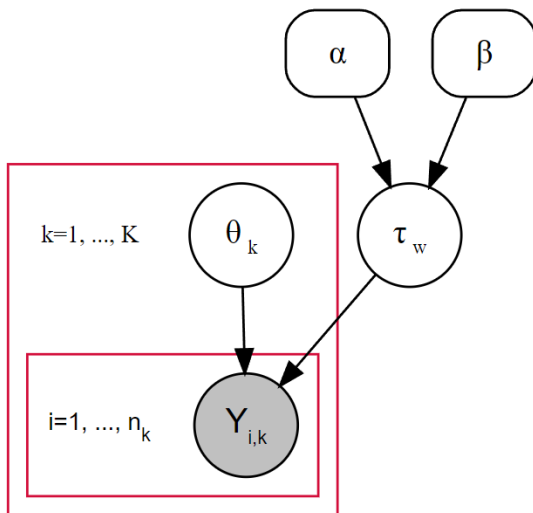
$$p(\theta_k|y, \mu, \tau_b, \tau_w) \propto \prod_{i=1}^{n_k} p(y_{ik}|\theta_k\tau_w)p(\theta_k|\mu, \tau_b)$$

- ▶ where  $p(y_{ik}|\theta_k\tau_w)$  and  $p(\theta_k|\mu, \tau_b)$  are both normally distributed, for all  $i$  and  $k$ .
- ▶ Then  $\theta_k|y, \mu, \tau_b, \tau_w \sim \mathcal{N}(\lambda_{n_k}, 1/\gamma_{n_k})$ , with

$$\lambda_{n_k} = \frac{n_k\tau_w\bar{y}_k + \tau_b\mu}{n_k\tau_w + \tau_b};$$

$$\gamma_{n_k} = n_k\tau_w + \tau_b.$$

## Conditional distributions – $\tau_w$



## Conditional distributions – $\tau_w$

- ▶ The conditional distribution for within group precision  $\tau_w$  is:

$$p(\tau_w|\theta, \alpha, \beta) \propto \prod_{k=1}^K \prod_{i=1}^{n_k} p(y_{i,k}|\theta_k, \tau_w) p(\tau_w|\alpha, \beta),$$

- ▶  $p(y_{i,k}|\theta_k, \tau_w)$  is normal for all  $i, k$ , and  $p(\tau_w|\alpha, \beta)$  is Gamma
  - ▶  $\Rightarrow \tau_w|\theta, \alpha, \beta \sim \mathcal{G}(\alpha_n, \beta_n)$ , with

$$\alpha_n = \sum_{k=1}^K \frac{n_k}{2} + \alpha;$$

$$b_K = 1/2 \left\{ \sum_{k=1}^K \sum_{i=1}^{n_k} (y_{i,k} - \theta_k)^2 \right\} + b.$$

## Inference for model

- ▶ Because the conditional distributions are available in all cases, we can update the model using a Gibbs sampler.
- ▶ Even if we use a different approach, the conditional distributions of the parameters give us an insight into what we learn from the data.
- ▶ The hierarchical structure of the model means that the group means  $\theta_1, \dots, \theta_K$  behave like the “data” for  $\mu$  and  $\tau_b$ ;
- ▶ Conversely,  $\mu$  and  $\tau_b$  behave like hyperparameters for each  $\theta_k$ , even though they are estimated from the data.



# Interpreting the model

- ▶ Note that, when estimating  $\theta_k$ , the conditional mean  $\lambda_{n_k}$  has the form

$$\lambda_{n_k} = \frac{n_k \tau_w \bar{y}_k + \tau_b \mu}{n_k \tau_w + \tau_b}$$

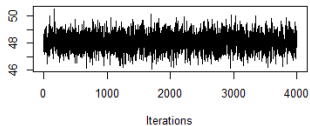
- ▶ So the estimate for  $\theta_k$  is effectively a weighted combination of elements
  - ▶  $\bar{y}_k$  estimated directly from the data and relating explicitly to group  $k$ ;
  - ▶  $\mu$  indirectly estimated from data across all groups.
- ▶ We are **pooling** information across all samples when assessing individual groups.

# Interpreting the model

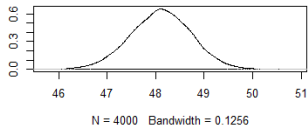
- ▶ In this context,  $\mu$  can be called a **shrinkage** factor.
- ▶ In a full hierarchical model, we “borrow” information between all groups when estimating the parameters of individual groups.
- ▶ This is useful when the sample size for some groups is small.

# Schools data - global parameters

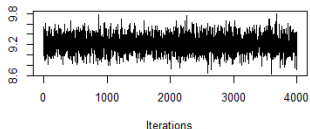
Trace of  $\mu$



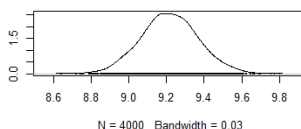
Density of  $\mu$



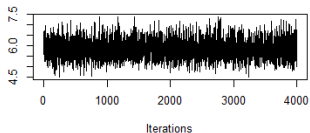
Trace of  $sd_w$



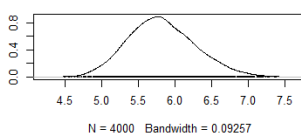
Density of  $sd_w$



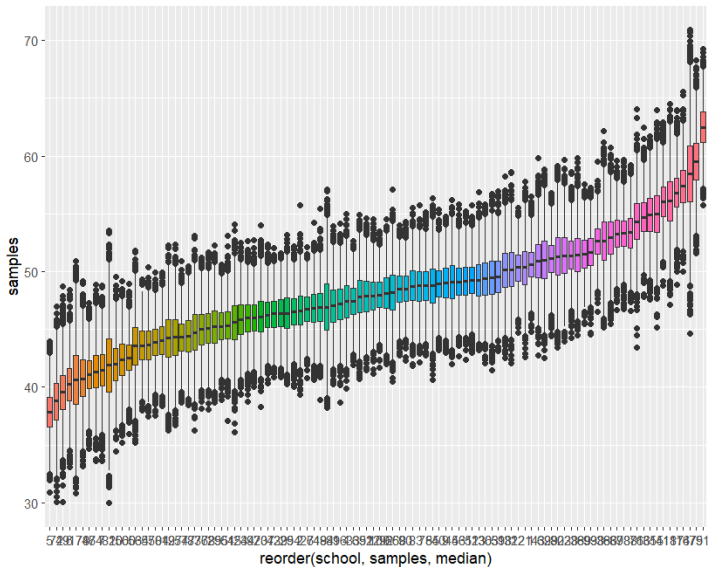
Trace of  $sd_b$



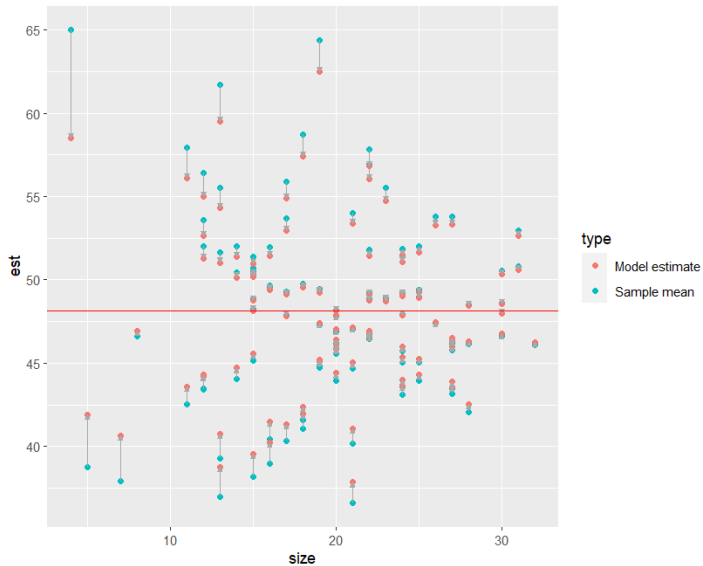
Density of  $sd_b$



# Schools data - schools means $\theta$



# Schools data - $\theta$ vs $\bar{y}$



# Conclusion

- ▶ We have looked at a hierarchical model to compare data from  $K$  groups.
  - ▶ This is our first “proper” statistical model.
- ▶ We estimated the model using Gibbs samplers automatically using Stan:
  - ▶ The structure of our model meant that inference was straightforward.
  - ▶ This structure also lets us pool information across several groups
  - ▶ This is helpful when information is limited (sample size is small) for some groups
- ▶ You should be clear on what each parameter in the models represents, and how to interpret the model output.
  - ▶ When a model has a lot of parameters we sometimes need to be creative in how we interpret this output.

# Extensions

- ▶ It is possible to extend the hierarchy of the model further if the data set has richer structure, for example:
  - ▶ schools in regions
  - ▶ students in classes
- ▶ Inference in this case would be very similar to the models we have examined.
- ▶ In fact, in many cases Bayesian inference can be thought of as “mechanical”, and there is dedicated software to implement such models as simply as possible.