



# A journey through Bayesian hierarchical models: Analyzing income and education in Thailand

---

Dr. Irving Gómez Méndez

April 2024

This work was done in collaboration with Dr. Chainarong Amornbunchornvej, **NECTEC**



Hierarchical models with one cluster: income per region





## What is a Bayesian model?

### Input:

We model a random variable using the **likelihood**  $Y \sim p(Y|\theta)$ .

We model the uncertainty on  $\theta$  using a **prior distribution**  $p(\theta)$ .

### Output:

We update our uncertainty on  $\theta$  through Bayes rule, getting the **posterior distribution**  $p(\theta|\mathbf{Y})$ .

We capture the total uncertainty on  $Y$  using the **posterior predictive distribution**  $p(Y|\mathbf{Y})$ .



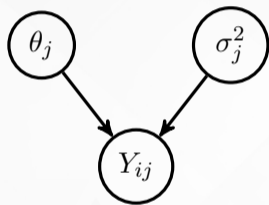
There are 76 provinces grouped in 6 different regions.

Let be  $Y_{ij}$  the income in province  $i$  belonging to region  $j$ .

We assume that  $Y_{ij}|\theta_j, \sigma_j^2 \sim \mathcal{N}(\theta_j, \sigma_j^2)$  (**Likelihood**).



## No pooling model



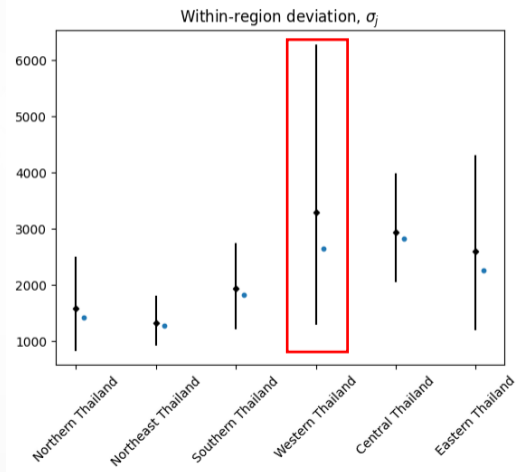
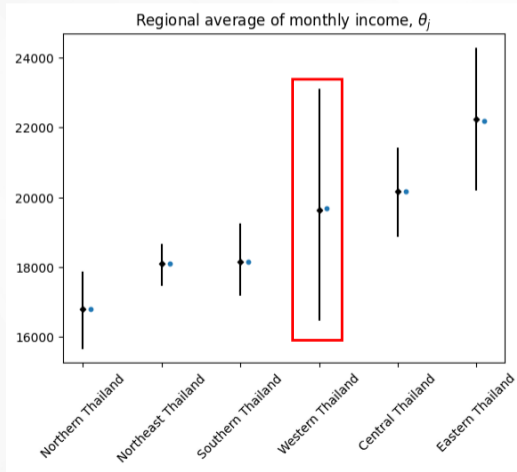
Each region has an independent model.

**Prior distribution**

$$p(\boldsymbol{\theta}, \boldsymbol{\sigma}^2) \propto \prod_{j=1}^J \frac{1}{\sigma_j^2} \mathbb{1}_{\mathbb{R}}(\theta_j) \mathbb{1}_{(0, \infty)}(\sigma_j^2)$$

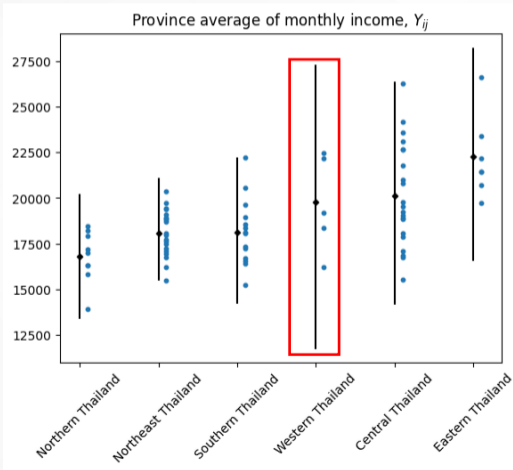


## Posterior distributions





## Posterior predictive distribution



Large intervals, because regions don't share information.

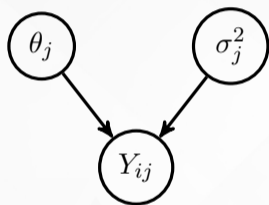
Intervals overlap for most of the regions.

All the parameters might be estimating the same quantity.

It is highly unlikely that the regions are independent between them.



## Complete pooling model



Same model for all the regions.

$\theta_j = \theta, \sigma_j = \sigma$ , for all  $j$ .

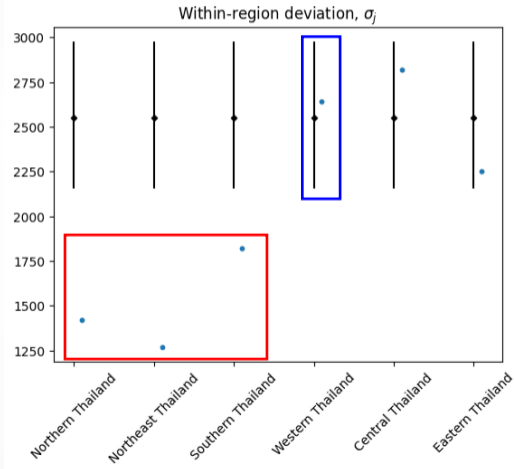
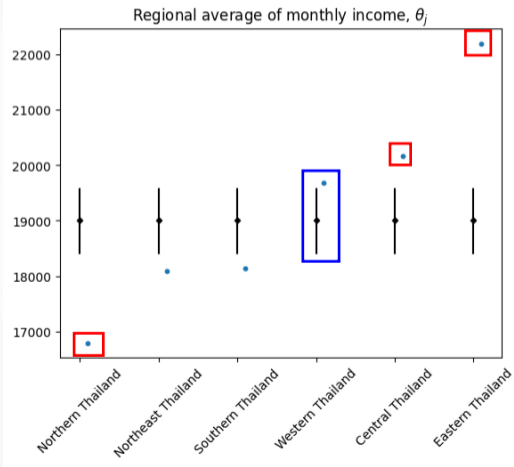
**Prior distribution**

$$p(\theta, \sigma^2) \propto \frac{1}{\sigma^2} \mathbb{1}_{\mathbb{R}}(\theta) \mathbb{1}_{(0, \infty)}(\sigma^2)$$



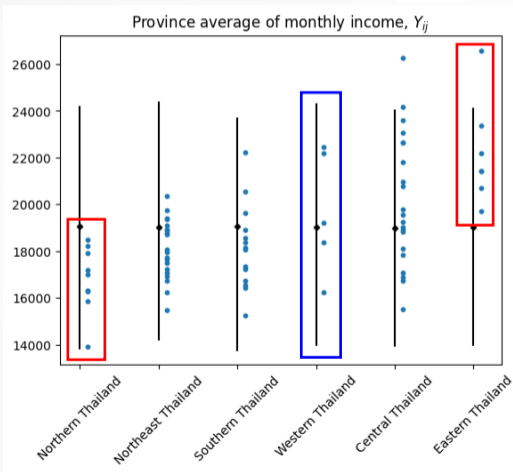


## Posterior distributions





## Posterior predictive distribution



Narrow intervals.

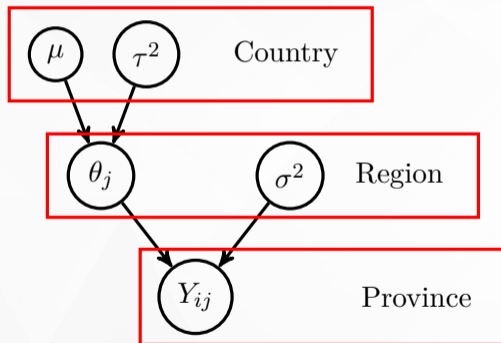
Common mean  $\theta$  can barely explain the mean of a few regions.

Overestimated  $\sigma$ .

Income in Northern Thailand is overestimated. Income in Eastern Thailand is underestimated.



## Hierarchical model with common within-cluster variance



### Likelihood

$$Y_{ij} | \theta_j, \sigma^2 \sim \mathcal{N}(\theta_j, \sigma^2)$$

### Prior distributions

$$\theta_j | \mu, \tau^2 \sim \mathcal{N}(\mu, \tau^2)$$

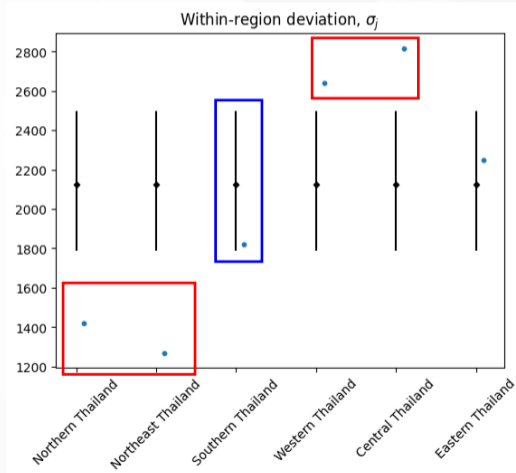
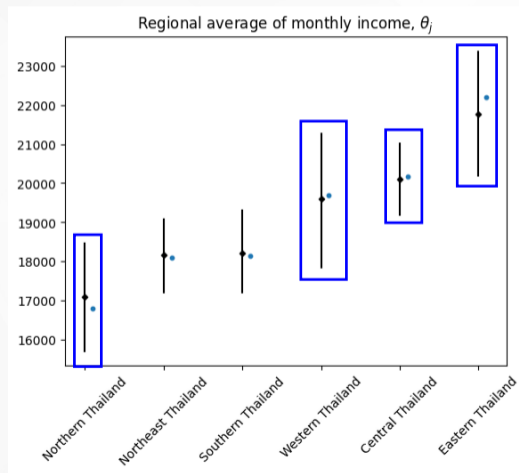
$$p(\mu) \propto \mathbb{1}_{\mathbb{R}}(\mu)$$

$$p(\tau^2) \propto \mathbb{1}_{(0, \infty)}(\tau^2)$$

$$p(\sigma^2) \propto \frac{1}{\sigma^2} \mathbb{1}_{(0, \infty)}(\sigma^2)$$

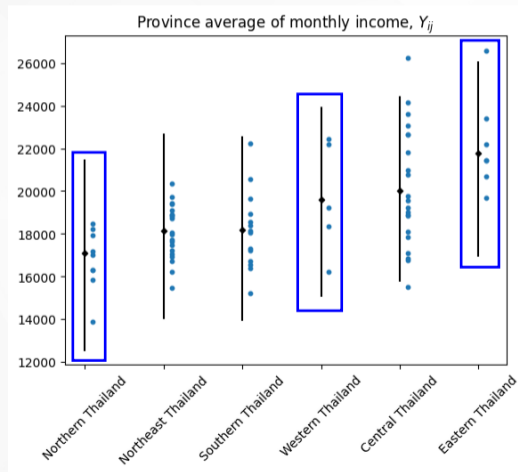


## Posterior distributions





## Posterior predictive distribution



Less uncertainty than independent models.

Better performance than complete pooling model.

The common variance is not overestimated.

A common variance cannot explain the observed variability.



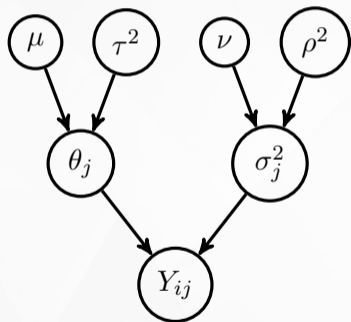
## Drawbacks of Bayesian models:

**We must ensure the existence of the posterior distribution.** This requires that **we analyze** the mathematical properties of the model.

**The model is not going to do it for us! That's the hatch!**



## Hierarchical model varying within-cluster variance



$$Y_{ij} | \theta_j, \sigma_j^2 \sim \mathcal{N}(\theta_j, \sigma_j^2)$$

$$\theta_j | \mu, \tau^2 \sim \mathcal{N}(\mu, \tau^2)$$

$$p(\mu) \propto \mathbb{1}_{\mathbb{R}}(\mu)$$

$$p(\tau^2) \propto \mathbb{1}_{(0, \infty)}(\tau^2)$$

$$\sigma_j^2 | \nu, \rho^2 \sim \text{Inverse-}\chi^2(\nu, \rho^2)$$

$$p(\rho^2) \propto \frac{1}{\rho^2} \mathbb{1}_{(0, \infty)}(\rho^2)$$

$$\nu \sim ?$$



It remains challenging to propose priors for the degrees of freedom of a distribution.

We consider three approaches:

- ▶ Set  $\nu$  equal to some estimated value  $\hat{\nu}$ . For example, using the method of moments

$$\hat{\nu} = \frac{2(E_{s^2})^2}{V_{s^2}} + 4.$$

- ▶ Use a vague prior, for example  $p(\nu) \propto \nu^{-h} \mathbb{1}_{(0,\infty)}(\nu)$ , for some  $h \in [0, \infty)$ .
- ▶ Use a regularizing prior, for example  $\nu \sim \text{Exponential}(1/\hat{\nu})$ .





## Cons:

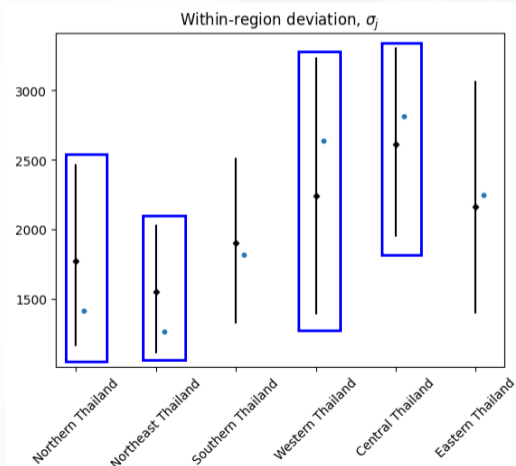
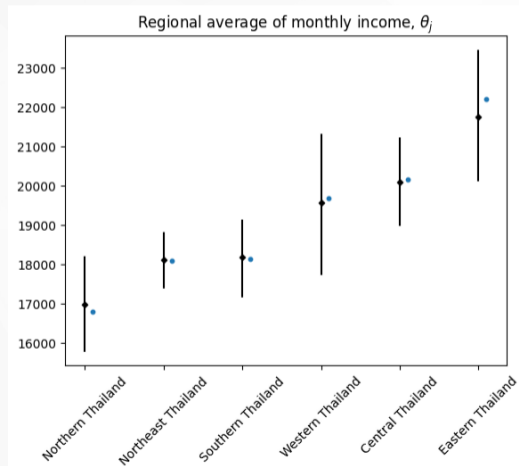
- ▶ Set  $\nu = \hat{\nu}$ : Eliminates the uncertainty on  $\nu$ .
- ▶ Improper distribution: Does not give guarantee for the existence of the posterior distribution.
- ▶ Regularizing prior: Which distribution to use?

## Pros:

- ▶ Set  $\nu = \hat{\nu}$ : We don't have to worry about the posterior.
- ▶ Improper distribution: There's nothing to estimate.
- ▶ Regularizing prior: Maintain the uncertainty and guarantee the existence of the posterior distribution.

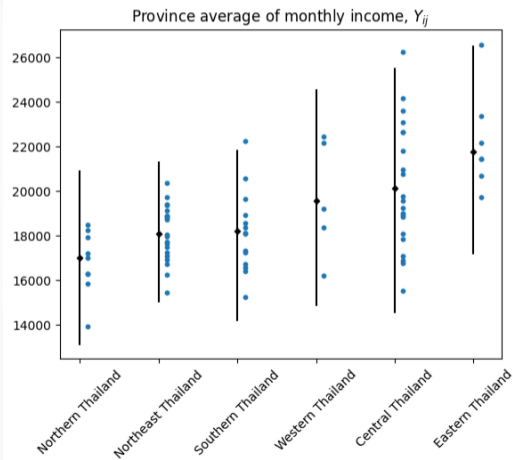


## Posterior distributions

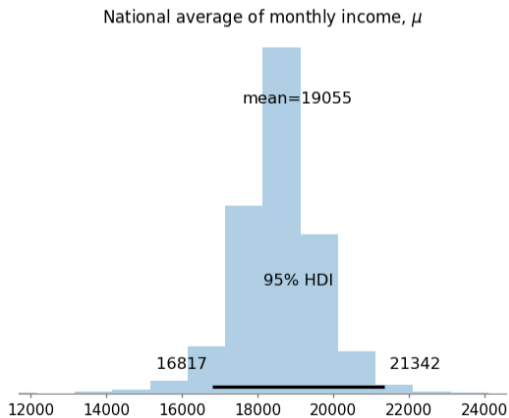




## Posterior predictive distribution

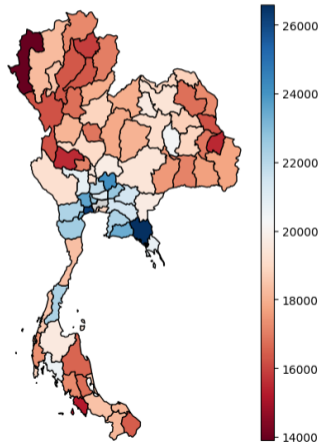


## Posterior distribution

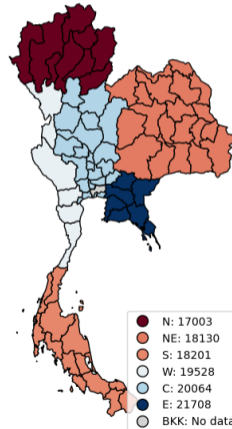


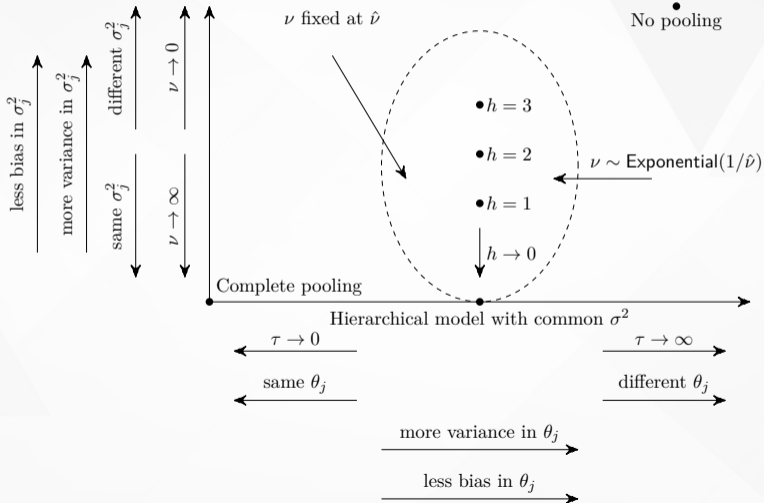


Province average of monthly income



Regional average of monthly income







Hierarchical model with two non-nested clusters: income per region and education level



We define the next rule:

<b>Education</b>	<b>Years of education</b>	<b>Education level</b>
Uneducated	0	Low
Kindergarten	0	
Pre-elementary school	3	
Elementary school	6	
Junior high school	9	Mid
Senior high school	12	
Vocational degree	14	
Bachelor degree	16	High
Post-graduate	19	



$$Y_{ijk} | \theta_j, \lambda_k, \sigma_{jk}^2 \sim \mathcal{N}(\theta_j + \lambda_k, \sigma_{jk}^2)$$

$$\theta_j | \mu, \tau^2 \sim \mathcal{N}(\mu, \tau^2)$$

$$\mu \sim \mathcal{N}(\hat{\mu}, \hat{\sigma}_\mu^2)$$

$$\tau^2 \sim \text{Exponential}(1/\hat{\tau}^2)$$

$$\lambda_k | \xi^2 \sim \mathcal{N}(0, \xi^2)$$

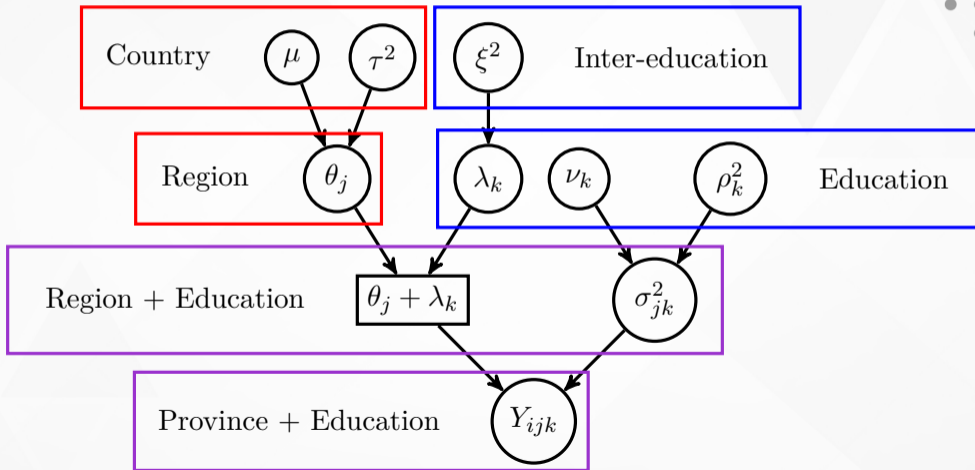
$$\xi^2 \sim \text{Exponential}(1/\hat{\xi}^2)$$

$$\sigma_{jk}^2 | \nu_k, \rho_k^2 \sim \text{Inverse-}\chi^2(\nu_k, \rho_k^2)$$

$$\nu_k^2 \sim \text{Exponential}(1/\hat{\nu}_k^2)$$

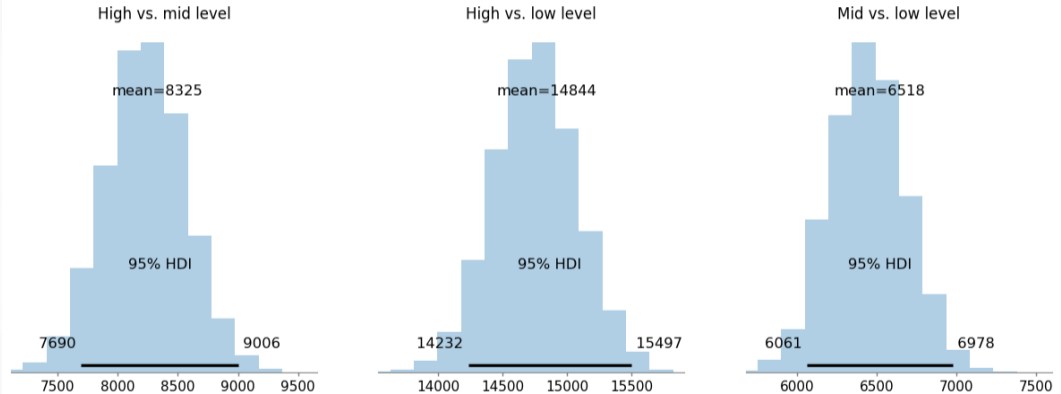
$$p(\rho_k^2) \propto \frac{1}{\rho_k^2} \mathbb{1}_{(0, \infty)}(\rho_k^2)$$





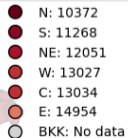
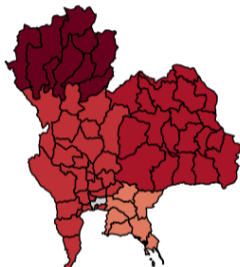


Difference in monthly income between different education levels,  $\lambda_k - \lambda_{k'}$

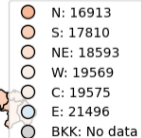
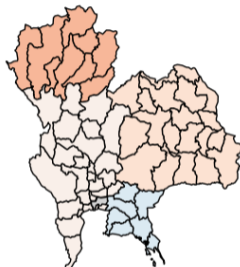




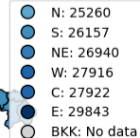
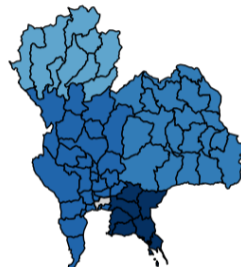
Education Level: Low



Education Level: Mid



Education Level: High





# Bayesian hierarchical regression: income considering years of formal education



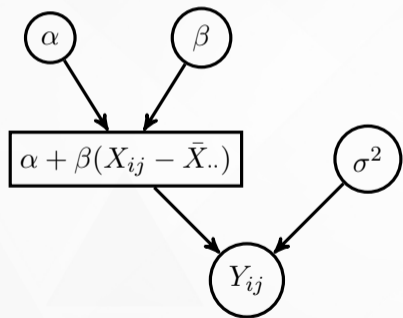
We define the next rule:

<b>Education</b>	<b>Years of education</b>	<b>Education level</b>
Uneducated	0	Low
Kindergarten	0	
Pre-elementary school	3	
Elementary school	6	
Junior high school	9	Mid
Senior high school	12	
Vocational degree	14	
Bachelor degree	16	High
Post-graduate	19	



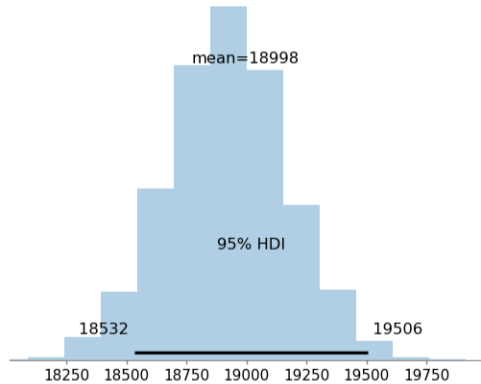
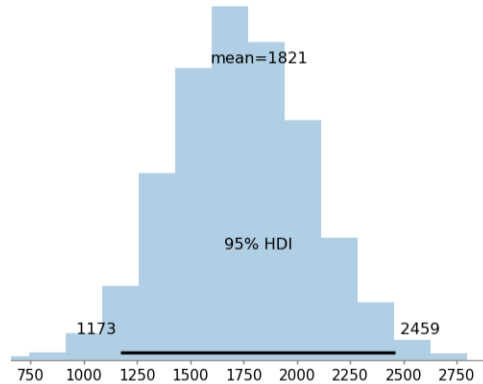
## National model

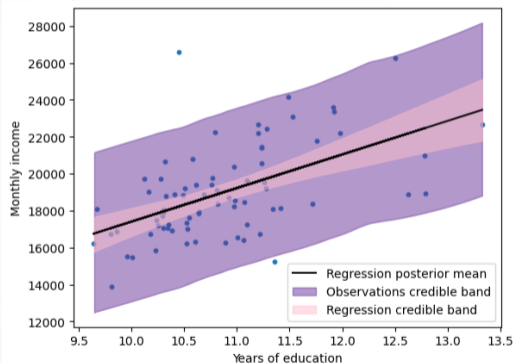
Let be  $X_{ij}$  the average years of formal education in the province  $i$  belonging to region  $j$ .



$$Y_{ij} | \alpha, \beta, \sigma^2 \sim \mathcal{N}(\alpha + \beta(X_{ij} - \bar{X}_{..}), \sigma^2)$$

$$p(\alpha, \beta, \sigma^2) \propto \frac{1}{\sigma^2} \mathbb{1}_{\mathbb{R}}(\alpha) \mathbb{1}_{\mathbb{R}}(\beta) \mathbb{1}_{(0, \infty)}(\sigma^2)$$

National average of monthly income,  $\alpha$ National ratio of income per year-of-education,  $\beta$ 



Suitable to explain national behavior.

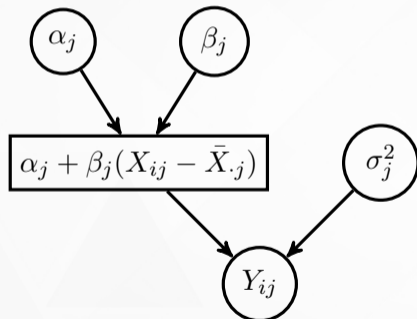
We cannot say anything at a regional level.





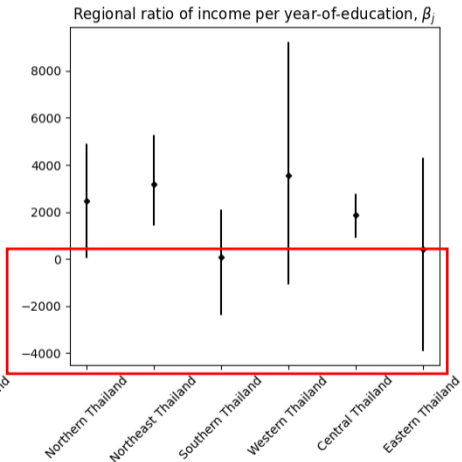
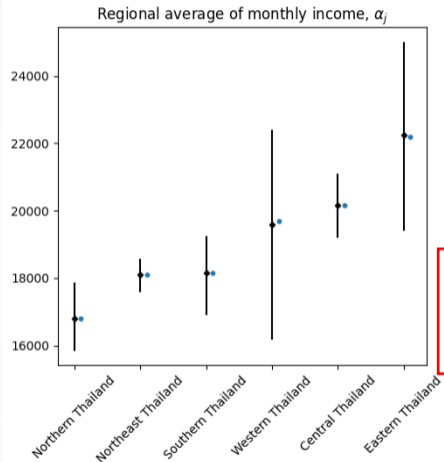
## Separate models

Let be  $X_{ij}$  the average years of formal education in the province  $i$  belonging to region  $j$ .



$$Y_{ij} | \alpha, \beta, \sigma^2 \sim \mathcal{N}(\alpha_j + \beta_j(X_{ij} - \bar{X}_{.j}), \sigma^2)$$

$$p(\alpha_j, \beta_j, \sigma_j^2) \propto \frac{1}{\sigma_j^2} \mathbb{1}_{\mathbb{R}}(\alpha_j) \mathbb{1}_{\mathbb{R}}(\beta_j) \mathbb{1}_{(0, \infty)}(\sigma_j^2)$$





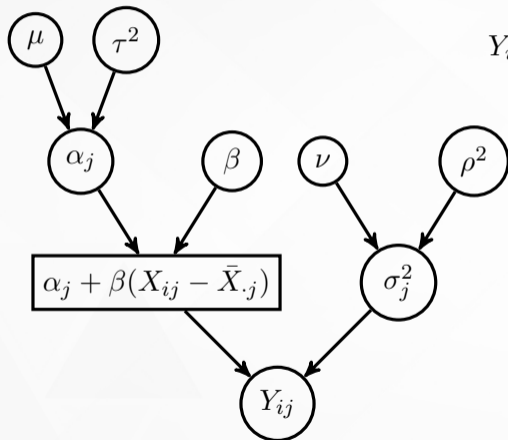
Large credible intervals.

The intervals include negative values for  $\beta_j$ , which seems implausible.

Pretending that each region is independent for the others seems unrealistic.



## Bayesian hierarchical regression varying intercepts



$$Y_{ij} | \alpha_j, \beta_j, \sigma_j^2 \sim \mathcal{N}(\alpha_j + \beta_j(X_{ij} - \bar{X}_{.j}), \sigma_j^2)$$

$$\alpha_j | \mu, \tau^2 \sim \mathcal{N}(\mu, \tau^2)$$

$$\mu \sim \mathcal{N}(\hat{\mu}, \hat{\sigma}_\mu^2)$$

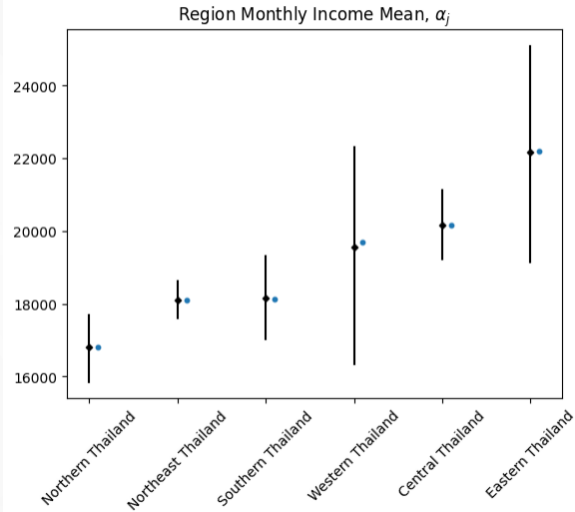
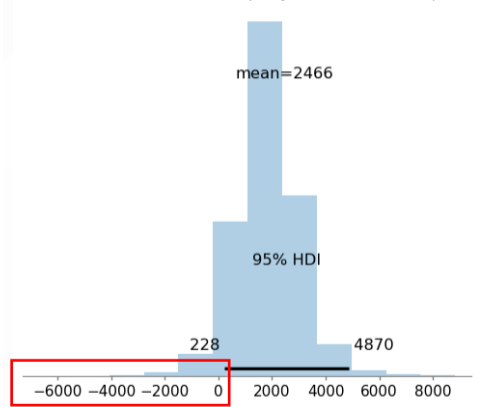
$$\tau^2 \sim \text{Exponential}(1/\hat{\tau}^2)$$

$$p(\beta) \propto \mathbb{1}_{\mathbb{R}}(\beta)$$

$$\sigma_j^2 | \nu, \rho^2 \sim \text{Inverse-}\chi^2(\nu, \rho^2)$$

$$\nu^2 \sim \text{Exponential}(1/\hat{\nu}^2)$$

$$p(\rho^2) \propto \frac{1}{\rho^2} \mathbb{1}_{(0, \infty)}(\rho^2)$$

National ratio of income per year-of-education,  $\beta$ 



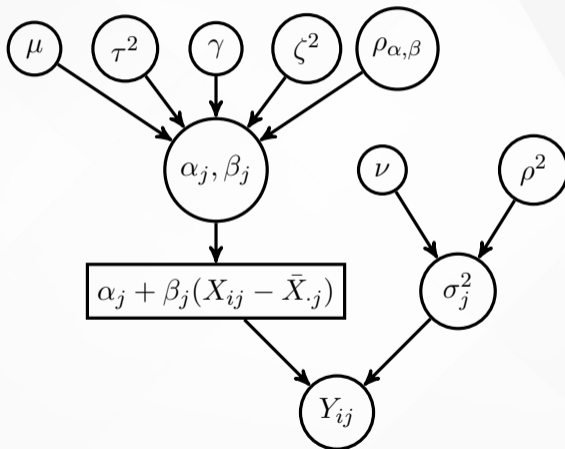
The distribution of  $\beta$  includes negative values.

Different slopes could be preferred.





## Bayesian hierarchical regression varying intercepts and slopes





$\alpha_j$  and  $\beta_j$  will follow a multivariate normal distribution with mean  $(\mu, \gamma)$  and a matrix of variances and covariances

$$S = \begin{pmatrix} \tau^2 & \tau\zeta\rho_{\alpha,\beta} \\ \tau\zeta\rho_{\alpha,\beta} & \zeta^2 \end{pmatrix},$$

which can be written as

$$S = \begin{pmatrix} \tau & 0 \\ 0 & \zeta \end{pmatrix} R \begin{pmatrix} \tau & 0 \\ 0 & \zeta \end{pmatrix},$$

where

$$R = \begin{pmatrix} 1 & \rho_{\alpha,\beta} \\ \rho_{\alpha,\beta} & 1 \end{pmatrix}$$

is the correlation matrix.





$$Y_{ij} | \alpha_j, \beta_j, \sigma_j \sim \text{Laplace}(\alpha_j + \beta_j(X_{ij} - \bar{X}_{\cdot j}), \sigma_j)$$

$$\alpha_j, \beta_j | \mu, \tau^2, \gamma, \zeta^2, \rho_{\alpha, \beta} \sim \text{MVN} \left( \begin{bmatrix} \mu \\ \gamma \end{bmatrix}, S \right)$$

$$\mu \sim \mathcal{N}(\hat{\mu}, \hat{\sigma}_\mu^2)$$

$$\tau^2 \sim \text{Exponential}(1/\hat{\tau}^2)$$

$$\gamma \sim \mathcal{N}(\hat{\gamma}, \hat{\sigma}_\gamma^2)$$

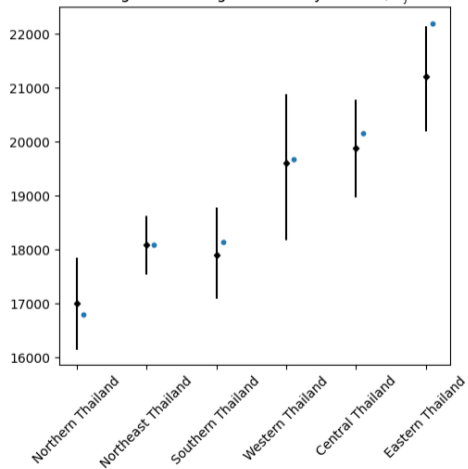
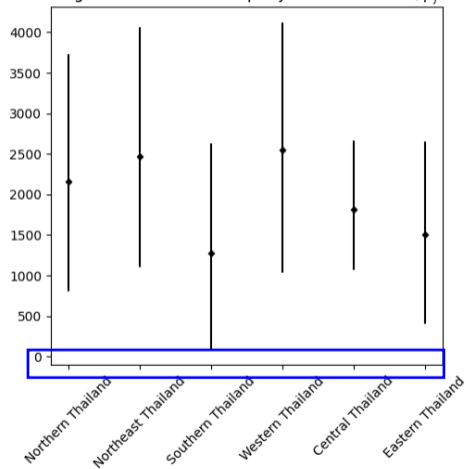
$$\zeta^2 \sim \text{Exponential}(1/\hat{\zeta}^2)$$

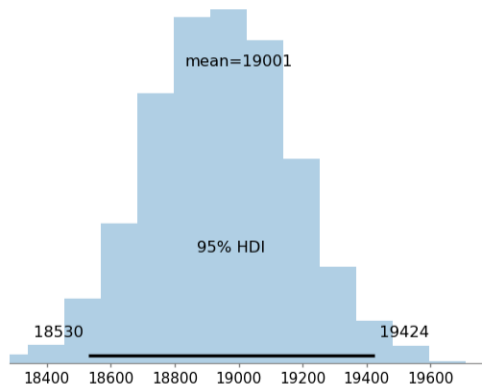
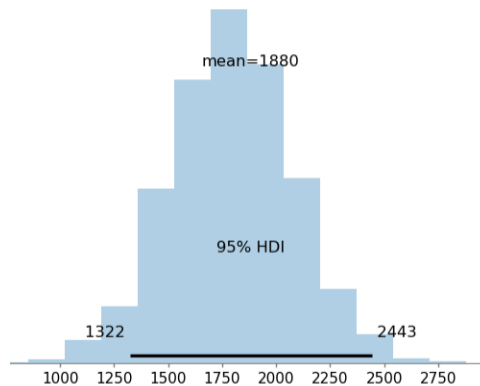
$$R \sim \text{LKJ}(2)$$

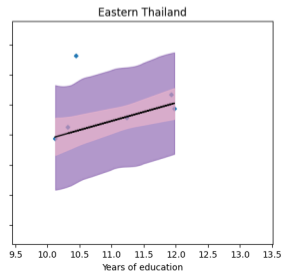
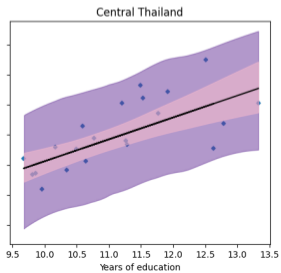
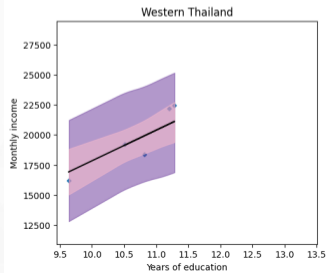
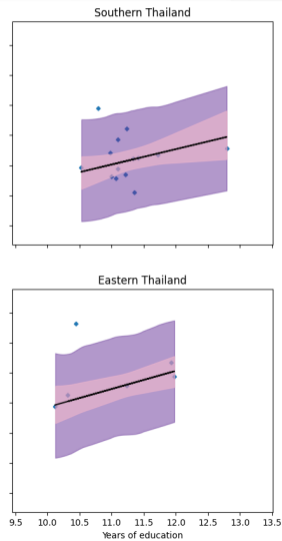
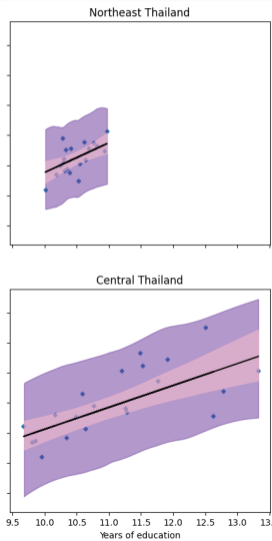
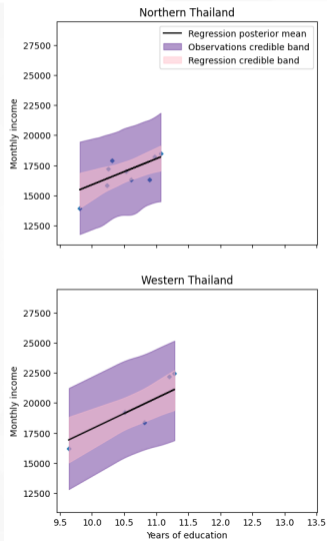
$$\sigma_j^2 | \nu, \rho^2 \sim \text{Inverse-}\chi^2(\nu, \rho^2)$$

$$\nu^2 \sim \text{Exponential}(1/\hat{\nu}^2)$$

$$p(\rho^2) \propto \frac{1}{\rho^2} \mathbb{1}_{(0, \infty)}(\rho^2)$$

Regional average of monthly income,  $\alpha_j$ Regional ratio of income per year-of-education,  $\beta_j$ 

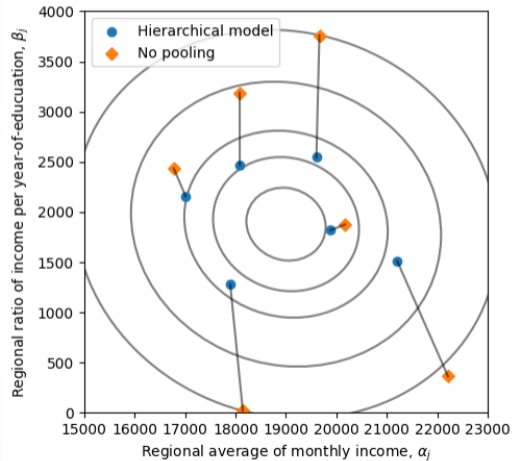
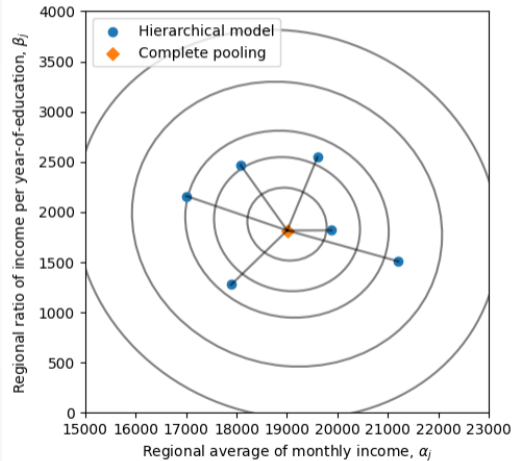
National average of monthly income,  $\mu$ National ratio of income per year-of-education,  $\gamma$ 





We can explain the data at a national and regional level.

There are no negative values for  $\beta_j$  nor  $\gamma$ .





*If you feel at all confused, it is only because you are paying attention.*