# GEOG0125
## ADVANCED TOPICS IN SOCIAL AND GEOGRAPHIC DATA SCIENCE
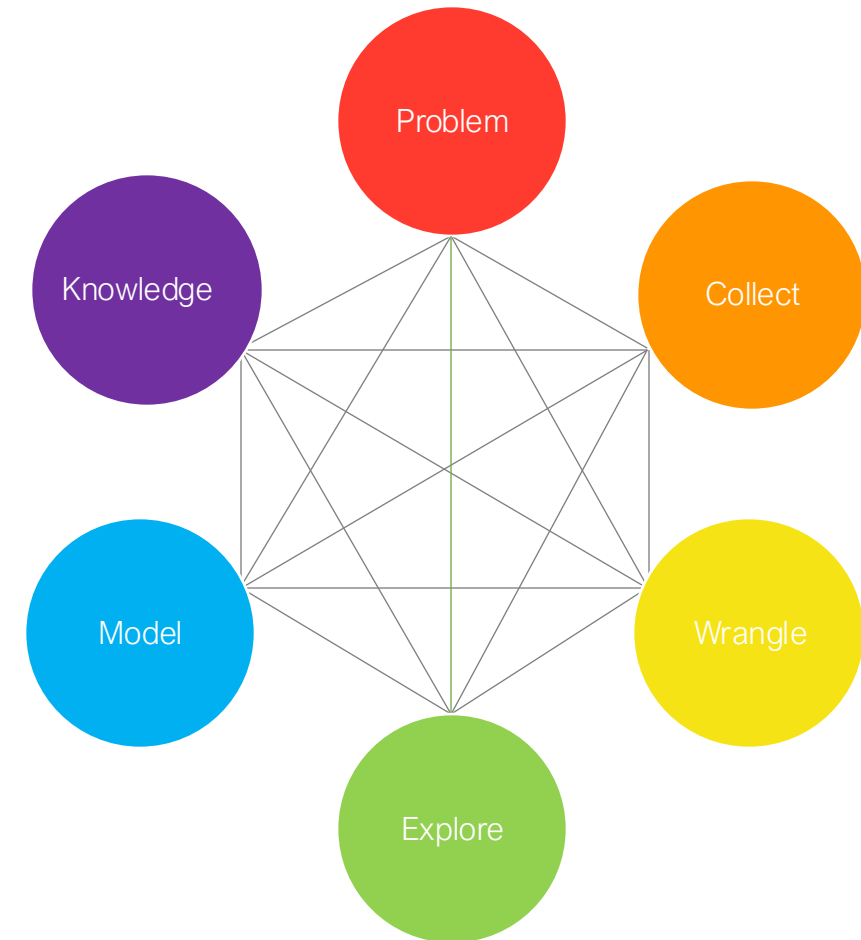
# BAYESIAN HIERARCHICAL REGRESSION MODELS

Dr Anwar Musah (a.musah@ucl.ac.uk)
Lecturer in Social and Geographic Data Science
UCL Geography

# Contents

- **What are Hierarchical Regression Models?**
  - Usage for exploring relationships with complex data structures (hierarchies, nesting, repeated measures, temporality etc.,)
  - Importance and why we use them to account for a specific data artefact
  - Wide applications of this specialized technique to other scientific domains

- **Components of a hierarchical model**
  - Statistical formulation
  - Random intercepts and Random slopes
  - Model types: Intercept-only, Random slope-only or both

- **Model Specification from a Bayesian Framework**

- **Examples and interpretation**

# What are Hierarchical Regression Models?

# Recall in Week 2 and 3, we have extensively covered these various types of regression models

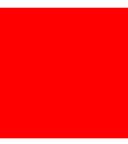| Distribution of dependent variable | Suitable Model (GLM or GAM) |
|---|---|
| Continuous measures: e.g., average income in postcode (£); concentrations of ambient particular matter (PM2.5); Normalised Vegetative Difference Index (NDVI) etc., | Linear regression |
| Binary measures (1 = "present" or 0 = "absent"): e.g., Person's voting for a candidate, Lung cancer risk, house infested with rodents etc., | Logistic Regression |
| Binomial measure (or proportion): e.g., prevalence of houses in a postcode infested with rodents, percentage of people in a village infected with intestinal parasitic worms, prevalence of household on a street segment victimised by crime etc., | Logistic Regression |
| Counts or discrete measures: e.g., number of reported burglaries on a street segment, number of riots in a county etc., | Poisson Regression |

Note 1: The models listed in this table assumes there is a "linear" relationship between the outcome and some independent variable(s).

These models are typically used within a generalised linear modelling framework.

Note 2: The models listed in this table have alternative versions; when the relationship between the outcome and some independent variable(s) are "non-linear", they can be used within a generalised additive modelling framework.

- Typically, the data structure or scenario we have been applying these models to are single row records or unit observations (i.e., for an individual, or a geographical unit etc.,)

- What about data structures with repeated measurements, or unit observation within a group, or temporal datasets?

# Data structures [1]

Imagine we have some dataset containing information on an individual-level.

| ID | Name | Maths Performance | Maths TSR | OFSTED Grade |
|---|---|---|---|---|
| SCH01 | Acton High School | 29 | 20.9 | 5 (Worst) |
| SCH02 | Brentside High School | 40 | 18.6 | 5 (Worst) |
| SCH03 | Greenford High School | 51 | 11.7 | 3 (Below average) |
| SCH04 | Northolt High School | 60 | 9.9 | 2 (Good) |
| SCH05 | Ellen Wilkinson School | 88 | 14.6 | 0 (Excellent) |
| SCH06 | Twyford Church of England | 76 | 6.3 | 1 (Very Good) |
| SCH07 | Featherstone High School | 73 | 5.3 | 1 (Very Good) |
| SCH08 | Drayton Manor High School | 80 | 12.9 | 0 (Excellent) |
| SCH09 | Dormers Wells High School | 67 | 16.5 | 2 (Good) |

Note 1: This dataset contain details for individual schools in Ealing Borough (inside London). Information on the overall maths performance of a school and the maths teacher-student ratio in a class.

We want to understand what historical and sociodemographic factors have an impact on a school's performance when it comes to mathematics.

Note 2: We would typically fit a linear regression model if we wanted to see how just **Maths TSR** and **OFSTED Grade** are linked with **Maths Performance** variable.

# Data structures [2]

Suppose we want to consider broader risk factors, not measured at an individual-level but at a group-level...

| ID | LSOA | Name | Maths Performance | Maths TSR | OFSTED Grade | LSOA IMD Resources | LSOA IMD Income |
|---|---|---|---|---|---|---|---|
| SCH01 | LSOA01 | Acton High School*** | 29 | 20.9 | 5 (Worst) | 5.2305 | 6.4734 |
| SCH02 | LSOA01 | Brentside High School | 40 | 18.6 | 5 (Worst) | 5.2305 | 6.4734 |
| SCH03 | LSOA01 | Greenford High School | 51 | 11.7 | 3 (Below average) | 5.2305 | 6.4734 |
| SCH04 | LSOA02 | Northolt High School | 60 | 9.9 | 2 (Good) | 1.2353 | 0.3491 |
| SCH05 | LSOA02 | Ellen Wilkinson School | 88 | 14.6 | 0 (Excellent) | 1.2353 | 0.3491 |
| SCH06 | LSOA03 | Twyford Church of England | 76 | 6.3 | 1 (Very Good) | 0.2396 | 1.9843 |
| SCH07 | LSOA03 | Featherstone High School | 73 | 5.3 | 1 (Very Good) | 0.2396 | 1.9843 |
| SCH08 | LSOA03 | Drayton Manor High School | 80 | 12.9 | 0 (Excellent) | 0.2396 | 1.9843 |
| SCH09 | LSOA04 | Dormers Wells High School | 67 | 16.5 | 2 (Good) | 3.1435 | 2.3679 |

For instance, other broader factors that might either be on an environmental, geopolitical, societal-level e.g., LSOA IMD public resource allocation for schools and average income scores. We have altered the structure of our dataset and made it far more complex...

# Data structures [3]

A hierarchical or multi-level structure in the dataset is formed

| ID | LSOA | Name | Maths Performance | Maths TSR | OFSTED Grade | LSOA IMD Resources | LSOA IMD Income |
|---|---|---|---|---|---|---|---|
| SCH01 | LSOA01 | Acton High School*** | 29 | 20.9 | 5 (Worst) | 5.2305 | 6.4734 |
| SCH02 | LSOA01 | Brentside High School | 40 | 18.6 | 5 (Worst) | 5.2305 | 6.4734 |
| SCH03 | LSOA01 | Greenford High School | 51 | 11.7 | 3 (Below average) | 5.2305 | 6.4734 |
| SCH04 | LSOA02 | Northolt High School | 60 | 9.9 | 2 (Good) | 1.2353 | 0.3491 |
| SCH05 | LSOA02 | Ellen Wilkinson School | 88 | 14.6 | 0 (Excellent) | 1.2353 | 0.3491 |
| SCH06 | LSOA03 | Twyford Church of England | 76 | 6.3 | 1 (Very Good) | 0.2396 | 1.9843 |
| SCH07 | LSOA03 | Featherstone High School | 73 | 5.3 | 1 (Very Good) | 0.2396 | 1.9843 |
| SCH08 | LSOA03 | Drayton Manor High School | 80 | 12.9 | 0 (Excellent) | 0.2396 | 1.9843 |
| SCH09 | LSOA04 | Dormers Wells High School | 67 | 16.5 | 2 (Good) | 3.1435 | 2.3679 |

We have 9 records but we created an hierarchy…
3 school records nested in LSOA01;
2 school records nested in LSOA02;
3 school records nested in LSOA03;
1 school record nested in LSOA04

# Data structures [4]

## Individual-level data

| ID | Name | Maths Performance | Maths TSR | OFSTED Grade |
|----|------|-------------------|-----------|--------------|
| SCH01 | Acton High School | 29 | 20.9 | 5 (Worst) |
| SCH02 | Brentside High School | 40 | 18.6 | 5 (Worst) |
| SCH03 | Greenford High School | 51 | 11.7 | 3 (Below average) |
| SCH04 | Northolt High School | 60 | 9.9 | 2 (Good) |
| SCH05 | Ellen Wilkinson School | 88 | 14.6 | 0 (Excellent) |
| SCH06 | Twyford Church of England | 76 | 6.3 | 1 (Very Good) |
| SCH07 | Featherstone High School | 73 | 5.3 | 1 (Very Good) |
| SCH08 | Drayton Manor High School | 80 | 12.9 | 0 (Excellent) |
| SCH09 | Dormers Wells High School | 67 | 16.5 | 2 (Good) |

## Group-level data

| LSOA | LSOA IMD Resources | LSOA IMD Income |
|------|--------------------|-----------------|
| LSOA01 | 5.2305 | 6.4734 |
| LSOA02 | 1.2353 | 0.3491 |
| LSOA03 | 0.2396 | 1.9843 |
| LSOA04 | 3.1435 | 2.3679 |

We have 9 records, but we created a hierarchy…

3 school records nested in LSOA01;
2 school records nested in LSOA02;
3 school records nested in LSOA03;
1 school record nested in LSOA04

A typical linear regression model would be severely inadequate for this problem due to the hierarchical structure that is formed in this dataset.

We would need a model that not only takes into account how the records are nested within a group; but one that would allow us to model the "within-group" variations formed in each group, as well as the "across-group" variations. Lastly, we will need a model that will let us fitted both individual-level and group-level variables, this is called a **hierarchical regression model**.
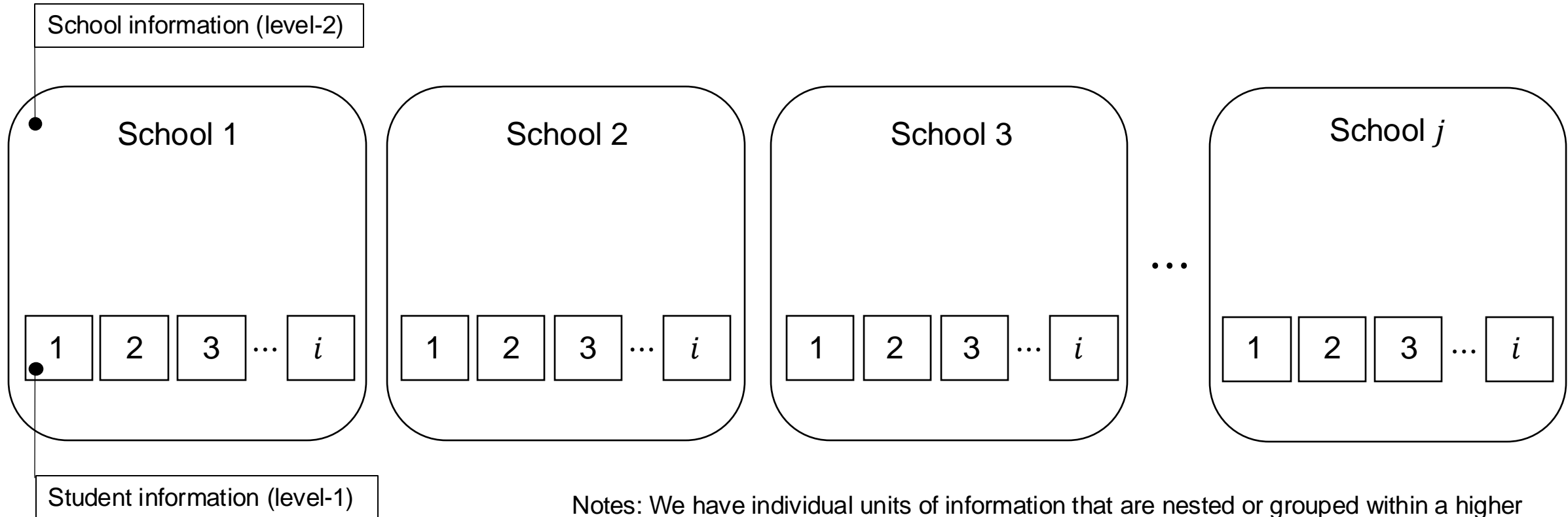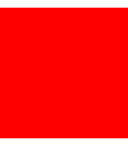
# Definition:

A **hierarchical regression model**, are a specialised group of regression-based models that are able to recognise the existence of hierarchies within a data structure and account for them. It is a statistical model used for exploring the relationship between a dependent variable with one or more independent variables while accounting for these hierarchical structures.

Key characteristics of the hierarchical regression model:

- While it is commonly known as **hierarchical models**, it is also commonly interchangeable with the terms: **Multilevel models**; **Mixed-effect models**, **Nested data model** or even **Random-effects models**.

- The hierarchies formed by the natural structure of the dataset are treated as **levels** in the hierarchical model. There. can be more than one-level formed in the hierarchical regression model. A **two- or three-level hierarchical regression models** are often used a lot in research; however, more and more levels beyond 3 makes the regression incredibly complexed.

- The model structure is based on **levels** – the lowest level always correspond to individual units; while higher levels are the groupings. For example, a survey of a set of $i$ number of students for their academic performance in $j$ number of schools, across a set of years $t$. The students are level-1 (individual-level); schools are level-2 (grouping); and years are the level-3 (grouping or repeated measurements).
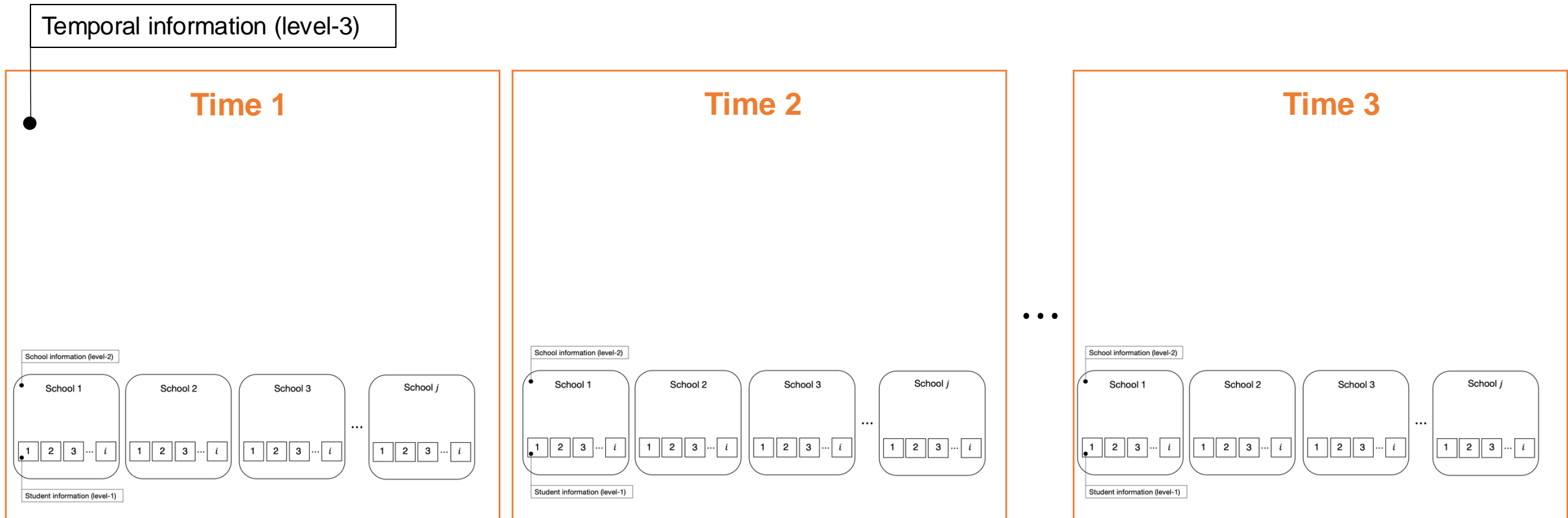
9

# We are Illustrating concisely what we mean by two- or three-level model structure [1]

School information (level-2)

| School 1 | School 2 | School 3 | ... | School $j$ |
|----------|----------|----------|-----|------------|
| 1 2 3 ... $i$ | 1 2 3 ... $i$ | 1 2 3 ... $i$ | | 1 2 3 ... $i$ |

Student information (level-1)

Notes: We have individual units of information that are nested or grouped within a higher measure. This is typically a **two-level structure,** and a **two-level hierarchical regression** model must be used for this scenario.

Building on the example highlighted in point 3 (see slide 9): A survey of a set of $i$ number of students for their academic performance in $j$ number of schools, across a set of years $t$. The students are level-1 (individual-level); schools are level-2 (grouping); and years are the level-3 (grouping or repeated measurements).

# We are Illustrating concisely what we mean by two- or three-level model structure [2]

Temporal information (level-3)



**Time 1**

**Time 2**

**Time 3**

Notes: We have individual units of information that are nested or grouped within a higher measure, whereby the same individuals (from the same units) are repeated (i.e., longitudinal). This is typically a **three-level structure** and so a **three-level hierarchical regression** model must be used for this scenario.

Building on the example highlighted in point 3 (see slide 9): A survey of a set of $i$ number of students for their academic performance in $j$ number of schools, across a set of years $t$. The students are level-1 (individual-level); schools are level-2 (grouping); and years are the level-3 (grouping or repeated measurements).

# Definition:

A **hierarchical regression model**, are a specialised group of regression-based models that are able to recognise the existence of hierarchies within a data structure and account for them. It is a statistical model used for exploring the relationship between a dependent variable with one or more independent variables while accounting for these hierarchical structures.

## Why are hierarchical regression models important:

- It is an elegant way to model datasets that have varying scales in their measurements ( - this artefact is caused by the multilevel or hierarchical structure in the dataset)

- It is a robust approach for accounting for **variations across individual units**, and at the same time, the "**within-group variations**" among groupings

- When we are modelling the direct relationship between the level-1 independent variables against the dependent variable, we can allow for direct interactions between level-1 and higher-level independent variables that were measured at a group-level

- We can quantify group-specific differences as well as group-specific coefficients through the usage of "**varying-slopes**" or "**varying-coefficients**"

Zoology

Food security

Environmental Criminology

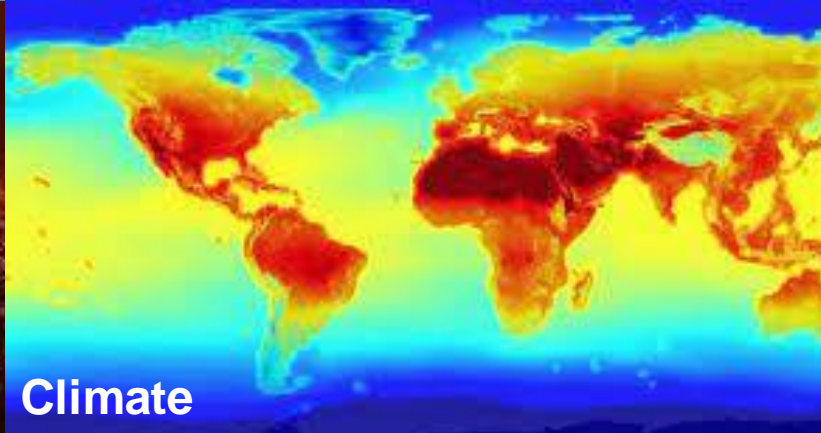Environmental & Spatial Epidemiology
Vector-borne disease

Palaeontogy and

Landscape ecology

Natural Disaster Science

Climate

Humanitarian crisis

13

# Components of a Hierarchical Regression Model

# Recall the base model formula for a GLM

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \varepsilon$$

See Week 2 and 3 notes

**Variables**

- $y$ is the dependent variable
- $x_1, x_2, x_3, \ldots, x_k$ are the independent variables (which we have **k** number of them)

**Parameters**

- $\beta_0$ is the intercept
- $\beta_1, \beta_2, \beta_3, \ldots, \beta_k$ are the slopes (or coefficients) for the corresponding variables $x_1, x_2, x_3, \ldots, x_k$
- $\varepsilon$ is the error term

Let's extend the above model into a hierarchical framework

15

# Mathematical reformulation of the base GLM regression model using indexes

- When there is a hierarchical structure in the dataset, the base form of the GLM can be explicitly reformulated to show the hierarchies with indexes. For instance:

  ❖ We let $i$ represent each individual unit or observation
  ❖ We let $j$ represent a group or cluster which an individual unit or observation $i$ is from.

- Mathematical formulation of such scenario will be as follows:

$$y_{i,j} = \beta_{0,j} + \beta_{1,j}x_{1,i,j} + \beta_{2,j}x_{2,i,j} + \cdots + \beta_{k,j}x_{k,i,j} + \varepsilon_{i,j}$$
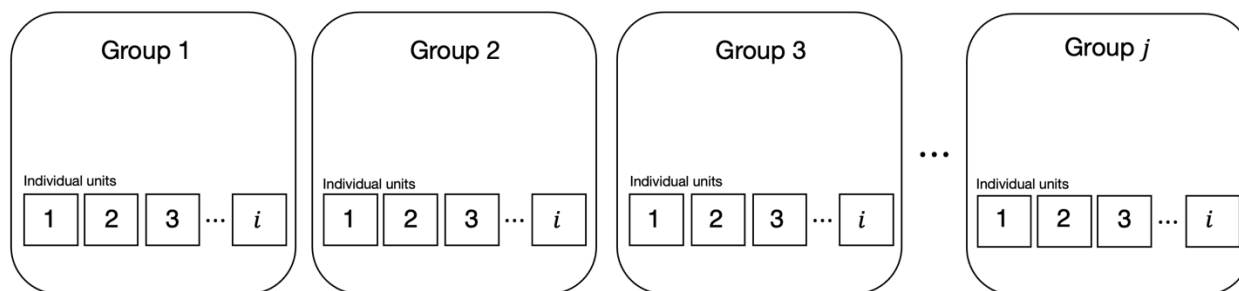
**Breakdown of the above statistical model**

**[1] Variables**

- $y_{i,j}$ is the dependent variable. Is the observed outcome $i$ in group $j$

- $x_1, x_2, x_3, \ldots, x_k$ are the $k$ number independent variables

- Notation $x_{k,i,j}$ is the actual observation. It means that its the $i$ observation in group $j$ for the variable $k$

**[2] Parameters**

- $\beta_{0,j}$ is the intercept

- $\beta_{1,j}, \beta_{2,j}, \beta_{3,j}, \ldots, \beta_{k,J}$ are the coefficients corresponding to $x_1, x_2, x_3, \ldots, x_k$

- $\varepsilon_{i,j}$ is an error term

2-level hierarchical data drawn in picture form



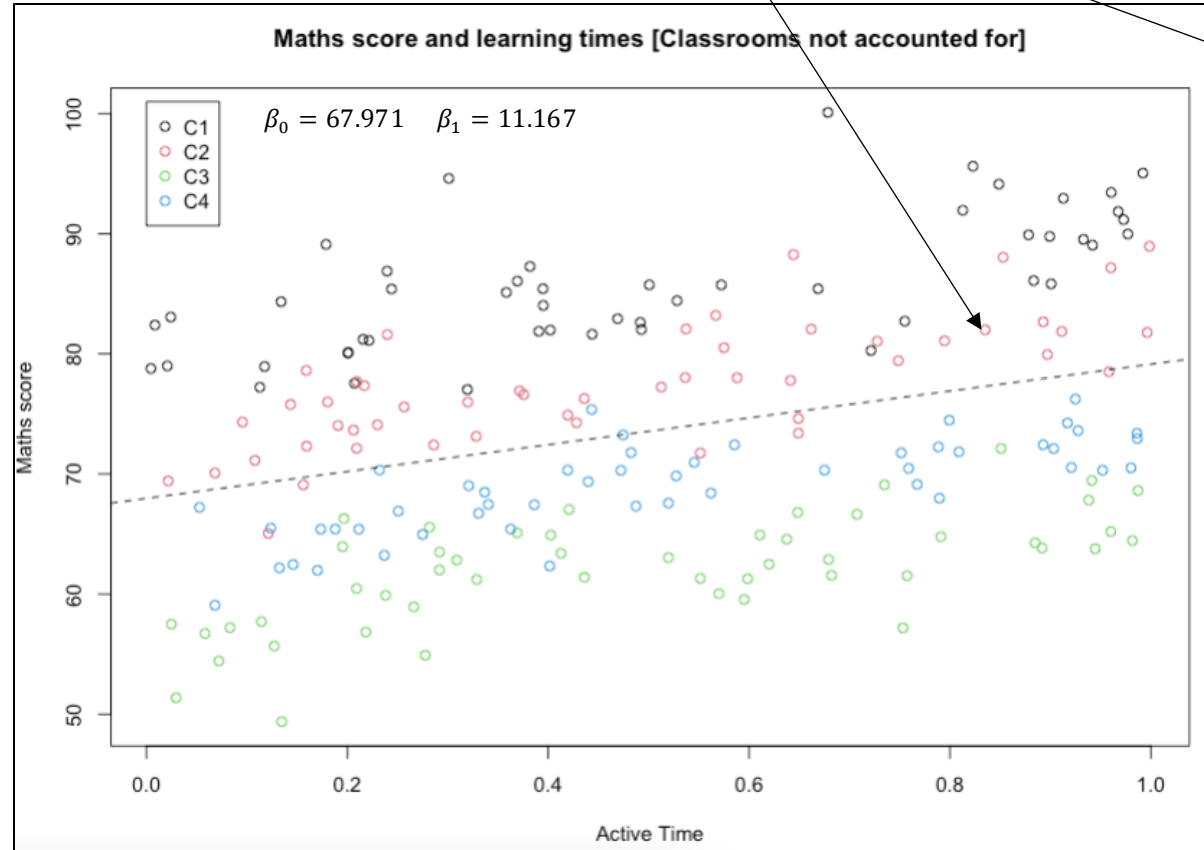2-level hierarchical data frame written in matrix algebraic form

| $i$ | $j$ | $y_{i,j}$ | $x_1$ | $x_2$ | $x_3$ | $\cdots$ | $x_k$ |
|---|---|---|---|---|---|---|---|
| 1 | 1 | $y_{1,1}$ | $x_{1,1,1}$ | $x_{2,1,1}$ | $x_{3,1,1}$ | $\cdots$ | $x_{k,1,1}$ |
| 2 | 1 | $y_{2,1}$ | $x_{1,2,1}$ | $x_{2,2,1}$ | $x_{3,2,1}$ | $\cdots$ | $x_{k,2,1}$ |
| 3 | 1 | $y_{3,1}$ | $x_{1,3,1}$ | $x_{2,3,1}$ | $x_{3,3,1}$ | $\cdots$ | $x_{k,3,1}$ |
| 1 | 2 | $y_{1,2}$ | $x_{1,1,2}$ | $x_{2,1,2}$ | $x_{3,1,2}$ | $\cdots$ | $x_{k,1,2}$ |
| 2 | 2 | $y_{2,2}$ | $x_{1,2,2}$ | $x_{2,2,2}$ | $x_{3,2,2}$ | $\cdots$ | $x_{k,2,2}$ |
| 1 | 3 | $y_{1,3}$ | $y_{1,1,3}$ | $x_{2,1,3}$ | $x_{3,1,3}$ | $\cdots$ | $x_{k,1,3}$ |
| 2 | 3 | $y_{2,3}$ | $y_{1,2,3}$ | $x_{2,2,3}$ | $x_{3,2,3}$ | $\cdots$ | $x_{k,2,3}$ |
| 3 | 3 | $y_{3,3}$ | $y_{1,3,3}$ | $x_{3,3,3}$ | $x_{3,3,3}$ | $\cdots$ | $x_{k,3,3}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\cdots$ | $\vdots$ |
| $i$ | $j$ | $y_{i,j}$ | $x_{1,i,j}$ | $x_{2,i,j}$ | $x_{3,i,j}$ | $\cdots$ | $x_{k,i,j}$ |

# Notation for the intercept and coefficient i.e., $\beta_{0,j}$ and $\beta_{k,j}$ - what are they?

- Let us consider the following scenarios: 50 students in 4 classes, we are interested to know how active learning time impacts maths score
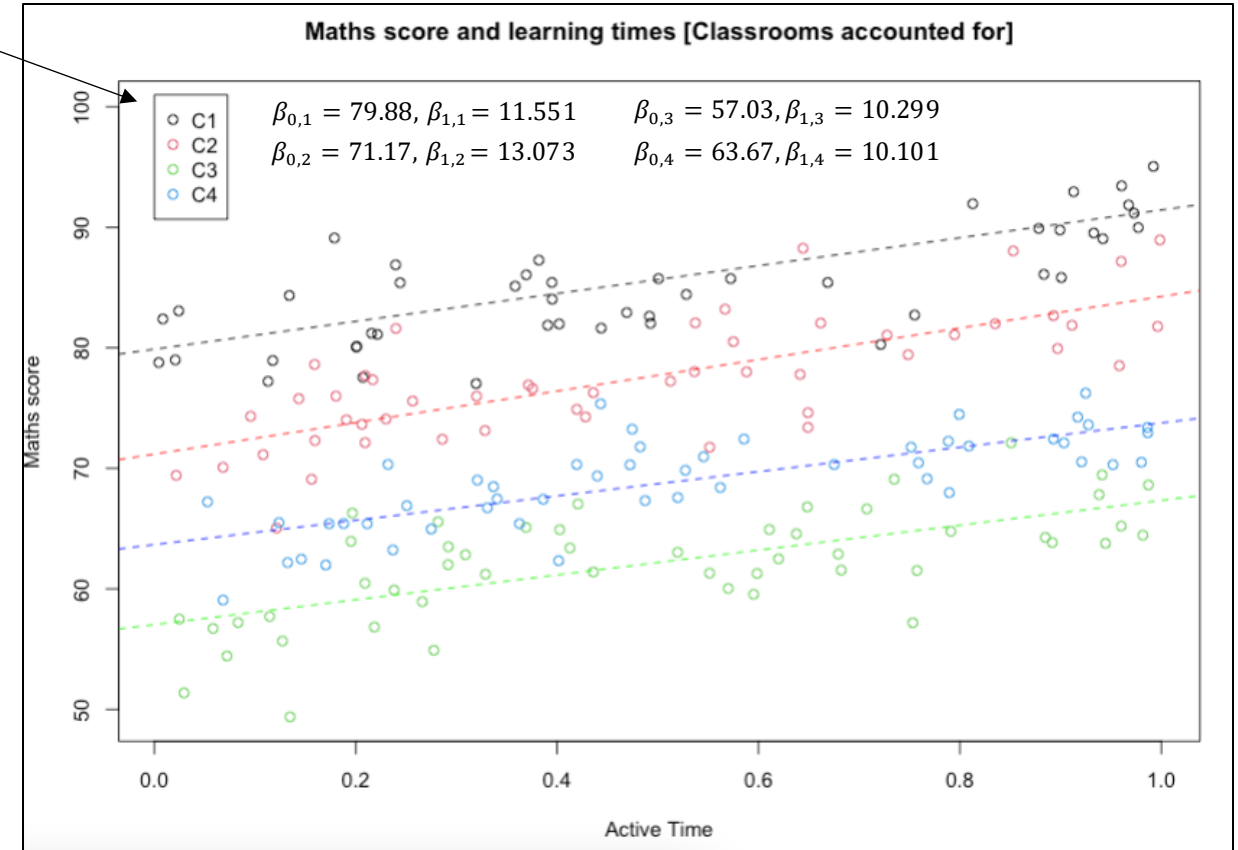
GLM model (not-indexed)

$$y = \beta_0 + \beta_1 x_1 + \varepsilon$$

GLM model (indexed)

$$y_{i,j} = \beta_{0,j} + \beta_{1,j} x_{1,i,j} + \varepsilon_{i,j}$$



Maths score and learning times [Classrooms not accounted for]

$\beta_0 = 67.971 \quad \beta_1 = 11.167$



Maths score and learning times [Classrooms accounted for]

$\beta_{0,1} = 79.88, \beta_{1,1} = 11.551 \quad \beta_{0,3} = 57.03, \beta_{1,3} = 10.299$

$\beta_{0,2} = 71.17, \beta_{1,2} = 13.073 \quad \beta_{0,4} = 63.67, \beta_{1,4} = 10.101$

Here, we can see that if we use a statistical model to analyse this data without regards for the group structure. We get a single intercept and a single coefficient. Here, we are assuming that this relationship between active times and maths are similar across all the 200 students regardless of the classrooms they are in. This is what we term as Fixed Effects scenario

However, we can see that this panel shows something different. Accounting for the classroom groups, **by fitting separate linear models** we get different intercepts (i.e., global mean specific to a group) with different slope (or slope variation). There is an indication that some variation within the groups are causing this pattern. This variation is known as a **Random Effects**, and it acting on our **intercepts and slopes**. This random effect must be accounted for and so doing this would mean reformulating the **indexed model with the random effects to show the equation in its true hierarchical form!**

17

# Mathematical formulation for hierarchical regression model (full form) [1]

- Remember, we are **strictly** using a simple case of the 2-level model scenario

    ❖ We let $i$ represent each individual unit or observation
    ❖ We let $j$ represent a group or cluster which an individual unit or observation $i$ is from.

- Mathematical formulation of such scenario will be as follows:

$$y_{i,j} = \beta_{0,j} + \beta_{1,j}x_{1,i,j} + \beta_{2,j}x_{2,i,j} + \cdots + \beta_{k,j}x_{k,i,j} + \varepsilon_{i,j}$$   Level 1 Equation

- For $\beta_{0,j}$ and for some $\beta_{k,j}$. Let us introduce some random effect (or random deviation) $u_{k,j}$ which causes this global intercept and slope coefficient to vary across some groups $j$ and incorporate them to the indexed model. We would have these new equations specifically for the intercept $\beta_{0,j}$ and for the coefficients $\beta_{k,j}$ from the indexed model:

$$\beta_{0,j} = \gamma_{00} + u_{0,j}$$
$$\beta_{1,j} = \gamma_{10} + u_{1,j}$$
$$\beta_{2,j} = \gamma_{20} + u_{2,j}$$
$$\vdots \qquad \vdots \qquad \vdots$$
$$\beta_{k,j} = \gamma_{k0} + u_{k,j}$$

Level 2 Equations

Note 1: The above level 1 equation is a normal regression that we know. We makes it a hierarchical regression model is when we incorporate these random effects by nesting these level 2 equations to the intercept and coefficients in the level 1 equation.

Note 2: Notice how the intercept and slopes in level 1 equation are indeed a function of the components of the level 2 counterparts? To get the hierarchical regression in its full form, we simply substitute the level 2 equation into level 1
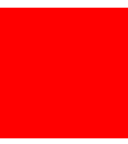
[1] Breakdown of components for the level 1 equation (individual units)
- $y_{i,j}$ is the dependent variable. Is the observed outcome $i$ in group $j$
- $\beta_{0,j}$ is the intercept
- $\beta_{1,j}, \beta_{2,j}, \beta_{3,j}, \ldots, \beta_{k,J}$ are the coefficients that correspond to $x_1, x_2, x_3, \ldots, x_k$
- $\varepsilon_{i,j}$ is residual term

[2] Breakdown of the components of level 2 equation (group units)
- $\gamma_{00}$ is a **fixed effect** (i.e., a constant term) associated with the intercept $\beta_{0,j}$
- $\gamma_{10}, \gamma_{20}, \cdots, \gamma_{k0}$ are **fixed effects** (i.e., constant terms) for the associated coefficients $\beta_{1,j}, \beta_{2,j}, \beta_{3,j}, \ldots, \beta_{k,J}$
- $u_{0,j}$ is the **random effects** (random variation caused by the groupings) on the intercept $\beta_{0,j}$
- $u_{1,j}, u_{2,j}, \cdots, u_{k,j}$ are **random effects** (i.e., random variation caused by the groupings) on the coefficients $\beta_{1,j}, \beta_{2,j}, \beta_{3,j}, \ldots, \beta_{k,J}$

# Mathematical formulation for hierarchical regression model (full form) [2]

$$y_{i,j} = \beta_{0,j} + \beta_{1,j}x_{1,i,j} + \beta_{2,j}x_{2,i,j} + \cdots + \beta_{k,j}x_{k,i,j} + \varepsilon_{i,j}$$

**Level 1 Equation**

$$\beta_{0,j} = \gamma_{00} + u_{0,j}$$
$$\beta_{1,j} = \gamma_{10} + u_{1,j}$$
$$\beta_{2,j} = \gamma_{20} + u_{2,j}$$
$$\vdots \quad \vdots \quad \vdots$$
$$\beta_{k,j} = \gamma_{k0} + u_{k,j}$$

- 1st equation is a random-intercept
- 2nd, 3rd and 4th and so on equations are random-slopes
- Note that these equation does not have a two-level independent variable that impacts the outcome

**Level 2 Equations**

- Substitute the level 2 model equations into the level 1 model equation:

$$\Rightarrow y_{i,j} = (\gamma_{00} + u_{0,j}) + (\gamma_{10} + u_{1,j})x_{1,i,j} + (\gamma_{20} + u_{2,j})x_{2,i,j} + \cdots + (\gamma_{k0} + u_{1,j})x_{k,i,j} + \varepsilon_{i,j}$$

- After substitution, we expanding the expression and rearrange as follows:

$$\Rightarrow y_{i,j} = \underbrace{\gamma_{00} + \gamma_{10}x_{1,i,j} + \gamma_{20}x_{2,i,j} + \cdots + \gamma_{k0}x_{k,i,j}}_{\text{Fixed part}} + \underbrace{u_{0,j} + u_{1,j}x_{1,i,j} + u_{2,j}x_{2,i,j} + \cdots + u_{k,j}x_{k,i,j} + \varepsilon_{i,j}}_{\text{Random part}}$$

Model's true form

Note: There are model scenarios

- $\gamma_{00}$ is the global intercept from the fixed part of the model we want to report

- $\gamma_{10}, \gamma_{20}, \ldots$ and $\gamma_{k0}$ are the coefficients from the fixed part of the model we want to report now

- $u_{0,j}, u_{1,j}, u_{2,j}, \ldots$ and $u_{k,j}$ as well as $\varepsilon_{i,j}$ they have variances for random part of the model we want to report
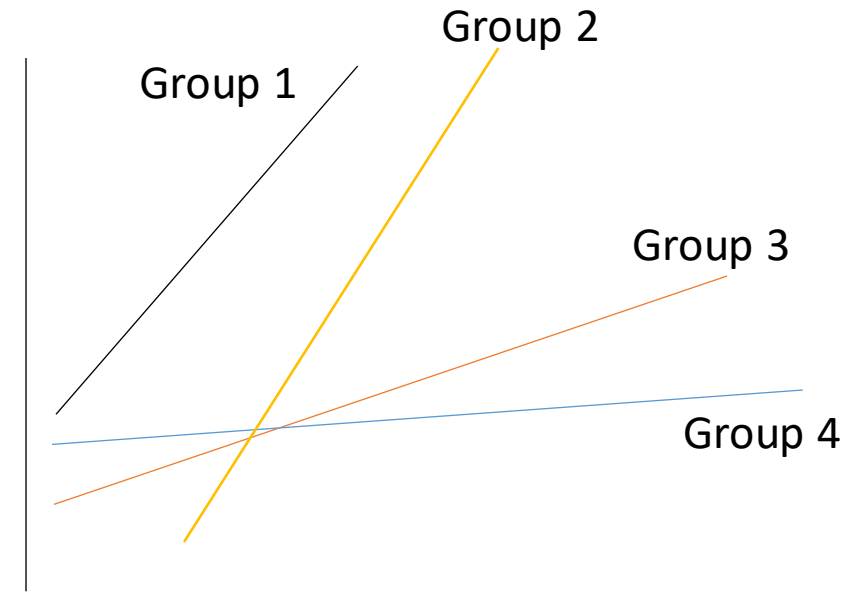
$$y_{i,j} = \beta_{0,j} + \beta_{1,j}x_{1,i,j} + \beta_{2,j}x_{2,i,j} + \cdots + \beta_{k,j}x_{k,i,j} + \varepsilon_{i,j}$$

Level 1 Equation

$$\beta_{0,j} = \gamma_{00} + u_{0,j}$$
$$\beta_{1,j} = \gamma_{10} + u_{1,j}$$
$$\beta_{2,j} = \gamma_{20} + u_{2,j}$$
$$\vdots \quad \vdots \quad \vdots$$
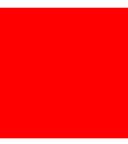$$\beta_{k,j} = \gamma_{k0} + u_{k,j}$$

Level 2 Equations

This is an example of a **random-slope model** which includes both a random-intercept and random-slopes. This means their group structures causes variation in the means across groups (i.e., intercepts) and slopes

$$y_{i,j} = \gamma_{00} + \gamma_{10}x_{1,i,j} + \gamma_{20}x_{2,i,j} + \cdots + \gamma_{k0}x_{k,i,j} + u_{0,j} + u_{1,j}x_{1,i,j} + u_{2,j}x_{2,i,j} + \cdots + u_{k,j}x_{k,i,j} + \varepsilon_{i,j}$$

Fixed part

Random part

Model's true form

# Random-intercept-only, Random-slopes & Random coefficient scenarios [2]

$$y_{i,j} = \beta_{0,j} + \beta_{1,j}x_{1,i,j} + \beta_{2,j}x_{2,i,j} + \cdots + \beta_{k,j}x_{k,i,j} + \varepsilon_{i,j}$$

Level 1 Equation

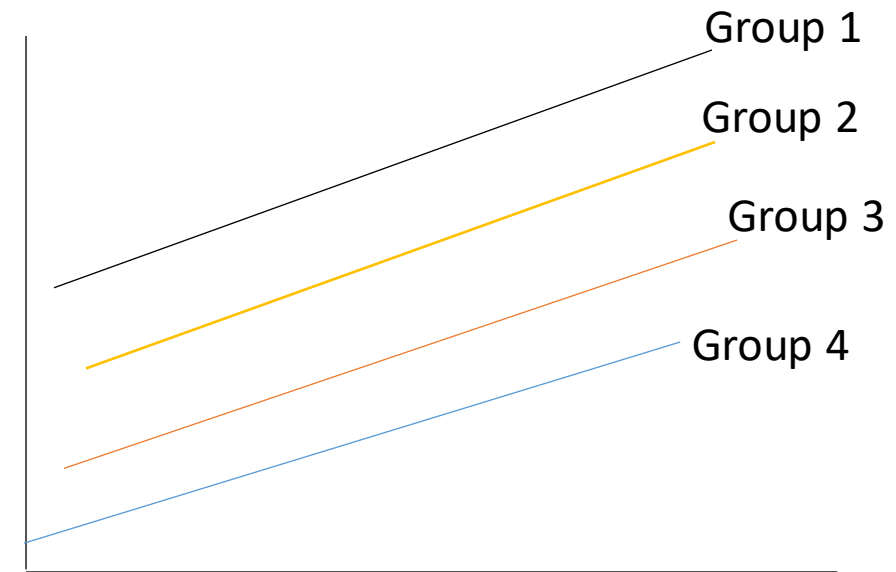$$\beta_{0,j} = \gamma_{00} + u_{0,j}$$

Level 2 Equation

Here, the model is much simpler:

$$y_{i,j} = \underbrace{\gamma_{00} + \beta_{1,j}x_{1,i,j} + \beta_{2,j}x_{2,i,j} + \cdots + \beta_{k,j}x_{k,i,j}}_{\text{Fixed part}} + \underbrace{u_{0,j} + \varepsilon_{i,j}}_{\text{Random part}}$$

Model's true form

This is an example of a **random-intercept-only model** which only includes a random-intercept and excludes the random-slopes. This means that the group structure causes variation on the means (i.e., group-specific intercepts) but not on slopes



Group 1

Group 2

Group 3

Group 4

# Random-intercept-only, Random-slopes & Random coefficient scenarios [1]

$$y_{i,j} = \beta_{0,j} + \beta_{1,j}x_{1,i,j} + \beta_{2,j}x_{2,i,j} + \cdots + \beta_{k,j}x_{k,i,j} + \varepsilon_{i,j}$$

Level 1 Equation

$$\beta_{0,j} = \gamma_{00} + \gamma_{01}Z_1 + u_{0,j}$$
$$\beta_{1,j} = \gamma_{10} + \gamma_{11}Z_1 + u_{1,j}$$
$$\beta_{2,j} = \gamma_{20} + \gamma_{21}Z_1 + u_{2,j}$$
$$\vdots \quad \vdots \quad \vdots \quad \vdots$$
$$\beta_{k,j} = \gamma_{k0} + \gamma_{k1}Z_1 + u_{k,j}$$

Level 2 Equations

Suppose we have an independent variable measure on the group-level impacting our outcome on the individual-level.

- Substitute the level 2 model equations with the variables into the level 1 model equation:

$$\Rightarrow y_{i,j} = (\gamma_{00} + \gamma_{01}Z_I + u_{0,j}) + (\gamma_{10} + \gamma_{11}Z_I + u_{1,j})x_{1,i,j} + (\gamma_{20} + \gamma_{21}Z_I + u_{2,j})x_{2,i,j} + \cdots + (\gamma_{k0} + \gamma_{k1}Z_I + u_{1,j})x_{k,i,j} + \varepsilon_{i,j}$$

- After substitution, we expanding the expression and rearrange as follows:

$$\Rightarrow y_{i,j} = \gamma_{00} + \gamma_{01}Z_1 +$$
$$\gamma_{10}x_{1,i,j} + \gamma_{20}x_{2,i,j} + \cdots + \gamma_{k0}x_{k,i,j} +$$
$$\gamma_{11}Z_1 x_{1,i,j} + \gamma_{20}Z_1 x_{2,i,j} + \cdots + \gamma_{k0}Z_1 x_{k,i,j} +$$
$$u_{0,j} + u_{1,j}x_{1,i,j} + u_{2,j}x_{2,i,j} + \cdots + u_{k,j}x_{k,i,j} + \varepsilon_{i,j}$$

Model's true form

$\gamma_{00}$ is the global or population mean

$\gamma_{01}$ is the random coefficient for $Z_1$

These are fixed effects coefficients for the variables in the level 1 equation

These are random coefficients for the interacting variables from the level 1 & 2 equation

These are the random effects

**Advice – make life easy for yourself and use the random-intercept-only model. If you have a level-2 variable as you won't have to deal with any interactions!**
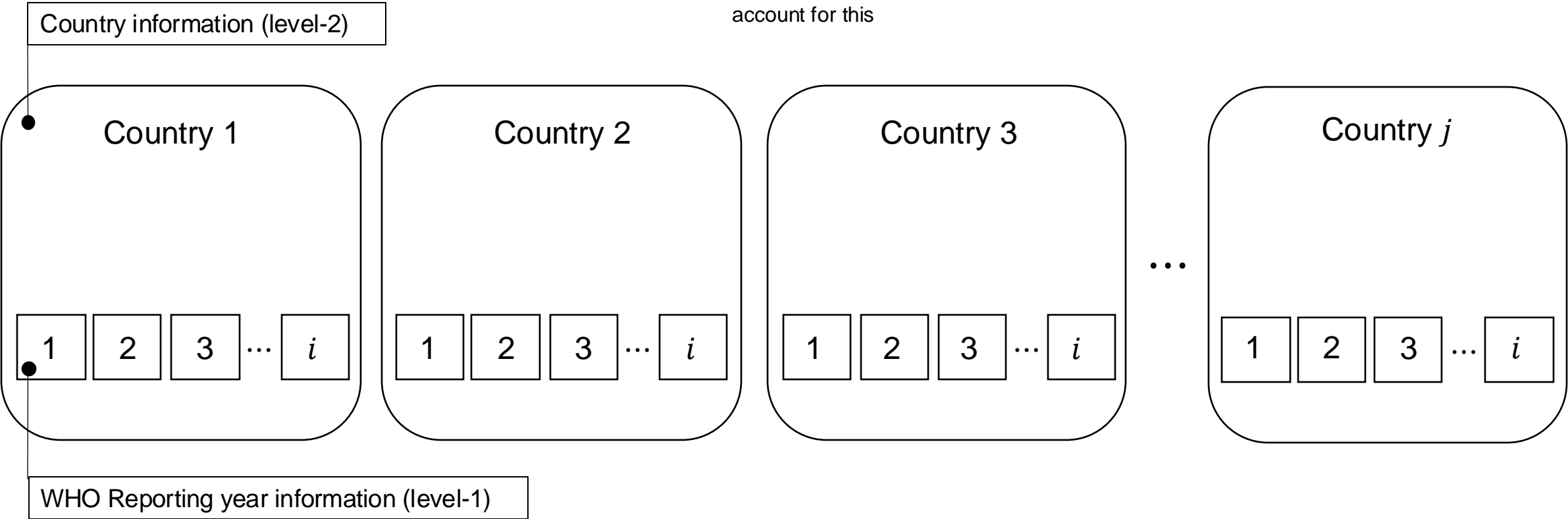
# An example and Interpretation

# Example: Assessing the impact of water and sanitation provision on Cholera burden in Sub-Saharan Africa [1]

**GOAL: Find the association between limited water and sanitation services (%) with incident Cholera (2000-2017) across 13 countries**

$y_{i,j}$ = Incident Cholera reported in the $i$ year in country $j$
$x_{1,i,j}$ = The proportion or coverage without water service in year $i$ in country $j$
$x_{2,i,j}$ = The proportion or coverage without sanitation services in year $i$ in country $j$

The 17 reporting years are clustered into 13 different countries. We want to know two things: the overall association between the cholera and these two variables. But we want the risk to varying across countries. Hence, we will use a **random-intercept and slope model** to account for this

Country information (level-2)

| Country 1 | Country 2 | Country 3 | ... | Country $j$ |
|---|---|---|---|---|
| 1  2  3  ...  $i$ | 1  2  3  ...  $i$ | 1  2  3  ...  $i$ | | 1  2  3  ...  $i$ |

WHO Reporting year information (level-1)

# Information

**Dependent variable**
- $y_{i,j}$ = Incident Cholera reported in the $i$ year in country $j$

**Primary independent variables**
- $x_{1,i,j}$ = The proportion or coverage without water service in year $i$ in country $j$
- $x_{2,i,j}$ = The proportion or coverage without sanitation services in year $i$ in country $j$

The 17 reporting years are clustered into 13 different countries. We want to know two things: the overall association between the cholera and these two variables. But we want the risk to varying across countries. Hence, we will use a **random-intercept and slope model** to account for this. Note that GDP, temperature and rainfall variables were included as apriori confounding variables.

# Model formulation

- Using a 2-level hierarchical model (random-intercept-only

$$y_{i,j} = \beta_{0,j} + \beta_{1,j}x_{1,i,j} + \beta_{2,j}x_{2,i,j} + \beta_3 x_{3,i,j} + \beta_4 x_{4,i,j} + \beta_5 x_{5,i,j} + \log(P_{i,j}) + \varepsilon_{i,j}$$

$\beta_{0,j} = \gamma_{00} + u_{0,j}$    allows the intercept to vary across countries (level-2)

$\beta_{1,j} = \gamma_{01} + u_{1,j}$    allows the 'no water service' variable to vary across countries (level-2)

$\beta_{2,j} = \gamma_{02} + u_{2,j}$    allows the "no sanitation service' variable to vary across countries (level-2)

- Specify likelihood function. The outcome is count outcome, thus it is Poisson model with overdispersion

$$y_{i,j} \sim \text{negbin}(\lambda_{i,j}, \phi)$$
$$\lambda_{i,j} = \beta_{0,j} + \beta_{1,j}x_{1,i,j} + \beta_{2,j}x_{2,i,j} + \beta_3 x_{3,i,j} + \beta_4 x_{4,i,j} + \beta_5 x_{5,i,j} + \log(P_{i,j})$$

- Define the priors for the intercept, fixed and random effects

$\gamma_{00} \sim \text{normal}(0, 1)$
$\gamma_{01} \sim \text{normal}(0, 1)$
$\gamma_{02} \sim \text{normal}(0, 1)$
$u_{0,\text{intercept}} \sim \text{normal}(0, \sigma_{\text{intercept}})$
$\sigma_{\text{intercept}} \sim \text{cauchy}(0, 0.5)$
$u_{0,\text{water}} \sim \text{normal}(0, \sigma_{\text{water}})$
$\sigma_{\text{water}} \sim \text{cauchy}(0, 0.5)$
$u_{0,\text{sanitation}} \sim \text{normal}(0, \sigma_{\text{sanitation}})$
$\sigma_{\text{sanitation}} \sim \text{cauchy}(0, 0.5)$
$\phi \sim \text{cauchy}(0, 0.4)$
$\beta_3 \sim \text{normal}(0, 1)$
$\beta_4 \sim \text{normal}(0, 1)$
$\beta_5 \sim \text{normal}(0, 1)$

Define the **data** block

```
data {
  int<lower=1> N;                                    // Number of observations
  int<lower=1> Country;                              // Number of countries
  int<lower=1, upper=Country> CountryID[N];          // Country IDs (Aligns the groupings to the units)
  int<lower=0> Cholera[N];                            // Cholera cases (Dependent variable)
  real Water[N];                                     // Water access variable
  real Sanitation[N];                                // Sanitation variable
  real GDP[N];
  real Rainfall[N];
  real Temperature[N];
  real Population[N];                                // Population (used as offset)
  real Kappa;                                        // Overdispersion prior (0.4)
}
```

Define the **parameters** block

```
parameters {
  real gamma00;                              // Overall intercept
  real gamma01;                              // Overall fixed effect for Water
  real gamma02;                              // Overall fixed effect for Sanitation
  real beta3;                                // Overall relationship for GDP
  real beta4;                                // Overall relationship for rainfall
  real beta5;                                // Overall relationship for temperature
  real random_intercept[Country];           // Country-specific random intercepts
  real random_slope_water[Country];          // Country-specific random slopes for Water
  real random_slope_sanitation[Country];    // Country-specific random slopes for Sanitation
  real<lower=0> group_intercept_sd;         // SD of random intercepts
  real<lower=0> group_slope_water_sd;       // SD of random slopes for Water
  real<lower=0> group_slope_sanitation_sd;   // SD of random slopes for Sanitation
  real<lower=0> phi;                         // Overdispersion parameter
}
```

Define the **transformed parameters** block

**transformed parameters** {
  // Here, we build the sub-equations that will allow the parameters to vary across groups

  real beta00[Country];
  real beta01[Country];
  real beta02[Country];

  for (j in 1:Country) {
    beta00[j] = gamma00 + random_intercept[j];      // Random intercept per country
    beta01[j] = gamma01 + random_slope_water[j];    // Random slope for Water per country
    beta02[j] = gamma02 + random_slope_sanitation[j];    // Random slope for Sanitation per country
  }
}

# Define the **model** block

```
model {
 // Priors for fixed effects
 gamma00 ~ normal(0, 1);
 gamma01 ~ normal(0, 1);
 gamma02 ~ normal(0, 1);
 beta3 ~ normal(0, 1);
 beta4 ~ normal(0, 1);
 beta5 ~ normal(0, 1);

 // Priors for random effects
 random_intercept ~ normal(0, group_intercept_sd);
 random_slope_water ~ normal(0, group_slope_water_sd);
 random_slope_sanitation ~ normal(0, group_slope_sanitation_sd);

 // Priors for standard deviations of random effects
 group_intercept_sd ~ cauchy(0, 0.5);
 group_slope_water_sd ~ cauchy(0, 0.5);
 group_slope_sanitation_sd ~ cauchy(0, 0.5);

 // Prior for overdispersion parameter
 phi ~ cauchy(0, Kappa);

 // Likelihood: Negative Binomial Poisson Regression
 for (i in 1:N) {
   Cholera[i] ~ neg_binomial_2_log(beta00[CountryID[i]] + beta01[CountryID[i]]*Water[i] + beta02[CountryID[i]]*Sanitation[i] + beta3*GDP[i] + beta4*Rainfall[i] +
beta5*Temperature[i] + log(Population[i]), phi);
 }

}
```

# Define the **generate quantities** block

```
generated quantities {
 // report the coefficients as relative risk ratios
  real gamma00_RR;
  real gamma01_RR;
  real gamma02_RR;

  gamma00_RR = exp(gamma00);
  gamma01_RR = exp(gamma01);
  gamma02_RR = exp(gamma02);

  real beta3_RR;
  real beta4_RR;
  real beta5_RR;

  beta3_RR = exp(beta3);
  beta4_RR = exp(beta4);
  beta5_RR = exp(beta5);

  // report the varying slopes as relative risk ratios
  vector[13] beta01_RR;
  beta01_RR = exp(beta01);

  vector[13] beta02_RR;
  beta02_RR[1] = exp(beta02[1]);
}
```

**Table 2.** Using a 2-level multilevel negative binomial Poisson regression model within a Bayesian framework to quantify the overall associated risk between basic water and sanitation services and incident Cholera (adjusted for GDP, rainfall, and temperature) for 13 countries in Sub Saharan Africa.

| Variables | Posterior RR (95 CrI) | P(RR > 1.00) | ESS ($\hat{R}$) | |
|---|---|---|---|---|
| Without Basic Water Services (%) | 1.70 (95% CrI: 0.90 to 2.88) | 0.95 | 25,747 ($\hat{R}$ = 1.00 < 1.05) | $\gamma_{01}$ |
| Without Basic Sanitation Services (%) | 1.18 (95% CrI: 0.73 to 1.85) | 0.74 | 21,269 ($\hat{R}$ = 1.00 < 1.05) | $\gamma_{02}$ |

1.) RR: Relative risks; 2.) 95% CrI: 95% Credibility Intervals; 3.) P(RR > 1.00): Interpretated as the exceedance probability, meaning the probability that the relative risks for a particular variable exceeding the null value of 1.00. This should be interpretated as the probability of cholera risk being excessively high in relation to that independent variable; 4.) ESS: Effective sample size; 5.) $\hat{R}$ less than 1.05 means the parameter estimates are valid and convergence was achieved. 6.) GDP: Gross Domestic Product

**Table 3.** Using an extended 2-level multilevel negative binomial Poisson regression model with varying slopes within a Bayesian framework to quantify the associated risk between basic water & sanitation services with incident Cholera (adjusted for GDP, rainfall, and temperature) for each of the 13 countries in Sub Saharan Africa.

| Country | Coverage with no basic water services Posterior RR (95% CrI) [P(RR > 1.00)] | Coverage with no sanitation services Posterior RR (95% CrI) [P(RR > 1.00)] |
|---|---|---|
| Benin | 1.63 (95% CrI: 0.67 to 3.32) [0.84] | 1.11 (95% CrI: 0.61 to 1.96) [0.57] |
| Burundi | 3.92 (95% CrI: 1.51 to 8.53) [1.00] | 2.69 (95% CrI: 1.22 to 5.13) [0.99] |
| DRC | 2.03 (95% CrI: 1.23 to 3.14) [1.00] | 1.45 (95% CrI: 0.81 to 2.40) [0.89] |
| Ghana | 2.02 (95% CrI: 0.72 to 4.62) [0.90] | 1.38 (95% CrI: 0.64 to 2.84) [0.74] |
| Ivory Coast | 5.54 (95% CrI: 1.36 to 16.58) [0.99] | 4.02 (95% CrI: 0.93 to 12.6) [0.97] |
| Kenya | 1.24 (95% CrI: 0.09 to 3.80) [0.51] | 0.86 (95% CrI: 0.07 to 2.59) [0.31] |
| Malawi | 2.40 (95% CrI: 0.78 to 6.06) [0.93] | 1.67 (95% CrI: 0.61 to 4.08) [0.79] |
| Mozambique | 1.61 (95% CrI: 1.11 to 2.29) [1.00] | 1.14 (95% CrI: 0.75 to 1.68) [0.71] |
| Niger | 0.93 (95% CrI: 0.50 to 1.58) [0.34] | 0.64 (95% CrI: 0.42 to 0.94) [0.01] |
| Nigeria | 0.66 (95% CrI: 0.19 to 1.61) [0.16] | 0.44 (95% CrI: 0.17 to 0.91) [0.01] |
| Somalia | 1.85 (95% CrI: 1.05 to 3.06) [0.98] | 1.31 (95% CrI: 0.76 to 2.15) [0.81] |
| Tanzania | 1.19 (95% CrI: 0.80 to 1.75) [0.78] | 0.84 (95% CrI: 0.57 to 1.23) [0.16] |
| Togo | 1.55 (95% CrI: 0.84 to 2.64) [0.91] | 1.08 (95% CrI: 0.70 to 1.64) [0.58] |

$\beta_{1,j}$ $\beta_{2,j}$

# Any questions?