

Centre for Multilevel Modelling

Random intercept models

A transcript of random intercept models presentation, by Rebecca Pillinger

- To watch the presentation go to [Random intercepts models - listen to voice-over with slides and subtitles](#) (If you experience problems accessing any videos, please email info-cmm@bristol.ac.uk)

1. [Random intercept models: What are they and why use them?](#)

- [Explanatory variables](#)
- [Using a single-level regression model](#)
- [Solution: Random Intercept Model](#)
- [Fixed and random part](#)
- [What does the model look like?](#)

2. [Random intercept models: Research questions and interpretation](#)

- [So how do we interpret the parameters?](#)
- [Substantive interest vs. nuisance](#)
- [Examples of research questions](#)

3. [Random intercept models: Adding more explanatory variables](#)

- [Example](#)
- [When is a variable a level?](#)
- [Guidelines](#)
 - [Examples](#)

4. [Random intercept models: Hypothesis Testing](#)

- [Likelihood ratio test](#)

5. [Random intercept models: Variance partitioning coefficients](#)

- [ρ and clustering](#)
- [Interpreting the value of ρ](#)
- [Clustering in the model](#)

6. [Random intercept models: the correlation matrix](#)

- [Assumptions of the random part](#)
- [V, the correlation matrix](#)
- [Covariance matrix for a single level model](#)
- [Covariance matrix for a random intercepts model](#)
- [Calculating the covariances](#)
- [V for random intercepts model](#)

7. [Random intercepts models: Residuals](#)

- [Multilevel residuals](#)
- [Why do we shrink?](#)
- [Points to note about shrinkage](#)

8. [Random intercepts models: Predictions](#)

- [Visualising the model](#)

1) Random intercept models: What are they and why use them?

Explanatory variables

What about explanatory variables?

- We've seen how to fit a variance components model
- This lets us see how much of the variance in our response is at each level
- But what if we want to look at the effects of explanatory variables?

Example

Suppose we have data on exam results for pupils within schools

- We fit a variance components model and we find 20% of the variance is at the school level.
- But can we interpret this as "20% of the variance in exam scores is caused by schools"?
- Schools differ in their intake policy and in the pupils who apply
- These differences also contribute to school-level variance
- So we would like to control for previous exam score

We've seen how to fit a variance components model and that lets us see how much of the variance in our response is at each level. But what if we want to look at the effects of explanatory variables?

So for example, suppose that we have data on exam results of pupils within schools and we fit a variance components model, and find that 20% of the variance is at the school level. But can we really interpret that as meaning "20% of the variance in exam scores is caused by schools"? Because schools differ in their intake policy and in the pupils who apply. So this 20% variance at the school level could actually be caused, partly or wholly, by the fact that the pupils were actually different before they entered the school. So we'd actually like to control for the previous exam score so that we can try and look at just the variance that's due to things that have happened whilst those pupils are at that school.

Using a single-level regression model

Using a single-level regression model

- We usually do this by fitting a regression model:
 $y_i = \beta_0 + \beta_1 x_i + e_i \quad e_i \sim N(0, \sigma_e^2)$
- When we have clustered data, using this model causes problems.

Clustered data
Data where observations in the same group are related, e.g.

- exam results of pupils within schools
- heights of children within families

Problem 1
If we fit this model to clustered data we get the wrong answers

Problem 2

- This model doesn't show us how much variation is at each level
- So we can't find out by using this model how much of an effect school has on exam score after controlling for intake score
- This is what we're interested in → problem

Well usually, if we want to control for something, we do that by fitting a regression model, like this, but when we have clustered data, using this single level model causes problems, as we've seen. And clustered data are data where observations in the same group are related, so, for example, exam results for pupils within schools, or heights of children within families.

So the first problem, of using this model with clustered data, is that we'll actually get the wrong answers, so for the standard errors, our estimates will actually be wrong. And we saw that a bit in the first presentation about different analysis strategies.

The second problem, is that that model doesn't actually show us how much variation is at the school level, and how much is at the pupil level. And so if we fit that model, we won't actually find out how much of an effect school has on exam score after controlling for intake score and since that's what we're interested in, that's a problem.

Solution: Random Intercept Model

Solution: Random Intercept Model

We combine the variance components and the regression models

```

graph TD
    VC["Variance components model  
y_ij = beta_0 + u_j + e_ij  
e_ij ~ N(0, sigma_e^2)  
u_j ~ N(0, sigma_u^2)"]
    SLR["Single level regression model  
y_i = beta_0 + beta_1 x_i + e_i  
e_i ~ N(0, sigma_e^2)"]
    RI["Random intercept model  
y_ij = beta_0 + beta_1 x_ij + u_j + e_ij  
e_ij ~ N(0, sigma_e^2)  
u_j ~ N(0, sigma_u^2)"]

    VC -- "+" --> RI
    SLR -- "+" --> RI
  
```

So what we do is we combine the variance components and the single level regression model and we get a random intercept model. So the random intercept model has got 2 random terms, just like the variance components model so we've got a variance of the level 1 random

term here $e_{ij} \sim N(0, \sigma_e^2)$

$$N(0, \sigma_e^2)$$

...a variance of the level 2 random term here

Fixed and random part

Fixed and random part

The random intercept model has two parts:

- a "fixed part"
- and a "random part"

$$y_{ij} = \underbrace{\beta_0 + \beta_1 x_{ij}}_{\text{fixed part}} + \underbrace{u_j + e_{ij}}_{\text{random part}}$$

Fixed part

- Parameters that we estimate are the coefficients

$$\beta_0, \beta_1, \dots$$

Random part

- Parameters that we estimate are the variances

$$\sigma_u^2 \text{ and } \sigma_e^2$$

- The "random part" is random in the same way that the error term e_i of the single level regression model is random:

- the u_j and e_{ij} are allowed to vary
- some unmeasured processes are generating the u_j and e_{ij}

The random intercept model has two parts. It's got a fixed part (which is the intercept and the coefficient of the explanatory variable times the explanatory variable) and it's got a random part, so that's this $u_j + e_{ij}$ at the end.

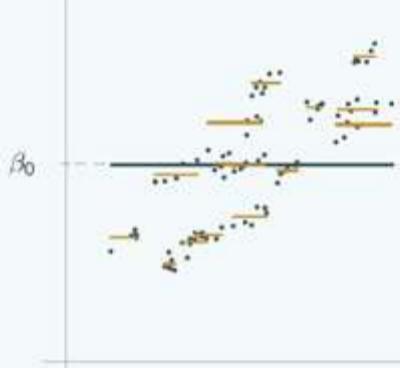
So the parameters that we estimate for the fixed part are the coefficients β_0, β_1 and so on and the parameters that we estimate for the random part are the variances, σ_u^2 and σ_e^2 .

The random part is random in the same way that the error term of the single level regression model is random. All that means is that the u_j and the e_{ij} are allowed to vary so that you can think of it as being that some unmeasured processes are generating the u_j and the e_{ij} .

What does the model look like?

So what does the random intercept model actually look like? Well first of all let's remind ourselves of what the variance components and the single level regression model look like. So the variance components model...

Variance components model



has a line for each group and all of those lines are parallel to the overall average line, and that's just a flat line, because here we plotted our data against our explanatory variable x but x hasn't actually been involved in our model so although you can see the underlying data, you can see that the underlying data have a relationship between x and y . In terms of our model which is the lines that we fitted here, there isn't any relationship between x and y because we didn't include one.

For the single level regression model, we only have one line, just one overall line, but that line isn't just flat, that line is showing the relationship between x and y .

Single level regression model

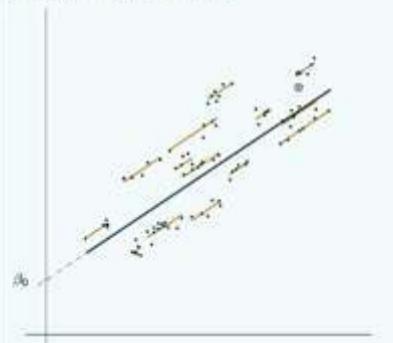


And we can colour in those graphs according to which group the points have come from.

So our random intercept model now:

What does the model look like?

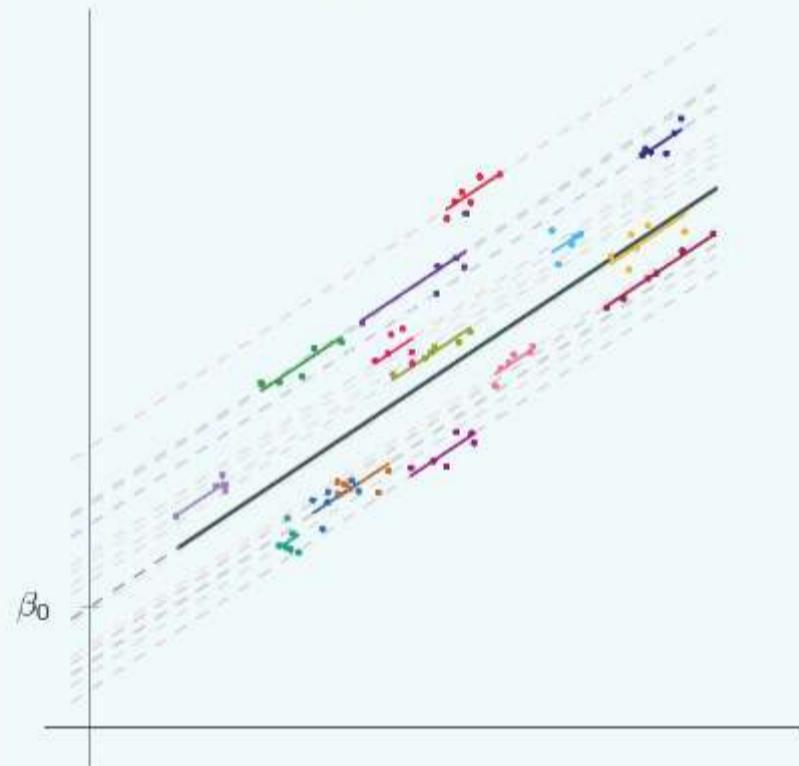
Random intercept model



Well, like the variance components model, our random intercept model has one line for each group, and, again, they're parallel, these lines, to the overall line. Like the single level regression model, the overall line is not just flat, the overall line slopes to fit our data points so it is showing a relationship between x and y and again we can colour in the points according to which group they come from.

What does the model look like?

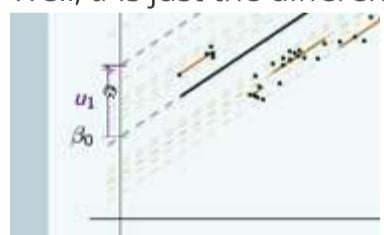
Random intercept model



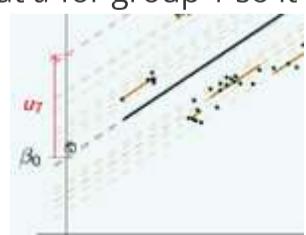
We can extend the group lines because our model equation doesn't specify the length of those lines. In theory they go on forever, but often we won't actually draw them fully extended, we'll only draw them for the range of our data.

So what does our u actually look like for this model?

Well, u is just the difference between the intercept of the overall line and the intercept of the group line so it's just this difference here:

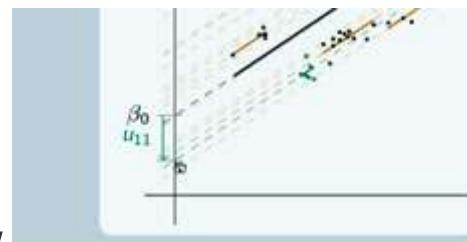


and in this case we're looking at u for group 1 so it's the difference between the intercept of the overall line and the



intercept for group 1. If we look now at group 7, u_7 here

is the difference between the intercept for the overall line, and



And if we look at group 11 now

For this group, the intercept of the group line is below the intercept of the overall line so u_{11} will be negative.

What does the model look like?

Overall line Like the single level regression model, the overall average line has equation $\beta_0 + \beta_1 x_{ij}$	Group lines: Like the variance components model, each group has its own line, parallel to the overall average line
---	--

The 'random intercept'

- For the single level regression model, the intercept is just β_0
- This is a parameter from the fixed part of the model
- For the random intercept model, the intercept for the overall regression line is still β_0
- For each group line the intercept is $\beta_0 + u_j$

Just to recap that, like the single level regression model, the overall line for the random intercept model has the equation $\beta_0 + \beta_1 x_{ij}$ and like the variance components model, each group has its own line, and those lines are parallel to the overall average line.

So what's this random intercept? Why do we call it a random intercept? Well, for the single level regression model, the intercept is just β_0 , and that's a parameter from the fixed part of the model. For the random intercept model, the intercept for the overall regression line is still β_0 but for each group line the intercept is $\beta_0 + u_j$ and you can see that if we go back to the graphs.

What does the model look like?

Random intercept model

So here we have group 1, so if we want to find the intercept for group 1 we take β_0 , the intercept for the overall line, and we add on u_1 , so the intercept is $\beta_0 + u_1$. It's true as well when we have the negative u , so the group 11 - here's the intercept of the group line



and what we do to get there is we take β_0 and we add on this negative value u_{11} so the intercept is still $\beta_0 + u_{11}$.

So this actually involves a parameter from the random part, this u_j , and so that means that this intercept is allowed to vary so we call it a random intercept. It's worth noting as well that the terminology is slightly confusing here. Sometimes we call this, $\beta_0 + u_j$, the random intercept, and sometimes we just call u_j the random intercept so you have to work out from the context which of those two is meant. But often, actually, statements about one of them are true for the other as well so it doesn't make much of a difference.

2) Random intercept models: Research questions and interpretation

What questions can we answer?

We can use the model to answer two kinds of question:

About variables <ul style="list-style-type: none"> ■ e.g. "What is the relationship between exam score at age 11 and exam score at age 16?" ■ We can answer this question because the model has allowed for the clustering in our data ■ This is a question about means ■ It is answered using the fixed part: the slope β_1 of the overall regression line 	About levels <ul style="list-style-type: none"> ■ e.g. "How much variation in pupils' progress between age 11 and 16 is at the school level?" ■ We can answer this question because the model has allowed for differences between schools in intake ■ This is a question about variances ■ It is answered using the random part: the level 2 variance σ_u^2
---	--

[What] questions can we answer using the random intercept model? Well we can use the model to answer two kinds of question.

1. Questions about variables: So for example "What is the relationship between exam score at age 11 and exam score at age 16?" So normally we would fit just the single level regression model but we can't because we have clustering in our data. But the model, the random intercept model, has allowed for the clustering in our data, so we can now answer the question. And it's a question about means, and we answer it using the fixed part, so in this case, the slope β_1 of the overall regression line. So, just in the same way that if we had fitted a single level regression model, we'd use our estimate of β_1 to answer the question.
2. We can also answer questions about levels. So for example "How much variation in pupils' progress between age 11 and 16 is at the school level?" And that's a question that we would answer with a variance components model, but we can't do that because we need to control for the differences in intake. The random intercepts model does control for those differences, and it's a question about variances, and it's answered using the random part, we use the level 2 variance, σ^2_u

So how do we interpret the parameters?

Interpreting the parameters	
Fixed part We can interpret the parameters as for a single level regression model <ul style="list-style-type: none">• β_1 is the increase in the response for a 1 unit increase in x	Random part We can interpret the parameters as for a variance components model Note that again the parameters we estimate are σ^2_u and σ^2_e , not u_j and e_{ij} <ul style="list-style-type: none">• σ^2_u is the unexplained variation at level 2 after controlling for the explanatory variables• σ^2_e is the unexplained variation at level 1 after controlling for the explanatory variables

Well, for the fixed part we can interpret the parameters just the same as for a single level regression model. So β_1 is the increase in the response for a 1 unit increase in x .

For the random part, we interpret the parameters just as for the variance components model, and again note that the parameters that we estimate are σ^2_u and σ^2_e , not u_j and e_{ij} , so we're interpreting the variances, not the individual school effects, just the same as for the variance components model. So σ^2_u is the unexplained variation at level 2 after we control for the explanatory variables. And σ^2_e is the unexplained variation at level 1 after we control for the explanatory variables.

Substantive Interest vs. nuisance

Substantive interest vs. nuisance	
Fixed part of interest <ul style="list-style-type: none">• β_1 may be the only thing of interest• We want to know what the effect of the explanatory variable is on the response• The clustering in our data is a nuisance• It prevents us fitting a single level regression model• σ^2_u is a nuisance parameter	Random part of interest <ul style="list-style-type: none">• σ^2_u and σ^2_e may be the only things of interest• We want to know how much variation is at each level• The relation between the response and the explanatory variable is a nuisance• It prevents us fitting a variance components model• β_1 is a nuisance parameter <ul style="list-style-type: none">• Sometimes both β_1 and σ^2_u & σ^2_e may be of interest• Which part is of interest all depends on our research question

Sometimes, when we come to interpret the estimates for our random intercepts model, we may only be interested in the fixed part of the model. So β_1 , for example, may be the only thing of interest. And that's when we want to know what the effect of the explanatory variable is on the response and the clustering in our data is just a nuisance which prevents us from fitting a single-level regression model. And so in that case we can look at σ^2_u as being just a nuisance parameter. So we don't look at our estimate for σ^2_u and interpret it.

On the other hand, sometimes only the random part is of interest, so in that case σ^2_u and σ^2_e are the only things that we want to look at. And that's when we want to know how much variation is at each level and in that case the relation between the response and the explanatory variable is a nuisance. So that's like our example - when we wanted to control for intake score when looking at school exam results. So this relation between the response and the explanatory variable prevents us from fitting a variance components model. So, in this case, we might look at β_1 as being a nuisance parameter.

Sometimes, both β_1 and σ^2_u and σ^2_e may be of interest. And which part is of interest, all depends on our research question.

So here are some -

Examples of research questions

Muijs (2003)

Do pupils make more progress in maths when receiving support from numeracy support assistants?

Levels: 2 school Answer: No
1 pupil

So this study [Muijs, 2003] looked at whether pupils make more progress in maths when they receive support from numeracy support assistants. So they had school at level 2 and pupil at level 1, and of course they expected that pupils within the same school would be more similar than pupils from different schools. So they had to fit the random intercept model, they couldn't fit a single level regression model, and their answer to that question was: No, pupils don't make any more progress in maths when they receive support from numeracy support assistants.

Southwell (2005)

Does memory of health campaign advertisements increase as extent of past drug use increases?

Levels: 2 advert Answer: No
1 individual

This study [Southwell, 2005] looked at whether memory of health campaign advertisements increases as extent of past drug use increases, and they were actually interested in whether the impact of the advertisements increases, and memory was what they were using to measure impact. They had advert at level 2 and individual at level 1 because each advert was seen by more than one person, and of course they expected memory of the same advert to be more similar than memory of different adverts, and their answer to the question, again, was: No, the memory of health campaign advertisements doesn't increase as the extent of past drug use increases. So they interpreted that as meaning, for individuals with more past drug use, an advertisement doesn't have any more impact than for individuals with less past drug use.

Judge et al. (2006)

Is job satisfaction negatively related to workplace deviance (misbehaviour)?

Levels: 2 individual Answer: Yes
1 occasion

This study [Judge et al., 2006] looked at whether job satisfaction is negatively related to something they called 'workplace deviance' which is misbehaviour at work - and they were looking at individuals over time. Each individual was measured as to their workplace deviance and their job satisfaction on several days, and of course they expected that for the same person, both the job satisfaction and the misbehaviour would be more similar than for a different person so they had to fit a random intercept model, they couldn't just fit a single-level regression model. And their answer was: Yes, job satisfaction is negatively related to workplace deviance, so people who are less satisfied at work misbehave more, or, people who [mis]behave less at work are more satisfied in their job.

Buckley et al. (2008)

Is flowering stem length linked to vegetative stem length in St John's wort?

Levels: 2 quadrat Answer: Yes
1 plant

This study [Buckley et al., 2008] looked at whether flowering stem length was related to vegetative stem length in St John's Wort. They had quadrat at level 2 and plant at level 1. A quadrat is a sort of sampling frame, maybe a metre square and they take every plant which falls inside that sampling frame. So obviously you'd expect plants in the same sampling frame to be more similar than plants in different sampling frames. That's partly because the conditions would be more similar within that one square metre, it's more likely to have similar levels of shade and similar levels of rainfall and so on. But it's also actually because St John's Wort reproduces asexually as well as sexually, so they didn't know whether plants within the same quadrat were actually clones of each other so they had to allow for the clustering at the

quadrat level, partly due to similar conditions but also due to probably similar genes in their plants. And their answer was: Yes, that flowering stem length is linked to vegetative stem length.

Goldstein et al. (2007)

Does which school a pupil attends affect their progress in maths between KS1 and KS2?

Levels: 2 school Answer: Yes
1 pupil

And finally, all these other 4 studies have been looking at the fixed part of the model so that for them the clustering has been a nuisance and they would have liked to have fitted a single level regression model. This study [Goldstein et al., 2007] is an example of a study for which the relationship between the explanatory variable and the response is a nuisance and they fit the random intercept model to allow for that relationship. Again, it's a case of pupils in schools and basically this is a paper giving the example that we had earlier with the response being exam score, this Key Stage 2 test, and the explanatory variable being the previous exam score, Key Stage 1. So, in that situation where the response is exam score and the explanatory variable is previous exam score, actually they're looking at progress and so they're looking at the variation in progress. So they have schools at level 2, pupils at level 1, and the answer is: Yes, which school a pupil attends does affect their progress in maths.

And here are the details of all those studies:

- Buckley, Y., Briese, D. and Rees, M. (2008) *Demography and management of the invasive plant species Hypericum perforatum. I. Using multi-level mixed-effects models for characterizing growth, survival and fecundity in a long-term data set* Journal of Applied Ecology, 40 pp 481 - 493
- Goldstein, H., Burgess, S. and McConnell, B. (2007) *Modelling the effect of pupil mobility on school differences in educational achievement* Journal of the Royal Statistical Society Series A, 170 4 pp 941 - 954
- Judge, T., Scott, B. and Ilies, R. (2006) Hostility, job attitudes and workplace deviance: Test of a multilevel model Journal of Applied Psychology 91 1 pp 126 - 138
- Muijs, D. (2003) *The effectiveness of learning support assistants in improving the mathematics achievement of low achieving pupils in primary school* Educational Research 45:3, pp 219 - 230
- Southwell, B. (2005) *Between messages and people: A multilevel model of memory for television content* Communication Research 32 1 112 - 140

3) Random intercept models: Adding more explanatory variables

Adding more explanatory variables

- We can easily add in more explanatory variables, as for a single level regression model
$$y_{ij} = \beta_0 + \beta_1 x_{1ij} + \beta_2 x_{2ij} + \beta_3 x_{3ij} + \beta_4 x_{2ij}x_{3ij} + \dots + u_j + e_{ij}$$
- Notice that we can also include interactions in the usual way
- Variables can be continuous or categorical
- Variables can be defined at a higher level (e.g. school mean intake score)
 - Mathematically they're the same
 - but the interpretation is different
 - so we'll leave them for now

Level 2 variance can increase

- When we add in a (level 1) variable, the variation at level 2 may decrease or increase (or stay the same)
- However, the level 1 variation and the total residual variation will both decrease (or stay the same)

So far we've looked at examples of random intercept models with only one explanatory variable but in fact we can easily add in more explanatory variables, just in the same way as for a single level regression model. So, we've got an example of that here where we've got another couple of explanatory variables added in. And interactions, we can also add in, in exactly the same way as for a single level model. Variables can be continuous or categorical. Of course, if we have a categorical variable we enter that as a set of dummy variables, just as for a single level regression model. And variables can be defined at a higher level. So, for example, school mean intake score, we say that's defined at a higher level because the school mean intake score will take the same value for every pupil within that school. Pupils in the same school will have the same value of this variable. Now mathematically, that's exactly the same. It goes into the model exactly the same and it goes into MLwiN in exactly the same way but the interpretation of these variables is different and so we're going to leave those for now and only look at variables defined at level 1.

Now when we put in explanatory variables, the level 2 variance can increase. Now we're used to the idea with single level regression models that when you put in an explanatory variable the variance goes down. For a random intercept model it's true that the level 1 variation and the total residual variation will both decrease when we put in an explanatory variable. But the level 2 variance can increase, and that's true

not only when we add in a second explanatory variable, but actually also when we add in the first explanatory variable, to go from the variance components model to the random intercepts model.

Example

Adding more explanatory variables

Example

This also applies when we add a variable to a variance components model to get a random intercept model. An example is in modelling house prices (with area as level 2 and house as level 1).

- We start with a variance components model
- We add house size as an explanatory variable
- More expensive areas tend to have smaller houses
- So before we add size, house prices appear more similar across areas
- The variance components model has less variation at level 2 than the random intercept model

Level 2 variance can increase

- When we add in a (level 1) variable, the variation at level 2 may decrease or increase (or stay the same)
- However, the level 1 variation and the total residual variation will both decrease (or stay the same)

So as an example of that let's look at modelling house prices. So in this case, area is level 2 and house is level 1. So we start with a variance components model and then we add in house size as an explanatory variable. And actually more expensive areas tend to have smaller houses. For example, in London the expensive areas tend to be in the centre of the town where there's not so much room for larger houses. Because the more expensive areas have smaller houses, that means that in the variance components model, before we add in house size, the fact that the houses are smaller is bringing down the average price for those more expensive areas so they don't appear as much more expensive. And similarly for the less expensive areas, which tend to have larger houses, before we add in house size, the average price for that area is brought up by the fact that the houses in that area are larger. So that means that before we add size, the house prices will appear more similar across areas and after we add size, the average price for the expensive areas is no longer being brought down because we've adjusted for the size and the house price for the less expensive areas is no longer being brought up, because we've adjusted for the size.

So the variance components model will have less variation at level 2 because the areas appear more similar than the random intercept model.

When is a variable a level?

When is a variable a level?

Not an easy question...

- Sometimes it's very clear-cut
 - we would not consider using a variable as a level
 - or else we would definitely want to put it in as a level
- In other cases, it's less obvious what we should do
 - often when we have a relatively small number of units
- Basically, two things govern the decision:
 - number of units and
 - exchangeability
- The units should be a representative draw from a population
- Exchangeability is a tricky concept, so it helps to draw up some guidelines

Exchangeability

The units are exchangeable if we could randomly reassign their codes without losing any information

So when we come to adding in more explanatory variables, we may wonder: When should we add in an explanatory variable as an explanatory variable and when instead should we put it in as a level? Because, for example, when we have our example of pupils within schools, we could have put in school as a fixed explanatory variable, and allowed each school to have a different effect in the fixed part of the model but we decided not to do that, we decided to put school in as a level, but, you know, the question comes up with many other explanatory variables: Should we put it in as an explanatory variable or should we put it in as a level? Well, it's not an easy question. Sometimes, it is very clear cut. Sometimes we wouldn't consider using a variable as a level, we would definitely put it in as an explanatory variable, and sometimes we definitely would want to put it in as a level, we definitely would not put it in as an explanatory variable. But often, it's less obvious what we should do, and it's often when we have a relatively small number of units, which would be at that level if we put it in, when it's less obvious what we should do. So, for example, if we are considering: Should we put in country as a level or as an explanatory variable and we might only have a few countries, so is that enough? Basically two things govern our decision about what to do, one is the number of units and the other is exchangeability.

Exchangeability: The units are exchangeable if we could randomly reassign their codes without losing any information. So that's quite a complicated concept to understand, so we're going to give some guidelines.

The units basically should be a representative draw from the population, so that's why these 2 things - number of units and exchangeability, are important, because obviously if we don't have enough units, then they will not be a representative draw, and if we don't have

exchangeability, that means that we're not drawing from a population. And, as we've said, exchangeability is a tricky concept, so we're going to give some guidelines.

Guidelines

Guidelines	
Probably a level <ul style="list-style-type: none">■ (If you had never heard of multilevel modelling) you might use it as the units of a single level regression model■ Nominal variable whose categories have no special meanings■ Before we run the analysis, we would not predict different results for any particular unit■ Large number of categories■ It's a particular kind of variable that we know can be used as a level	Probably a variable <ul style="list-style-type: none">■ (Even if you had never heard of multilevel modelling) you wouldn't use it as the units of a single level regression model■ Ordinal or continuous variable■ Categories have special meanings■ Before we run the analysis, we would predict different results for some categories■ Small number of categories (ok for lower level)

These basically take into account both those points about number of units and about exchangeability. So a kind of overall guideline is: If you had never heard of multilevel modelling, might you use this variable as the units of a single level regression model? So in our example of schools, if you had never heard of multilevel modelling, you might decide that the best way to analyse your data for pupils within schools would be to take averages for each school of both your response and your explanatory variables, and then have those averages as your data points and fit a single level regression model to those. Of course now we know that this is not a sensible thing to do because we have been learning about multilevel modelling and we know that that will give you the wrong standard errors. But since, if you hadn't heard of multilevel modelling you might do that, that implies that school is probably a level, you probably can use it as the level 2 units in a multilevel model. On the other hand, if you ask that question for something like gender, then, even if you had never heard of multilevel modelling, you would not take the average for men and the average for women and then fit a single level regression model to those 2 data points. So that suggests that, for gender, we don't want to put that in as a level, we should put it in as a variable.

Here are some other guidelines which actually explain why we might answer either way to these questions, that kind of break this down a bit.

If we have a nominal variable whose categories have no special meanings, so that's like school for example. School 63 doesn't have a special meaning compared to school 111, it really doesn't matter which school gets which of those categories. So in that case, the variable can probably be put in as a level. But on the other hand, if we have an ordinal or continuous variable, so something like income bands, or just income measured as a continuous variable, then we'd probably have to put that in as a variable, not as a level. Or if our categories have special meanings, like with gender, male and female have special meanings and if we swap those labels around, then, when we try to interpret our estimates, we'll be looking at the wrong thing. In that case, the categories do have special meanings and so if that's true, we probably don't want to put that in as a level, we want to put it in as a variable.

Also, if, before we run the analysis, we wouldn't predict different results for any particular unit, then we can probably put it in as a level. So with the schools for example, before we run the analysis we have no reason to believe that any particular school will have a greater effect than any other. If we did, then that would mean that we would probably have to take that school out and put it in the fixed part of the model. If, before we run the analysis, we would predict different results for some categories, so like with gender, we'd predict different results for men and women, and we can probably guess which way round, whether we'll have a greater estimate for women or for men, we probably can't put that in as a level, we probably have to put that in as a variable.

And then there's the number of categories, so for school we have a large number of categories so it's probably OK to put that in as a level, but for gender, we have a small number of categories, just two categories, so we probably can't put that in - at level 2. It is important to note that the number of categories doesn't matter for lower levels, it's only the highest level, that the number of categories matters in this decision.

And finally, if it's a particular kind of variable that we know can be used as a level, then, even if it appears to violate some of these conditions, then it's still OK. Examples would be repeated measures, so when we've got several time occasions, that would be an example of a variable that is ordinal, but we can find out from textbooks and courses that we can use that variable as a level so we can put that in.

So, some **examples** to understand how to do this:

Example

Hospital is probably a level If we have data on treatment outcomes for patients in 100 UK hospitals.	Ethnicity is not a level If we have data on height for people of a variety of ethnicities:
<ul style="list-style-type: none">■ You might naively use hospital as the units in a single level model (if you didn't know about multilevel modelling)■ Hospital ID is a nominal variable whose categories have no special meanings■ We would not expect any particular hospital to give different results beforehand	<ul style="list-style-type: none">■ Even if we could measure ethnicity very finely using many categories, we would not use it as the units in a single level model■ The categories of ethnicity have a special meaning■ We expect different results for different categories■ We probably don't have many different categories in any case

Well, hospital is probably a level. Suppose that we have data on treatment outcomes for patients in 100 UK hospitals. Now, if we didn't know about multilevel modelling we might naively use the hospitals as units in a single level model. We might take the average for each hospital of the response and the explanatory variables and fit a single level model to those data points. Also, hospital ID is a nominal variable whose categories have no special meanings. Hospital 28 is no different from hospital 163, as a category name, it doesn't mean anything special, and we wouldn't expect any particular hospital to give different results, before we run the analysis.

Ethnicity, on the other hand, is not a level, if we have data on height for people of a variety of ethnicities. Even if we could measure ethnicity very finely, using many categories, even if we could measure ethnicity using 50 categories, say, we still wouldn't use it as the units in a single level model. We wouldn't take the average for each ethnicity and fit a single level regression model to those points and that's because the categories of ethnicity have a special meaning. Black has a special meaning compared to Asian, and we do expect different results for different categories. In the case of height we might expect Asian (and if we have sub-categories of Asian those sub-categories) to have smaller values of height than some of the other ethnicities. And in any case, the idea that we might be able to measure ethnicity by 50 categories is probably unrealistic. Probably we don't have very many ethnicity categories, we may only have 6 categories, say. So for these reasons we can't put ethnicity in as a level.

4) Random intercept models: Hypothesis Testing

Hypothesis testing	
<ul style="list-style-type: none">■ Hypothesis testing is an important part of interpretation■ We don't only want to know the size of the fixed effects and the amount of variance at each level■ We also want to know whether the fixed effects are significant■ and whether there is a significant amount of variance at level 2 <p>Fixed part</p> <ul style="list-style-type: none">■ Divide the coefficient by its standard error to get $z = \frac{\beta_1}{\text{s.e.}(\beta_1)}$ <ul style="list-style-type: none">■ If $z \geq 1.96$ (or informally if $z \geq 2$), then β_1 is significant at the 5% level	<p>Random part</p> <ul style="list-style-type: none">■ We CAN'T just divide σ_u^2 by $\text{s.e.}(\sigma_u^2)$ and compare the modulus with 1.96■ Instead we have to fit the model with and without u_j and do a likelihood ratio test to see whether σ_u^2 is significant

[Hypothesis] testing is of course an important part of interpretation because we don't only want to know the size of the fixed effects and the amount of variance at each level, we also want to know whether the fixed effects are significant and whether there's a significant amount of variance at level 2.

For the fixed part, hypothesis testing is just the same as for a single level model. We just divide the coefficient by its standard error, to get z, and then we take the modulus of z, and if that's bigger than 1.96 (or informally we can use 2) then β_1 is significant at the 5% level.

For the random part we can't just divide σ_u^2 by its standard error and compare the modulus with 1.96. Instead, we have to fit the model with and without u_j and do a likelihood ratio test comparing those 2 models, to see whether σ_u^2 is significant.

Likelihood ratio test

Likelihood ratio test

- We fit $y_{ij} = \beta_0 + \beta_1 x_{1ij} + \beta_2 x_{2ij} + \dots + u_j + e_{ij}$ ①
and $y_{ij} = \beta_0 + \beta_1 x_{1ij} + \beta_2 x_{2ij} + \dots + e_{ij}$ ②
and note the likelihoods
- The test statistic is $2(\log(\text{likelihood}(①)) - \log(\text{likelihood}(②)))$
 - MLwiN gives $-2 \times \log(\text{likelihood})$ in the **Equations** window
 - So we just take $(\text{MLwiN's value for } ②) - (\text{MLwiN's value for } ①)$
- The null hypothesis is that $\sigma_u^2 = 0$ and so we don't need u_j in the model
- We compare the test statistic to the $\chi^2_{(1)}$ distribution, then divide the corresponding *p*-value by 2 (since $\sigma_u^2 \geq 0$)
- The degrees of freedom are 1 because there is one more parameter, σ_u^2 , in ① compared to ②

So here are the two models that we fit, This one has u_j :

$$y_{ij} = \beta_0 + \beta_1 x_{1ij} + \beta_2 x_{2ij} + \dots + u_j + e_{ij} \quad ①$$

and this one doesn't have u_j :

$$y_{ij} = \beta_0 + \beta_1 x_{1ij} + \beta_2 x_{2ij} + \dots + e_{ij} \quad ②$$

So this is our random intercept model:

$$y_{ij} = \beta_0 + \beta_1 x_{1ij} + \beta_2 x_{2ij} + \dots + u_j + e_{ij} \quad ①$$

and this is a single level regression model:

$$y_{ij} = \beta_0 + \beta_1 x_{1ij} + \beta_2 x_{2ij} + \dots + e_{ij} \quad ②$$

which is exactly the same in terms of the explanatory variables. The only difference is that it doesn't have u_j . So we fit those and we note the likelihoods. The test statistic is 2 times the log of the likelihood of the random intercept model minus the log of the likelihood of the single level regression model:

$$2(\log(\text{likelihood}(①)) - \log(\text{likelihood}(②)))$$

And MLwiN, in fact, gives minus 2 times the log likelihood in the Equations window so if we're using MLwiN, then we can just take MLwiN's value for the single level model and subtract MLwiN's value for the random intercepts model:

- MLwiN gives $-2 \times \log(\text{likelihood})$ in the **Equations** window
- So we just take $(\text{MLwiN's value for } ②) - (\text{MLwiN's value for } ①)$

The null hypothesis is that $\sigma_u^2 = 0$. and if that's true then we don't need u_j in the model so we can just fit the single level regression model .

So we compare the test statistic to the $\chi^2_{(1)}$ distribution with 1 degree of freedom. And then we can divide the corresponding *p*-value by 2 since σ_u^2 has to be greater than or equal to 0.

The null hypothesis is that $\sigma_u^2 = 0$ and so we don't need u_j in the model

We compare the test statistic to the $\chi^2_{(1)}$ distribution, then divide the corresponding *p*-value by 2 (since $\sigma_u^2 \geq 0$)

There is one degree of freedom because there is one more parameter, σ_u^2 , in the random intercept model, compared to the single level regression model.

The degrees of freedom are 1 because there is one more parameter, σ_u^2 , in ① compared to ②

5) Random intercept models: Variance partitioning coefficients

Variance partitioning coefficients	
<ul style="list-style-type: none"> For variance components models, we saw that the VPC is a useful way to see how the variance divides up This is even more true for random intercept models, since the total amount of variance may change as we add explanatory variables, making comparison hard <p>Note We most often use 'Level 1 variance' to mean 'Residual variance at level 1' – not 'Variance of y at level 1'. Similarly for 'Level 2 variance'.</p>	<p>Calculating the VPC</p> <ul style="list-style-type: none"> Recall that $\rho = \frac{\text{Level 2 variance}}{\text{Total residual variance}}$ <ul style="list-style-type: none"> For a random intercept model, <ul style="list-style-type: none"> Level 1 variance = σ_e^2 Level 2 variance = σ_u^2 Total residual variance = $\sigma_e^2 + \sigma_u^2$ So $\rho = \frac{\sigma_u^2}{\sigma_u^2 + \sigma_e^2}$ This is just the same as for the variance components model

So when we were looking at variance components models, we found that the variance partitioning coefficient is a useful way to see how the variance divides up between levels. And it's actually even more true for random intercept models because, as we saw earlier, the total amount of variance may change as we add explanatory variables, which can make comparison hard.

It's important to note, in what we're about to say, that we most often use level 1 variance to mean residual variance at level 1, not the variance of y at level 1 and the same thing for level 2 variance.

So how do we calculate the variance partitioning coefficient? Well, remember that the variance partitioning coefficient is the level 2 variance divided by the total residual variance. And for a random intercept model, our level 1 variance is σ_e^2 , our level 2 variance is σ_u^2

and the total residual variance is $\sigma_e^2 + \sigma_u^2$. So our variance partitioning coefficient is σ_u^2 over $\sigma_u^2 + \sigma_e^2$ and that's just exactly the same as for the variance components model.

ρ and clustering

ρ and clustering	
<p>Another way to think of ρ is that it measures the clustering</p> <p>Large ρ</p> <ul style="list-style-type: none"> When ρ is large, a lot of the variance is at level 2 so units within each group are quite similar but there is a lot of difference between groups Values of the response are largely determined by which group the unit belongs to. So the data are very clustered <p>Large $\rho \Rightarrow$ a lot of clustering</p>	<p>Small ρ</p> <ul style="list-style-type: none"> When ρ is small, not much variance is at level 2 so units within each group may be quite dissimilar but there is not much difference between groups Which group a unit belongs to does not have much impact on the response So the data are not very clustered <p>Small $\rho \Rightarrow$ little clustering</p>

Another way to think of ρ is that it measures the clustering.

If we have a large value of ρ then a lot of the variance is at level 2 and that means that units within each group are quite similar but that there's a lot of difference between groups and so values of the response are largely determined by which group the unit belongs to, and so the data are very clustered.

If we have a small value of ρ then not much variance is at level 2, and so units within each group may be quite dissimilar but there's not much difference between groups. So which group a unit belongs to, doesn't have much of an impact on the response and the data are not very clustered.

So here's an example of a graph for a particular value of ρ :

A small value of ρ



Is this a graph where we have a large value of ρ or where we have a small value of ρ ? Well, this is a small value of ρ because in this case, there's not much variation in the group lines, but there's quite a bit of variation in the individual data points around their group lines.

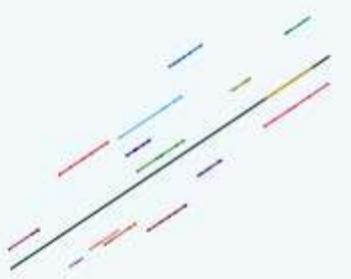
Here's an example of a different value of ρ :

A large value of ρ



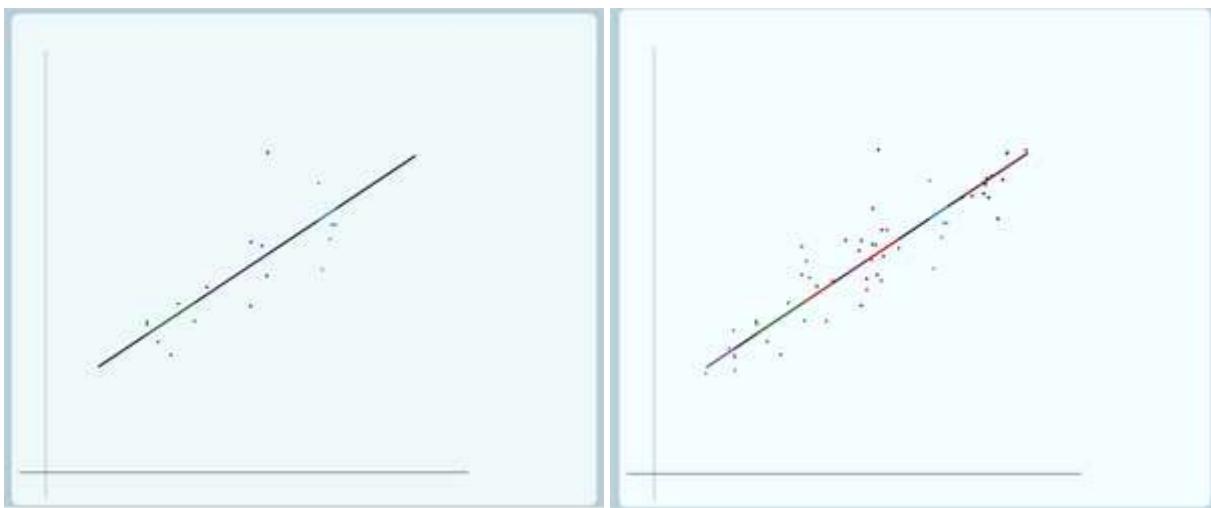
Is this a small value of ρ or a large value of ρ ? Well, this is a large value of ρ because in this case there's a lot of variation of the group lines around the overall line but not much variation of the individual data points around their group lines.

So in theory, ρ can actually be as large as 1 or as small as 0 - so what do those values actually look like? Well here's one of them: Is this $\rho = 1$ or $\rho = 0$?



Well this is $\rho = 1$. In this case, there's no variation at level 1. There's variation at level 2 because the group lines vary round the overall regression line but there's no variation at level 1 because the individual data points don't vary around their group lines. So in this case, we have complete dependency. If we know the value of the explanatory variable and we know which group a data point belongs to, then we know the value of the response for that data point. And obviously, that doesn't really tend to happen in practice in the social sciences.

Now if we want to look at $\rho = 0$ (we'll build it up gradually so it's easier to see):



In this case, there's no variation at level 2 so the group lines don't vary around the overall regression line, all of the group lines lie on the overall regression line. But the individual data points do still vary around their group line, still have level 1 variance.

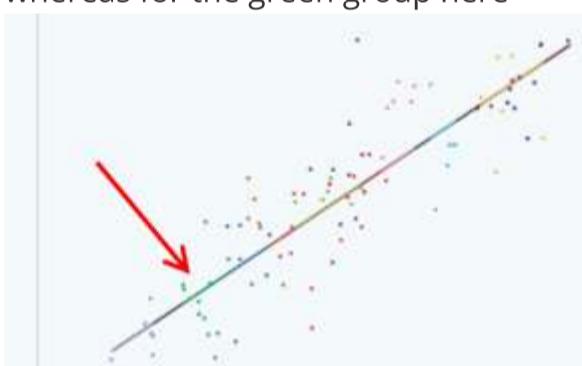
Interpreting the value of ρ

If we look at the formula for ρ , as we said, we can see that in theory the smallest it can be is 0. So in this case when ρ is 0 we don't have any dependency. If we know the value of the explanatory variable for a data point then knowing which group it comes from doesn't tell us anything more about its value of the response. In this particular example, we do have dependency in terms of x . You can see that for the yellow group here,



all of the data points have quite high values of x

whereas for the green group here



all of the data points have quite low values of x

So in terms of x there is dependency. But we're interested in whether there's dependency in y after controlling for x and there isn't any dependency there. So in this case, we could actually use a single level regression model.

Interpreting the value of ρ	
Theoretical limits for ρ:	<ul style="list-style-type: none"> • Looking at the formula for ρ, we can see that in theory the smallest it can be is 0 and the largest it can be is 1 <ul style="list-style-type: none"> * 1 indicates maximum clustering * 0 indicates no clustering (single level data structure)
In practice:	<ul style="list-style-type: none"> • We never expect to see a value of 0 or 1 for ρ • If ρ is small enough we can use a single level model • But we don't make that decision by looking at the value of ρ; we use the likelihood ratio test described earlier
What is a large value for ρ?	<ul style="list-style-type: none"> • It depends on the subject area and what the units of each level are. <ul style="list-style-type: none"> • We expect more clustering for observations on occasions within individuals than observations on people within families • and more clustering for observations on people within families than pupils within schools, for example

So as we just said, if we look at the formula for ρ we can see that, in theory, the smallest ρ can be is 0 and the largest it can be is 1 and 1 indicates maximum clustering and 0 indicates no clustering. So it's basically a single level data structure. In practice, though, we never expect to see a value of 0 or 1 for ρ . If ρ is small enough, we can use a single level model but ρ doesn't have to be 0, just small enough. But we don't make that decision by looking at the value of ρ , we do that by testing the significance of σ^2_u in our random intercepts model using the likelihood ratio test which we described in the section on [hypothesis testing](#).

So, OK, in practice we're not going to see a value of 1 for ρ . What is a large value? Well, that actually depends on the subject area, and on what the units of each level are. If we have observations on occasions within individuals we expect quite a lot of clustering but we expect less

clustering for observations on people within families and we expect less clustering again for observations on pupils within schools. And for a more detailed presentation of ρ and dependency you can see the audio presentation [Measuring Dependency](#) which is also on our web site.

Clustering in the model

Clustering in the model	
Clustering We've talked a lot about clustering: <ul style="list-style-type: none"> ■ Clustering is the reason we can't use a single level regression model ■ Clustering is why we have variance at both level 1 and level 2- why our response is determined at two levels ■ We've seen how to measure this clustering ■ and how to interpret it 	Incorporating the clustering <ul style="list-style-type: none"> ■ We haven't seen yet how the clustering is incorporated into the model: ■ how does the random intercepts model allow for similarities between different observations from the same group? ■ To discover this, we need to look at the correlation matrix V. ■ And to do that we need to return to the technicalities of the model

So we've talked a lot about clustering in the previous presentations. Clustering was the reason why we can't use a single level regression model. It's why we have variance at both level 1 and level 2 - so why our response is actually determined at two levels. We've seen how to measure the clustering and how to interpret it.

But what we haven't seen yet, is how the clustering is incorporated into the model, why it is that the random intercept model can allow for the clustering but the single level model can't. So how does the random intercepts model actually allow for similarities between different observations from the same group? Well, to discover this we need to look at the correlation matrix, V and in order to do that we're going to need to return to the technicalities of the model.

6) Random intercept models: the correlation matrix

Assumptions of the random part

Assumptions of random part	
Let's first recall the assumptions for a single level model and for a variance components model: Single level model $y_i = \beta_0 + \beta_1 x_i + e_i \quad e_i \sim N(0, \sigma_e^2)$ $\text{Cov}(e_i, x_i) = 0 \quad \text{Cov}(e_i, e_j) = 0$	
Variance components model $y_{ij} = \beta_0 + u_j + e_{ij} \quad u_j \sim N(0, \sigma_u^2)$ $e_{ij} \sim N(0, \sigma_e^2)$ $\text{Cov}(u_{j1}, u_{j2}) = 0 \quad \text{Cov}(u_{j1}, e_{i1j}) = 0 \quad \text{Cov}(e_{i1j}, e_{i2j}) = 0$ $\text{Cov}(u_{j1}, e_{i2j}) = 0 \quad \text{Cov}(e_{i1j}, e_{i2j}) = 0$	

So, first of all, let's recall the assumptions for a single level model and for a variance components model.

For the single level model we assume that the residual and the covariates are uncorrelated and that 2 different residuals are uncorrelated.

For the variance components model we assume that the level 2 residuals for 2 different groups are uncorrelated, that the level 1 and the level 2 residuals are uncorrelated, whether from the same or a different group and that 2 different level 1 residuals are uncorrelated, again whether from the same or from a different group.

Assumptions of random part	
Assumptions for a random intercept model are a mixture: Random intercept model $y_{ij} = \beta_0 + \beta_1 x_{ij} + u_j + e_{ij} \quad u_j \sim N(0, \sigma_u^2)$ $e_{ij} \sim N(0, \sigma_e^2)$ $\text{Cov}(u_{j1}, u_{j2}) = 0 \quad \text{Cov}(u_{j1}, e_{i1j}) = 0 \quad \text{Cov}(e_{i1j}, e_{i2j}) = 0$ $\text{Cov}(u_{j1}, e_{i2j}) = 0 \quad \text{Cov}(e_{i1j}, e_{i2j}) = 0$ $\text{Cov}(u_j, x_{ij}) = 0 \quad \text{Cov}(e_{ij}, x_{ij}) = 0$ <ul style="list-style-type: none"> ■ Level 2 residuals for different groups are uncorrelated ■ Level 1 residuals for different observations are uncorrelated ■ Level 2 and level 1 residuals are uncorrelated ■ Residuals and covariates are uncorrelated 	

For the random intercept model, we have a mixture of assumptions from the single level and from the variance components model. We again assume that the level 2 residuals for different groups are uncorrelated and we assume that the level 2 and the level 1 residuals are uncorrelated. We assume that the level 1 residuals for different observations are uncorrelated and we assume that residuals and covariates are uncorrelated.

V, the correlation matrix

V, the correlation matrix				
The correlation matrix gives the correlation between every pair of level 1 units in our dataset after controlling for the explanatory variables				
Single level model				
1	1	2	3	...
2	$\text{Cor}(y_1 - \hat{y}_1, y_2 - \hat{y}_1)$	$\text{Cor}(y_2 - \hat{y}_2, y_1 - \hat{y}_1)$	$\text{Cor}(y_3 - \hat{y}_3, y_1 - \hat{y}_1)$...
3	$\text{Cor}(y_1 - \hat{y}_1, y_3 - \hat{y}_1)$	$\text{Cor}(y_2 - \hat{y}_2, y_3 - \hat{y}_1)$	$\text{Cor}(y_3 - \hat{y}_3, y_3 - \hat{y}_1)$...
4	$\text{Cor}(y_1 - \hat{y}_1, y_4 - \hat{y}_1)$	$\text{Cor}(y_2 - \hat{y}_2, y_4 - \hat{y}_1)$	$\text{Cor}(y_3 - \hat{y}_3, y_4 - \hat{y}_1)$...
5	$\text{Cor}(y_1 - \hat{y}_1, y_5 - \hat{y}_1)$	$\text{Cor}(y_2 - \hat{y}_2, y_5 - \hat{y}_1)$	$\text{Cor}(y_3 - \hat{y}_3, y_5 - \hat{y}_1)$...
6	$\text{Cor}(y_1 - \hat{y}_1, y_6 - \hat{y}_1)$	$\text{Cor}(y_2 - \hat{y}_2, y_6 - \hat{y}_1)$	$\text{Cor}(y_3 - \hat{y}_3, y_6 - \hat{y}_1)$...
7	$\text{Cor}(y_1 - \hat{y}_1, y_7 - \hat{y}_1)$	$\text{Cor}(y_2 - \hat{y}_2, y_7 - \hat{y}_1)$	$\text{Cor}(y_3 - \hat{y}_3, y_7 - \hat{y}_1)$...
8	$\text{Cor}(y_1 - \hat{y}_1, y_8 - \hat{y}_1)$	$\text{Cor}(y_2 - \hat{y}_2, y_8 - \hat{y}_1)$	$\text{Cor}(y_3 - \hat{y}_3, y_8 - \hat{y}_1)$...
9	$\text{Cor}(y_1 - \hat{y}_1, y_9 - \hat{y}_1)$	$\text{Cor}(y_2 - \hat{y}_2, y_9 - \hat{y}_1)$	$\text{Cor}(y_3 - \hat{y}_3, y_9 - \hat{y}_1)$...
...

The correlation matrix gives the correlation between every pair of level 1 units in our dataset after controlling for the explanatory variables. The columns - each column is one individual in our dataset and each row is also one individual in our dataset.

So here's the correlation of individual 1 with themselves:

Single level model	
	1
1	$\text{Cor}(y_1 - \hat{y}_1, y_1 - \hat{y}_1)$

Here's the correlation between individual 1 and individual 2:

Single level model	
	1
1	$\text{Cor}(y_1 - \hat{y}_1, y_2 - \hat{y}_1)$
2	$\text{Cor}(y_1 - \hat{y}_1, y_2 - \hat{y}_2)$

Here's the correlation between individual 2 and individual 9:

	1	2
1	$\text{Cor}(y_1 - \hat{y}_1, y_1 - \hat{y}_1)$	$\text{Cor}(y_2 - \hat{y}_2, y_1 - \hat{y}_1)$
2	$\text{Cor}(y_1 - \hat{y}_1, y_2 - \hat{y}_2)$	$\text{Cor}(y_2 - \hat{y}_2, y_2 - \hat{y}_2)$
3	$\text{Cor}(y_1 - \hat{y}_1, y_3 - \hat{y}_1)$	$\text{Cor}(y_2 - \hat{y}_2, y_3 - \hat{y}_1)$
4	$\text{Cor}(y_1 - \hat{y}_1, y_4 - \hat{y}_1)$	$\text{Cor}(y_2 - \hat{y}_2, y_4 - \hat{y}_1)$
5	$\text{Cor}(y_1 - \hat{y}_1, y_5 - \hat{y}_1)$	$\text{Cor}(y_2 - \hat{y}_2, y_5 - \hat{y}_1)$
6	$\text{Cor}(y_1 - \hat{y}_1, y_6 - \hat{y}_1)$	$\text{Cor}(y_2 - \hat{y}_2, y_6 - \hat{y}_1)$
7	$\text{Cor}(y_1 - \hat{y}_1, y_7 - \hat{y}_1)$	$\text{Cor}(y_2 - \hat{y}_2, y_7 - \hat{y}_1)$
8	$\text{Cor}(y_1 - \hat{y}_1, y_8 - \hat{y}_1)$	$\text{Cor}(y_2 - \hat{y}_2, y_8 - \hat{y}_1)$
9	$\text{Cor}(y_1 - \hat{y}_1, y_9 - \hat{y}_1)$	$\text{Cor}(y_2 - \hat{y}_2, y_9 - \hat{y}_1)$

And so this matrix extends towards the right, it extends downwards as well, it basically covers our entire dataset. The correlations that we're taking, what we do is we take away the fixed part of our model and we take the correlation between what's left so the correlation between the random parts of the model, because that will be the correlation after controlling for the explanatory variables.

So how do we actually work out what these correlations come to?

Calculating the entries of V

How do we calculate the entries $\text{Cor}(y_i - \hat{y}_i, y_j - \hat{y}_j)$?

- First note that for a single level model, $y_i - \hat{y}_i = e_i$
- So we're trying to calculate $\text{Cor}(e_i, e_j)$
- Next, recall that we get the correlation by dividing the covariance by the total variance
- So we calculate the covariance matrix then divide all terms by the total variance to get the correlation matrix

How do we calculate the covariances?

We need two things:

- the handy rule

$$\text{Cov}(a, a) = \text{Var}(a)$$

 for any a
- and the assumptions of the model which we just saw

Well first of all, note that, for the single level model, this here - $y_i - \hat{y}_i$ - is just going to be just e_i , the one random term. So what the correlations were actually trying to calculate are just these - $\text{Cor}(e_i, e_j)$ - the correlation between the level 1 residual for each pair of individuals. The other thing to remember is that the correlation is the covariance divided by the total variance so the easiest way to work out the entries of this matrix is to calculate the covariance matrix and then divide all of the terms by the total variance which will give us the correlation matrix.

So to calculate the covariances, we need two things, one of them is this handy rule - $\text{Cov}(a, a) = \text{Var}(a)$ - which says that the covariance of something with itself is just the variance of that thing. And the other is all the assumptions that we've just seen.

Covariance matrix for a single level model

Covariance matrix for single level model

	1	2	3	4	5	6	...
1	$\text{Cov}(e_1, e_1)$	$\text{Cov}(e_1, e_2)$	$\text{Cov}(e_1, e_3)$	$\text{Cov}(e_1, e_4)$	$\text{Cov}(e_1, e_5)$	$\text{Cov}(e_1, e_6)$...
2	$\text{Cov}(e_2, e_1)$	$\text{Cov}(e_2, e_2)$	$\text{Cov}(e_2, e_3)$	$\text{Cov}(e_2, e_4)$	$\text{Cov}(e_2, e_5)$	$\text{Cov}(e_2, e_6)$...
3	$\text{Cov}(e_3, e_1)$	$\text{Cov}(e_3, e_2)$	$\text{Cov}(e_3, e_3)$	$\text{Cov}(e_3, e_4)$	$\text{Cov}(e_3, e_5)$	$\text{Cov}(e_3, e_6)$...
4	$\text{Cov}(e_4, e_1)$	$\text{Cov}(e_4, e_2)$	$\text{Cov}(e_4, e_3)$	$\text{Cov}(e_4, e_4)$	$\text{Cov}(e_4, e_5)$	$\text{Cov}(e_4, e_6)$...
5	$\text{Cov}(e_5, e_1)$	$\text{Cov}(e_5, e_2)$	$\text{Cov}(e_5, e_3)$	$\text{Cov}(e_5, e_4)$	$\text{Cov}(e_5, e_5)$	$\text{Cov}(e_5, e_6)$...
6	$\text{Cov}(e_6, e_1)$	$\text{Cov}(e_6, e_2)$	$\text{Cov}(e_6, e_3)$	$\text{Cov}(e_6, e_4)$	$\text{Cov}(e_6, e_5)$	$\text{Cov}(e_6, e_6)$...
7	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\ddots

So the covariance matrix for the single level model is this, basically we're just taking the covariance of these terms.

Calculating the covariances

The same observation (diagonal terms)

$$\begin{aligned}\text{Cov}(e_{ii}, e_{ii}) &= \text{Var}(e_{ii}) \\ &= \sigma_e^2\end{aligned}$$

(by definition)

Two different observations

$$\text{Cov}(e_{ii}, e_{jj}) = 0$$

(by assumption)

For the same observation - that's all the diagonal terms - the covariance of each level 1 residual with itself is just the variance and we define that to be σ_e^2 so all the diagonal terms are

$$\sigma_e^2.$$

Two different observations, the covariance is 0 because that was one of our assumptions.

So here's the covariance matrix for a single level model:

Covariance matrix for a single level model

	1	2	3	4	5	6	7	8	9	10	11	12	13	...
1	σ_e^2	0	0	0	0	0	0	0	0	0	0	0	0	...
2	0	σ_e^2	0	0	0	0	0	0	0	0	0	0	0	...
3	0	0	σ_e^2	0	0	0	0	0	0	0	0	0	0	...
4	0	0	0	σ_e^2	0	0	0	0	0	0	0	0	0	...
5	0	0	0	0	σ_e^2	0	0	0	0	0	0	0	0	...
6	0	0	0	0	0	σ_e^2	0	0	0	0	0	0	0	...
7	0	0	0	0	0	0	σ_e^2	0	0	0	0	0	0	...
8	0	0	0	0	0	0	0	σ_e^2	0	0	0	0	0	...
9	0	0	0	0	0	0	0	0	σ_e^2	0	0	0	0	...
10	0	0	0	0	0	0	0	0	0	σ_e^2	0	0	0	...
11	0	0	0	0	0	0	0	0	0	0	σ_e^2	0	0	...
12	0	0	0	0	0	0	0	0	0	0	0	σ_e^2	0	...
13	0	0	0	0	0	0	0	0	0	0	0	0	σ_e^2	...
...

σ_e^2 down the diagonal, and 0 everywhere else and the total variance for the single level model is just σ_e^2 so we divide by that and get 1s down the diagonal and 0s everywhere else. This shows us that for this single level model, of course each individual - the correlation with themselves is 1 - but any pair of different individuals, their random terms are uncorrelated. So, after we've controlled for the explanatory variables, there's no correlation between a pair of different individuals if we use the single level model.

Covariance matrix for a random intercepts model

Covariance matrix: random intercepts model														
L2	1	1	2	2	3	...								
L1	1	2	1	2	2	1								
1	$\text{Cov}(\eta_{1,1}, \eta_{1,1})$	$\text{Cov}(\eta_{1,2}, \eta_{1,1})$	$\text{Cov}(\eta_{2,1}, \eta_{1,1})$	$\text{Cov}(\eta_{2,2}, \eta_{1,1})$	$\text{Cov}(\eta_{3,1}, \eta_{1,1})$	$\text{Cov}(\eta_{1,2}, \eta_{1,2})$	$\text{Cov}(\eta_{1,3}, \eta_{1,2})$	$\text{Cov}(\eta_{2,2}, \eta_{1,2})$	$\text{Cov}(\eta_{2,3}, \eta_{1,2})$	$\text{Cov}(\eta_{3,2}, \eta_{1,2})$	$\text{Cov}(\eta_{1,3}, \eta_{1,3})$	$\text{Cov}(\eta_{2,3}, \eta_{1,3})$	$\text{Cov}(\eta_{3,3}, \eta_{1,3})$...
1	$\text{Cov}(\eta_{1,1}, \eta_{1,2})$	$\text{Cov}(\eta_{1,2}, \eta_{1,2})$	$\text{Cov}(\eta_{2,1}, \eta_{1,2})$	$\text{Cov}(\eta_{2,2}, \eta_{1,2})$	$\text{Cov}(\eta_{3,1}, \eta_{1,2})$	$\text{Cov}(\eta_{1,3}, \eta_{1,2})$	$\text{Cov}(\eta_{1,2}, \eta_{1,3})$	$\text{Cov}(\eta_{2,2}, \eta_{1,3})$	$\text{Cov}(\eta_{2,3}, \eta_{1,3})$	$\text{Cov}(\eta_{3,2}, \eta_{1,3})$	$\text{Cov}(\eta_{1,3}, \eta_{1,3})$	$\text{Cov}(\eta_{2,3}, \eta_{1,3})$	$\text{Cov}(\eta_{3,3}, \eta_{1,3})$...
2	$\text{Cov}(\eta_{1,1}, \eta_{1,2})$	$\text{Cov}(\eta_{1,2}, \eta_{1,2})$	$\text{Cov}(\eta_{2,1}, \eta_{1,2})$	$\text{Cov}(\eta_{2,2}, \eta_{1,2})$	$\text{Cov}(\eta_{3,1}, \eta_{1,2})$	$\text{Cov}(\eta_{1,3}, \eta_{1,2})$	$\text{Cov}(\eta_{1,2}, \eta_{1,3})$	$\text{Cov}(\eta_{2,2}, \eta_{1,3})$	$\text{Cov}(\eta_{2,3}, \eta_{1,3})$	$\text{Cov}(\eta_{3,2}, \eta_{1,3})$	$\text{Cov}(\eta_{1,3}, \eta_{1,3})$	$\text{Cov}(\eta_{2,3}, \eta_{1,3})$	$\text{Cov}(\eta_{3,3}, \eta_{1,3})$...
2	$\text{Cov}(\eta_{1,1}, \eta_{1,2})$	$\text{Cov}(\eta_{1,2}, \eta_{1,2})$	$\text{Cov}(\eta_{2,1}, \eta_{1,2})$	$\text{Cov}(\eta_{2,2}, \eta_{1,2})$	$\text{Cov}(\eta_{3,1}, \eta_{1,2})$	$\text{Cov}(\eta_{1,3}, \eta_{1,2})$	$\text{Cov}(\eta_{1,2}, \eta_{1,3})$	$\text{Cov}(\eta_{2,2}, \eta_{1,3})$	$\text{Cov}(\eta_{2,3}, \eta_{1,3})$	$\text{Cov}(\eta_{3,2}, \eta_{1,3})$	$\text{Cov}(\eta_{1,3}, \eta_{1,3})$	$\text{Cov}(\eta_{2,3}, \eta_{1,3})$	$\text{Cov}(\eta_{3,3}, \eta_{1,3})$...
3	$\text{Cov}(\eta_{1,1}, \eta_{1,2})$	$\text{Cov}(\eta_{1,2}, \eta_{1,2})$	$\text{Cov}(\eta_{2,1}, \eta_{1,2})$	$\text{Cov}(\eta_{2,2}, \eta_{1,2})$	$\text{Cov}(\eta_{3,1}, \eta_{1,2})$	$\text{Cov}(\eta_{1,3}, \eta_{1,2})$	$\text{Cov}(\eta_{1,2}, \eta_{1,3})$	$\text{Cov}(\eta_{2,2}, \eta_{1,3})$	$\text{Cov}(\eta_{2,3}, \eta_{1,3})$	$\text{Cov}(\eta_{3,2}, \eta_{1,3})$	$\text{Cov}(\eta_{1,3}, \eta_{1,3})$	$\text{Cov}(\eta_{2,3}, \eta_{1,3})$	$\text{Cov}(\eta_{3,3}, \eta_{1,3})$...
3	$\text{Cov}(\eta_{1,1}, \eta_{1,2})$	$\text{Cov}(\eta_{1,2}, \eta_{1,2})$	$\text{Cov}(\eta_{2,1}, \eta_{1,2})$	$\text{Cov}(\eta_{2,2}, \eta_{1,2})$	$\text{Cov}(\eta_{3,1}, \eta_{1,2})$	$\text{Cov}(\eta_{1,3}, \eta_{1,2})$	$\text{Cov}(\eta_{1,2}, \eta_{1,3})$	$\text{Cov}(\eta_{2,2}, \eta_{1,3})$	$\text{Cov}(\eta_{2,3}, \eta_{1,3})$	$\text{Cov}(\eta_{3,2}, \eta_{1,3})$	$\text{Cov}(\eta_{1,3}, \eta_{1,3})$	$\text{Cov}(\eta_{2,3}, \eta_{1,3})$	$\text{Cov}(\eta_{3,3}, \eta_{1,3})$...
3	$\text{Cov}(\eta_{1,1}, \eta_{1,2})$	$\text{Cov}(\eta_{1,2}, \eta_{1,2})$	$\text{Cov}(\eta_{2,1}, \eta_{1,2})$	$\text{Cov}(\eta_{2,2}, \eta_{1,2})$	$\text{Cov}(\eta_{3,1}, \eta_{1,2})$	$\text{Cov}(\eta_{1,3}, \eta_{1,2})$	$\text{Cov}(\eta_{1,2}, \eta_{1,3})$	$\text{Cov}(\eta_{2,2}, \eta_{1,3})$	$\text{Cov}(\eta_{2,3}, \eta_{1,3})$	$\text{Cov}(\eta_{3,2}, \eta_{1,3})$	$\text{Cov}(\eta_{1,3}, \eta_{1,3})$	$\text{Cov}(\eta_{2,3}, \eta_{1,3})$	$\text{Cov}(\eta_{3,3}, \eta_{1,3})$...
...

- Now $y_{ij} - \hat{y}_{ij} = u_j + e_{ij}$
- This is exactly the same as for a variance components model
- Recall when we were calculating residuals for a variance components model we defined r_{ij} to be $u_j + e_{ij} = y_{ij} - \hat{y}_{ij}$
- We'll use that notation again here to save space.

So, moving onto the random intercepts model now - now we've got a more complicated data structure we've got actually level 1 units within level 2 units so, for example, pupils within schools.

L2	1	1	2	2	3	...
L1	1	2	1	2	1	...
↑	↑	↑	↑	↑	↑	
So now this is pupil 1 within school 1	this is pupil 2 within school 1	Here's pupil 1 within school 2	pupil 2 within school 2	pupil 1 within school 3		

(L2 = 1, L1 = 1) So now this is pupil 1 within school 1, (L2 = 1, L1 = 2) this is pupil 2 within school 1. (L2 = 2, L1 = 1) Here's pupil 1 within school 2, (L2 = 2, L1 = 2) pupil 2 within school 2, (L2 = 3, L1 = 1) pupil 1 within school 3.

L2	L1	
1	1	pupil 1 within school 1
1	2	← pupil 2 within school 1
2	1	← pupil 1 within school 2

And the same for the rows

(L2 = 1, L1 = 1) pupil 1 within school 1, (L2 = 1, L1 = 2) pupil 2 within school 1, (L2 = 2, L1 = 1) pupil 1 within school 2, ...and so on.

Obviously this is just a reduced example so it will all fit on the screen, obviously really we'd expect to have more pupils per school.

For the random intercept model, this thing that we're taking the covariance of, is just $u_j + e_{ij}$ and we've actually written this here as r_{ij} because, if you remember, in the variance components model, when we were calculating residuals we actually defined r_{ij} to be just $u_j + e_{ij}$. So we've written that here because it takes less space basically. And this covariance matrix is just the same as for a variance components model, if we were taking the covariance matrix for a variance components model we'd have exactly the same thing to calculate.

Calculating the covariances

The same observation

$$\begin{aligned} &= \text{Cov}(u_{j_1} + e_{ij_1}, u_{j_1} + e_{ij_1}) \\ &= \text{Cov}(u_{j_1}, u_{j_1}) + 2\text{Cov}(u_{j_1}, e_{ij_1}) \\ &\quad + \text{Cov}(e_{ij_1}, e_{ij_1}) \\ &= \text{Var}(u_{j_1}) + 0 + \text{Var}(e_{ij_1}) \\ &= \sigma_u^2 + \sigma_e^2 \end{aligned}$$

Different observations, same group

$$\begin{aligned} &= \text{Cov}(u_{j_1} + e_{ij_1}, u_{j_2} + e_{ij_2}) \\ &= \text{Cov}(u_{j_1}, u_{j_1}) + \text{Cov}(u_{j_1}, e_{ij_2}) \\ &\quad + \text{Cov}(u_{j_2}, e_{ij_1}) + \text{Cov}(e_{ij_1}, e_{ij_2}) \\ &= 0 + 0 + 0 + 0 \\ &= 0 \end{aligned}$$

Different groups

$$\begin{aligned} &= \text{Cov}(u_{j_1} + e_{ij_1}, u_{j_2} + e_{ij_2}) \\ &= \text{Cov}(u_{j_1}, u_{j_1}) + \text{Cov}(u_{j_1}, e_{ij_2}) \\ &\quad + \text{Cov}(u_{j_2}, e_{ij_1}) + \text{Cov}(e_{ij_1}, e_{ij_2}) \\ &= 0 + 0 + 0 + 0 \\ &= 0 \end{aligned}$$

- We defined σ_u^2 to be the level 2 variance
- Now we can see it is also the covariance between level 1 units from the same group

Calculating the covariances

When we're calculating these covariances, for the same observation, this is what we get, and we won't go into the details but that comes out as the Variance of u_j + 0 + Variance of e_{ij} so these are just - the variance of u_j is just defined to be σ_u^2 and the Var of e_{ij} is defined to be σ_e^2 . For the same individual (on the diagonal terms) you have $\sigma_u^2 + \sigma_e^2$. For 2 individuals from different groups, this covariance comes out to be this:

Different groups

$$\begin{aligned} &= \text{Cov}(u_{j_1} + e_{ij_1}, u_{j_2} + e_{ij_2}) \\ &= \text{Cov}(u_{j_1}, u_{j_2}) + \text{Cov}(u_{j_1}, e_{ij_2}) \\ &\quad + \text{Cov}(u_{j_2}, e_{ij_1}) + \text{Cov}(e_{ij_1}, e_{ij_2}) \\ &= 0 + 0 + 0 + 0 \\ &= 0 \end{aligned}$$

and all of these are things that we assumed were 0 when we were looking at the assumptions of the model. So for 2 individuals from different groups, we just have a covariance of 0.

For 2 different individuals from the same group, the covariance comes out to be this:

Different observations, same group

$$\begin{aligned} &= \text{Cov}(u_{j_1} + e_{ij_1}, u_{j_1} + e_{ij_1}) \\ &= \text{Cov}(u_{j_1}, u_{j_1}) + \text{Cov}(u_{j_1}, e_{ij_1}) \\ &\quad + \text{Cov}(u_{j_1}, e_{ij_1}) + \text{Cov}(e_{ij_1}, e_{ij_1}) \\ &= \text{Var}(u_{j_1}) + 0 + 0 + 0 \\ &= \sigma_u^2 \end{aligned}$$

$$= \text{Cov}(u_{j_1}, u_{j_1}) + \text{Cov}(u_{j_1}, e_{ij_1}) \\ + \text{Cov}(u_{j_1}, e_{ij_1}) + \text{Cov}(e_{ij_1}, e_{ij_1})$$

and these 3 terms we assumed were 0 but this one -

$$\begin{aligned} &= \text{Cov}(u_{j_1}, u_{j_1}) + \text{Cov}(u_{j_1}, e_{ij_1}) \\ &\quad + \text{Cov}(u_{j_1}, e_{ij_1}) + \text{Cov}(e_{ij_1}, e_{ij_1}) \end{aligned}$$

- is the variance of u_j and we defined that to be σ_u^2 so for 2 different observations from the same group we have σ_u^2 as the covariance.

Now it might seem a bit strange - we introduced σ_u^2 at the beginning when we were looking at the variance components model for the first time - we introduced that as the level 2 variance and now we're saying - it's the covariance between level 1 units from the same group. Actually if you think about it, it's not that strange that those two things would be equal because the differences between groups are kind of the similarities within groups, it's just two different ways of looking at the same thing really.

Covariance matrix for random intercept model

L2	1	1	1	1	2	2	3	...
L1	1	2	3	4	1	2	1	...
1	$\sigma_u^2 + \sigma_e^2$	$\rho\sigma_u^2 + \sigma_e^2$	$\rho\sigma_u^2 + \sigma_e^2$	$\rho\sigma_u^2 + \sigma_e^2$	0	0	0	...
2	$\rho\sigma_u^2 + \sigma_e^2$	$\sigma_u^2 + \sigma_e^2$	$\rho\sigma_u^2 + \sigma_e^2$	$\rho\sigma_u^2 + \sigma_e^2$	0	0	0	...
3	$\rho\sigma_u^2 + \sigma_e^2$	$\rho\sigma_u^2 + \sigma_e^2$	$\sigma_u^2 + \sigma_e^2$	$\rho\sigma_u^2 + \sigma_e^2$	0	0	0	...
4	$\rho\sigma_u^2 + \sigma_e^2$	$\rho\sigma_u^2 + \sigma_e^2$	$\rho\sigma_u^2 + \sigma_e^2$	$\sigma_u^2 + \sigma_e^2$	0	0	0	...
2	0	0	0	0	$\sigma_u^2 + \sigma_e^2$	$\rho\sigma_u^2 + \sigma_e^2$	0	...
2	0	0	0	0	$\rho\sigma_u^2 + \sigma_e^2$	$\sigma_u^2 + \sigma_e^2$	0	...
3	0	0	0	0	0	0	$\sigma_u^2 + \sigma_e^2$...
3	0	0	0	0	0	0	$\rho\sigma_u^2 + \sigma_e^2$...
3	0	0	0	0	0	0	$\rho\sigma_u^2 + \sigma_e^2$...
3	0	0	0	0	0	0	$\rho\sigma_u^2 + \sigma_e^2$...

Divide through by total variance, $\sigma_u^2 + \sigma_e^2$

$$\frac{\sigma_u^2 + \sigma_e^2}{\sigma_u^2 + \sigma_e^2} = 1$$

$$\frac{0}{\sigma_u^2 + \sigma_e^2} = 0$$

$$\frac{\rho\sigma_u^2 + \sigma_e^2}{\sigma_u^2 + \sigma_e^2} = \rho$$

Here's our covariance matrix with all the terms filled in, so we've got $\sigma^2_u + \sigma^2_e$ down the diagonal so for each individual the covariance with themselves is the total variance and for 2 different individuals from the same group we've got σ^2_u and for 2 individuals from different groups we've got zero covariance.

So we now divide through by the total variance $\sigma^2_u + \sigma^2_e$ so we can get the correlation matrix so for the diagonal terms that will just become 1. For 2 individuals from different groups that will be 0 still and for 2 individuals from the same group, 2 different individuals from the same group, that will come out to be the variance partition coefficient, ρ . So that's why, as we saw with the variance components model, the intra-class correlation is the same as the variance partitioning coefficient.

V for random intercepts model

V for random intercepts model										...	
L2	1	1	1	1	2	2	3	3	4	...	
L1	1	2	3	4	1	2	1	2	3	4	...
1	1	ρ	ρ	ρ	0	0	0	0	0	0	...
1	2	ρ	1	ρ	0	0	0	0	0	0	...
1	3	ρ	ρ	1	0	0	0	0	0	0	...
1	4	ρ	ρ	ρ	1	0	0	0	0	0	...
2	1	0	0	0	0	1	ρ	ρ	ρ	ρ	...
2	2	0	0	0	0	ρ	2	0	0	0	...
3	1	0	0	0	0	0	ρ	1	ρ	ρ	...
3	2	0	0	0	0	0	ρ	1	2	0	...
3	3	0	0	0	0	0	0	ρ	1	ρ	...
3	4	0	0	0	0	0	0	0	ρ	1	...

- The correlation matrix is identical to the matrix for the variance components model
- As expected, observations within the same group are correlated but observations from different groups are uncorrelated

The correlation matrix shows us that for the random intercepts model and the variance components model (because that has the same correlation matrix) 2 individuals from different groups are uncorrelated after we control for the explanatory variables but 2 individuals from the same group, after we control for the explanatory variables, have correlation ρ . So that shows how the random intercepts model is allowing for that correlation between different observations from the same group.

7) Random intercepts models: Residuals

So, having looked at the correlation between the residuals, let's look at the residuals themselves in a bit more detail.

Residuals									
Having looked at the correlation between residuals, let's look at the residuals themselves in more detail.									
Residuals are estimates for the random part									
For a single level model, we often say 'estimates for the error term'									
Reminder									
For a single level model, the residual for an observation is an estimate for e_i									
$y_i = \beta_0 + \beta_1 x_i + e_i$									
If we write $\hat{y}_i = \beta_0 + \beta_1 x_i$ then									
$\hat{e}_i = y_i - \hat{y}_i$: the observed value - the value predicted by the regression line									



The residuals are estimates for the random part. For a single level model we often say *estimates for the error term*. And just a quick reminder - for a single level model - the residual for an observation is an estimate for e_i so we've got here an estimate for \hat{e}_{54} and an estimate for \hat{e}_{44} here and if we write - $\hat{y}_i = \beta_0 + \beta_1 x_i$ - then our estimate for our residual is the observed value of y minus \hat{y}_i . It's the observed value minus the value predicted by the regression line. So that's what we can see here, the estimate for the residual is the distance between the data point and the overall regression line.

Why are we actually interested in the residuals?

Why are we interested in the residuals?

Often we're not, but they can be useful in some cases:

Diagnostics

- We can plot the residuals to check their normality
- This is part of checking how well the model fits

Rankings

- We can rank level 2 units by u_j
- e.g. school league tables

Interest in a unit

- We can find out how a particular unit compares to the average

Prediction/ visualisation

- The level 2 residuals are needed to make predictions for individuals in a particular level 2 unit
- We need them to graph the group lines

Well often we're not, but they can be useful in some cases.

So one example is diagnostics. We can plot the residuals to check their normality and that's part of checking how well the model fits.

We can also rank level 2 units by their level 2 residuals and an example of that is school league tables.

Sometimes we're interested in a particular unit, we want to find out how that particular unit compares to the average.

And then we need the level 2 residuals in order to make predictions for individuals in a particular level 2 unit and we need them in fact to graph the group lines which are the predicted lines for each group.

Why are we interested in the residuals?

Often we're not, but they can be useful in some cases:

Diagnostics

- 'Was our normality assumption justified?'

Interest in a unit

- 'How is Hospital 18 doing?'
- 'How is Pupil 6 doing compared to the rest of their school?'

Rankings

- 'After accounting for intake, which school performs the best?'

Prediction/ visualisation

- 'What is the expected weight of a salmon from Fish Farm 28?'
- 'What does our model look like?'

A few examples of those uses:

For the use of diagnostics we might be asking "Was our normality assumption justified?"

For rankings we might be asking "After accounting for intake, which school performs the best?"

For interest in a unit, we might be asking "How is Hospital 18 doing?" or: "How is Pupil 6 doing compared to the rest of their school?"

And in terms of prediction and visualisation, we might be asking "What is the expected weight of salmon from Fish Farm number 28?" or "What does our model look like?"

Multilevel residuals

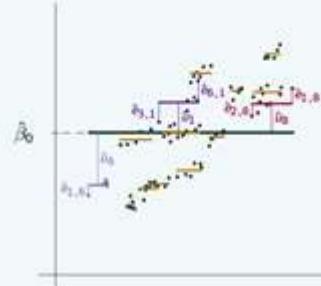
Multilevel residuals

Variance components model

$$y_{ij} = \beta_0 + u_j + e_{ij}$$

Recall that now that we have 2 random terms, we have 2 kinds of residual:

- the level 2 residual, an estimate for u_j
- the level 1 residual, an estimate for e_{ij}



So just recall briefly that for the variance components model, because we have 2 random terms, we have 2 kinds of residual. We've got the level 2 residual, which is an estimate for u_j , and the level 1 residual which is an estimate for e_{ij} .

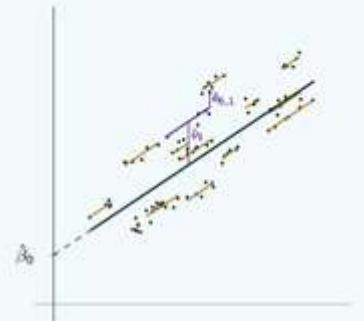
Multilevel residuals

Random intercept model

$$y_{ij} = \beta_0 + \beta_1 x_{ij} + u_j + e_{ij}$$

Again now that we have 2 random terms, we have 2 kinds of residual:

- the level 2 residual, an estimate for u_j
- the level 1 residual, an estimate for e_{ij}



And for the random intercept model, again, now that we have 2 random terms, we have 2 kinds of residual, again an estimate for u_j and an estimate for e_{ij} .

And the calculation is very similar too for the variance components model:

Calculation of residuals

Raw mean residual	Shrinkage factor
<ul style="list-style-type: none"> • Recall $r_{ij} = y_{ij} - \hat{y}_{ij}$ and • the raw mean residual \bar{r}_j is the mean of r_{ij} for group j 	$k = \frac{\hat{\sigma}_u^2}{\hat{\sigma}_u^2 + \hat{\sigma}_e^2/n_j}$ Be careful! The shrinkage factor is similar to the VPC ρ .
Level 2 residual	
<ul style="list-style-type: none"> • Just as for the variance components model, we shrink the raw residuals towards the overall mean 	
	$\hat{u}_j = \bar{r}_j \times k = \bar{r}_j \times \frac{\hat{\sigma}_u^2}{\hat{\sigma}_u^2 + \hat{\sigma}_e^2/n_j}$
Level 1 residual	
<ul style="list-style-type: none"> • Once again, the fixed effects model uses the raw residual r_j 	
	$\hat{e}_{ij} = y_{ij} - \hat{y}_{ij} - \hat{u}_j = r_{ij} - \hat{u}_j$

Remember that we have the raw mean residual, which is the mean of r_{ij} , for group j , and r_{ij} is just the datapoint minus the value predicted by the overall regression line.

And we have the shrinkage factor again with the same formula as for the variance components model - and again we need to be careful because the shrinkage factor is quite similar to the VPC, the variance partitioning coefficient p , so it's important not to get those confused.

So the level 2 residual is calculated in just the same way as for the variance components model. We take the raw mean residual and multiply by the shrinkage factor. Once again, if we have the fixed effects model then we wouldn't shrink the residual, we'd just use the raw residual, but for the multilevel models, for the variance components model and for the random intercept model, we do shrink the residual.

Level 1 residual is simpler, we just take the observed value and subtract the value predicted by the overall regression line and the level 2 residual, so that's just the raw mean residual minus the level 2 residual.

So why do we actually shrink?

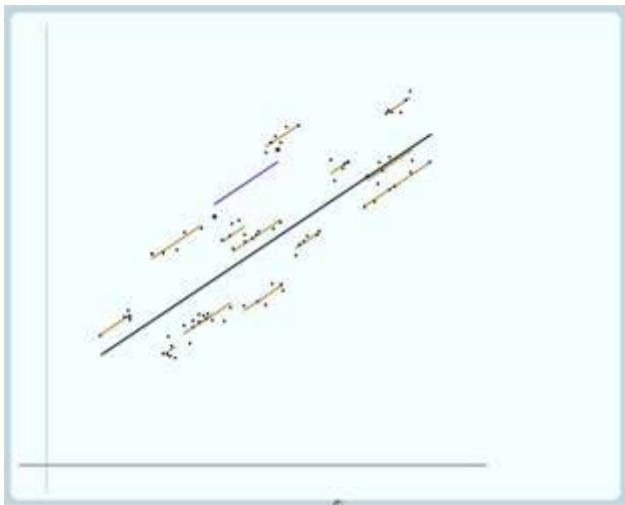
Why do we shrink? A thought experiment

The situation
<ul style="list-style-type: none"> • Suppose we have exam results for 6 pupils each from a number of schools • Now suppose we drop 4 pupils from one school, just to see what will happen • How close will the school line using the remaining 2 pupils be to the "correct" school line using all 6 pupils?
How close is the line?
<ul style="list-style-type: none"> • It depends how "typical" those 2 pupils are of the full set of 6. • Since we're picking 2 pupils from 6, it's quite likely that we might pick 2 untypical pupils • Then the school line drawn using those 2 will be quite far from the school line using all 6 pupils- as happens in this example

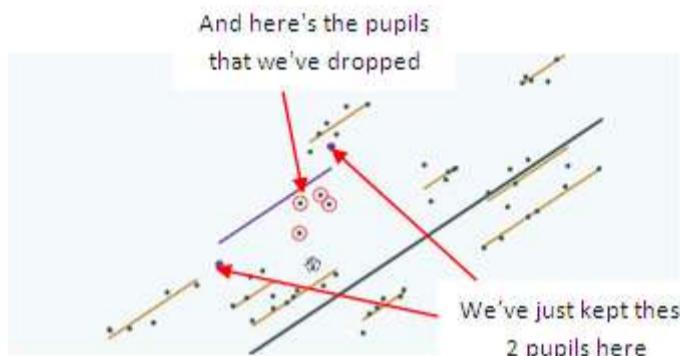
Well, in order to explain that, let's do a thought experiment. So imagine that we have exam results for 6 pupils each from a number of schools and suppose that we drop 4 pupils from 1 school just to see what will happen. If we do that, how close will the school line using the remaining 2 pupils be to the correct line using all 6 pupils?

Well, that will depend on how typical those 2 pupils are of the full set of 6. Because we are picking 2 pupils from 6 it is quite likely that we might pick 2 untypical pupils and then the school line drawn using those 2 will be quite far from the school line using all 6 pupils.

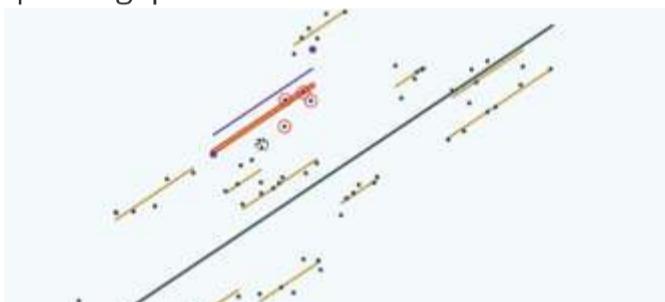
And if we draw a graph of our particular example because that happens in this case:



Here's the school where we have dropped 4 pupils,



We've drawn the line here using just the 2 pupils that we kept and we can put on, again, the line drawn using all 6 pupils, you can see there's quite a gap between those 2 lines:



Well, we're interested, in fact, in drawing the line for all 6 pupils -

Thought experiment cont.

Drawing the group line

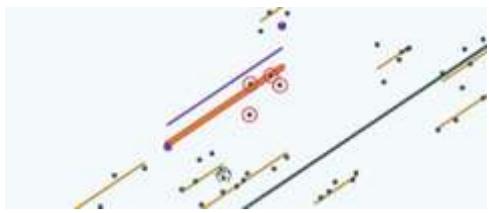
- We're interested in the line for all 6 pupils
- We want our line when we have 2 pupils
 - to stay close to the line for all 6 pupils
 - not to move around a lot according to which 2 pupils we pick
- How can we draw a line that does those things?

Information about the dropped points

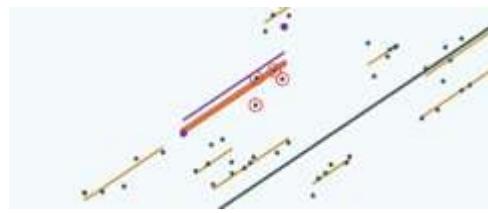
- To draw a line like that, we need some information on the 4 pupils we dropped
- It seems we don't have any, but actually that's not true
- The position of the other school lines tells us something about where the 4 pupils and the line for this school are likely to be:
 - they are more likely to be closer to the overall average line.
 - So we can improve our positioning of the line by shrinking it in towards the overall average

it's much more interesting to know where the line for all 6 pupils is than the line for just those 2 pupils because we're interested in the school effect rather than the school effect for those particular 2 pupils. So we want our line - when we have 2 pupils - to stay close to the line for all 6 pupils and we don't want it to move around a lot depending on which 2 pupils we pick. How can we actually draw a line that does those things?

We need some information about the 4 pupils that we dropped and it seems as though we don't have any but actually that's not true. The position of the other school lines tells us something about where the 4 pupils and the line for this school are likely to be because they're more likely to be closer to the overall average line. If those 4 pupils that we dropped were further away from the overall average line than the pupils we kept they'd be going towards the direction of being outliers. But, if they're closer to the overall average line, then they're more in line with the whole distribution of schools so it's more likely that they're closer to the overall average line. So we can improve our positioning of the line by shrinking it in towards the overall average. So let's see that:



So here's the line as we draw it just using those two pupils



and if we now shrink it towards the overall average, this is where it ends up

...Let's do that again, so here's where we draw it using just those two pupils

...and if we shrink it towards the overall average here's where it ends up

So obviously, in our thought experiment situation, we don't know where those four pupils were, we don't know where that group line is. So if we want to do that in terms of how it looks to us, given the information available to us, here's the group line drawn just using those two pupils...



...before we shrink



and after we shrink here it is

...here it is before we shrink

...and here it is after we shrink.

So you might wonder: Is shrinking actually always better?

Thought experiment cont.

Is shrinkage always better?

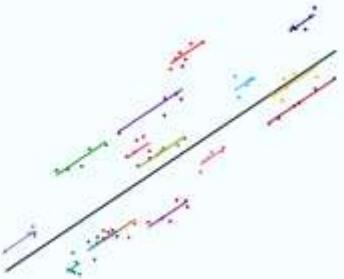
- It is possible that shrinking the line will move it away from the line using all pupils
- If we don't know where the other 4 pupils are we don't know if shrinking moves the line closer or further away
- But we do know that shrinking will move the line closer more often than it will move it further away
- So shrinking the line is always our best chance to get the line closer

How does this generalise?

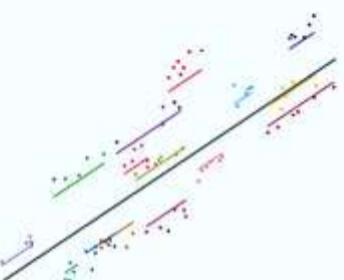
- We are in the same situation with our dataset as a whole
- We have just 6 pupils from each school
- There are really many more
- The position of the other school lines gives us information about the likely position of the pupils not included in the dataset
- Again, we get a better estimate of the group line by shrinking it in towards the overall average

Well it is possible, that when we shrink, we will move the line away from the line using all pupils because those other 4 pupils could have been further from the overall average line so it's possible when we move towards the overall average line we are moving further away from them. But we actually don't know where the other 4 pupils are so we don't know whether shrinking will move the line closer or further away. However, we do know that shrinking will move the line closer more often than it will move it further away. And that means that shrinking the line is always our best chance to get the line closer. Shrinking the line is our best estimate. It might put our line further from the group line using all the points but our best estimate is to shrink the line and move it closer to the overall average line.

So OK, that was our thought experiment, but how does it actually generalise? Well we're actually in the same situation with our dataset as a whole. We have just 6 pupils from each school, and there are really many more. There might be 500 or 1,000 pupils in the entire school. The position of the other school lines gives us information about the likely position of the pupils not included in the dataset. So that means that, again, we'll get a better estimate of the group line by shrinking it in towards the overall average.



So here are our group lines drawn using just the information from the pupils in each school and if we now use the information from the pupils in all schools to draw each line by shrinking towards the overall average - here's what we get:



Let's see that again, so here's the lines before we shrink



...and here's the lines after we shrink

A couple of important points to note about shrinkage

Points to note about shrinkage	
When do we shrink? Always! <ul style="list-style-type: none"> • We always shrink the residuals because we always have a sample from each level 2 unit. <ul style="list-style-type: none"> * even if we have 499 out of the 500 pupils attending a school * even if we have all the pupils attending all the schools in the UK • Even if our dataset contains the whole population, we regard that as a sample from the superpopulation and shrink 	How much do we shrink by? Recall that: <ul style="list-style-type: none"> • We don't shrink by a fixed amount <ul style="list-style-type: none"> * If we have 500 pupils from a school we shrink less than if we have 7 • The amount we shrink by depends on the absolute number of level 1 units, not the proportion of the total for that level 2 unit • We can also see that the amount of shrinkage depends on the variances σ_u^2 and σ_e^2

When do we shrink? Well, we always shrink, and that's because we always have a sample from each level 2 unit. Even if we have 499 out of 500 pupils attending the school, that's still a sample. Even if we have all the pupils attending all the schools in the UK, we can still think of that as a sample from the super-population, our actual population of UK schools is a realisation from the super-population of all possible UK schools.

How much do we actually shrink by? Well, remember from the variance components model, we don't shrink by a fixed amount. If we have 500 pupils from a school we shrink less than if we have 7. Remember that the amount we shrink by depends on the absolute number of level 1 units, not the proportion of the total for that level 2 unit. If we have 10 pupils from a school then we shrink by the same amount, whether that school actually only has 10 pupils, or whether that school really has 500 pupils. We can also see that the amount of shrinkage will depend on the variances σ_u^2 and σ_e^2 .

And for more about the residuals and shrinkage - for a more detailed examination of the same topic - you can listen to our [audio-visual presentation on residuals](#) which is also available from the web site.

8) Random intercepts models: Predictions

Let's look at making predictions for the random intercept model now.

Predictions

There are several reasons for making predictions:

- Model testing
To see how close predictions from the model are to the values we observe
- Model visualisation
To try to understand what happens to our response when we change the values of the explanatory variables, using graphs
- Estimates for units not in the dataset
To obtain an estimate for the value of the response for units not in the dataset
 - that you have values of the explanatory variables for
 - existing units or hypothetical ones

We focus on the second use

And there are several reasons why we might want to do that.

One of them is model testing. So if we want to see how close predictions from the model are to values that we actually observe, we might want to make some predictions.

Another is model visualisation. If we want to try and understand what happens to our response when we change the values of the explanatory variables then we might want to draw graphs of predicted values.

And the third reason is estimates for units not in our dataset. This would be estimates for units for which we have values of the explanatory variables and it can either be existing units that we actually have measurements on or it can be hypothetical units. And we're actually going to focus just on the second use in this section, we're only going to consider model visualisation.

Visualising the model

Visualising the model

This should already be familiar from single level models

Visualising a single level model

$$y_i = \beta_0 + \beta_1 x_i + e_i$$

- We plot $y = \hat{\beta}_0 + \hat{\beta}_1 x$ to get our graph
- $\hat{\beta}_0 + \hat{\beta}_1 x$ is actually \hat{y} , our predicted value.
- We can add on the actual data points

This idea should already be familiar from single level regression models because when we visualise a single level model we just plot $\hat{\beta}_0 + \hat{\beta}_1 x$ to get our graph of the overall regression line and this is just plotting the fixed part prediction, just our predicted value for the single level model. And we can add on the actual data points as well.

Visualising the random intercepts model

Overall regression line

- Prediction from the fixed part gives the overall regression line
- Prediction:
 $\hat{y}_{ij} = \hat{\beta}_0 + \hat{\beta}_1 x_{1ij}$
- This is just the same as the line for the single level model
- The value of \hat{y}_{ij} does not depend on the group j in this case, only the explanatory variables
- So this prediction only produces one line
- This is what we would predict if we didn't know which group a data point belonged to

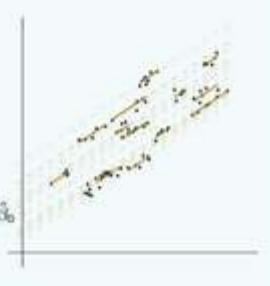
So for the random intercepts model, if we predict from the fixed part, again we get the overall regression line. So the prediction is just the same as for the single level regression model and the graph looks the same. The value of our prediction doesn't depend on group in this

case. There's no difference in what this comes to depending on which group we're looking at so we only have one line produced by this prediction. And this is what we'd predict if we didn't know which group a data point belonged to. Again we can put the data points on.

Visualising the random intercepts model

Group lines

- Adding in the group residual u_j gives the group lines
- Prediction:
 $\hat{y}_j = \hat{\beta}_0 + \hat{\beta}_1 x_j + u_j$
- Now the value of \hat{y}_j depends on the group j as well as the explanatory variables
- So there is a different line for each group
- The group line is what we would predict if we did know which group a data point belonged to



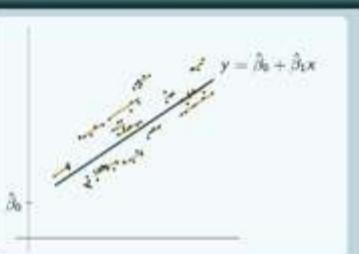
The group lines - we get by adding in the group residual to the fixed part prediction so now this is our prediction. And the value of this does depend on which group because this estimated u_j is different for every group - and incidentally you can see here how the residuals come into the prediction - because this is our level 2 residual.

So we've got a different line for each group because our predicted value depends on our group, and the group line is what we'd predict if we knew which group a data point belongs to. And again we can put on the data points.

We can combine the predictions from the fixed and random part into one graph to get a complete visualisation of the model and we can add on the actual data points for comparison. And usually we only draw our group lines for the range of values that we have in our dataset.

So here's a picture of our random intercept model:

Visualising the random intercepts model



Complete model

- We can combine the predictions from the fixed and random part in one graph to get a complete visualisation of the model
- and we can add in the actual data points for comparison
- Usually we only plot predictions for the range of values we have in our dataset