

1

Overview

- Goal: Learn to run basic (multilevel/hierarchical) regression in R.
 - Also focus on analysis of Intensive Longitudinal Data (ILD).
- Structure:
 - 3 days.
 - Alternating between lectures/discussion and practicals.

Tilburg University logo and 'tesc' logo at the bottom.

2

Topics


- Day 1:
 - (Multiple) Regression (the way I think about it).
 - How to approach a regression analysis....any analysis really.
 - Hierarchical data (what it is, what it implies, and what to do with it).
- Day 2 (common issues and special flavors of multilevel models).
- Day 3 (longitudinal and ILD data and whatever comes up during day 1 and 2).

Tilburg University logo and 'tesc' logo at the bottom.


3

The "Drill Sergeants"

Leonie, what do you think this course is going to be like?



Wrong on so many levels!



Tilburg University logo and 'tesc' logo at the bottom.

4

Regression

TILBURG UNIVERSITY



5

Regression

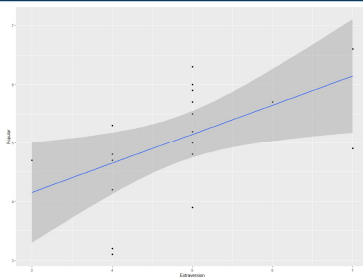
- What can you tell me about multiple regression?
- Describe what it does.

TILBURG UNIVERSITY



6

Regression

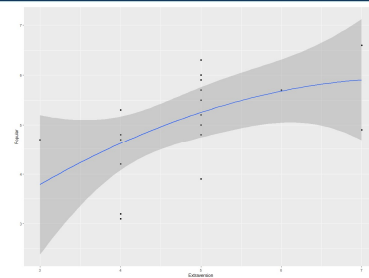


TILBURG UNIVERSITY



7

Regression

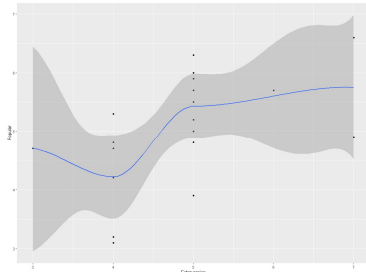


TILBURG UNIVERSITY



8

Regression

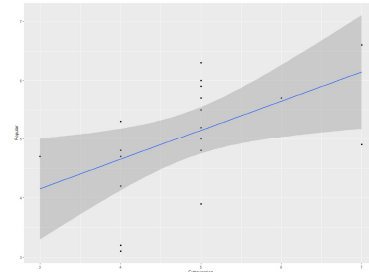


TILBURG UNIVERSITY

tesc

9

Regression → Curve Fitting!!

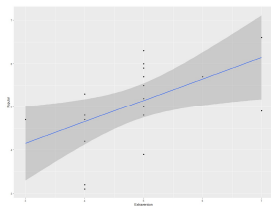


TILBURG UNIVERSITY

tesc

10

Curve Fitting



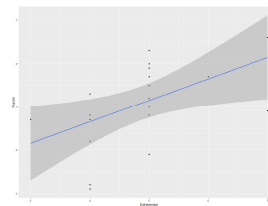
- Ok, so what are you saying with this model about your data?

TILBURG UNIVERSITY

tesc

11

Curve Fitting



- Ok, so what are you saying with this model about your data?

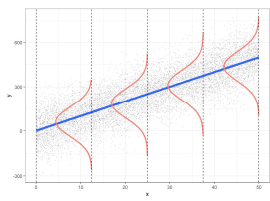
- **L**: There is a linear relationship between the mean response (Y) and the explanatory variable (X),
- **I**: The errors are independent—there's no connection between how far any two points lie from the regression line,
- **N**: The responses are normally distributed at each level of X , and
- **E**: The variance or, equivalently, the standard deviation of the responses is equal for all levels of X .

TILBURG UNIVERSITY

tesc

12

Curve Fitting

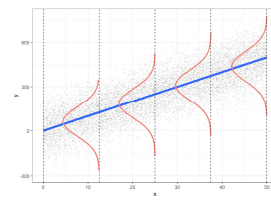


- Ok, so what are you saying with this model about your data?
- **L:** There is a linear relationship between the mean response (Y) and the explanatory variable (X).
- **I:** The errors are independent—there's no connection between how far any two points lie from the regression line.
- **N:** The responses are normally distributed at each level of X, and
- **E:** The variance or, equivalently, the standard deviation of the responses is equal for all levels of X.

TILBURG UNIVERSITY **t e s c**

13

Curve Fitting



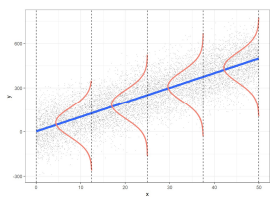
- **L:** There is a linear relationship between the mean response (Y) and the explanatory variable (X).
- **I:** The errors are independent—there's no connection between how far any two points lie from the regression line.
- **N:** The responses are normally distributed at each level of X, and
- **E:** The variance or, equivalently, the standard deviation of the responses is equal for all levels of X.

What if one or more don't hold?

TILBURG UNIVERSITY **t e s c**

14

Curve Fitting



- Take home:

Regression ≠ confined to linear, normal, homoscedastic....or

independent

TILBURG UNIVERSITY **t e s c**

15

Visualize, Visualize, Visualize

- Visualization is like exercise; we know we should but typically don't.
- This is a bigger issue than we realize.
- We analyze data through modeling, but the result of that modeling is only useful if the model matches the data (or vice versa).
- You should spend several hours just visualizing and looking at your data!
 - And do more than just check scatterplots for linearity and outliers ;).

TILBURG UNIVERSITY **t e s c**

16

Current Practices

- Usually only provide descriptions of our data through tables with summary statistics (means, sd's).
- These summaries are only useful with symmetric distributions!
- Could provide info on median, mode, kurtosis, skew, but....who really get's that anyway?
 - Even more cognitively taxing when comparing multiple groups or multiple levels.
 - Two groups can have different means but the same mode, different modes but the same mean, or the same mean and standard deviation but a meaningful skew (Heino et al., 2019)

Example

Description:

Accelerometer data indicated, that girls were as active as boys (mean 65 vs. 67 min).

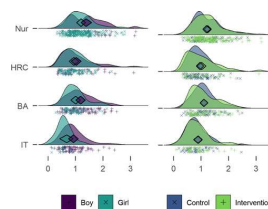
	IT	BA	HRC	Nur	Control	Intervention	Boy	Girl	Full sample
n	163	282	213	402	528	638	471	643	1166
Mean daily accelerometer wear time hours	839.4 (67.9)	848.3 (65.3)	839.4 (68.4)	848.0 (66.0)	834.7 (69.2)	861.0 (67.1)	839.9 (71.3)	842.9 (67.0)	848.6 (69.3)
Mean daily breaks in sitting	21.2 (6.8)	25.9 (7.2)	24.7 (6.9)	28.6 (7.8)	24.8 (7.4)	27.0 (7.8)	23.2 (6.8)	27.9 (7.8)	26.0 (7.7)
Mean daily hours spent sitting or lying down	617.5 (92.1)	535.7 (100.2)	514.9 (103.6)	499.2 (84.9)	519.6 (108.4)	534.8 (97.9)	570.5 (102.8)	499.1 (92.4)	527.7 (103.2)
Mean daily MVPA hours	51.2 (24.8)	66.4 (27.2)	58.0 (27.7)	74.4 (30.7)	63.3 (27.8)	68.0 (31.3)	66.9 (32.3)	64.9 (27.6)	65.8 (29.6)

Example

	Nur	HRC	BA	IT
g ² control	$\chi^2 = 13.0, df = 3, p = 0.004$	$\chi^2 = 17.5, df = 3, p = 0.000$	$\chi^2 = 40.3, df = 3, p < 0.001$	$\chi^2 = 24.2, df = 3, p < 0.001$
g ² intervention	$\chi^2 = 71.8, df = 3, p < 0.001$	$\chi^2 = 50.5, df = 3, p < 0.001$	$\chi^2 = 96.7, df = 3, p < 0.001$	$\chi^2 = 36.1, df = 3, p < 0.001$
g ² total	$\chi^2 = 84.8, df = 3, p < 0.001$	$\chi^2 = 68.0, df = 3, p < 0.001$	$\chi^2 = 136.9, df = 3, p < 0.001$	$\chi^2 = 60.3, df = 3, p < 0.001$
g ² interaction	$\chi^2 = 58.8, df = 3, p < 0.001$	$\chi^2 = 33.0, df = 3, p < 0.001$	$\chi^2 = 56.4, df = 3, p < 0.001$	$\chi^2 = 11.9, df = 3, p = 0.008$

Table 2. Means, standard deviations and some distributional properties of a single variable in different educational tracks the participants were nested in. Nur = Practical nurse, HRC = Hotel, restaurant and catering studies, BA = Business and administration, IT = Business information technology

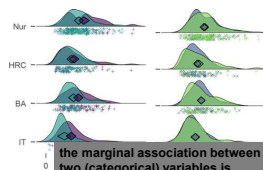
Example



Notice:

- boys did more moderate-to-vigorous physical activity in every educational track.
- In spite of this, girls appeared more active when combining the educational tracks
- Reason: much more people in the practical nurse track, as well as those people being mostly girls.
- This is also known as Simpson's paradox, and is best investigated by visualizing data.

Example



Notice:

- boys did more moderate-to-vigorous physical activity in every educational track.
- In spite of this, girls appeared more active when combining the educational tracks
- Reason: much more people in the practical nurse track, as well as those people being mostly girls.
- This is also known as Simpson's paradox, and is best investigated by visualizing data.

the marginal association between two (categorical) variables is qualitatively different from the partial association between the same two variables after controlling for one or more other variables.

21

Current Practices

- Data visualizations are crucial supplements to large numerical tables of descriptive statistics (Tay, Parrigon, Huang, & LeBreton, 2016).
- They are more straightforward ways to provide lots of important information (including uncertainty).
 - Require much less statistical/mathematical knowledge from the reader.
- Also, providing extensive visualizations (including the raw data) will aid open-science and replication/reproducibility.

22

Practical Data Example

- Study on the popularity of highschool students.
 - Total of 246 students from 12 different classes.
 - Determined how extraversion, gender, and teacher experience influenced a student's popularity.
- List of all the variables:
 - pupil: pupil identification variable, not needed in the analysis
 - class: class identification variable, the linking variable to define the 2 - level structure
 - student-level independent variables: extraversion (continuous; higher scores mean higher extraversion) and gender (dichotomous; 0=male, 1=female)
 - class-level independent variables: teacher experience (in years)
 - outcome variable: popular (continuous outcome variable at the student-level, higher scores indicate higher popularity)

23

Getting to know your data.

Variable	n	mean	sd	min	max	range	se
Extraversion	20	4.85	0.99	3.0	7.0	4.0	0.22
Popular	20	5.08	0.95	3.1	6.6	3.5	0.21

24

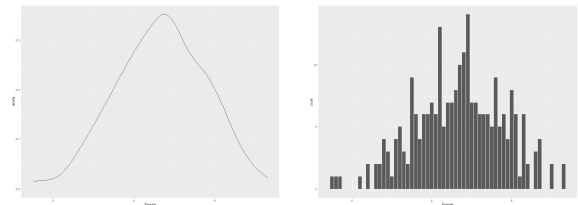
Getting to know your data.

Variable	n	mean	sd	min	max	range	se
Extraversion	20	4.85	0.99	3.0	7.0	4.0	0.22
Popular	20	5.08	0.95	3.1	6.6	3.5	0.21

Variable	n	mean	sd	min	max	range	skew	kurtosis	se
Extraversion	20	4.85	0.99	3.0	7.0	4.0	0.59	0.11	0.22
Popular	20	5.08	0.95	3.1	6.6	3.5	-0.50	-0.58	0.21

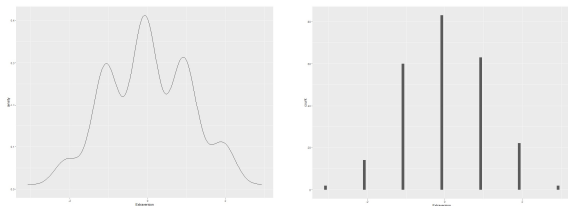
25

Getting to know your data.



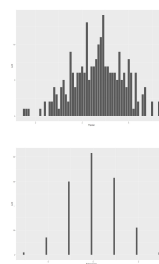
26

Getting to know your data.



27

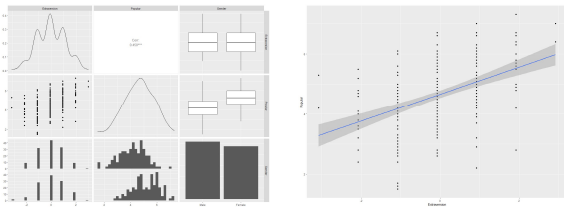
Getting to know your data.



- Any issues?
- Can these pictures tell me anything I need to know for modeling?
- Did we actually test any of the **LINE** assumptions?

28

Getting to know your data.

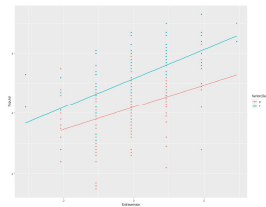


TILBURG UNIVERSITY



29

Getting to know your data.



TILBURG UNIVERSITY



30

- Conclusions:
 - Popularity and Extraversion not-normal, but that is not necessarily a bad thing.
 - Extraversion might be more ordinal than continuous (how we model depends on theory and measurement).
 - Proportion of boys and girls pretty similar (important for interpretation of effects and interactions)
 - No clear indications of interaction effect.

Intermediate conclusion

$$y = b_0 + b_1X_1 + b_2X_2 + e$$

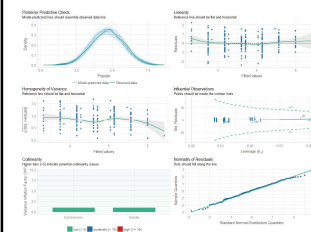
$$Popular = b_0 + b_1Extraversion + b_2Gender + e$$

TILBURG UNIVERSITY



31

Getting to know your data.

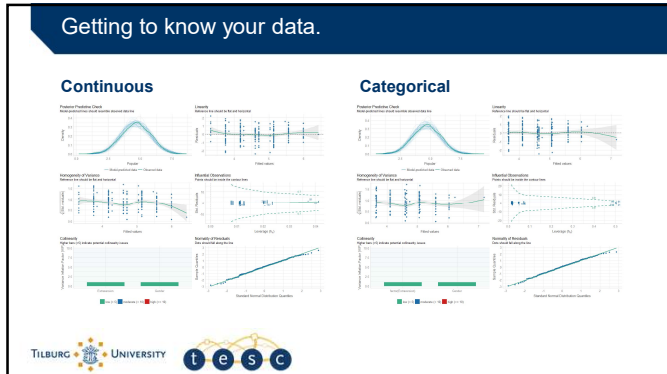


TILBURG UNIVERSITY

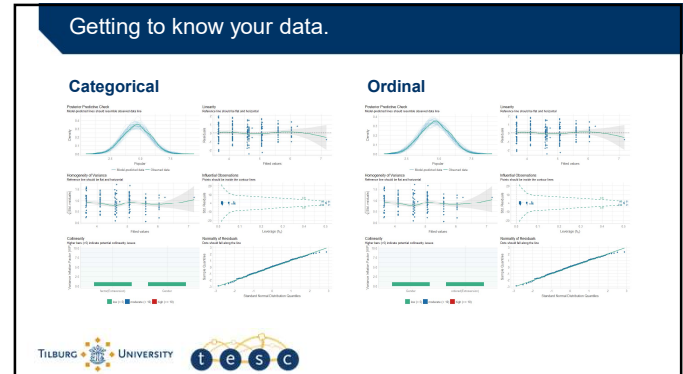


32

- Conclusions:
 - Homogeneity of Variance not ideal, but all else ok.
- Now what?



33



34

Robustness

- If you really worry about something in your data not matching your model....change the model.
- Heteroscedasticity? → Sandwich estimator, MELSM
- Non-normality → Choose different distribution
- Outliers → Choose t-distributed residuals
- Non-linearity → GAMM/GP/Polynomials
- Dependence → Model the dependence structure (e.g. with multilevel)

TILBURG UNIVERSITY

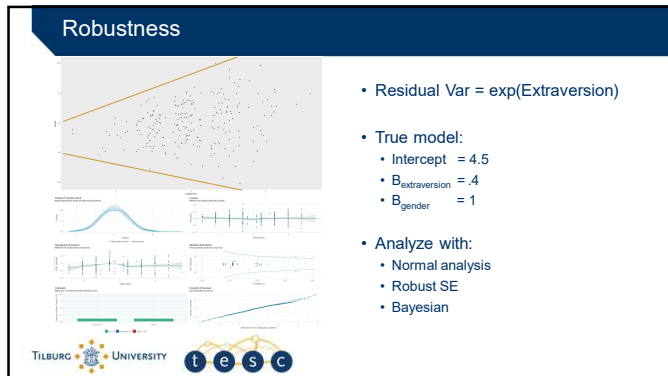
35

Robustness

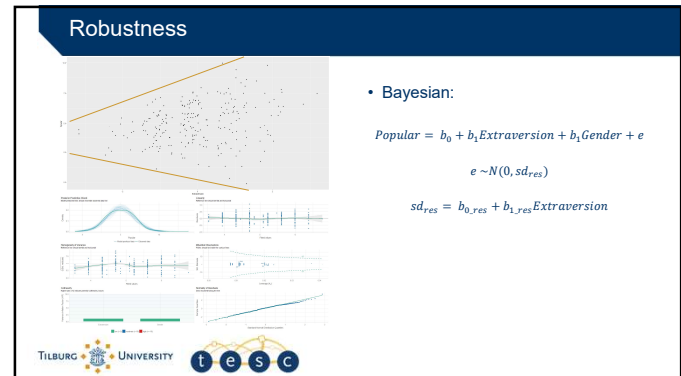
- Never forget...you don't have to model all aspect of your sample data!
- Only what is relevant!

TILBURG UNIVERSITY

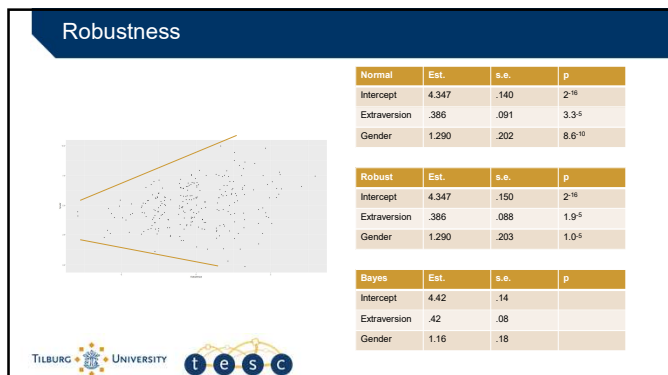
36



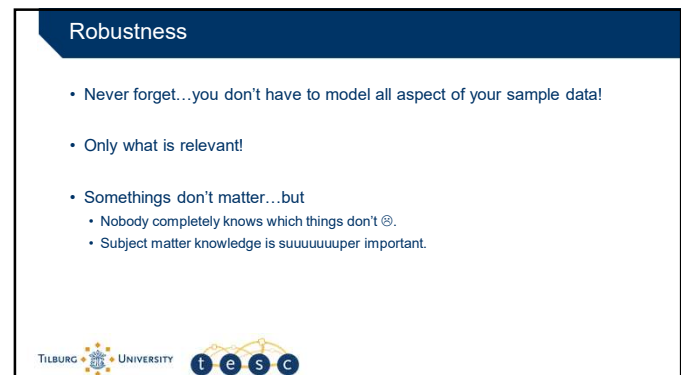
37



38



39



40

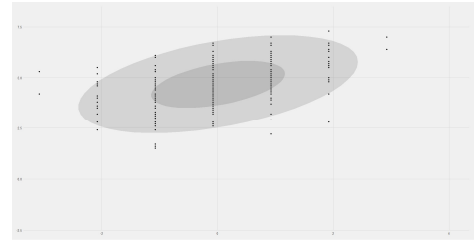
Robustness

- Heteroscedasticity? → Sandwich estimator, MELSM
- Non-normality → Choose different distribution
- Outliers → Choose t-distributed residuals
- Non-linearity → GAMM/GP/Polynomials
- Dependence → Model the dependence structure (e.g. with multilevel)

THESE ARE REASONS TO USE BAYESIAN STATISTICS/THE BRMS PACKAGE.

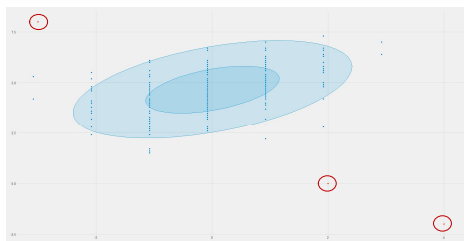
41

Example 2: Outliers



42

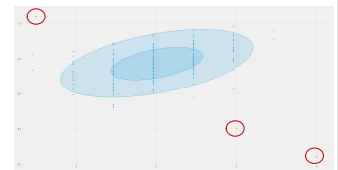
Example 2: Outliers



43

Example 2: Outliers

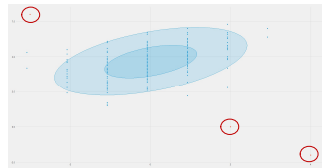
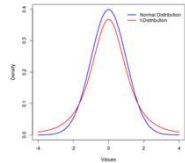
- The problem is that observations are further away than normal distribution "allows".



44

Example 2: Outliers

- The problem is that observations are further away than normal distribution "allows".



Example 2: Outliers

Original	Est.	s.e.	p
Intercept	4.206	.077	2 ⁻¹⁶
Extraversion	.430	.050	1.0 ⁻¹⁶
Gender	0.961	.111	6.6 ⁻¹⁶

Normal	Est.	s.e.	p
Intercept	4.160	.097	2 ⁻¹⁶
Extraversion	.255	.061	3.8 ⁻⁵
Gender	.964	.140	5.2 ⁻¹¹

Robust	Est.	s.e.	p
Intercept	4.160	.091	2 ⁻¹⁶
Extraversion	.255	.125	.043
Gender	.964	.140	4.9 ⁻¹¹

Bayes (t-dist)	Est.	s.e.	p
Intercept	4.22	.08	
Extraversion	.42	.05	
Gender	.96	.12	

Conclusion 1

- Data analysis is all about modeling!
- Your model needs to represent the underlying structure in your data (as close as possible)
- So, think and look(!!) at your data...forget what different 'types' of analyses you know at first.
- So basically forget the structure of every stats course ever and forget SPSS 😊

Conclusion 1 - Addendum

- Never forget...you don't have to model all aspect of your sample data!
- Only what is relevant!
- Somethings don't matter...but
 - Nobody completely knows which things don't ☹.
 - Subject matter knowledge is suuuuuuper important.

Practical 1

49

Multilevel Analysis

50

Multilevel Hierarchical Analyses

51

Recap

- Analyzing your data is all about modeling.
- If your data is linear, model it as such!
 - If it isn't...don't.

52

Recap

- Analyzing your data is all about modeling.
- If your data is linear, model it as such!
 - If it isn't...don't.
- If the residual variance is the same "across" the board, model it as such!
 - If it isn't...don't

53

Recap

- Analyzing your data is all about modeling.
- If your data is linear, model it as such!
 - If it isn't...don't.
- If the residual variance is the same "across" the board, model it as such!
 - If it isn't...don't
- If scores are normally distributed around the mean estimates....you get the idea

54

Dependent Data

- Hierarchical data, or data with dependence between the observations, is no different than outliers etc.
- If dependence is a characteristic of your data, model it.
- Maybe do it with "multilevel analysis" if needed/appropriate(!).

55

Example

- Let's start with an example,
- I have 58 pupils from three different classes
- Interested whether being more extraverted makes you more popular, and if there is a gender difference
- List of all the variables:
 - pupil*: pupil identification variable
 - class*: class identification variable
 - student-level independent variables: *extraversion* (continuous; higher scores mean higher extraversion) and *gender* (dichotomous; 0=male, 1=female)
 - popular*: continuous outcome variable at the student-level (higher scores indicate higher popularity).

56

Example

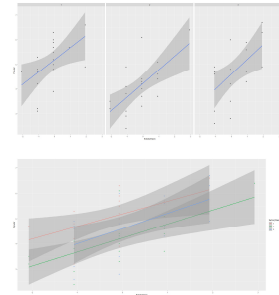
- How would you analyze this?
- Just mention all the different ways you can think of!
- No wrong answers!

Pupil	Class	Extraversion	Gender	Popular
1	1	-0,08	1	6,30
2	1	1,92	0	4,90
3	1	-1,08	1	5,30
4	1	-2,08	1	4,70
5	1	-0,08	1	6,00
6	1	-1,08	0	4,70
7	1	-0,08	0	5,90
8	1	-1,08	0	4,20
9	1	-0,08	0	5,20
10	1	-0,08	0	3,90

57

Some Options: Option 1

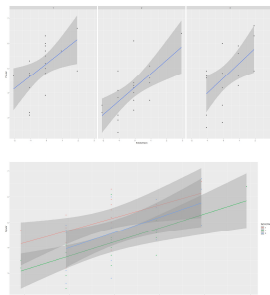
- Analyze all classes separately



58

Some Options : Option 1

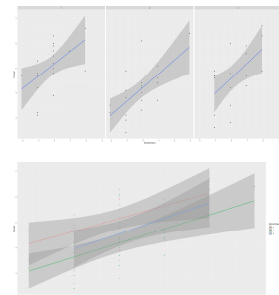
- Analyze all classes separately
- Benefits?
- Disadvantages?



59

Some Options : Option 1

- Analyze all classes separately
- Benefits?
- Disadvantages?
- How did we deal with the dependence in our data?



60

Richard McElreath

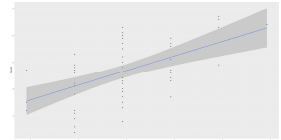
- *Many statistical models also have anterograde amnesia. As the models move from one cluster—individual, group, location—in the data to another, estimating parameters for each cluster, they forget everything about the previous clusters. They behave this way, because the assumptions force them to. Any of the models from previous chapters that used dummy variables to handle categories are programmed for amnesia. These models implicitly assume that nothing learned about any one category informs estimates for the other categories—the parameters are independent of one another and learn from completely separate portions of the data. This would be like forgetting you had ever been in a café, each time you go to a new café. Cafés do differ, but they are also alike*



61

Some Options: Option 2

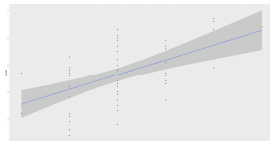
- Ok...so...analyze together then?



62

Some Options: Option 2

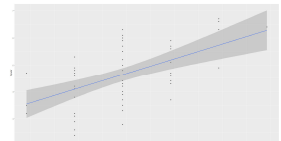
- Ok...so...analyze together then?
- What is the main issue in doing so?



63

Some Options: Option 2

- Ok...so...analyze together then?
- What is the main issue in doing so?
- Hint: How much data do we have?



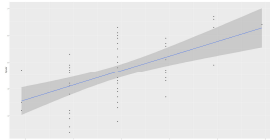
64

Some Options: Option 2

- Hint: How much data do we have?

$$n_{eff} = \frac{n}{1 + (n_{clus} - 1)\rho}$$

- Ok, so we have less data than data points.
- Now what?



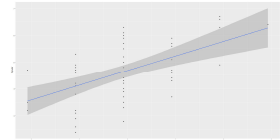
65

Some Options: Option 2

- Adjust se's

$$v_{eff} = v(1 + (n_{clus} - 1)\rho)$$

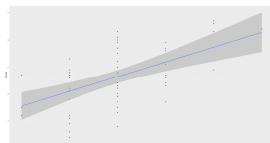
- Fortunately, smart ways to do this.
- Cluster robust s.e's



66

Some Options: Option 2

- Advantages?
- Disadvantages?



67

Some Options: Option 2

- Pro-tip: Use the formula below in combination with G*power for power analyses!

$$n_{eff} = \frac{n}{1 + (n_{clus} - 1)\rho} \longrightarrow n = n_{eff} * (1 + (n_{clus} - 1)\rho)$$

- This way you can do a power analysis even if you don't have all the information needed for a simulation study.
- G*power gives you the n_{eff} you need

68

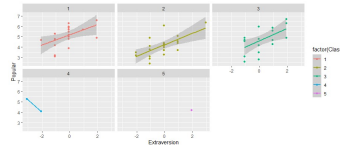
Some Options: Option 3 ---- Pool(ing) Party

- Who ever heard of partial-pooling?
- I'll add two classes to our dataset to illustrate it better
 - One class with two pupils
 - One class with one pupil

69

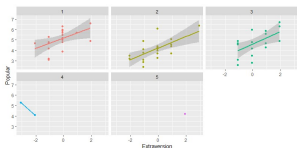
Pool(ing) Party

- What would be a problem if I analyzed everyone as N=1?



70

Pool(ing) Party

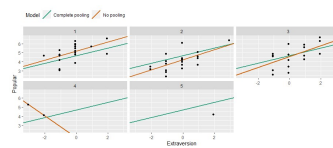


- Each panel shows an independently estimated regression line.
- This approach of fitting a separate line for each class is sometimes called the no pooling model
- Information from different classes is NOT combined or pooled together.

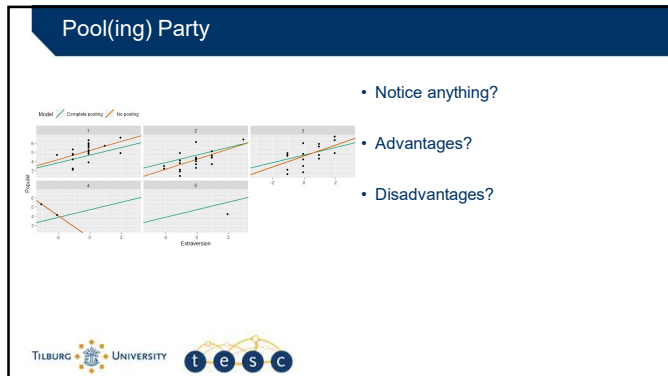
71

Pool(ing) Party

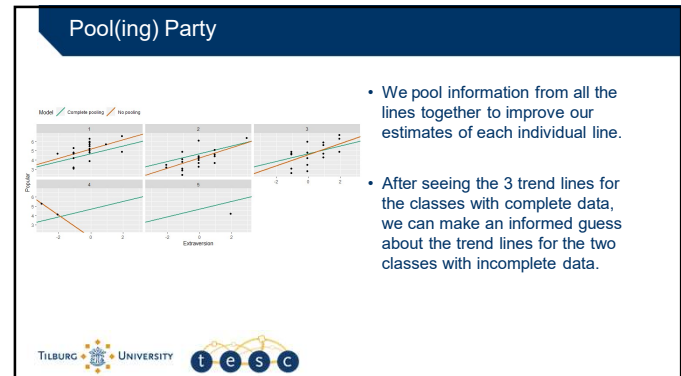
- Could also combine information from all classes (like we saw earlier).
- Fit a single line for the combined data set, unaware that the data came from different participants.



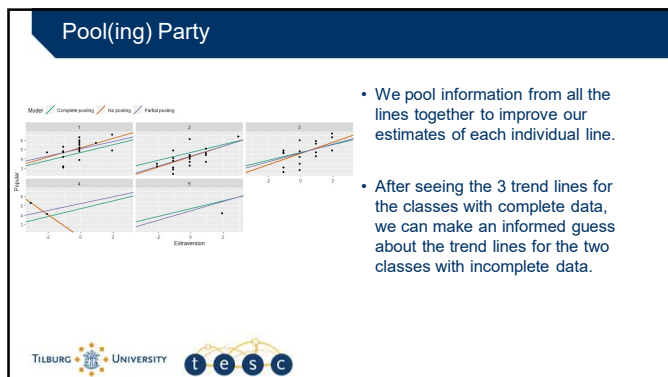
72



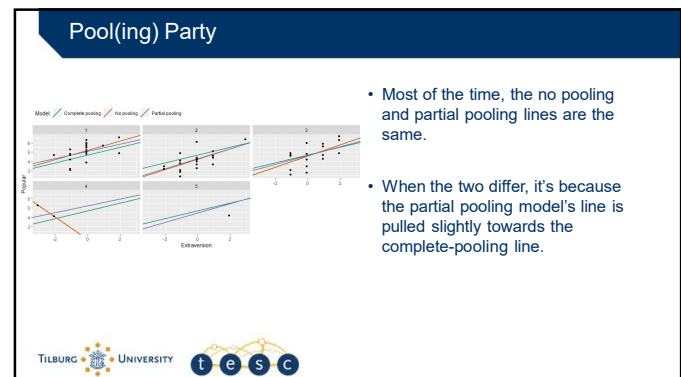
73



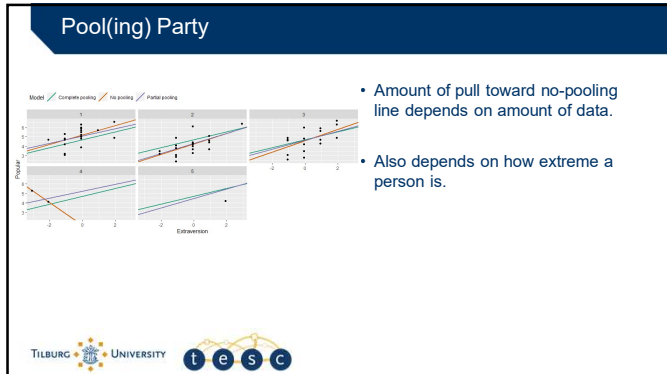
74



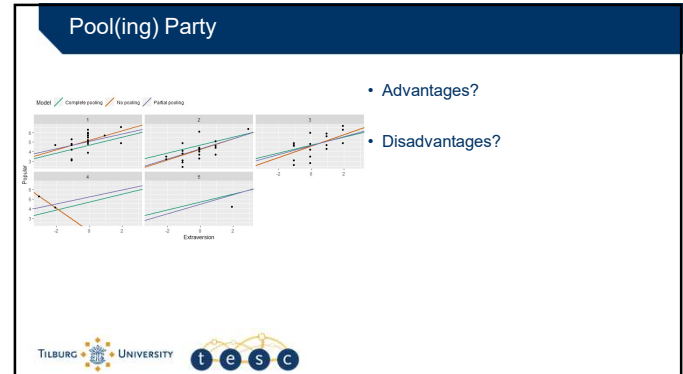
75



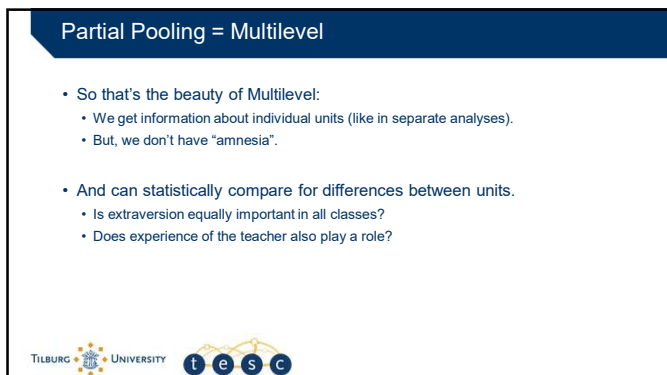
76



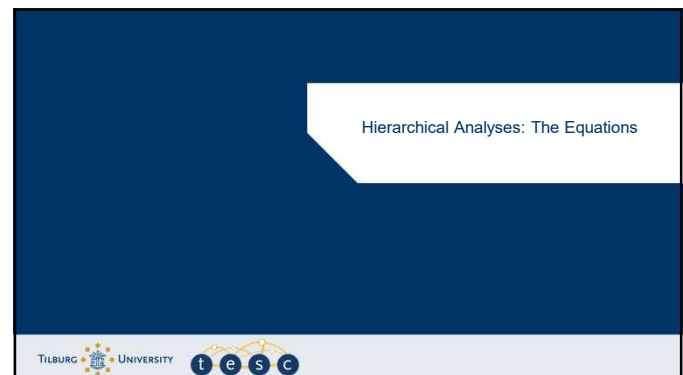
77



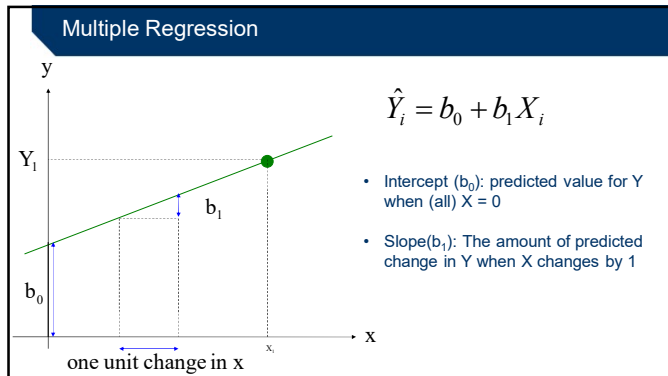
78



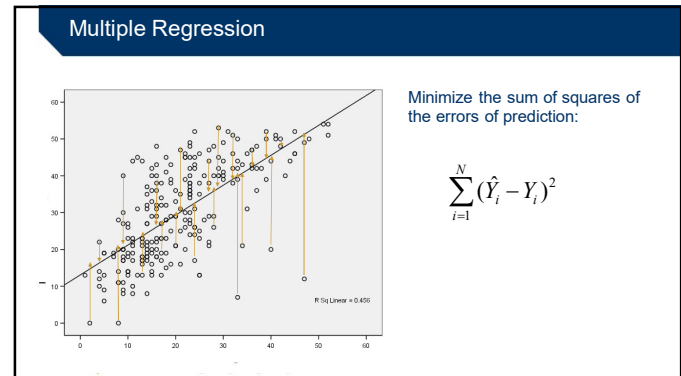
79



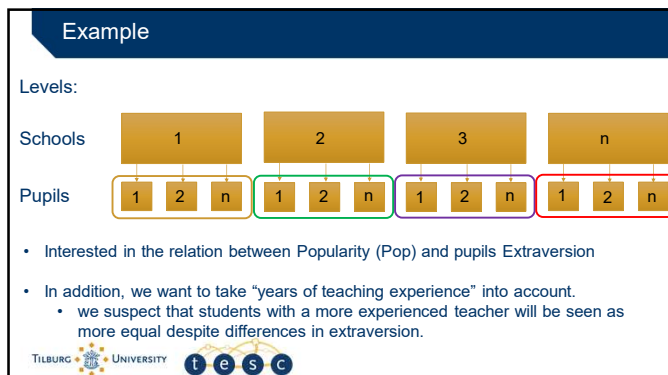
80



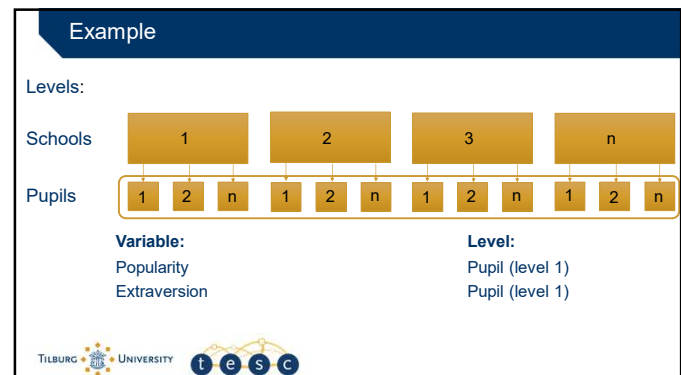
81



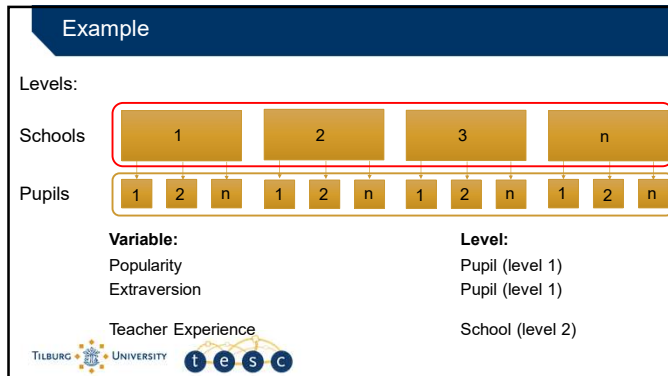
82



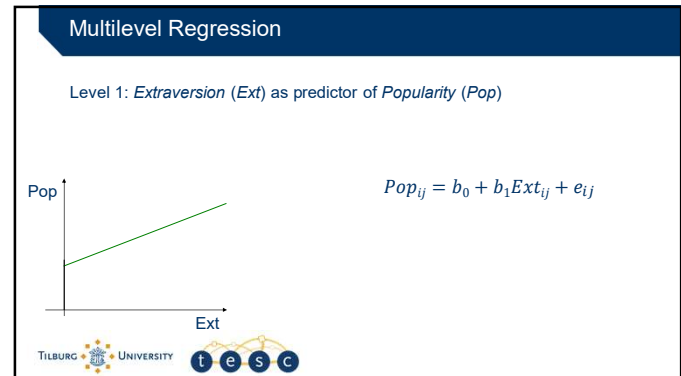
83



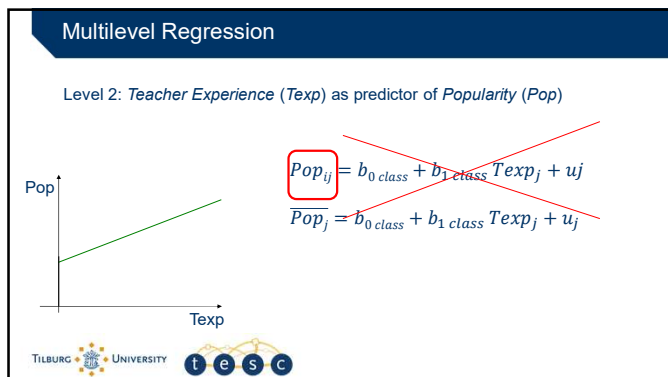
84



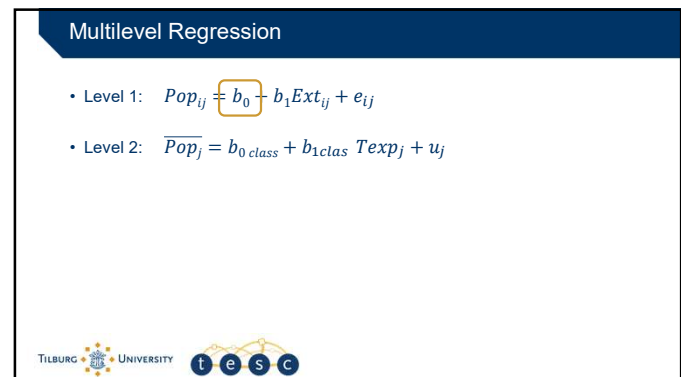
85



86



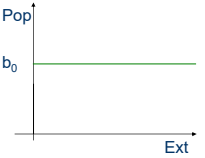
87



88


Multilevel Regression

- Level 1:



$$Pop_{ij} = b_0 + e_{ij}$$

$$b_0 = \overline{Pop}$$


TILBURG UNIVERSITY 

89

Multilevel Regression

Level 1: $Pop_{ij} = b_0 + b_1 Ext_{ij} + e_{ij}$

Level 2: $\overline{Pop}_j = b_{0\ class} + b_{1\ class} Texp_j + u_j$


TILBURG UNIVERSITY 

90

Multilevel Regression

Level 1: $Pop_{ij} = b_{0j} + b_1 Ext_{ij} + e_{ij}$

Level 2: $b_{0j} = b_{0\ class} + b_{1\ class} Texp_j + u_j$

TILBURG UNIVERSITY 

91

Multilevel Regression


Level 1: $Pop_{ij} = b_{0j} + b_1 Ext_{ij} + e_{ij}$

Level 2: $b_{0j} = b_{0\ class} + b_{1\ class} Texp_j + u_j$

Combined:

$$Pop_{ij} = b_{0j} + b_1 Ext_{ij} + e_{ij}$$

$$= b_{0\ class} + b_1 Ext_{ij} + b_{1\ class} Texp_j + e_{ij} + u_j$$

TILBURG UNIVERSITY 

92

Multilevel Regression

Level 1: $Pop_{ij} = b_{0j} + b_{1j}Ext_{ij} + e_{ij}$

Level 2: $b_{0j} = \gamma_{00} + \gamma_{01}Texp_j + u_{0j}$

Combined:

$$Pop_{ij} = b_{0j} + b_{1j}Ext_{ij} + e_{ij}$$

$$= \gamma_{00} + b_{1j}Ext_{ij} + \gamma_{01}Texp_j + e_{ij} + u_{0j}$$

93

Multilevel Regression

Level 1: $Pop_{ij} = b_{0j} + b_{1j}Ext_{ij} + e_{ij}$

Level 2: $b_{0j} = \gamma_{00} + \gamma_{01}Texp_j + u_{0j}$

$$b_{1j} = \gamma_{01} + u_{1j}$$

Combined:

$$Pop_{ij} = b_{0j} + b_{1j}Ext_{ij} + e_{ij}$$

$$= \gamma_{00} + \gamma_{01}Ext_{ij} + \gamma_{01}Texp_j + e_{ij} + u_{0j} + u_{1j}Ext_{ij}$$

94

Multilevel Regression

- Can you relate these equations back to partial pooling?

95

Practical 2

96