

SOC542 Statistical Methods in Sociology II

Data structures

Thomas Davidson

Rutgers University

April 25, 2022

Course updates

- ▶ Presentations next week
 - ▶ 10 minutes to present replication project
 - ▶ Introduction
 - ▶ Main replication
 - ▶ Sensitivity checks
 - ▶ Bayesian replication
 - ▶ Conclusions

Plan

- ▶ Violations of regression assumptions
- ▶ Robust and clustered standard errors
- ▶ Fixed effects
- ▶ Random effects
- ▶ Space, time and social structure

Violations of regression assumptions

IID and heteroskedasticity

- ▶ Our approach to regression modeling has been based on the assumption that our data are independently and identically distributed
 - ▶ e.g. Random samples from a known population
- ▶ In practice, this assumption is often violated
 - ▶ Groups with different distributions
 - ▶ Non-independent observations
- ▶ OLS assumes that residuals are homoskedastic, but observed data often have heteroskedastic structures.
 - ▶ This is particularly common if data are sampled from different groups with variation in the underlying data generation process.

Violations of regression assumptions

Impact on standard errors

- ▶ Confidence intervals that are too narrow too narrow
- ▶ More likely to commit Type I errors (false positives) by incorrectly rejecting the null hypothesis
- ▶ Inaccurate description of a plausible range of effect sizes

Violations of regression assumptions

Robust and clustered standard errors

- ▶ Adjust standard errors to account for violations of assumptions
 - ▶ **Robust/Heteroskedasticity consistent** standard errors
 - ▶ **Clustered standard errors** can be used to account for particular types of grouping

Robust and clustered standard errors

Intuition

- ▶ Variance component of the model is *inconsistent* due to heteroskedasticity or other model misspecification
 - ▶ This implies that we will not converge on the true population parameter, even with large samples.
- ▶ Corrections can be applied to variance components using a **sandwich** estimator.¹

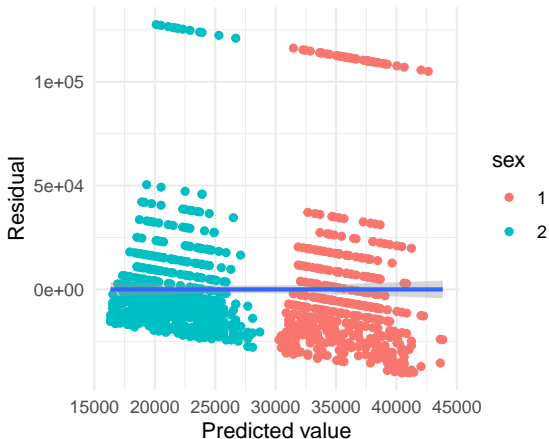
¹See King and Roberts 2015 for further technical discussion.

Robust and clustered standard errors

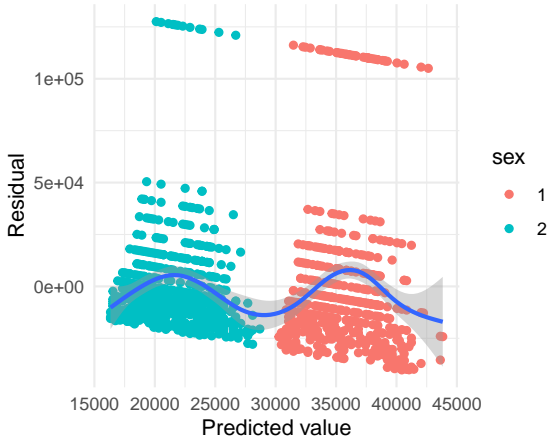
Estimating a simple model: Income as a function of age and sex (GSS 2020)

```
m <- lm(realrinc ~ age + sex, data = gss2020)
```


Heteroskedastic residuals



Heteroskedastic residuals



Robust and clustered standard errors

Calculation in R

There is no need to re-estimate the model. Robust standard errors can be calculated using `sandwich::vcocHC`. The `lmtest::coeftest` function allows us to easily apply the function and format the adjusted model for presentation.²

²[Grant McDermott's blog](#) has an excellent walkthrough of standard error adjustments using this function.

Robust and clustered standard errors

Clustering by sex

We can use the same function to apply other kinds of standard error correction. For example, we could cluster the errors by sex (although this is not warranted in this case).

```
m.r.g <- coeftest(m, vcov = vcovCL(m, cluster = ~ sex))
```

Robust and clustered standard errors

	OLS	OLS (robust)	OLS (clustered)
(Intercept)	26285.153*** (3458.317)	26285.153*** (3063.128)	26285.153*** (1929.282)
age	199.354** (66.361)	199.354** (61.976)	199.354*** (40.463)
sex2	-13940.750*** (1915.327)	-13940.750*** (1961.312)	-13940.750*** (68.045)
Num.Obs.	1077	1077	1077
Log.Lik.	-12675.353	-12675.353	-12675.353

+ $p < 0.1$, * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Robust and clustered standard errors

Caveats

- ▶ Robust and clustered standard errors have become extremely popular and are often the default approach in applied econometrics
 - ▶ Stata makes it particularly easy to specify them: `reg x y, robust`
- ▶ But standard error corrections are not a panacea and do not address underlying issues with model misspecification, as King and Roberts (2015) demonstrate.

Fixed effects

- ▶ **Fixed effects** are a useful tool for dealing with unobservables and reducing the threat of omitted variable bias when data have a grouping structure.
 - ▶ We previously used fixed-effects to account for any unexplained village-level factors when using the Diffusion of Microfinance dataset.
- ▶ Adding dummy variables for each group “soaks up” *between-group* variation, allowing us to estimate *within-group* effects.

Fixed effects

- ▶ A fixed-effects model can be written like a standard regression model. γ is a vector of coefficients, one dummy variable for each group.

$$y_i = \beta_1 x_i + \gamma_j + u_i$$

- ▶ It is common to drop the intercept from these models, although it is not required.

Fixed effects

Pooling

- ▶ **Pooling** refers to how observations are pooled together to estimate averages.
- ▶ Considering data with a grouped structure,
 - ▶ Standard regression approaches imply **complete pooling** since all available data to estimate a population mean.
 - ▶ Any variation between groups is effectively ignored.
 - ▶ Fixed-effects regression implies **no pooling** as a separate mean is estimated for each group.
 - ▶ No information is shared across groups. Assumption that variation between groups is effectively infinite.

Data and Methodology

- ▶ GSS panel data
 - ▶ Sample of 2016 and 2018 respondents were re-interviewed in 2020 (online)
 - ▶ Formatted data so each observation is one person-year

Data and Methodology

- ▶ Outcome
 - ▶ natcrime: Are we spending too much, about right, or too little on halting the rising crime rate?
 - ▶ Dichotomized (1 = too little, 0 = too much/about right)
- ▶ Predictors
 - ▶ Sex, age, race, political ideology
- ▶ Fixed effects
 - ▶ Survey year
 - ▶ Region
- ▶ Model
 - ▶ Linear probability model, listwise deletion to drop missing observations

Fixed effects

Implementation in R

We can easily specify fixed effects models using the `fixest` package.³

```
library(fixest)
ols <- lm(natcrime ~ sex + race + age + polviews, data = gss.new)
fe.r <- feols(natcrime ~ sex + race + age + polviews | region, data = gss)
fe.y <- feols(natcrime ~ sex + race + age + polviews | year, data = gss)
fe.ry <- feols(natcrime ~ sex + race + age + polviews | region + year,
```

³By default this model removes the main intercept from models with fixed effects.

Fixed effects

	Pooled	Region FE	Year FE	Both FE
(Intercept)	0.379*** (0.060)			
sex2	0.137*** (0.032)	0.147* (0.052)	0.139+ (0.039)	0.149* (0.052)
race2	0.143** (0.046)	0.117+ (0.057)	0.140+ (0.037)	0.113+ (0.057)
race3	0.190*** (0.057)	0.179** (0.044)	0.189+ (0.062)	0.177** (0.042)
age	0.004** (0.001)	0.004* (0.002)	0.004 (0.001)	0.004* (0.002)
polviewsConservative	0.077+ (0.045)	0.078* (0.032)	0.076+ (0.018)	0.077* (0.032)
polviewsLiberal	-0.119** (0.038)	-0.104+ (0.046)	-0.119* (0.016)	-0.104+ (0.046)
Num.Obs.	855	855	855	855
Log.Lik.	-550.539	-540.979	-549.875	-540.292

Fixed effects

Interpretation

- ▶ The fixed effects have accounted for unexplained variation between regions and over time, allowing us to measure the aggregate effect of our predictors on the dependent variable.
- ▶ In this case, there appears to be more regional variation than temporal variation
 - ▶ Perhaps this is not surprising since we are considering the same people over time

Fixed effects

Limitations of fixed effects

- ▶ No pooling
 - ▶ No information sharing across groups, only within-group variation analyzed
- ▶ Perfect multicollinearity
 - ▶ Time-invariant group-level variables are perfectly correlated with fixed effects and dropped from the model

Random effects

Comparing fixed and random effects

- ▶ Consider case where we observe random variables y and x , where observations belong to j groups.
- ▶ The fixed-effects formulation is given by

$$y_i = \beta_1 x_i + \gamma_j + u_i, u_i \sim N(0, \sigma_u^2)$$

- ▶ A random-intercepts model takes a more complex formula, where each element of γ_j is drawn from a distribution:

$$y_i = \beta_0 + \beta_1 x_i + \gamma_j + u_i$$

$$u_i \sim N(0, \sigma_u^2)$$

$$\gamma_j \sim N(0, \sigma_\gamma^2)$$

Random effects

Partial pooling and shrinkage

- ▶ The RE model considers the groups as related through a common distribution, whereas the entities in an FE model are unconnected.
- ▶ Random effects models are characterized by **partial pooling**
 - ▶ Information is shared among groups as intercepts are drawn from a common distribution.
- ▶ This tends to reduce overfitting compared to no pooling, since information in each group helps to improve estimates for every other group.
- ▶ **Shrinkage** describes how group-level estimates are pushed towards a common mean.
 - ▶ This is particularly helpful if there are small groups, where group means might be inaccurately estimated with a fixed effects model.

Random effects

Nesting

- ▶ Random effects models allow us to directly model more complex nested data structures
 - ▶ e.g. Education researchers might want to consider Level 1 (student), Level 2 (classroom), Level 3 (school), Level 4 (district)
- ▶ Unlike fixed effects, where all variance is explained by the fixed effect, variables can be incorporated at different levels
- ▶ Shrinkage/partial pooling helps to prevent overfitting

Random effects

A note on terminology

- ▶ These models are referred to using a range of different names including mixed effects, random effects, and hierarchical models. Moreover, the term “fixed effects” is also used in different ways, adding to the confusion.
- ▶ The “fixed part” or “population” component of a random effects model is the part that does not vary across groups.
 - ▶ e.g. $y_i = \beta_0 + \beta_1 x_i$
- ▶ The “random part” varies across groups
 - ▶ e.g. γ_i

Random effects

Estimation in R

The `lme4` package can be used to estimate Maximum Likelihood random effects models in R. `lmer` function can fit a standard model; `glmer` generalizes to other link functions. The random part of the model - in this case a random intercept for region - is specified in parentheses.

```
library(lme4)
re.r <- lmer(natcrime ~ sex + race + age + polviews + (1|region),
            data = gss.new)
re.r.logit <- glmer(natcrime ~ sex + race + age + polviews + (1|region)
                  data = gss.new, family = binomial)
```

Random effects

Estimating random coefficient models

Since we have multiple observations for each respondent, we could allow each respondent to have their own intercept. This kind of model would not be possible if we used fixed-effects.

```
re.r.id.logit <- glmer(natcrime ~ sex + race + age + polviews +  
                      (1 | region) + (1 | id),  
                      data = gss.new, family = binomial)
```

Random effects

	Region FE	Region RE	Logit	Logit + Resp
sex2	0.147 (0.052)	0.143 (0.033)	1.928 (0.297)	2.238 (0.481)
race2	0.117 (0.057)	0.130 (0.047)	1.920 (0.461)	2.399 (0.794)
race3	0.179 (0.044)	0.188 (0.057)	2.685 (0.838)	3.581 (1.476)
age	0.004 (0.002)	0.004 (0.001)	1.017 (0.006)	1.022 (0.008)
polviewsConservative	0.078 (0.032)	0.076 (0.045)	1.455 (0.328)	1.648 (0.479)
polviewsLiberal	-0.104 (0.046)	-0.113 (0.038)	0.602 (0.105)	0.531 (0.123)
(Intercept)		0.365 (0.063)	0.515 (0.151)	0.452 (0.176)
SD (Intercept)		0.060	1.250	1.139
SD (sex2)		0.060	1.250	3.159

View random intercepts

The random component of the model can be extracted using the `ranef` function. This shows the point estimates for the region level deviations from the population intercept.

```
ranef(re.r)
```

```
## $region
##      (Intercept)
## 1 -0.074394507
## 2  0.016011669
## 3  0.048694832
## 4 -0.009759633
## 5  0.050229324
## 6 -0.018106992
## 7  0.035204832
## 8 -0.058255524
## 9  0.010375999
##
## with conditional variances for "region"
```

Plot random intercepts

We can get more information by using the broom.mixed package:

```
library(broom.mixed)
broom.mixed::tidy(re.r, effects = "ran_vals", conf.int = TRUE) %>%
  mutate(across(where(is.numeric), round, 3)) %>%
  select(level, term, estimate, conf.low, conf.high) %>%
  head(5) %>% kable()
```

level	term	estimate	conf.low	conf.high
1	(Intercept)	-0.074	-0.165	0.016
2	(Intercept)	0.016	-0.057	0.089
3	(Intercept)	0.049	-0.011	0.109
4	(Intercept)	-0.010	-0.096	0.077
5	(Intercept)	0.050	-0.007	0.107

Random effects

Random coefficients

- ▶ In addition to random intercepts, we can also allow the slopes to vary by group.
- ▶ For example, we might want to see whether the effect of sex on attitudes varies across regions.
- ▶ Such a model includes the population coefficient, β_{sex} and a group-level deviation $\gamma_{j,sex}$.

Random effects

Estimating random coefficient models

We can easily modify the formula to include random slopes. In this case, we allow the slope of race to vary according to the region.

The control argument is included due to estimation issues.⁴

```
rc.logit <- glmer(natcrime ~ sex + race + age + polviews + (1 + sex|region,
                  data = gss.new, family = binomial,
                  control= glmerControl(optimizer="bobyqa", optCtrl=list(
```

⁴Warnings suggest potential problems with the model fit that require more detailed exploration.

Random effects

	Region RE	Region RE & Race RC
(Intercept)	0.515* (0.151)	0.554* (0.164)
sex2	1.928*** (0.297)	1.866* (0.455)
race2	1.920** (0.461)	1.973** (0.474)
race3	2.685** (0.838)	2.822*** (0.885)
age	1.017** (0.006)	1.017** (0.006)
polviewsConservative	1.455+ (0.328)	1.442 (0.327)
polviewsLiberal	0.602** (0.105)	0.612** (0.108)
Num.Obs.	855	855
Log.Lik.	-522.744	-520.561

Random effects

	Region RE	Region RE & Race RC
(Intercept)	0.515* (0.151)	0.554* (0.164)
sex2	1.928*** (0.297)	1.866* (0.455)
sd__(Intercept)	1.250	1.198
cor__(Intercept).sex2		0.509
sd__sex2		1.674

+ $p < 0.1$, * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Random effects

```
sex.slopes <- broom.mixed::tidy(rc.logit, effects = "ran_vals", conf.in  
  mutate(across(where(is.numeric), round, 3)) %>%  
  filter(term == "sex2") %>%  
  select(estimate, conf.low, conf.high)  
sex.slopes %>% head(5) %>% kable()
```

estimate	conf.low	conf.high
-0.728	-1.435	-0.020
0.225	-0.399	0.850
0.415	-0.152	0.983
0.465	-0.311	1.240
0.195	-0.302	0.693

Random effects

Extracting random slopes

```
fixed.coef <- broom.mixed::tidy(rc.logit, effects = "fixed", conf.int =  
  mutate(across(where(is.numeric), round, 3)) %>%  
  filter(term == "sex2") %>%  
  select(estimate)  
sex.slopes %>% mutate(sex_region = estimate + fixed.coef$estimate) %>%  
  select(sex_region) %>% head(5) %>% kable()
```

sex_region
-0.104
0.849
1.039
1.089
0.819

Random effects

Model selection

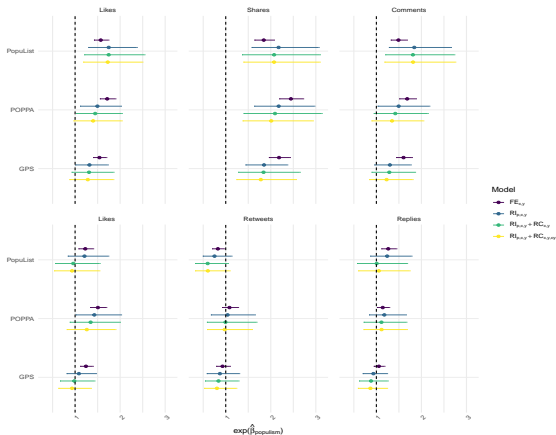
- ▶ There is some debate over how to decide between fixed and random effects specifications.
- ▶ One convention is to use a test to identify whether random effects should be included over fixed effects, but this has been questioned (Bell and Jones 2019).
- ▶ Random effects are a more flexible approach to capture complex structures (McElreath 2020), but are not as parsimonious as fixed effects specifications.

Random effects

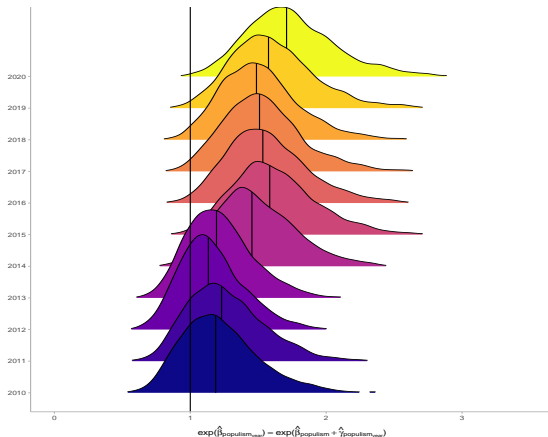
Advanced multilevel modeling

- ▶ Cross-level interactions can reveal relationships between different levels
 - ▶ e.g. In a model to predict child's test scores, one could interact child-level and school-level variables
- ▶ The “within-between” decomposition approach allows effects to be disentangled within and between units (see Bell and Jones 2019)
- ▶ Bayesian hierarchical modeling offers a more stable approach to complex models than MLE
 - ▶ `brms` uses the same syntax as `lme4` for model specification

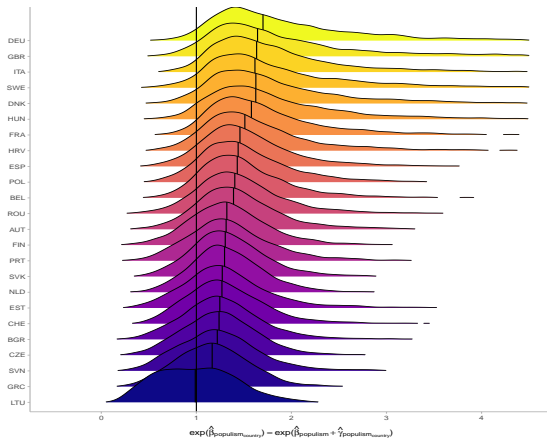
Example: Populism and social media engagement



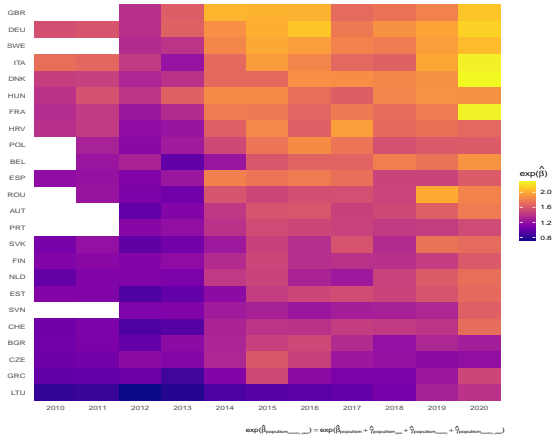
Example: Populism and social media engagement



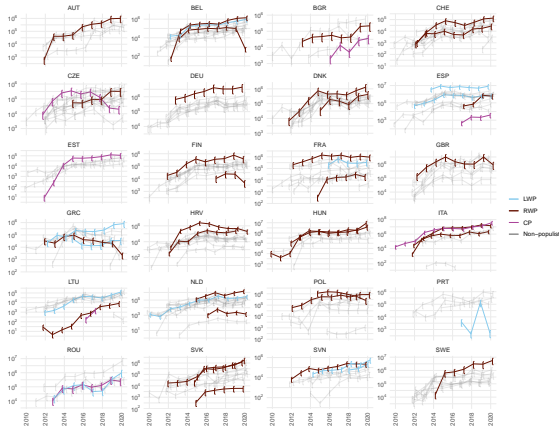
Example: Populism and social media engagement



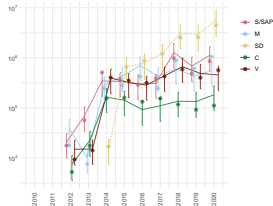
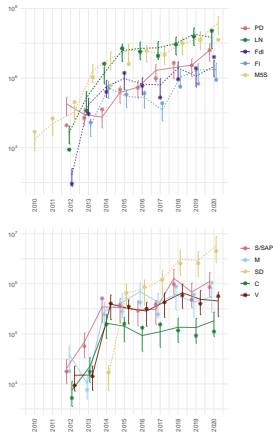
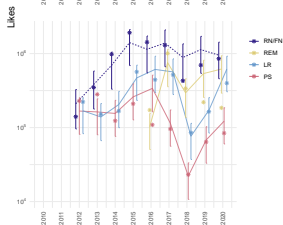
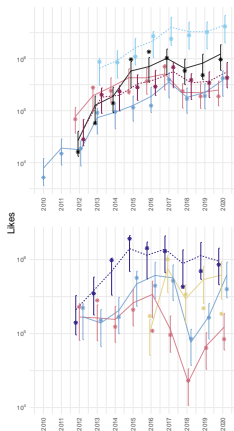
Example: Populism and social media engagement



Example: Populism and social media engagement



Example: Populism and social media engagement



Space, time and social structure

Autocorrelation

- ▶ Autocorrelation implies that something is correlated with itself
- ▶ Violation of IDD assumption
- ▶ Unlikely to be an issue when using randomly sampled cross-sectional data, but a problem in many applied settings

Space, time and social structure

Types of autocorrelation

- ▶ Temporal autocorrelation is the most typical case, where measurements are correlated with time
 - ▶ e.g. Given quarterly GDP, we expect high correlation between GDP_t and GDP_{t-1}
- ▶ Spatial autocorrelation implies that measurements are correlated with spatial proximity
 - ▶ e.g. County-level GDP more similar between proximate counties than distant ones.
- ▶ Network autocorrelation implies that measurements are correlated with network position
 - ▶ This is typically a problem if we want to sample measurements from individuals who have some relationship with one another
 - ▶ e.g. Children in a classroom who are friends are more likely to have similar interests than children who are not friends

Space, time and social structure

When to be concerned?

- ▶ Repeated measurements
 - ▶ Temporal autocorrelation
- ▶ Spatial structure to measurements
 - ▶ Spatial autocorrelation
- ▶ Non-random or network sampling
 - ▶ Network autocorrelation

Space, time and social structure

When to be (even more) concerned?

- ▶ Repeated measurements of spatial units
 - ▶ Spatial and temporal autocorrelation
- ▶ Repeated measurements of social networks
 - ▶ Spatial and network autocorrelation

Space, time and social structure

Solutions

- ▶ Standard error corrections
 - ▶ Appropriate error structures
- ▶ Fixed and random effects
 - ▶ Directly model data structure
- ▶ Data processing
 - ▶ e.g. De-trending and de-seasoning time series variables
- ▶ Model specification
 - ▶ e.g. Lagged variables, differences
- ▶ More advanced approaches
 - ▶ ERGM and SAOM models for networks

Space, time and social structure

Takeaways

- ▶ Standard GLMs alone are often insufficient to account for the way data are structured
- ▶ Standard error corrections are often necessary, but not a panacea
- ▶ Fixed effects and random effects models allow structure to be modeled in different ways
- ▶ More complex types of structure and dynamics should be directly modeled to avoid misleading inferences

Summary

- ▶ IID assumptions often violated and data structures must be accounted for
- ▶ A variety of statistical techniques can be used to explicitly model these structures
 - ▶ Standard error corrections
 - ▶ Fixed effects
 - ▶ Random effects
 - ▶ Advanced model specifications