

Overview of three big cohort projects

ELSA, TILDA, & SHARE



English Longitudinal Study of Ageing (ELSA)

- 20 years, 19k participants
- 10 waves, \pm every 2 years
- samples representative to people aged 50+ living in private households in England
- roughly every other wave limited medical assessment and collection of biological data (blood, saliva)



English Longitudinal Study of Ageing (ELSA)

Topics

- Demographic – age, marital stat, children, education
- Economic – income, pension, employment, housing
- Health – physical & mental, health behaviours, height + weight, COVID
- Psychosocial – social isolation, well-being, self-perception, cognitive functions



How to access?

- <https://www.elsa-project.ac.uk/accessing-elsa-data>
- (free) access through UK Data Service
- account needed, affiliation with UCD helps

The Irish Longitudinal Study on Ageing (TILDA)

- 15 years, 20k+ observations
- 6 waves, \pm every 2 years
- samples representative to Ireland population aged 50+; household response rate 62%
- medical assessment and collection of biological data in Waves 1&3



Staidéar Fadaimseartha na
hÉireann um Dhul in Aois

The Irish Longitudinal
Study on Ageing

The Irish Longitudinal Study on Ageing (TILDA)

Topics

- Demographic – migration
- Economic – similar to ELSA
- Health – cardiovascular, executive functions, gait, balance
- Psychosocial – similar to ELSA



Staidéar Fadaimseartha na
hÉireann um Dhul in Aois

The Irish Longitudinal
Study on Ageing

How to access?

- <https://tilda.tcd.ie/data/accessing-data/>
- on-site (Dublin Center) hot desk with full data
 - €10k fee for full data access if you have available budget
- online access with partial data through Irish Social Science Data Archive (ISSDA)
 - signed data access form must be submitted by email

Survey of Health, Ageing and Retirement in Europe (SHARE)

- 20 years, 28 European countries and Israel, 600k+ observations
- 9 waves, \pm every 2 years
- samples theoretically representative to country population aged 50+; different sampling methods
- country-specific drop-off questionnaires with unique variables



Survey of Health, Ageing and Retirement in Europe (SHARE)

Topics

- largely overlapping with ELSA and TILDA
- specific SHARELIFE questionnaire in Wave 3 on timeline of close personal relationships, partnerships
- specific COVID questionnaires in Wave 8 and 9



How to access?

- <https://share-eric.eu/data/become-a-user>
- free for research purposes
- data access form must be submitted and approved to create an account
- conditions of use have extensive citation requirements
 - they change periodically, check both on start and end of your project

Watch out for variable availability throughout the projects



Gateway to Global Ageing Data

- <https://g2aging.org/hrd/overview>
- neat project, which records and harmonizes compatible variables across 11 big projects
- engine can search for comparable variables present in datasets
- very useful when working with data from multiple projects

Advanced statistical analyses

SURVIVAL ANALYSIS

LINEAR MIXED MODEL (MULTILEVEL)

A solid orange horizontal bar at the bottom of the slide.

Tips & tricks for working with large data

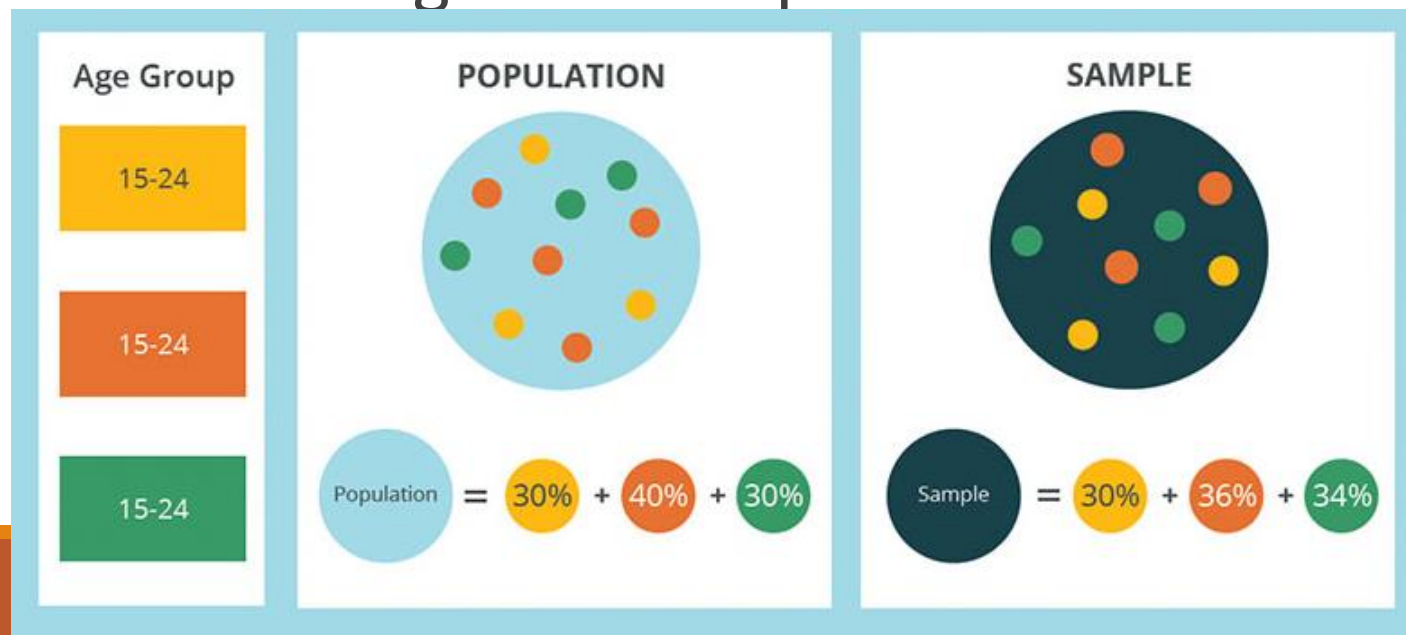
- always read the codebook / release guide first to get the sense of data
- do not rely on descriptions in data files
- often, recoded variables (e.g. total scale scores) are present

Table 9: Variables in the Technical Variables Module

Variable	Label
<i>fam_resp</i>	Family respondent
<i>fin_resp</i>	Financial respondent
<i>hou_resp</i>	Household respondent
<i>mn005_</i>	Single or couple interview
<i>mn016_</i>	Mother in household
<i>mn017_</i>	Father in household
<i>mn018_</i>	Mother-in-law in household
<i>mn019_</i>	Father-in-law in household
<i>mn024_</i>	Nursing home interview
<i>mn026_</i>	First respondent from couple or single
<i>mn028</i>	Eligible for dried blood spots collection (bs &gv_dbs)
<i>mn029</i>	Eligible for linkage
<i>mn030</i>	Eligible for social networks module (sn)
<i>mn031_</i>	Eligible for mini childhood module
<i>mn032_</i>	Eligible for social exclusion items
<i>mn038_</i>	Eligible for Accelerometry
<i>mn040_</i>	Need to ask consent question (ex123)
<i>mn041_</i>	Need to ask retirement info
<i>mn101_</i>	Questionnaire version (longitudinal vs. baseline)
<i>mn103_</i>	SHARELIFE life history interview (only w7)
<i>mn104_</i>	Household moved

Tips & tricks for working with large data

- have a strategy for dealing with participant missing data
- search for harmonized data if comparing countries
- weighting cases can support representativeness, but the logic of assigning your own weights is complex



Tips & tricks for working with large data

- learn how to merge datasets based on IDs and other variables
- learn how to switch wide and long data formats (Restructure in SPSS)

Wide Format

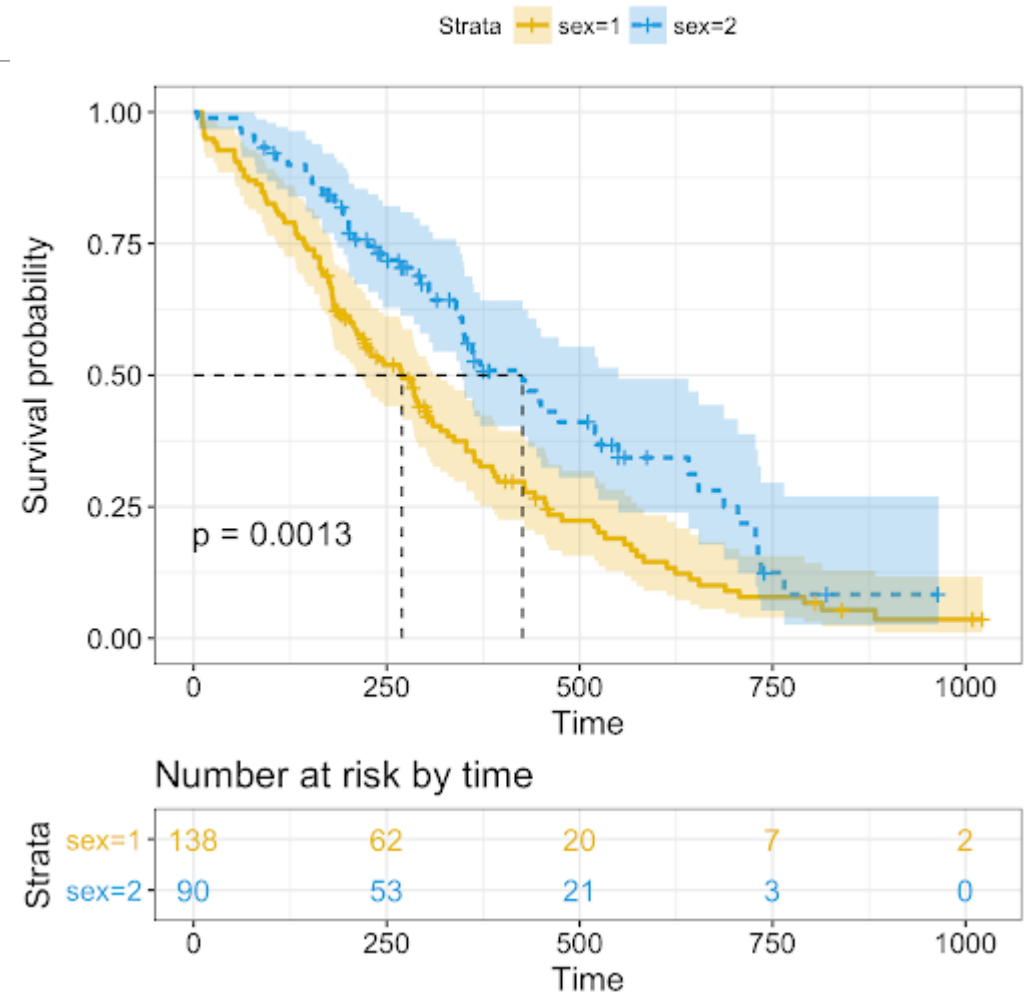
Team	Points	Assists	Rebounds
A	88	12	22
B	91	17	28
C	99	24	30
D	94	28	31

Long Format

Team	Variable	Value
A	Points	88
A	Assists	12
A	Rebounds	22
B	Points	91
B	Assists	17
B	Rebounds	28
C	Points	99
C	Assists	24
C	Rebounds	30
D	Points	94
D	Assists	28
D	Rebounds	31

Survival analysis

- Survival analysis operates with **time left** until a one-time event happens
- estimated odds are extrapolated across the whole observation period
- convertible to survival or hazard ratios
- Can handle multiple groups, or single-time covariates



Survival analysis

- similar to logistic regression

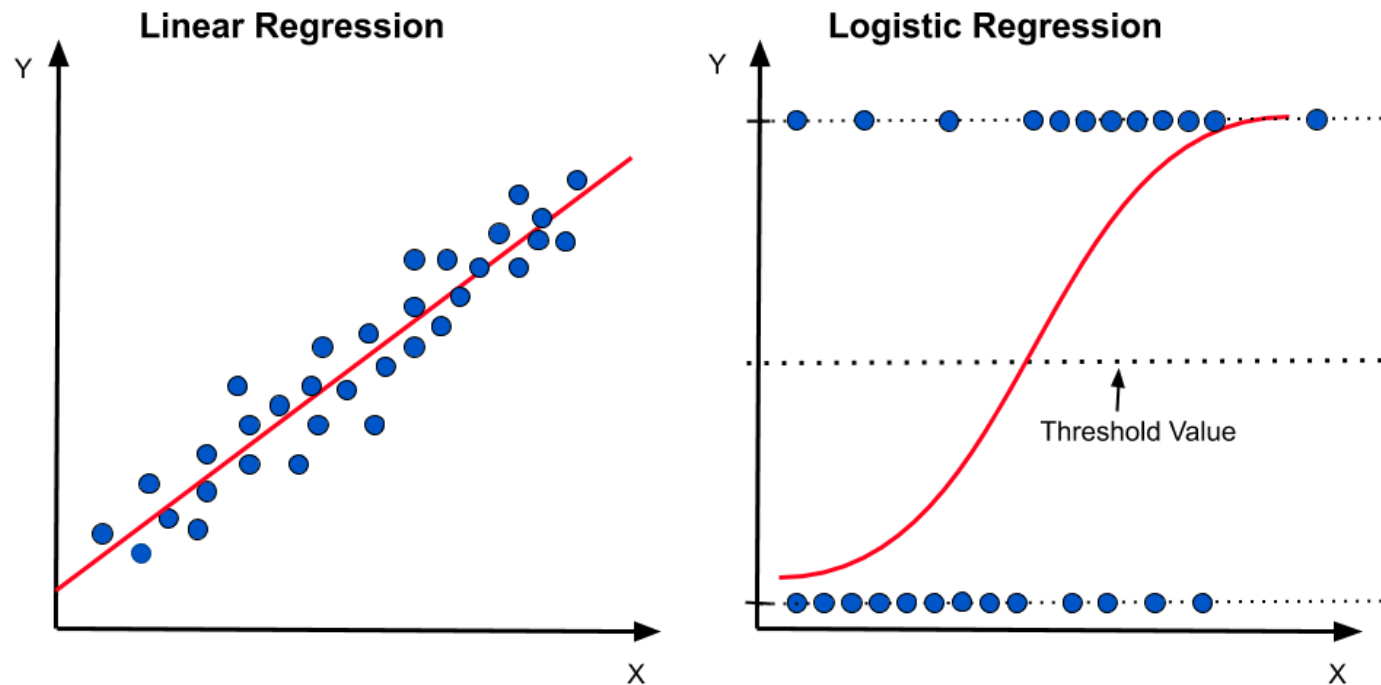


Image source:

<https://medium.com/@praveenraj.gowd/why-is-logistic-regression-called-logistic-regression-and-not-a-logistic-classification-5a418293040d>

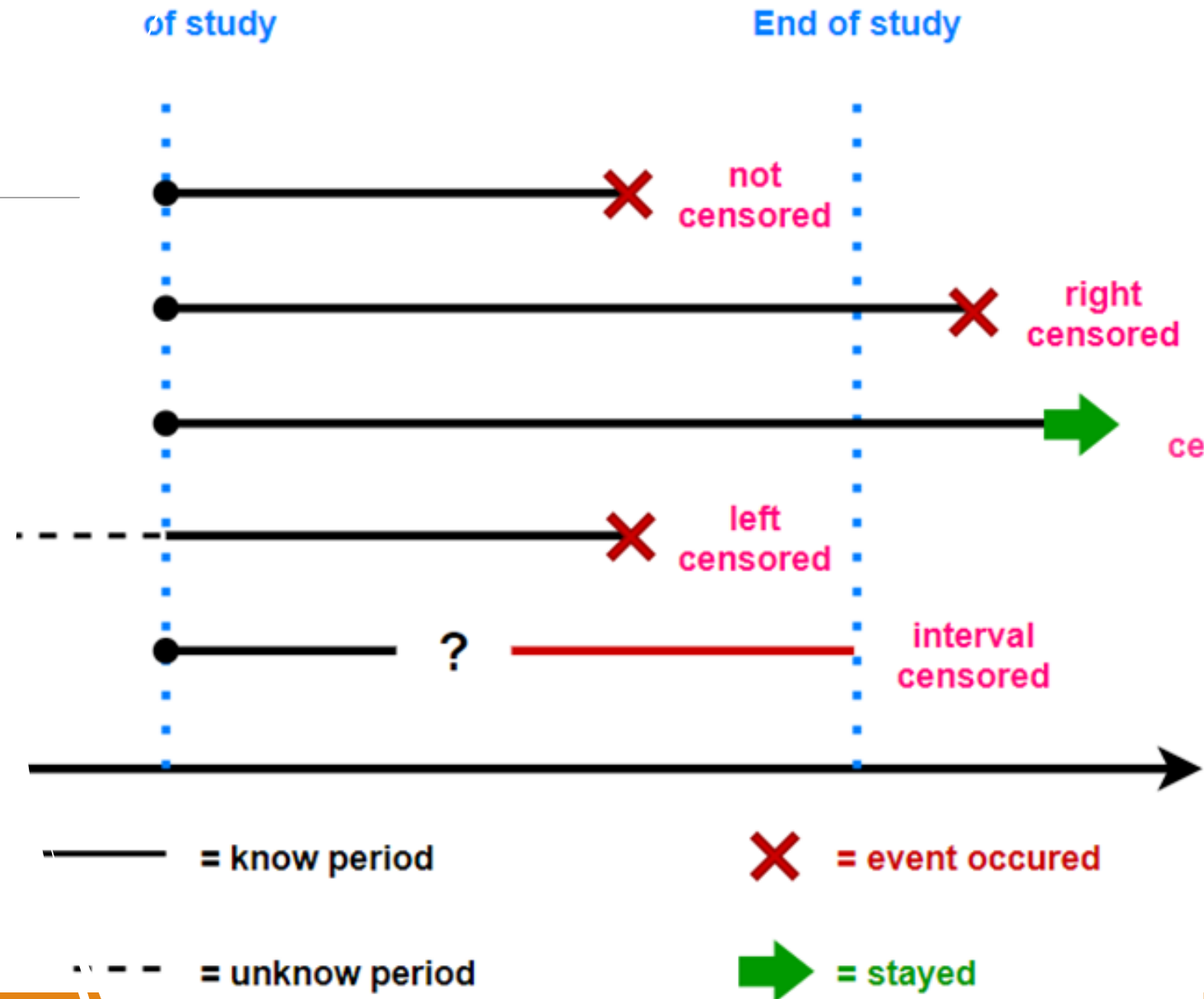
Survival analysis

- as in logreg, output are regression coefficients transformed into **Odds Ratios**
 - linreg coefficient(B) – if predictor value X increases by 1, by how many units does the Y change?
 - the coefficient value is additive; predicted value $Y = \text{intercept} + X \cdot (B)$
 - logreg coefficient(also B) – if predictor value X increases by 1, how do the odds of Y happening change?
 - the OR value is multiplicative; predicted odds $Y = \exp(\text{intercept}) * X \cdot \exp(B)$

Survival analysis - censorship

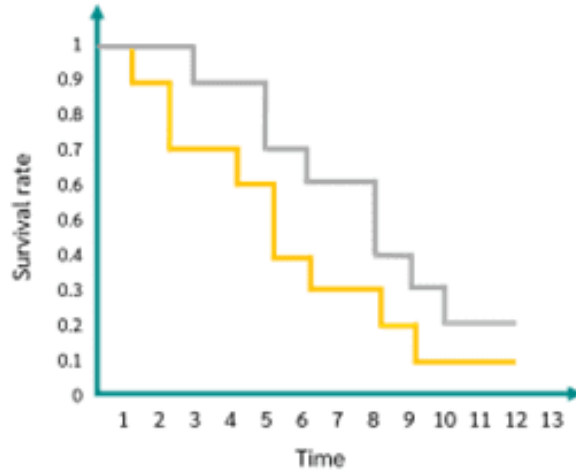
Ideally, we want clear
starting and ending
point of observing.

Problems with
generalizing out of
observed interval.

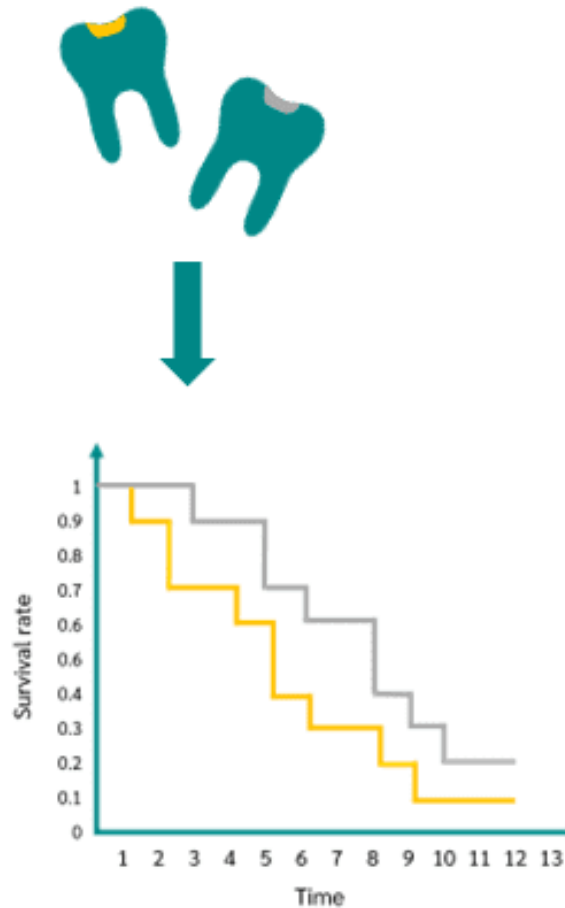


Survival analysis types

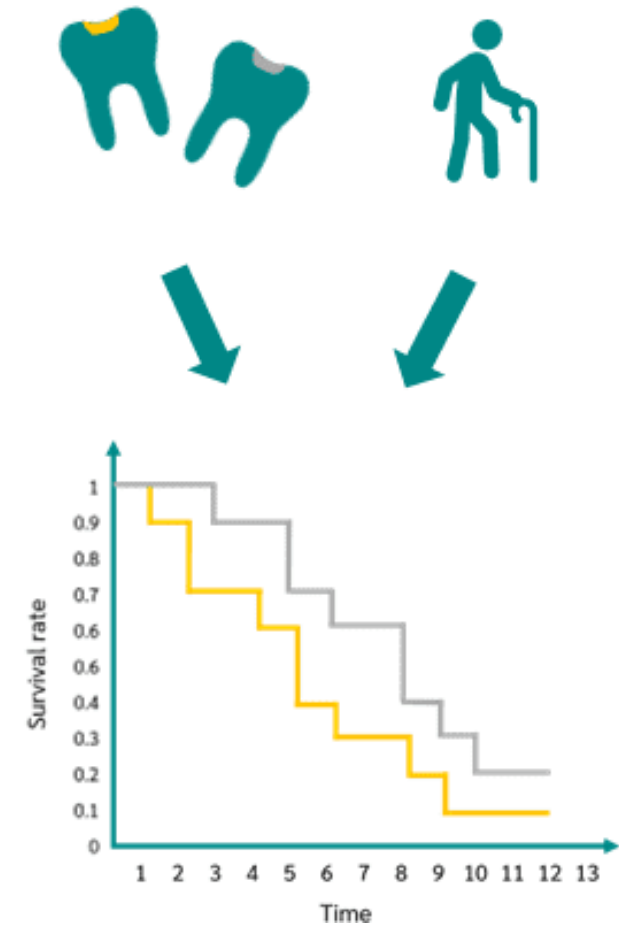
Kaplan Meier Curve



Log Rank Test



Cox Regression



Cox regression results example

Factor	Group	N	HR	90% CI		P value
Sex	Male	195	1.000 (Reference)			
	Female	960	1.115	0.975	1.276	0.1816
Age	age < 40	203	1.000 (Reference)			
	40 ≤ age < 60	815	1.043	0.912	1.193	0.6087
	age ≥ 60	137	1.066	0.879	1.294	0.5861
family history	Non history	115	1.000 (Reference)			
	HCC	596	0.915	0.814	1.029	0.2130
	Cirrhosis	341	0.971	0.812	1.162	0.7908
	HBV	103	0.952	0.801	1.133	0.6441
Compensation status	Decompensation	345	1.000 (Reference)			
	Compensation	810	1.114	0.998	1.244	0.1058
time lengths from cirrhosis treatment to HCC onset	Non-treated	627	1.000 (Reference)			
	Total	528	0.810	0.721	0.911	0.0004
	LAM	102	1.045	0.873	1.250	0.6877
	ADV	181	0.792	0.687	0.914	0.0072
	ETV	83	0.716	0.585	0.877	0.0068
	LAM+ADV	95	0.822	0.684	0.989	0.0808
	Others	67	0.731	0.589	0.906	0.0166

Source: Bi, Jingfeng & Zhang, Zheng & Qin, Enqiang & Hou, Jun & Liu, Shuiwen & Liu, Zengmin & Li, Shuo & Wei, Zhenman & Zhong, Yanwei. (2017). Nucleoside analogs treatment delay the onset of hepatocellular carcinoma in patients with HBV-related cirrhosis. Oncotarget. 8. 10.18632/oncotarget.18075.

Linear mixed (multilevel) model

- very flexible, allows for designs with multiple measurements per participant (e.g., longitudinal)
- works similarly to regression, yields regression coefficients with the same meaning
- can also be used for nonlinear associations or categorical dependent variables, those models are called generalized

Linear mixed (multilevel) model

Many statistical tests in psychology are simplified cases of complex mathematical modeling techniques



Generalized Linear
(Mixed) Model



Regression



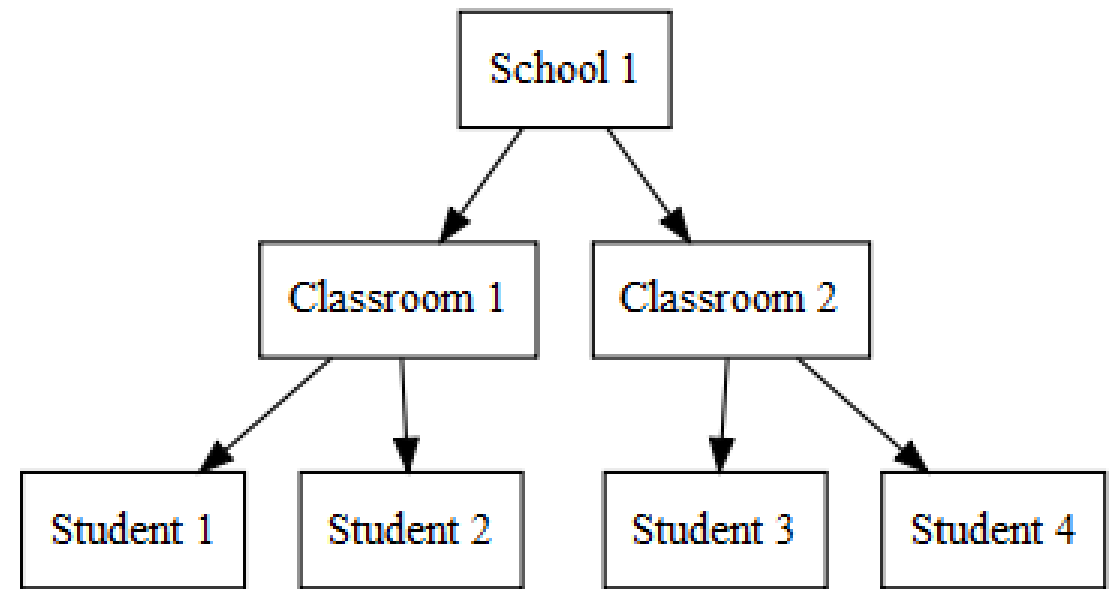
ANOVA



t-test

Multilevel models are similar to regression

- simple linear regression:
 - $Y \sim b_0 + b_1 * X$
- simple linear mixed model
 - $Y \sim b_0 + b_1 * X + (1 | Z)$
 - Z is random effect factor, across which we want to „average“ the results, usually grouping variable – participant ID, strata, etc.



Linear mixed model results example

```
Formula: testScore ~ bodyLength2 + (1 | mountainRange)
Data: dragons
```

Random effects:

Groups	Name	Variance	Std.Dev.
mountainRange	(Intercept)	339.7	18.43
	Residual	223.8	14.96

Number of obs: 480, groups: mountainRange, 8

The random effect part tells you how much variance you find among levels of your grouping factor(s), plus the residual variance

Fixed effects:

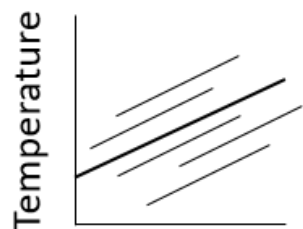
	Estimate	Std. Error	t value
(Intercept)	50.3860	6.5517	7.690
bodyLength2	0.5377	1.2750	0.422

The fixed effect part is very similar to a linear model output:

Intercept and error

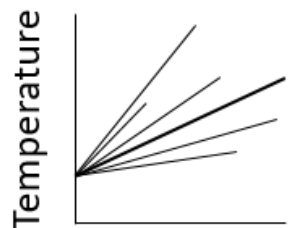
Slope estimate and error

Image source: <https://ourcodingclub.github.io/tutorials/mixed-models/>



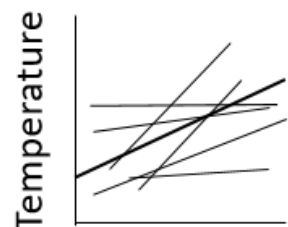
$$(1 | \text{SiteID})$$

- Different sites have different intercepts
- The slope of temperature over time is the same for all sites



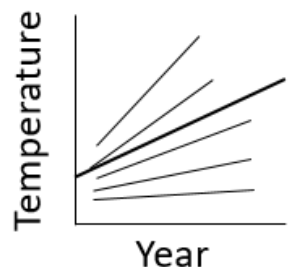
$$(0 + \text{Year} | \text{SiteID})$$

- The intercept is the same for all sites.
- The slope of temperature over time varies among sites



$$(0 + \text{Year} | \text{SiteID}) + (1 | \text{SiteID})$$

- No covariance between the intercept and the slope.
- Statistically easier to compute
- Is this realistic or not? If there is no reason to expect it then why model it? If you have a massive dataset then may as well model the covariance.



$$(1 + \text{Year} | \text{SiteID}) + (1 | \text{SiteID})$$

- If there is a positive covariance, then a lower intercept would mean that the slope for that site would be lower.
- Note: if the correlation value reported is 1 or -1 then it suggests there are issues calculating this in the model and it can't estimate it well.

Linear mixed (multilevel) model

- (RE)ML estimator, nested models easily comparable
- sensitive to misspecification of covariances
- able to construct complex models
- gives fixed effects, then the variability of dependent variable across groups (random effects)

Final Comparison

Survival Analysis

- **best** to model continuous time until the occurrence of a one-time event
- **bad** flexibility of the model, cannot handle groups within groups well or time-varying covariates; expects proportional hazard ratios
- requires specific data structure

Multilevel Analysis

- **great flexibility** in predicting dependent variable across multiple observations and/or groups; quite benevolent about missing data
- **needs clearly defined assumptions** regarding the covariance structure
- requires long data structure