

## Lecture 3: FITTING STATISTICAL MODELS TO DATA

### X X X X LECTURE 3: FITTING STATISTICAL MODELS TO DATA

QUESTION

WHAT DO WE MEAN BY FITTING MODELS TO DATA?

- Note that we're not fitting data to models but models to data.
- Our models must describe the relationship & dist. of the data.

Example: Test Performance & Age

Variables  
Dependent = Test performance  
Predictor = Age

Relationship stated:  $\text{perf} - \text{age}$  = curved relationship

Goals: (i) Estimate Marginal mean of performance across all ages

(ii) Est. Mean performance conditioned on age (Relationship of Age w/ Mean Perf)

(I) Model 1 = "Mean only" model for performance

2 Perform = Mean & variance; Normal dist.

(II) Model 2 = "Conditional" model for performance as a function of age

↳ Normal dist. Mean defined by a quad formula & variance ( $\sigma^2$ )

• The data, after examination, is normal  $\rightarrow$  that, Q-Q plot

• Viz relationship b/w Performance & Age  $\rightarrow$  shows inverse "U" shape

↳ scatterplot very useful

(III) "Mean Only" model = regression model  $\rightarrow$  higher test perf. on a single mean + error

↳  $\text{perf} = m + e$  | errors denote deviation of obs around the mean

• Check assumptions that errors are Normal. This is done by looking at residuals via hist. or Q-Q.

(IV) "Conditional" model

↳  $\text{perf} = a + b(\text{age}) + c(\text{age}^2) + e$

Params = a, b, c, d Regression coeff.

e = random error

↓  
describe relationship  
of age w/ performance

b  $\rightarrow$  estimate = linear portion of the equation (function)  
 C estimate  $\rightarrow$  describes acceleration / deceleration in function as a function of age

- Looking @ ratio of estimates to SE; all will seem to be non-zero

- After fitting model; remember to assess fit of model

| L Red dashed lines = predicted values based on the fitted quad. funct.

Assess fit by looking @ the residuals

| In the second plot, there a lot of residuals centered & spread along the predicted mean.

- Let's look @ a poor model.

| Predicted values of perf show a curvilinear relationship

| Note the red dashed lines || our assumption of a linear relationship

## WK1

### TYPES OF VARIABLES IN STATISTICAL MODELS

- There is a dichotomy of data when it comes to dependent variables & independent variables exists in data modelling (Model fitting)

(i) **DEPENDENT variables**  $\rightarrow$  They are defined by our research question

(ii) **INDEPENDENT variables**  $\rightarrow$  take are looking and relation with DV  $\in$  IVs

| prone to manipulation by an investigator or observed

- Observational studies are mostly focused on describing relationships, i.e. randomized experiments that enable us to make stat. inference.

using **CONTROL variables**  $\rightarrow$  Randomization helps deal w/ the confounding

| problem in observational studies

If we  $\infty$  an independent variable included into the model but is  $\perp$  theoretical perspective it

- helps us control its value when looking @ a relationship where the confounding variable is known

- Missing DATA  $\rightarrow$  We need to find out if the cases that are going to be dropped

| are systematically off from those ~~retained~~ retained

| imputation  $\rightarrow$  predict missing data as a function on known data.

WK 1

## DIFFERENT STUDY DESIGNS GENERATE DIFFERENT TYPES OF DATA: IMPLICATIONS FOR MODELING

- Like the title says: different study designs generate different types of data
- It's important to know where the data come from.

Goal = Estimate parameters that best describe the data & that vanish  
↳ there is a possibility that different values on a dependent variable of interest maybe correlate. Many times they might be introduced by our study design.

### WK 1 OBJECTIVES OF MODEL FITTING: INFERENCE vs PREDICTION

- In Reg. Objectives A model fitting
  - ① Making up about rel. between variables
  - ② Predictor of future outcomes
- Objective I → Making inference
  - Once we have estimates of param.  $\pm$  SE we can [hypothesis test]  
[conf. interval]
  - If  $c$  is significantly different from 0; there'd be evidence for linear relationship b/w performance & age.
- Objective II - Making predictions (regression model =  $5.11 + 0.2 \times \text{age}$ )
  - ↳ When making prediction from a fitted model, we also need  $-0.26(\text{age})^2 + e$  to describe the amount of uncertainty associated w/ the prediction.

### WK 1 → Learning Notes < BAYESIAN STATS >

- P-values, in frequentist stats, represent the probability of the data under the null hypothesis. That is not what scientists really want to know
- They want to know the opposite: the probability of the hypothesis given the data, this is what Bayesian stats helps us calculate

## Plot

### Poison Prediction & Prediction uncertainty

Watch out for only looking for point estimates when carrying out predictive modeling. This isn't all!

Always plot your uncertainty [the grey box around the line] u = slope  
b = intercept

Speaks of uncertainty of the slope

The ~~second~~ diagram has a lot of uncertainty about its parameters

2 ways to account for uncertainty

→ prediction bounds

(I) Plot data = To look @ variance

(II) Look @ standard errors = Speaks of uncertainty around the

Model 2 ~~needs~~ to work on ~~it~~ prediction estimates

caution -

## Wk 2

### Linear Regression - INTRODUCTION

→ quantitative

Primary outcome (dependent variable) = Cart Wheel Dist. (inches)

\* other variables in terms (correlation) → Height; Completed Cartwheels

Data shows reasonable normality. If you were asked to predict the distance for the next random adult, we could report the mean = 82.48 in.

But there are other variables that could influence what the CWDst could be.

So is there a relationship? We examine using a scatterplot.

When looking @ a scatterplot, we'd like to provide a summary that'll include

Form, slope, strength, Outliers

→ can be quantified by specifying  $r$  &  $r^2$

\* R-squared, the coeff of determination, shows how much Y is explained by X. So height only accounts for 11% of the variability in Cart wheel distance.

- Draw best line of fit; Estimate of  $\hat{y} = b_0 + b_1 x$   
Least Squares Line

$b_0$  = intercept,  $b_1$  = slope  $\rightarrow \Delta y \text{ over } \Delta x$ .

estimated response when  $x=0$ , not usually useful

### Ways of Convoy up with the Line & Best Fit

- ② Least Square Regress = Find the line that minimizes the total squared

$\langle \text{Actual} - \text{Predictor} \rangle$  (observed) error

→ intercept,  $b_0 = 7.55$ ; slope,  $b_1 = 1.11$

An adult who is 1 in  
than another adult would have  
a car that is 1.1 in  
longer on average

$$\text{Predicted CND} \approx 7.5518 + 1.1076(\text{Height})$$

- Extrapolation = any prediction that is outside the range & where we modeled the relationship
- Our regression model gives us an estimated mean response based on the value of  $x$
- Residuals = Observed  $\hat{x}$  - Predicted  $\hat{x}$  → Need to do some model checking.

## WK 2 LINEAR REGRESSION INFERENCE

Remember our question: "Is there a significant relationship b/w CND & age and height?"

- A slope of 0 doesn't help because we get same  $y$  on diff.  $x$
- If we try to plot the entire population on the scatterplot, it'll quickly become busy but we could come up with the data and come up w/ the true regression line
- There is a true intercept & slope. It is that  $\beta_i$  underlying true slope that we would want to assess i.e. if  $\beta_i > 0$

- Hypotheses,  $H_0$ : True Slope ( $\beta_i$ ) = 0

we have a  $b_i$   $\stackrel{\sim}{=} \beta_i$  and we need to find out how far away it is from 0.

- the standard err (0.670) estimates how far away  $b_i$  is from  $\beta_i$  on average
- the t-statistic (1.653) measures the distance of  $b_i$  from 0 in standard errors

- so we have a p-value of 0.112 for a two-tail test  $H_0: \beta_1 = 0$

Since we are dealing w/ a one-sided tail test for  $H_1: \beta_1 > 0$ ,

$$\text{the p-value} = 0.112/2 = 0.056 \quad (\text{marginally significant})$$

a 10% level  
significant @ 0.05  
but not @ a 5% level

- We can also look @ the ~~95%~~ 95% CI of the  $\beta_1 = [-0.278, 2.498]$

true slope it covers 0

- The confidence bands show the confidence interval bands around the mean response @  $x$  estimator,  $\hat{\beta}_1$ . Notice that the bands are curved. The intervals are narrow for values closer to our sample mean. [Mean response @ a given  $x$ ]

→ over  $\hat{x}$  resp.  $\mu$  effectively CI to the

because the mean must lie resp. to the pop mean

- The conf. interval for an individual height is wider than that of the mean.

- Assumption → errors need to have a mean of 0 and a std dev. independent of  $y$ ; the height

- We can't look @ the true errors; only residuals. We use that to check for our normality.

Looking @ the Residual Plot - we want to see a random scatter of pts around 0

whether points falling <sup>w/</sup> a constant horizontal band, i.e. the variability @ low heights & higher heights are similar.

- Adding another variable; "complete" | Pay attention to the interpretation of the regression coefficients.

## Wk 2 Ready: review of Regression Analysis + Linear Models

- The goal of any statistical approach to modeling data is to take a sample of data that represents a population, then use that sample to estimate some facet of the population.

For regression analysis, it's focus is on the conditional mean value of a single dependent variable ( $y$ ) corresponding to a given set of predictor variables  $(x_1, x_2, \dots, x_p)$

- The conditional mean is a function because it depends on the predictor variables.

- **regression analysis** = A set of techniques for taking a sample of a population, then using that data to estimate a conditional mean for that population.

→ Most widely used method = LINEAR LEAST SQUARES

Linear Least Squares represent the conditional mean function as a linear function of the predictor variables.

- The conditional mean function can be modeled as;  $E[y|x_1, x_2] = \beta_0 + \beta_1 x_1 + \beta_2 x_2$

$E[y|x_1, x_2]$  read as "the conditional mean of  $y$  given  $x_1$  and  $x_2$ "  
 ↳ "expectation" = mean or average value.

- $\beta_0, \beta_1, \beta_2$  = regression parameters → Numerical values that define the line

- Sometimes the intercept is not interpretable. How you would find Age & BMI of

- ① to predict the respective Blood Pressure?

However, there is a case when the intercept is interpretable, that is when all the predictor variables have been centered to have mean zero. In this case, the intercept is the cond. mean of the dependent var. when all the predictor vars are @ their mean values. The intercept  $\beta_0$ : Blood Pre. =  $\beta_0 + \beta_1(\text{Age}) + \beta_2(\frac{\text{BMI}}{\text{Height}})$

would be the average BP for someone who has avg Age & avg height

Representing a <sup>explic</sup> model in "generative form";  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$

$\epsilon$  = unexplained variation or error or noise

represent part of the data that cannot be explained using the <sup>Predictor</sup> variables

- Dependent variable → response or outcome variable

- Predictor variable → Independent variable, covariate or regressor

- Mean → Expectation, average

- Regression parameter → Slope, coefficient, effect

- Unexplained variation → Error, noise

- Conditional mean function → Regression function

- Constant variance → Homoscedasticity

- Nonconstant variance → Heteroscedasticity

linear models → linear regression model, linear model

Fitted values → Predicted values

R-squared = Proportion of explained variance  
or Coeff. of determination

The term "linear" in linear least squares is usually used to refer only to (ii) and not to sense (i)

where (i) & (ii) are the senses in which the conditional mean model used in a LLS analysis is linear (i) - It is linear in the predictor variables.

\* (ii) - It is linear in the regression parameters.

Note: It is important to understand that the CMF must be a linear function of its regression parameters if we plan to use linear least squares for estimation

### ... A NOTE ON CAUSALITY ...

When describing the regression model with CMF say,  $E[y|x_1, x_2] = 1 + 3x_1 - 2x_2$ , and saying "when  $x_1$  goes up by 1 unit and  $x_2$  is held fixed,  $y$  goes up by 3 units". Instead say; "when comparing individuals whose  $x_1$  values differ by 1 unit, and whose  $x_2$  values are the same, the value of  $y$  will differ by average by 3 units".  $\rightarrow$  Similarly, instead of saying: " $x_1$  affects  $y$ ", it is better to say that " $x_1$  is associated with  $y$ , after controlling for  $x_2$ ".

• Estimation variance in a regression model is sharply influenced by

3 other characteristics: (i) the level of conditional variance

The scatter in the conditional data around the conditional mean.  $\rightarrow CV \geq EV$

(ii) the variance of the predictor variables

(iii) the correlation among the predictor variables

this is bad variance because greater CV results in greater uncertainty about the regression parameter

- (i) refers to how dispersed the values of the diff predictor variables are w/in the dataset being used to fit the model. The main aim of a regression model

is to establish how diff in the values of a pred. variable relates to diff in the expected val. of the outcome. This is referred to as "good variance" since greater varia-

w in the predictor variable results in less uncertainty about the regression param.

(iii) Correlation among the predictor variables = CORLINEARITY. It makes it difficult for the model to properly attribute a trend to any of the correlated variables.

This doesn't mean that there should be 0 correlation b/w predictor variables.  
greater collinearity results in greater uncertainty in  
but we are saying that the greater the  
the param. estimation. This can be overcome by collecting more data.

Collinearity is viewed as "bad correlation" while correlating b/w predictor var. & the dependent var. is "good correlation" because they allow us to fit models that do a better job of explaining the variation in the dependent variable.

- EXPLAINED VARIATION / R-SQUARES → The  $\alpha$  the squared Pearson correlation coefficient between the fitted values and the observed value of the dependent variable

In general, a higher R-squared reflects a better fit of the model but its interpretation should be qualified in 2 ways:

(i) "Goodness of fit" refers to more than just the mean function. To have a model that fits well, we would like to capture the variance structure as well as the mean structure.

(ii) Higher R-squared can reflect "overfitting", in which the model fits the data in-hand better than it will fit equivalent data that we observe in the future.

As noted before, several factors affect parameter estimates: (i) sample size, (ii) conditional variance in the dependent variable (iii) variance of the predictor variable (iv) collinearity

Note that *i* & *ii* impact all parameters equally; *iii* & *iv* impact different parameters in a model to different extents.

### - Categorical Predictor Variables

A categorical variable is not numerical and hence can't be directly inserted into the linear predictor:  $B_0 + B_1x_1 + \dots + B_px_p$ . This issue is typically addressed by coding the categorical var. into one or more INDICATOR VARIABLES → These can be 0 or 1 for each level of a categorical var. It takes 1 for "yes" on a level & "0" for "no" at same level. For example, the indicator var for gender is Male = 1 Female = 0.

	Male	Female	Breakfast	Lunch	Dinner	
1	1	0	0	1	0	
2	0	1	1	0	0	
3	0	1	0	0	1	
4	1	0	1	0	0	

→ for get?

Meal eaten =  
Breakfast =  
Lunch =  
Dinner =  
Total meal of  
the day

When including indicator variables as predictors var in a regression model, one indicator derived from the parent categorical variable must be dropped.

This is because the sum of all indicator var. from the same parent var. is inherently equal to 1 and hence the same as the intercept. For the model to be estimable, we cannot include the same variable twice as a predictor variable. Omitting one indicator var. solves this issue, and the level that is dropped = the reference level.

- For example, we can ~~drop~~ the "Male" ind. var & only use the "female" indicator var. hence "Male" = the reference level.

The interpretation of a regression param. for an indicator var. should account for the reference level. The regression param. for any non-reference level is the contrast between units w/ that level of the parent variable, and units w/ the reference level of the parent variable. Example: The coeff. for the "female" variable would be interpreted as the difference b/w females & males. In a model for blood pressure, we might find that the coeff for the "female" var. = -4. This would mean that holding other factors fixed, females have on average 4 mm/Hg lower blood pressure than males.

- Residuals = the differences between each fitted value from its observed value. A scatter plot of the residuals against the fitted values tells us about the structure of the fitted model. In particular, if the degree of scatter in this plot increases or decreases from left to right, we have discovered a MEAN VARIANCE RELATIONSHIP. This implies that no only is there heteroscedasticity but also the cond. variances differ in a way that's predictable from the mean.

→ Also, knowing that, a  $\bar{x}/s$  relationship is present indicates that while the linear least squares est. of the regression param are meaningful, their SE may not be correct & we shouldn't rely on the results obtained from standard procedures for stat inf. in linear models.

- Scatter plots of the residuals on individual covariates can sometimes reveal nonlinear structure that the model failed to capture. If for example, there is a "U-shaped" relationship b/w the residuals and the values of a covariate (e.g.  $x_2$ ), then it's better to include  $x_2^2$  or some other transformed version of  $x_2$  in the model.

## W2 Logistic Regression - Methods

n=25

Binary variable of Inter. = Cartwheel completion

Research Q ≈ Is there relationship betw age ≈ P(Cart completion)

- Variable of Int → takes only 2 values

Balk of data b/w 22 &amp; 25

Linear regression line does quite fit here.

- We need to modify our response linear so that all our estimates should be between 0 & 1

Instead of predicting the completion status using a linear regression model, we'll

Predict a transformed version of the probability of success

→ Logit Transformation (uses the logit function or log odds)

$$\text{Proper Logit function} = \ln\left(\frac{P}{1-P}\right) ; P = \text{Probability of success}$$

odds =  $P(\text{success}) / P(\text{Failure})$

- The Logistic function of P is a symmetric distribution

Note: The higher the value on the logit function, the higher the prob. of success (vice versa)

- There are no maximums of the logit function. They extend to  $-\infty \approx \infty$

$$\text{Our model} \quad \text{logit}(\hat{y}) = b_0 + b_1 x$$

Scatter plot → The logistic reg line is a curved line because we're looking at the P(success) scale  
 + Logistic regression line

→ Don't extrapolate i.e. estimate y based on an x that isn't in the range of the data.

Python output

$$\text{Intercept} = -4.4213$$

$$\text{Age} = 0.2096$$

$$\text{logit}(\hat{y}) = -4.42 + 0.2096(\text{Age})$$

+ few attempts to interpretation  
 $1.23 = e^{0.2096}$

• Prediction Uncertainty = ~~standard~~ Bands are smaller close near the mean  $\hat{y}$ .  
but were there are fewer observations, the bands are wider.

• transform log odds to probabilities  $\hat{y} = \delta(\text{Covariates}) \rightarrow$  The predictor values  
mainly because our  $n$  is small  $\leftarrow$  bands are curved but  
not smooth

Assumptions = (i) The model;  $\text{logit}(y)$  has to be linearly associated with dependent variable  
(ii) residual plots aren't that helpful since  $y$  takes only 2 values. This  
can be solved by having more variation & also having more covariates

## WK 2 Logistic Regression LINE

• Confidence Interval (95%) for Slope: Sample slope  $\pm$  "a few"  $\cdot$  SE

from software result  $0.2096 \pm 1.96(0.171) \leftarrow b_1 \pm z \pm \text{SE}_{b_1}$

$$95\% \text{ CI} = (-0.126, 0.545) \rightarrow \text{NOT significant because contains } 0$$

• Hypothesis testing:  $H_0: \beta_1 = 0$ ;  $H_A: \beta_1 \neq 0$

$$\text{test stat} = \frac{b_1 - 0}{\text{SE}_{b_1}} = \frac{0.2096}{0.171} \approx 1.225 \rightarrow p\text{ value} = 0.221$$

## WK2 Reading: Overview of regression w/ BINARY outcomes

The linear approach is useful for many diff types of variables; especially quantitative variables. However, the need arises to conduct a regression analysis using a binary dependent variable.

Ex: Consider a research study on the US adult population and look @ whether a given person in this population has ever been a smoker in their lifetime. That is a categorical variable w/ 2 levels, so it doesn't have an expected mean [i.e. the mean of "smokers" and "non-smokers"]

the variable values, if coded as  $0 \pm 1$ , mean nothing. However we can choose one cate,

say the "smoker" cate. and aim to model the probability that the outcome occurs. Equivalently, we can code the outcome numerically as 1 (for "smoker") and 0 ("non-smoker") and then

the expected value of this ~~binary~~ equivalent. [Both approaches are equivalent]

A binary dependent variable such as smoking & fetus generally can't be modelled with linear regression because trying to model a binary dependent variable w/ linear regression

would often obtain fitted probabilities that are greater than 1 or less than 0.

Link function → used when working w/ binary dependent variables and is used to map the probability (or mean) to a value on the real no line which (in this case) models using the linear predictor. Logistic function or Log odds function ( $\log \left[ \frac{p}{1-p} \right]$ ) is the link function commonly used for regression w/ a binary outcome

- Odds → A quantity derived from a probability. It can take on any non-negative number.  
 $L = \frac{p}{1-p}$ . That odds of 3  $\equiv$  we are 3 times more likely to observe a 1 than to observe a 0.
- The odds of  $1/3 \equiv$  we are 3 times more likely to observe a 0 than to observe a 1.
- The odds  $= 1 \equiv$  neutral because we are equally likely to observe both outcomes.  $B(\frac{1/2}{1/2})$
- A further transformation that is often applied is to take the (natural) logarithm of the odds yielding the log odds  $\sim \log\left(\frac{p}{1-p}\right)$ . The neutral pt for the log odds  $= 0$ ; i.e.  $\log(1) \equiv p=1/2$ . The log odds is symmetric around 0 in the sense that a log odds of  $v$  and a log odds of  $-v$  convey the same strength of association in terms of the base relative to the neutral value. This symmetry makes log odds a good choice for modeling via the linear predictor.

### Logistic Regression

Here, we set  $\log\left[\frac{p}{1-p}\right] = b_0 + b_1 x_1 + \dots + b_p x_p$ . For example, if the outcome we are modeling is smoking history, and the covariates are age and gender (coded as female = 1 if the ref. level of this categorical var. is "male") then the model is

$$\Rightarrow \log\left[\frac{p}{1-p}\right] = b_0 + b_1 (\text{age}) + b_2 (\text{female})$$

- If  $b_1 > 0$ , then older people are more likely to be smokers than younger people. Whereas if  $b_1 < 0$ , older people are less likely to be smokers than younger people.
- If  $b_2 > 0$ , then age is unrelated to the probability that a person smokes (w/in a specific group).
- If  $b_2 > 0$ , then females are more likely to smoke than males and if  $b_2 < 0$ , then females are less likely to smoke than males.

### PARAMETER INTERPRETATION

The quantitative interpretation of the logistic regression param. are best explained:

- (i) additively in terms of the log odds or
- (ii) multiplicatively in terms of the odds.

Example: Age parameter 0.04; same gender comparison but ages differ by one year

(i) The older person has log odds for smoking that is 0.04 units greater than that of the younger person. This relationship is linear in age, so if we compare two people of the same gender whose age differ by 10 yrs, then the older person has 0.4 greater ~~log odds~~ <sup>odds</sup> for smoking than the younger person.

(ii) Note that when exponentiating the log odds to get the odds they are no longer linear but ~~rather~~ <sup>rather</sup> multiplicative. Since  $\exp(0.04) \approx 1.04$ , comparing two people of the same gender whose ages differ by one year, the older person has 1.04 greater odds of smoking than the younger person. When comparing people of the same gender whose ages differ by 10 yrs, the older person has  $\exp(0.4) \approx 1.49$  times greater odds of smoking than the younger person.

## WK 2 PROBLEMS $\rightarrow$ using R STANZA dataset

R-Squared can be defined as — (i) Squared Pearson corr. coeff between the covariates data for just 1 outcome & 1 covariate  $\leftarrow$  multiple covariates (ii) Squared Pearson corr. coeff b/w the fitted values & the true values

## WK 3 WHAT ARE MULTILEVEL MODELS & HOW DO WE FIT THEM?

• We will see looking @ appropriate stat. models for datasets that are generated by study design that introduce dependencies in the data.

• Now we're working with data where the variables are collinear

Multilevel models = used to model dependent data

Not every variable here, the regression coeff vary  $\left\{ \text{not fixed} \right\}$  across the higher level clusters.

→ Helps estimate the variability  $\leftarrow$  <sup>amongst</sup> subjects or clusters in terms of the coefficients of interest

• We estimate relationships & parameters that estimate variances

$$\text{Level 1: } Y_{ij} = \beta_{0j} + \beta_{1j} X_{1ij} + e_{ij}$$

Dependent variable:  $y$  which is defined for a given observation  $i$ , nested within cluster  $j$ .

$\beta_{0j}$  is called cluster-specific intercept

$\beta_{1j}$  is called cluster-specific slope

\* Notice how the  $\beta$  is a function of coeff which

have a subscript ( $j$ ) which means that the regression coeff are determined by what cluster,  $j$ , we are referring to

\* these coeffs are referred to as random coeffs and not params because they change w/ $j$ . They randomly vary depending on the cluster.

Level 1	Level 2:	$\beta_{0j} = \beta_0 + u_{0j}$	$\beta_{1j} = \beta_1 + u_{1j}$
Below level 1 & there are eqns for the random coefft		$\rightarrow$ the random effect for each cluster that makes the random coefft vary.	$\rightarrow$ the regression parameter

- \* Multilevel models allow us to decompose the unexplained variance in a given outcome into "between" & "within-cluster" variance
  - \* Random effects capture the level 1 error terms ( $\epsilon$ )

Key Question: How much of unexplained variance due to between-  
cluster variance arises in the intercepts or slopes of a given model?

WK3

# Multivariate Linear Regression Models

Model for a continuous dependent variable  $Y$ , measured on person in cluster  $j$

$$y_{ij} = \beta_0 + \beta_1 x_{1ij} + u_{0j} + u_{1j} x_{1ij} + e_{ij}$$

→ error term

Fixed part  
that capture the overall interest & model & overall relationship  
of  $x$  with  $y$

Random effect  
capture the cluster variance

→ captures what's not predicted by the regress coeffs & the rand. effects

- Fixed effects are used to find Mean of the dependent variable as it is not

$\text{EBLUPS} = \text{Empirical Best Linear Unbiased Predictions}$

**W3** | MULTILEVEL Lotusne Rekresjon Nedem

$$\ln \left[ \frac{P(y_{ij} = 1)}{1 - P(y_{ij} = 1)} \right] = \text{logit}[P(y_{ij} = 1)] = \beta_0 + \beta_1 x_{ij} + u_{0j} + u_{1j} x_{1r_j}$$

single well equation properties, band - effects

Q: Suppose that you examine the variability among randomly sampled higher-level clusters of observations (e.g., schools) in the proportions on a binary variable of interest, and you find visual evidence of significant variability in the proportions. You then fit a logistic regression model to these data in Python, but forgot to fit a multilevel model including random effects of the clusters, how would this affect your analysis?

- \* By omitting explicit random effects in the logistic regression model, the standard errors of our estimated fixed effects would likely be too low, and the estimates of the fixed effects would likely be incorrect because we're failing to explicitly adjust for the random effects of the higher level clusters. Omitting random effects when they are important is the same type of model specification error as omitting the fixed effect of an important predictor variable.

WK3

### WHAT ARE MULTILEVEL MODELS & WHY DO WE FIT THEM?

- We don't include random effects in marginal models because we don't care about the between-cluster variance rather we want to model the <sup>WITHIN</sup> ~~cluster~~ <sup>var</sup>.
- If we fail to account for the fact that observations are correlated within the same higher level cluster, we run the risk of underestimating our standard errors.

WK4

### SHOULD WE USE SURVEY WEIGHTS WHEN FITTING MODELS

#### MODELS

WK4

### BAYESIAN APPROACHES TO CLASSIFICATION & MODELLING

- Bayesian stats looks at belief & probability

