

R: Models for Clustered and Correlated Data

Mary Ryan

University of California, Irvine
Presented at NICAR 2019

March 7-10

What's wrong with linear regression?

"Essentially, all models are wrong, but some are useful"

— George Box, *Empirical Model-Building and Response Surfaces*, pg. 424

What's wrong with linear regression?

"Essentially, all models are wrong, but some are useful"

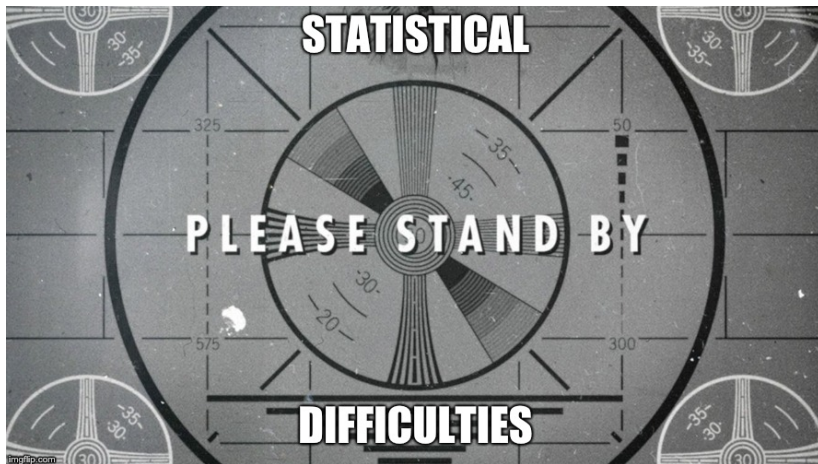
— George Box, *Empirical Model-Building and Response Surfaces*, pg. 424

Other iterations:

- ▶ “Remember that all models are wrong; the practical question is **how wrong do they have to be to not be useful**” (*Empirical Model-Building and Response Surfaces*, pg. 74)
- ▶ “The most that can be expected from any model is that it can supply **a useful approximation to reality**: All models are wrong; some models are useful” (*Statistics for Experimenters*, pg. 440)

What's wrong with linear regression?

- ▶ Regular linear regression assumes each datapoint is *independent* of each other
- ▶ What if we have multiple datapoints for each person/school/hospital/location we're measuring?



Ordinary Least Squares (OLS)

- ▶ If we have data \mathbf{X} and response \vec{Y} , we can find our regression coefficients through:

$$\mathbf{X}^T(\vec{Y} - \mathbf{X}\vec{\beta}) = 0$$

- ▶ Implement this using the `lm()` function

Iteratively Reweighted Least Squares (IRLS)

- ▶ Put a weight on the equation that estimates our coefficients
 - ▶ Accounts for the fact that we are no longer dealing with completely independent data points

$$\mathbf{X}^T \mathbf{W} (\vec{Y} - \mathbf{X} \vec{\beta}) = 0$$

- ▶ Implement this using the `gee()` function from the `gee` package
 - ▶ Can also use `lme()` from the `nlme` package, and `lmer()` from the `lme4` package



AND NOW BACK TO
OUR REGULARLY
SCHEDULED
PROGRAMMING

The `gee()` Function

- ▶ `gee` stands for Generalized Estimating Equation

Like with the `lm()` function, there are several arguments we need to fill in:

- ▶ **Formula:** this is the same formula you would plug in for `lm()`, of the form `response ~ variable1 + variable2 + ...`
- ▶ **id:** this is a variable in your dataframe that identifies your clusters. If I have 12 patients with 3 datapoints each, each datapoint needs to have something that tells us which patient it is coming from. Usually this is done as the very first column of your dataframe, where the id can be a number or a string.
- ▶ **data:** like with `lm()`, this is the name of your dataframe

The `gee()` Function

- ▶ **family**: the default for this argument is “gaussian”“, which just means Normal. We generally won’t put anything in for this argument unless we’re dealing with binary data (we’ll see this later).
- ▶ **corstr**: this tells the function how we want to do our weights.
 - ▶ There are 3 main options for this:
 1. “independence”: this is the default and will get us `lm()`
 2. “exchangeable”: this gives us random intercepts
 3. “AR-M”: this gives us random slopes and random intercepts. With this though, we also need to specify a “Mv” argument, which will be 1.

Side note: two other popular functions for modeling longitudinal data are `lme()` from the `nlme` package, and `lmer()` from the `lme4` package. These work similarly to the `gee()` function but have slightly different syntax and technically require stronger statistical assumptions to use. I generally stick to `gee()`.

Example 1: Childhood Wheezing and Maternal Smoking

- ▶ Ohio dataset
 - ▶ resp: an indicator of wheeze status (1=yes, 0=no)
 - ▶ id: a numeric vector for subject id
 - ▶ age: a numeric vector of age, 0 is 9 years old
 - ▶ smoke: an indicator of maternal smoking at the first year of the study
- ▶ 537 unique subjects
- ▶ 2148 total observations

Question: how does maternal smoking affect wheezing?

Example 2: Median Housing Value in Texas

- ▶ Housing dataset (select variables)
 - ▶ countyID: ID number for each county
 - ▶ countyName: Name of county
 - ▶ yrsSince2009: Number of years since 2009
 - ▶ Median: Median housing value in county
 - ▶ totPop18plus: Population 18 years or older
 - ▶ BAtotPop18plus: Population 18 years or older with at least a Bachelors
 - ▶ BApctTotPop18plus: Percentage of population 18 years or older with at least a Bachelors
- ▶ 53 unique counties
- ▶ 420 unique observations
- ▶ 8 years of data

Question: how does the percentage of Bachelors-holders in a district affect mean housing value?

Example 3: California API Scores

- ▶ API dataset (select variables)
 - ▶ CDS: County/District/School code
 - ▶ RTYPE: Record Type: (D=District, S=School, X=State)
 - ▶ DNAME: District name
 - ▶ API: Base API score
 - ▶ PCT_AA, PCT_AS, PCT_HI: Percentage of African American, Asian, and Hispanic students
 - ▶ charter_df: Indicator if record is direct-funded charter
 - ▶ charter_ndf: Indicator if record is non direct-funded charter
 - ▶ P_EL: Percent English learners
 - ▶ MEALS: Percentage of Students Tested that are eligible for Free or Reduced Price Lunch Program

Example 3: California API Scores

- ▶ 1047 unique school districts
- ▶ 7178 total observations
- ▶ 7 years of data

Question: how does charter status affect district API score?

Slide with R Output

```
summary(cars)
```

##	speed	dist
##	Min. : 4.0	Min. : 2.00
##	1st Qu.:12.0	1st Qu.: 26.00
##	Median :15.0	Median : 36.00
##	Mean :15.4	Mean : 42.98
##	3rd Qu.:19.0	3rd Qu.: 56.00
##	Max. :25.0	Max. :120.00

Slide with Plot

