

Correlated Data

Mine Dogucu

Noah Johnson

April 5, 2018

General Linear Models require responses to be approximately normally distributed and independent. Through generalized linear models, we learned about handling responses that are not normally distributed (e.g. Poisson, Binomial). From now on we learn about multilevel models / linear mixed effect models / hierarchical linear models that can model response that violate the independence assumption. These models have responses that are correlated thus not independent.

```
hsb <- read.csv('hsb.csv')
str(hsb)
```

```
## 'data.frame':    7185 obs. of  10 variables:
## $ schoolid: int  1224 1224 1224 1224 1224 1224 1224 1224 1224 1224 ...
## $ minority: int   0  0  0  0  0  0  0  0  0  0 ...
## $ female  : int   1  1  0  0  0  0  1  0  1  0 ...
## $ ses      : num  -1.528 -0.588 -0.528 -0.668 -0.158 ...
## $ mathach  : num   5.88 19.71 20.35  8.78 17.9 ...
## $ size     : int  842 842 842 842 842 842 842 842 842 ...
## $ sector   : int   0  0  0  0  0  0  0  0  0 ...
## $ pracad   : num   0.35 0.35 0.35 0.35 0.35 0.35 0.35 0.35 0.35 ...
## $ disclim  : num   1.6 1.6 1.6 1.6 1.6 ...
## $ himinty  : int   0  0  0  0  0  0  0  0  0 ...
```

In this class and in the next few, we will be using the hsb.csv dataset. We will try to understand math achievement (mathach) of students based on their socio-economic status (ses). We will use the notation Y_{ij} for math achievement of i th student in the j th school. There are **7185** students within **160** schools in this dataset.

Calculate the average math achievement score for all students in the dataset.

```
avg_mathach <- hsb %>% summarise(avg_mathach = mean(mathach))
```

$Y_{..} = 12.7478526$

Pick a random school (make sure each school has equal probability to be selected). For **only** the school you selected, find the mean math achievement.

```
random_school_id <- hsb$schoolid %>%
  unique() %>%
  sample(1)

print(random_school_id)
```

```
## [1] 3152
```

```
my_school <- hsb %>% filter(schoolid == random_school_id)
```

```
mean_school_math_score <- my_school %>%
  summarise(mathach = mean(mathach))
```

$Y_{.j} = 13.2090385$

For your school, fit a general linear model where math achievement is the response and ses is the predictor. Record the coefficients:

```
my_school_model <- lm(mathach ~ ses, data = my_school)
summary(my_school_model)

##
## Call:
## lm(formula = mathach ~ ses, data = my_school)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -16.1823  -2.9224  -0.2391   5.1961  10.1103
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   13.123      0.984   13.336 <2e-16 ***
## ses           2.768      1.172    2.362  0.0221 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.091 on 50 degrees of freedom
## Multiple R-squared:  0.1004, Adjusted R-squared:  0.08238
## F-statistic: 5.578 on 1 and 50 DF,  p-value: 0.02212
 $\hat{\beta}_{0j}$  : 13.1231163  $\hat{\beta}_{1j}$  : 2.7682496
```

For your school, fit a general linear model where math achievement is the response and grand mean centered ses is the predictor. Record the coefficients:

```
hsb$s ses.grandmean_c <- hsb$s ses - mean(hsb$s ses)
my_school <- hsb %>% filter(schoolid == random_school_id)

my_school_model.grand <- lm(mathach ~ ses.grandmean_c, data = my_school)
summary(my_school_model.grand)

##
## Call:
## lm(formula = mathach ~ ses.grandmean_c, data = my_school)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -16.1823  -2.9224  -0.2391   5.1961  10.1103
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   13.123      0.984   13.336 <2e-16 ***
## ses.grandmean_c  2.768      1.172    2.362  0.0221 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.091 on 50 degrees of freedom
## Multiple R-squared:  0.1004, Adjusted R-squared:  0.08238
## F-statistic: 5.578 on 1 and 50 DF,  p-value: 0.02212
 $\hat{\beta}_{0j}$  : 13.1235131  $\hat{\beta}_{1j}$  : 2.7682496
```

For your school, fit a general linear model where group mean centered math achievement is the response and ses is the predictor. Record the coefficients:

```
school_means <- hsb %>%
  group_by(schoolid) %>%
  summarise(ses.groupmean = mean(ses))

hsb <- merge(hsb, school_means, by = "schoolid")

hsb$ses.groupmean_c <- hsb$ses - hsb$ses.groupmean

my_school_model.group <- lm(ses.grandmean_c ~ ses, data = my_school)
summary(my_school_model.group)
```

```
## Warning in summary.lm(my_school_model.group): essentially perfect fit:
## summary may be unreliable

##
## Call:
## lm(formula = ses.grandmean_c ~ ses, data = my_school)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.207e-15  2.500e-18  2.856e-17  4.846e-17  1.175e-16
##
## Coefficients:
##              Estimate Std. Error    t value Pr(>|t|)
## (Intercept) -1.434e-04  2.470e-17 -5.804e+12  <2e-16 ***
## ses          1.000e+00  2.942e-17  3.399e+16  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.78e-16 on 50 degrees of freedom
## Multiple R-squared:      1, Adjusted R-squared:      1
## F-statistic: 1.155e+33 on 1 and 50 DF, p-value: < 2.2e-16
 $\hat{\beta}_{0j} : -1.4335421 \times 10^{-4}$   $\hat{\beta}_{1j} : 1$ 
```

We will use group-mean centering moving on.

Different effect types

Empty Model - One-Way Random-Effect ANOVA

```
install_load('lme4')

## Loading required package: Matrix
##
## Attaching package: 'Matrix'
## The following object is masked from 'package:tidyr':
##
##      expand
```

```
model.null <- lmer(mathach ~ 1 + (1|schoolid), data=hsb)
summary(model.null)
```

```
## Linear mixed model fit by REML ['lmerMod']
## Formula: mathach ~ 1 + (1 | schoolid)
## Data: hsb
##
## REML criterion at convergence: 47116.8
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -3.0631 -0.7539  0.0267  0.7606  2.7426
##
## Random effects:
##  Groups   Name      Variance Std.Dev.
## schoolid (Intercept)  8.614    2.935
## Residual                39.148    6.257
## Number of obs: 7185, groups: schoolid, 160
##
## Fixed effects:
##              Estimate Std. Error t value
## (Intercept)  12.6370    0.2444   51.71
```

Model Notation:

Parameter Estimates:

ICC