Residuals

Grinnell College

May 7, 2025

Grinnell College

STA 209 is the cat's meow

May 7, 2025

Review

Below is a model of chicken weight in days since birth according to one of four separate diets

```
1 > lm(weight ~ Time + Diet, ChickWeight) %>% summary()
2
3
            Estimate Std. Error t value
                                              Pr(>|t|)
4 (Intercept) 10.924 3.361 3.25
                                                0.0012 **
5 Time
         8.750 0.222 39.45 < 0.000000000000002 ***
6 Diet2
           16.166 4.086 3.96 0.00008556049098 ***
7 Diet3 36.499 4.086 8.93 < 0.000000000000002 ***
8 Diet4
      30.233 4.107 7.36 0.000000000064 ***
9
10 Multiple R-squared: 0.745, Adjusted R-squared: 0.744
11 F-statistic: 419 on 4 and 573 DF, p-value: <0.0000000000000002
```

- 1. Write out the equation for this linear model
- 2. What proportion of the total variability in chick weight is described by this model?
- 3. Which diet resulted in the heaviest weight? The lightest?
- 4. After 10 days, what is the predicted difference in weight between a chicken on Diet 2 and Diet 3?

Regression is a model posits linear relationship between dependent variable y and independent variable X of the form

$$y = \beta_0 + \beta_1 X + \epsilon$$

- Expand this to include combinations of independent variables, both qualitiative and quantitative
- Today our focus is on the error term ϵ (epsilon)



$$y = \beta_0 + \beta_1 X + \epsilon$$

Assumptions:

- Variability in error should be the same for all values of X, i.e., error same for all observations

Analyzing the error terms gives us a way to test the assumptions of our model



Fitted line with residual



Part 1: Checking Assumptions

Residuals and assumptions

Three common ways to investigate residuals visually:

- 1. Plot histogram of residuals (normality)
- 2. Plot residuals against covariate (linearity, constant variance)
- Plot residuals against new covariates (pattern identification)

Check how far away each point is from the line.



This distance is called residuals

Take the residuals and put them in a bag.

Are my errors normal? Is the variance the same at all levels of *X*?



Constant Variance



Tests of linearity



Tests of linearity

Sometimes a transformation of a variable (in this case, $\log({\sf weight}))$ can help correct trends



Fitted Line

Residual Plot

Part 2: Investigating Patterns

Correlated Covariates

Consider a simple linear model in which a covariate X is used to predict some value y

$$\hat{y} = \hat{\beta}_0 + X\hat{\beta}_1$$

The residuals associated with this describe the amount of variability that *is yet to be explained*

$$r = \hat{y} - y$$

The idea is to find new covariates *associated* with this residual, in effect "mopping up" the remaining uncertainty

Suppose I have:

- Quantitative outcome y
- Quantitative predictor X
- Categorical predictor indicating group



$$\hat{y} = \hat{\beta_0} + X\hat{\beta_1}$$





Grinnell College



We see there is a clear association between a missing variable and the residuals. So we add that variable!

Grinnell College

STA 209 is the cat's meow

$$\hat{y} = \hat{\beta}_0 + X\hat{\beta}_1 + \mathbb{1}_A\hat{\beta}_2$$



STA 209 is the cat's meow



Considering new covariates



Now consider a situation in which we wish to predict fuel economy with three separate models:

- 1. Using weight
- 2. Using weight and engine displacement
- 3. Using weight and quarter mile time

Starting with the first model, we can consider the relationship of the residuals with both engine displacement and quarter mile time

Considering new covariates (quantiative)



Considering new covariates (quantiative)

An interesting thing occurs when we try to create a regression model of the residuals with the original variable:

residuals
$$= \hat{eta}_0 + \hat{eta}_1$$
Weight $= 0$

Grinnell College

Considering new covariates (quantiative)

skip

When considering adding new variables to our regression model, we want to add those that will "mop up" the residuals that are left after considering weight

This brings us to the idea of **correlated variables**, or variables that have evidence of a *linear relationship* with one another

How correlated our variables may be will impact how much of the residual they are able to account for

Correlated Covariates

We can consider two extremes: if two quantitative variables are perfectly correlated, knowing the value of one variable means we also know the value of the other

This means that, in terms of predicting an outcome, adding a highly correlated variable to our model will contribute little new information and will not be very useful

By constrast, if two variables have perfectly uncorrelated, then knowing the value of one tells us nothing about the value of another

Adding an uncorrelated variable to our model thus offers more potential to "mop up" the variability that was not explained by the first variable

SKip

Correlated Covariates

multicolinearity: two or more independent variables are correlated and it causes problems because the share information



Residual Plots

Displacement has little association with the residuals, while quarter mile time has quite a bit of association. This suggests that adding QM time will "mop up" more unexplained variance



Correlated Covariates

```
1 > lm(mpg ~ wt, mtcars) %>% summary()
2
3 Estimate Std. Error t value Pr(>|t|)
4 (Intercept) 37.285 1.878 19.86 < 0.000002 ***
5 wt -5.344 0.559 -9.56 0.000013 ***
6 \text{ R-squared} = 0.75
7
8 > lm(mpg ~ wt + disp, mtcars) %>% summary()
9
10 Estimate Std. Error t value Pr(>|t|)
11 (Intercept) 34.96055 2.16454 16.15 0.000000049 ***
12 wt -3.35083 1.16413 -2.8 0.0074 **
13 disp -0.01772 0.00919 -1.93 0.0636.
                               adding displacement changed our weight estimate
14 R-squared = 0.78
15
                                and increased sd
16 > lm(mpg ~ wt + qsec, mtcars) %>% summary()
17
   Estimate Std. Error t value Pr(>|t|)
18
19 (Intercept) 19.746 5.252 3.76 0.00077 ***
20 wt -5.048 K → 0.484 -10.43 0.0000000025 ***
21 gsec 0.929 Adding 0.265 3.51 0.00150 **
22 R-squared = 0.82 /4 mile time backy charged our weight estimate and decrease so
       Grinnell College
                          STA 209 is the cat's meow
                                                  May 7, 2025
                                                              28 / 29
```

- 1. Number of assumptions for linear model
 - Linearity
 - Normal errors
 - Constant Variance
- 2. Need way to determine which new variables to add to model
- 3. Examining errors effective way to test assumptions and investigate new covariates
- 4. Relationship between correlation of predictors and residual analysis