### Problems with predictors

ST552 Lecture 18

Charlotte Wickham 2019-02-20

### Roadmap

#### Done:

- Regression model set up and assumptions
- Least squares estimates and properties
- Inference
- Diagnostics

### To Do:

- Specific problems that arise and some extensions
- Model Selection (week 8) <-- lasso</li>
- Some case studies (week 9)
- Non-linear, binary data (week 10) ← [asso

There will be 8 homeworks total, recall your lowest score is dropped.

Problems with predictors (Faraway 7)

- Collinearity
- Linear transformations of variables
- Errors in predictors

# data(seatpost, package = "faraway") ?seatpost

Car drivers like to adjust the seat position for their own comfort. Car designers would find it helpful to know where different drivers will position the seat depending on their size and age. Researchers at the HuMoSim laboratory at the University of Michigan collected data on 38 drivers.

The dataset contains the following variables:

Age, Age, years Weight, Weight, lbs HtShoes, Height in shoes, cm Ht, Height bare foot, cm Seated, Seated height, cm Arm, lower arm length, cm Thigh, Thigh length, cm Leg, Lower leg length, cm hipcenter, horizontal distance of the midpoint of the hips from a fixed location in the car, mm

(2)03

```
library(ggplot2)
ggplot(seatpos, aes(Ht, hipcenter)) +
  geom_point()
```



## lmod <- lm(hipcenter ~ ., data = seatpos) sumary(lmod)</pre>



Exact collinearity If X<sup>T</sup>X is singular, we say there is exact collinearity. There is at least one column that is a linear combination of the others. In R you will get NA for some estimates. Solution drop a column involved, or add constraints on the parameters

you problem

V.

Collinearity or multi-collinearity refers to the case where X<sup>T</sup>X is close to singular. There is at least one column that is almost a linear combination of the others. Or in other words one column is highly correlated with a combination of others. In practice this leads to imprecise estimates (i.e. estimates with large standard errors)

### Variance inflation factors

Let  $R_i^2$  be the  $R^2$  from the regression of the *j*th explanatory variable on all the other explanatory variables. That is, the proportion of the variation in the *j*th explanatory variable that is explained by the other explanatory variables.

### This is not a violation of the assumptions

- Multi-collinearity does not violate any regression assumption.
- Our t-tests, F-tests, confidence intervals and prediction intervals all behave as they should.
- The problem is the interpretation of individual parameter estimates. It no longer makes much sense to talk about "... the effect of X<sub>1</sub> holding other variables constant" because we have observed a relationship between X<sub>1</sub> and the other variables.
- We can't separate the effects of the variables that are collinear, and our standard errors reflect this accurately by being large.

In the seat example: large model  $R^2$  but nothing is individually significant. Large standard errors on terms that should be highly significant.

- Look at the correlation matrix of the explanatory variables. But, this will only identify pairs of explanatories that are correlated (not complicated relationships)
- 2. Regress  $X_i$  on other variables and look for high  $R^2$ , equivalently directly find variance inflation factors.
- 3. Look at the eigenvalues of  $X^{T}X$  and look for condition numbers

$$\kappa = \sqrt{rac{\lambda_1}{\lambda_p}} > 30$$

Go through example in  ${\sf R}$ 

### What to do about multicollinearity?

- Most importantly identify when it occurs, so you don't make stupid statements about individual parameter estimates.
- For prediction, it isn't a problem as long as future observations have the same structure in the explanatory variables.
- For explanation, we can't separate the effect of variables that measure the same thing, do joint tests instead.
- Dropping an offending variable is only necessary if you, for some reason, want a model with as few terms as possible. Do not conclude that a variable dropped due to multicollinearity isn't related to the response!

We assumed fixed X.

You can also use least squares if X is random before you observe it, and you want to do inference conditional on the observed X.

If, X is measured with error, i.e.  $X = X_a + \delta$   $Y = X_a \beta + \epsilon$   $Y = X_a \beta + \epsilon$   $X = X_b + \epsilon$  X

There are "errors in variables" estimation techniques.

### Linear transformations of predictors



 $X_j$  of one standard deviation is associated with a change in response of  $\beta_j \dots$ "

Also, can be useful to re-express a predictor in more reasonable units. For example, expressing income in \$1000s rather than \$1s.