

BIFX 553 - Discussion 5

Randy Johnson

February 16, 2017

Setup

```
library(missForest) # for imputation
library(tidyverse)
library(broom) # for tidy model display

theme_set(theme_classic() +
  theme(axis.line.x = element_line(color = 'black'),
        axis.line.y = element_line(color = 'black'),
        text = element_text(size = 15)))
```

Tests of Association

Single parameter tests

We've already seen and discussed p-values in regression output, but we will discuss them in more detail here. Our favorite model is currently

$$\log(nodes_i) = \beta_0 + age_i * \beta_1 + size_i * \beta_2 + grade_i * \beta_3 + \varepsilon_i.$$

Single parameter tests

Given this model, what are the statistically significantly associated predictors of the number of nodes? How would you describe these associations? Can we remove any variables from the model without losing information?

```
# revisit our model
full.model <- lm(lnodes ~ age + size + grade + meno + lpgr + ler + hormon,
                 data = gbsg)
tidy(full.model)
```

##	term	estimate	std.error	statistic	p.value
## 1	(Intercept)	0.503928657	0.365716045	1.3779233	1.686816e-01
## 2	age	-0.001092263	0.005326082	-0.2050782	8.375726e-01
## 3	size	0.019532987	0.002393120	8.1621425	1.603046e-15
## 4	grade	0.132535536	0.062808817	2.1101422	3.521185e-02
## 5	menopremenopausal	-0.060554085	0.109291246	-0.5540616	5.797194e-01
## 6	lpgr	-0.025311822	0.018632700	-1.3584624	1.747688e-01
## 7	ler	0.002606807	0.019428389	0.1341751	8.933039e-01
## 8	hormonno tamoxifen	-0.070174511	0.074036040	-0.9478426	3.435473e-01

Multiple degree of freedom tests

That last question is perhaps best answered with a multiple degree of freedom test. Lets say that we want to check our model against a model without age, grade, pgr, er and the size:grade interaction. We can do this with the `anova` function in R.

```
# this is our alternate model  
alt.model <- update(full.model,  
                    . ~ . - age - lpgr - ler - hormon)
```

Multiple degree of freedom tests

It appears as if we could trim the `full.model` down a bit in favor of the `alt.model`. Does this fit well with our disease model?

```
# check if we are loosing a significant amount of  
# information if we stick with the alternate model  
anova(full.model, alt.model)
```

```
## Analysis of Variance Table  
##  
## Model 1: lnodes ~ age + size + grade + meno + lpgr + ler + hormon  
## Model 2: lnodes ~ size + grade + meno  
##   Res.Df    RSS Df Sum of Sq      F Pr(>F)  
## 1      678 536.34  
## 2      682 539.27 -4    -2.9389 0.9288 0.4466
```

Multiple degree of freedom tests

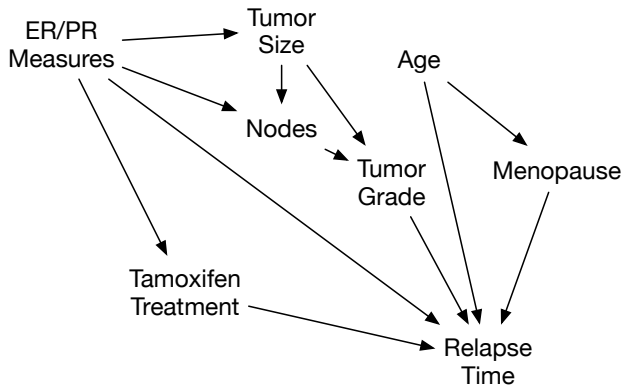


Figure 1: DAG summarizing our disease model.

Missing Data

Types of missing data

Missing completely at random: no factors relating to the samples (measured or not) influenced which data are missing.

- ▶ Example: A study ends prematurely, and some data are not able to be collected, independent of any sample characteristics.

Types of missing data

Missing at random: some important factors may have influenced which data are missing, but the probability of missingness is a function of variables that were measured.

- ▶ Example: Individuals with depression may be more likely to be lost to followup, resulting in missing data. As long as loss to followup isn't related to the variable that is missing (e.g. number of cigarettes smoked each week during the study), we can assume that the number of cigarettes smoked by individuals we did observe is representative of the number of cigarettes smoked by individuals we didn't observe.

Types of missing data

Informative missing: missing data are biased in some way by other confounding variables. This results in estimates that are higher or lower than they should be. This is often nearly impossible to detect without some outside information (e.g. experience with past studies or knowledge of the population under study).

- ▶ Example: Some unmeasured confounding variable (location: poor neighborhood) influences heavy smokers in the treatment group to drop out of the study at a higher rate than individuals who smoke less.

Problems stemming from missing data

Always characterize missingness of data. Ask questions like:

- ▶ What is the rate of missingness in each group?
- ▶ Are there any factors in our disease model that would cause an individual to have missing data?
- ▶ What other relationships between observed data and missingness exist?

Informative missingness can cause unexpected problems, including false associations.

- ▶ Example: In the previous presidential election, the “undecided” voters (i.e. likely voters with missing data) voted disproportionately for President Trump. This failed assumption resulted in unreliable polling leading up to the election.

Dealing with missing data: Ignore missing data

```
set.seed(239847)
n <- 100
# generate a dataset with a lot of missing data
dat <- data_frame(x1 = rnorm(n),
                  x2 = rnorm(n),
                  g = rbinom(n, 1, .5),
                  y = x1 + x2 + x1*x2 + (g == 1) + rnorm(n)) %>%
  prodNA()
dat
```

```
## # A tibble: 100 × 4
##       x1      x2      g      y
##   <dbl> <dbl> <int> <dbl>
## 1  0.43788108  2.0317248    NA  3.4269158
## 2  0.51206258 -0.4317878     1 -0.3195028
## 3  0.37422403  0.2143055     0  1.0164943
## 4  0.59381328  1.0086273     1  1.6125397
## 5  1.97367993  0.2218173     1  1.4531170
## 6 -0.73490075  0.3132395     0 -0.5675547
## 7  0.86301637 -0.2132996     0  0.7172030
## 8 -1.28152664 -0.2756235     0 -2.5823621
## 9  1.14205170  1.0586741     1  4.9751146
## 10 -0.09092382  1.8772961     1  1.0887537
## # ... with 90 more rows
```

Dealing with missing data: Ignore missing data

```
# look at the relationship between x and y by g  
lm(y ~ x1*x2 + g, data = dat) %>%  
  tidy()
```

##	term	estimate	std.error	statistic	p.value
## 1	(Intercept)	0.04631757	0.1777562	0.260568	7.952892e-01
## 2	x1	0.90845393	0.1310362	6.932848	2.809523e-09
## 3	x2	0.90704660	0.1392570	6.513474	1.483615e-08
## 4	g	0.70943829	0.2655294	2.671788	9.626272e-03
## 5	x1:x2	1.48408561	0.1522519	9.747566	3.969240e-14

Dealing with missing data: Replace missing data with the group mean

```
# replace
dat <- mutate(dat,
              mx1 = ifelse(is.na(x1), mean(dat$x1, na.rm = TRUE), x1),
              mx2 = ifelse(is.na(x2), mean(dat$x2, na.rm = TRUE), x2),
              mg = ifelse(is.na(g), mean(dat$g, na.rm = TRUE), g))

# look at the relationship between x and y by g
lm(y ~ mx1*mx2 + g, data = dat) %>%
  tidy()
```

##	term	estimate	std.error	statistic	p.value
## 1	(Intercept)	0.1157842	0.1746073	0.6631121	5.091650e-01
## 2	mx1	1.0477377	0.1377082	7.6083911	4.710378e-11
## 3	mx2	0.9235718	0.1485312	6.2180344	2.152961e-08
## 4	g	0.6426178	0.2562871	2.5074140	1.418613e-02
## 5	mx1:mx2	1.4630055	0.1647692	8.8791187	1.517980e-13

Dealing with missing data: Impute

```
options(warn = -1)
imp <- select(dat, -mx1, -mx2, -mg) %>%
  as.data.frame() %>% # won't accept a tibble
  missForest()
```

```
## missForest iteration 1 in progress...done!
## missForest iteration 2 in progress...done!
## missForest iteration 3 in progress...done!
## missForest iteration 4 in progress...done!
## missForest iteration 5 in progress...done!
```

```
lm(y ~ x1*x2 + g, data = imp$ximply) %>%
  tidy()
```

	term	estimate	std.error	statistic	p.value
## 1	(Intercept)	-0.003054707	0.13611685	-0.0224418	9.821426e-01
## 2	x1	0.963516376	0.10643253	9.0528370	1.748050e-14
## 3	x2	0.989975090	0.09921711	9.9778661	1.842074e-16
## 4	g	0.735507903	0.20356932	3.6130587	4.860852e-04
## 5	x1:x2	1.388411851	0.11147306	12.4551335	1.101380e-21