# IS5 in R: Comparing Groups (Chapter 17)

*Margaret Chien and Nicholas Horton (nhorton@amherst.edu)*

*July 16, 2018*

## Introduction and background

This document is intended to help describe how to undertake analyses introduced as examples in the Fifth Edition of *Intro Stats* (2018) by De Veaux, Velleman, and Bock. More information about the book can be found at http://wps.aw.com/aw_deveaux_stats_series. This file as well as the associated R Markdown reproducible analysis source file used to create it can be found at http://nhorton.people.amherst.edu/is5.

This work leverages initiatives undertaken by Project MOSAIC (http://www.mosaic-web.org), an NSF-funded effort to improve the teaching of statistics, calculus, science and computing in the undergraduate curriculum. In particular, we utilize the `mosaic` package, which was written to simplify the use of R for introductory statistics courses. A short summary of the R needed to teach introductory statistics can be found in the mosaic package vignettes (http://cran.r-project.org/web/packages/mosaic). A paper describing the mosaic approach was published in the *R Journal*: https://journal.r-project.org/archive/2017/RJ-2017-024.

## Chapter 17: Comparing Groups

```r
library(mosaic)
library(readr)
library(janitor)
```

### Section 17.1: A Confidence Interval for the Difference Between Two Proportions

```r
# Creating a data frame for Seatbelts
Seatbelts <- rbind(
  do(2777)        * data.frame(passenger = "F", belted = TRUE),
  do(4208 - 2777) * data.frame(passenger = "F", belted = FALSE),
  do(1363)        * data.frame(passenger = "M", belted = TRUE),
  do(2763 - 1363) * data.frame(passenger = "M", belted = FALSE)
) %>%
  select(passenger, belted)
```

Here, the `do()` function creates a number of rows for the data frame.

```r
set.seed(234)
numsim <- 10000

# What does do() do?
resample(Seatbelts) %>%
  group_by(passenger) %>%
  summarise(proportion = sum(belted)/n()) %>%
  summarise(diffprop = abs(diff(proportion))) # Difference of proportions from one random resample
```

```
## # A tibble: 1 x 1
##   diffprop
##      <dbl>
```

```
## 1     0.170
```

```r
resample(Seatbelts) %>%
  group_by(passenger) %>%
  summarise(proportion = sum(belted)/n()) %>%
  summarise(diffprop = abs(diff(proportion))) # Difference of proportions from another random resample
```
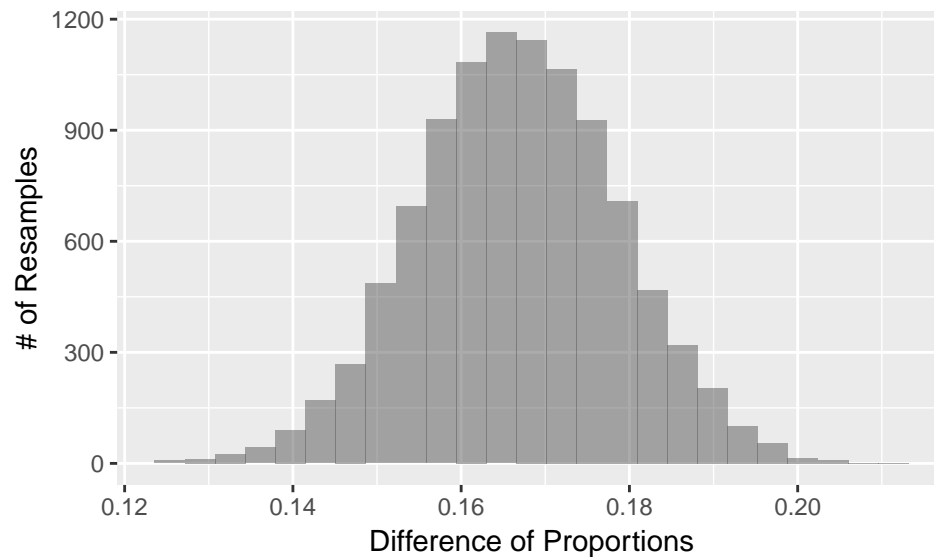
```
## # A tibble: 1 x 1
##   diffprop
##      <dbl>
## 1    0.178
```

```r
do(2) * resample(Seatbelts) %>%
  group_by(passenger) %>%
  summarise(proportion = sum(belted)/n()) %>%
  summarise(diffprop = abs(diff(proportion))) # Calculates two differences
```

```
##    diffprop
## 1 0.1573754
## 2 0.1548267
```

```r
# We need 10000 differences of proportions
seatbeltresamples <- do(numsim) * resample(Seatbelts) %>%
  group_by(passenger) %>%
  summarise(proportion = sum(belted)/n()) %>%
  summarise(diffprop = abs(diff(proportion)))

# Figure 17.1, page 542
gf_histogram(~ diffprop, data = seatbeltresamples) %>%
  gf_labs(x = "Difference of Proportions", y = "# of Resamples")
```



### Example 17.1: Finding the Standard Error of a Difference in Proportions

```r
# Creating the data set for online profiles
OnlineProf <- rbind(
  do(248 * .57)     * data.frame(gender = "M", profile = TRUE),
  do(248 * .43 + 1) * data.frame(gender = "M", profile = FALSE), # Add one for rounding errors
  do(256 * .70)     * data.frame(gender = "F", profile = TRUE),
```

```
  do(256 * .30 + 1) * data.frame(gender = "F", profile = FALSE)
)
tally(~ gender, data = OnlineProf)
```

```
## gender
##   M   F
## 248 256
```

```
OnlineProfM <- OnlineProf %>%
  filter(gender == "M")
sepboys <- ((mean(~ profile, data = OnlineProfM) * (1 - mean(~ profile, data = OnlineProfM)))/nrow(Onli
sepboys
```

```
## [1] 0.03145024
```

```
OnlineProfF <- OnlineProf %>%
  filter(gender == "F")
sepgirls <- ((mean(~ profile, data = OnlineProfF) * (1 - mean(~ profile, data = OnlineProfF)))/nrow(Onli
sepgirls
```

```
## [1] 0.02866236
```

```
sep <- (sepboys^2 + sepgirls^2)^.5
sep
```

```
## [1] 0.04255171
```

**Example 17.2: Finding a Two-Proportion $z$-Interval**

```
zstats <- qnorm(p = c(.025, .975))
(mean(~ profile, data = OnlineProfF) - mean(~ profile, data = OnlineProfM)) + zstats * sep
```

```
## [1] 0.04727054 0.21407019
```

**Section 17.2: Assumptions and Conditions for Comparing Proportions**

**Section 17.3: The Two-Sample $z$-Test: Testing for the Difference Between Proportions**

**Step-By-Step Example: A Two-Proportion $z$-Test**

```
# Create the data set
SleepHabits <- rbind(
  do(205)       * data.frame(gen = "GenY", internet = TRUE),
  do(293 - 205) * data.frame(gen = "GenY", internet = FALSE),
  do(235)       * data.frame(gen = "GenX", internet = TRUE),
  do(469 - 235) * data.frame(gen = "GenX", internet = FALSE)
)
```

```
# Mechanics
ngeny <- nrow(filter(SleepHabits, gen == "GenY"))
ngeny
```

```
## [1] 293
```

```
ygeny <- nrow(filter(SleepHabits, gen == "GenY" & internet == TRUE))
ygeny
```

```
## [1] 205
```

```
pgeny <- mean(~ internet, data = filter(SleepHabits, gen == "GenY"))
pgeny
```

```
## [1] 0.6996587
```

```
ngenx <- nrow(filter(SleepHabits, gen == "GenX"))
ngenx
```

```
## [1] 469
```

```
ygenx <- nrow(filter(SleepHabits, gen == "GenX"  & internet == TRUE))
ygenx
```

```
## [1] 235
```

```
pgenx <- mean(~ internet, data = filter(SleepHabits, gen == "GenX"))
pgenx
```

```
## [1] 0.5010661
```

```
sepgen <- ((pgeny * (1 - pgeny))/ngeny + (pgenx * (1 - pgenx))/ngenx)^.5
sepgen
```

```
## [1] 0.03535867
```

```
pdiff <- pgeny - pgenx
pdiff
```

```
## [1] 0.1985926
```

```
z <- (pdiff - 0)/sepgen
z
```

```
## [1] 5.616518
```

```
2 * pnorm(q = z, lower.tail = FALSE)
```

```
## [1] 1.948444e-08
```

**Section 17.4: A Confidence Interval for the Difference Between Two Means**

**Example 17.7: Finding a Confidence Interval for the Difference in Sample Means**
```
# Not sure if creating data set is really necessary
# page 555
nord <- 27
nref <- 27
yord <- 8.5
yref <- 14.7
sord <- 6.1
sref <- 8.4

seys <- 2.0
diffy <- yref - yord # 6.2
tstats <- qt(p = c(.025, .975), df = 47.46)
tstats
```

```
## [1] -2.011226  2.011226
```

4

```
me <- tstats * seys
me
```

```
## [1] -4.022452  4.022452
```

```
diffy + me
```

```
## [1]  2.177548 10.222452
```

**Section 17.5: The Two-Sample *t*-Test: Testing for the Difference Between Two Means**

**Step-By-Step Example: A Two-Sample *t*-Test for the Difference Between the Two Means**

```
# page 556
BuyingCam <- read_csv("http://nhorton.people.amherst.edu/is5/data/Buy_from_a_friend.csv")
```

```
## Parsed with column specification:
## cols(
##   Friend = col_integer(),
##   Stranger = col_integer()
## )
```
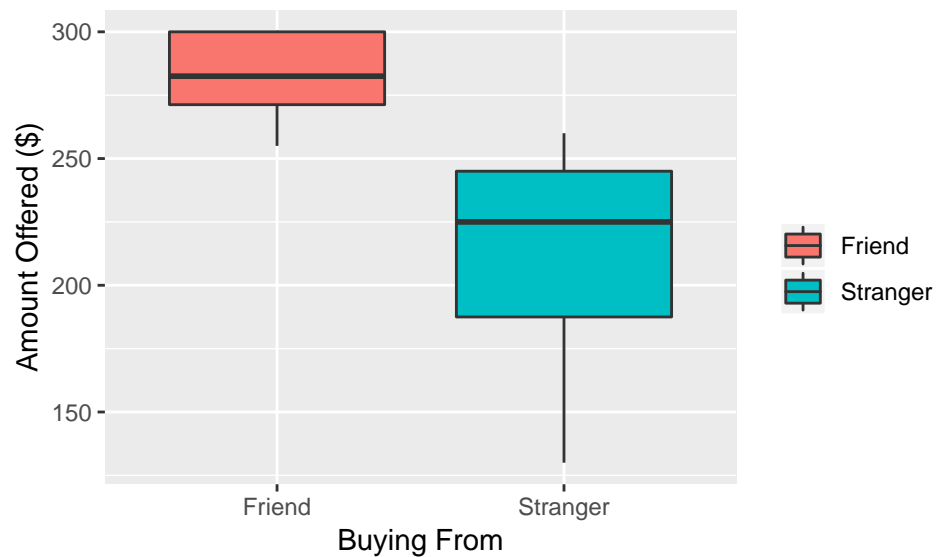
By default, `read_csv()` prints the variable names. These messages can be suppressed using the `message=FALSE` code chunk option to save space and improve readability.

```
library(tidyr) # for gather() function
```

```
##
## Attaching package: 'tidyr'
```

```
## The following object is masked from 'package:Matrix':
##
##     expand
```
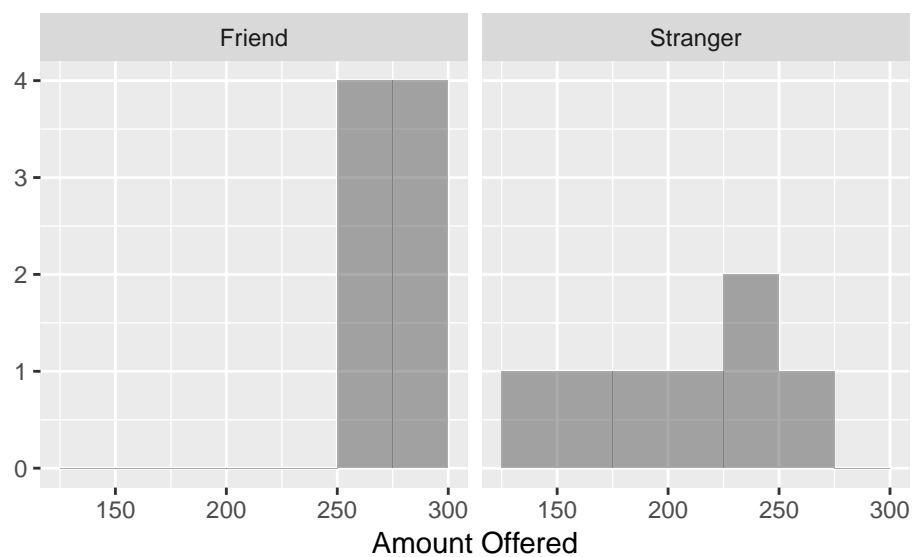
```
BuyingCam <- BuyingCam %>%
  gather(key = buying_type, value = amount_offered, Friend, Stranger)
# Model
gf_boxplot(amount_offered ~ buying_type, fill = ~ buying_type, data = BuyingCam) %>%
  gf_labs(x = "Buying From", y = "Amount Offered ($)", fill = "")
```

```
## Warning: Removed 1 rows containing non-finite values (stat_boxplot).
```

```r
gf_histogram(~ amount_offered, binwidth = 25, center = 12.5, data = BuyingCam) %>% # doesn't exactly ma
  gf_facet_wrap(buying_type ~ .) %>%
  gf_labs(x = "Amount Offered", y = "")
```

```
## Warning: Removed 1 rows containing non-finite values (stat_bin).
```



```r
# Mechanics
favstats(~ amount_offered | buying_type, data = BuyingCam)
```

```
##   buying_type min     Q1 median  Q3 max     mean       sd n missing
## 1      Friend 255 271.25  282.5 300 300 281.8750 18.31032 8       0
## 2    Stranger 130 187.50  225.0 245 260 211.4286 46.43223 7       1
```

**Section 17.6: Randomization Tests and Confidence Intervals for Two Means**

```r
Cars <- read_csv("http://nhorton.people.amherst.edu/is5/data/Car_speeds.csv")
```

```
## Parsed with column specification:
```

```
## cols(
##    direction = col_character(),
##    speed = col_double()
## )
```

```r
# Figure 17.2 (page 560) is the same as Figure 4.4 (page 102)
favstats(~ speed | direction, data = Cars)
```

```
##   direction   min      Q1 median      Q3   max     mean       sd   n
## 1      Down 10.27 20.4675 22.885 25.3525 32.95 22.71708 3.622006 250
## 2        Up 15.08 22.4975 25.155 28.1600 34.97 25.25172 3.856331 250
##   missing
## 1       0
## 2       0
```

**Section 17.7: Pooling**

**Section 17.8: The Standard Deviation of a Difference**