

IS5 in R: Displaying and Describing Data (Chapter 2)

Nicholas Horton (nhorton@amherst.edu)

December 17, 2020

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. 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 (<https://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 2: Displaying and Describing Data

Section 2.1: Summarizing and Displaying a Categorical Variable

```
library(mosaic)
library(readr)
library(janitor) # for variable names
options(digits = 3)
Titanic <- read_csv("http://nhorton.people.amherst.edu/is5/data/Titanic.csv")
```

By default, `read_csv()` prints the variable names. These messages were suppressed using the `message=FALSE` code chunk option to save space and improve readability.

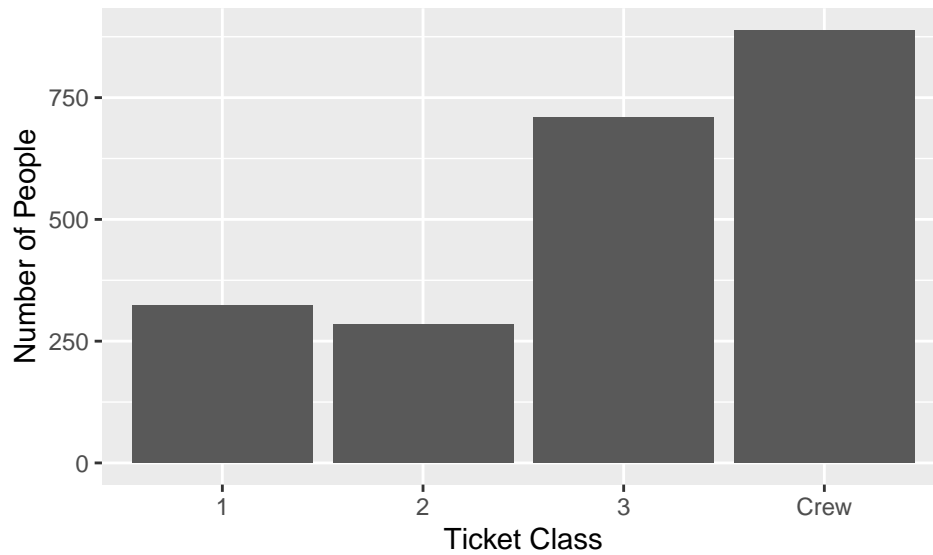
```
# Table 2.2, page 19
tally(~Class, data = Titanic)
```

```
## Class
##      1      2      3 Crew
## 324  285  710  889
```

```
# Table 2.3
tally(~Class, format = "percent", data = Titanic)
```

```
## Class
##      1      2      3 Crew
## 14.7 12.9 32.2 40.3
```

```
# Figure 2.2, page 19
gf_bar(~Class, data = Titanic) %>%
  gf_labs(x = "Ticket Class", y = "Number of People")
```



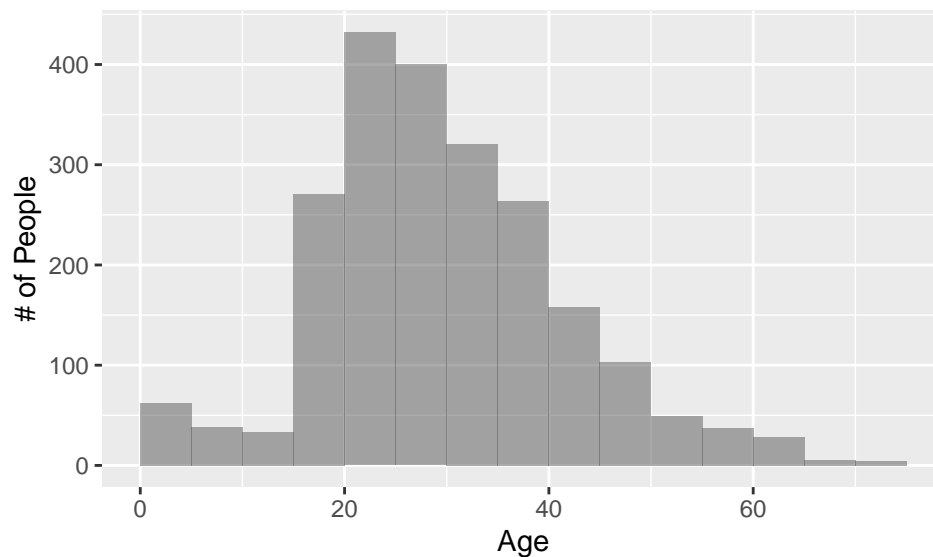
`GOAL(~ X)` is the general form of the modeling language for one variable in the `mosaic` package. We use `gf_bar()` to make a bar graph using the `ggformula` system, which is automatically downloaded with the `mosaic` package.

Section 2.2: Displaying a Quantitative Variable

```
# Figure 2.7, page 24
gf_histogram(~Age, data = Titanic, binwidth = 5, ylab = "# of People", center = 5 / 2)
```

Ages of Those Aboard the Titanic

```
## Warning: Removed 3 rows containing non-finite values (stat_bin).
```



The function generates a warning because three of the ages are missing; this output can be suppressed by adding `warning=FALSE` as an option in this code chunk.

```
# Example 2.3, page 25
```

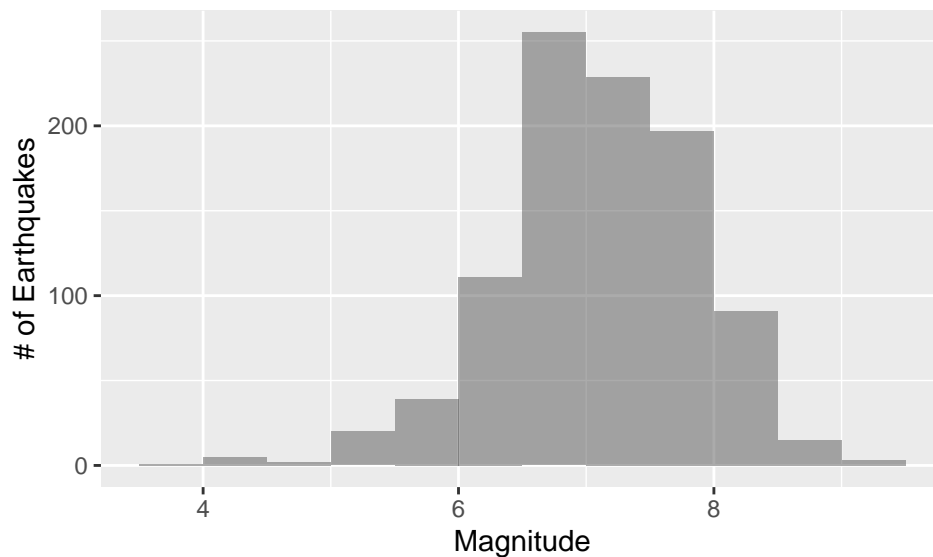
```
Earthquakes <- read_csv("http://nhorton.people.amherst.edu/is5/data/Tsunamis_2016.csv")
```

Earthquakes and Tsunamis

```
##
## -- Column specification -----
## cols(
##   Year = col_double(),
##   Focal_Depth = col_double(),
##   Primary_Magnitude = col_double(),
##   Country = col_character(),
##   Latitude = col_double(),
##   Longitude = col_double(),
##   Deaths = col_double(),
##   Missing = col_double(),
##   Injuries = col_double(),
##   `Damage($M)` = col_double()
## )

gf_histogram(~Primary_Magnitude,
  data = Earthquakes, binwidth = 0.5,
  ylab = "# of Earthquakes", xlab = "Magnitude", center = 0.25
)
```

```
## Warning: Removed 119 rows containing non-finite values (stat_bin).
```



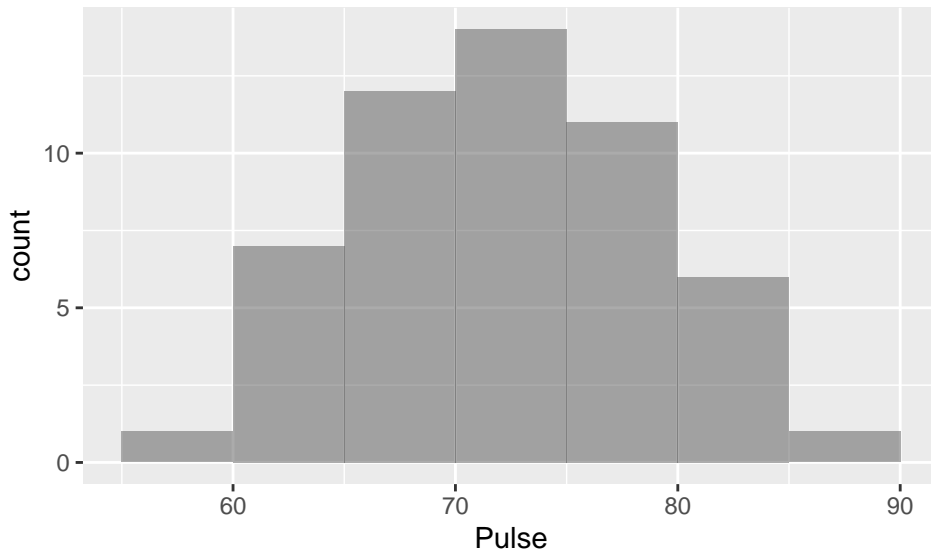
Stem-and-Leaf Displays See page 26.

Figure 2.8, page 26

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

```
##
## -- Column specification -----
## cols(
##   Pulse = col_double()
## )
```

```
gf_histogram(~Pulse, data = Pulse_rates, binwidth = 5, center = 5 / 2)
```



```
with(Pulse_rates, stem(Pulse))
```

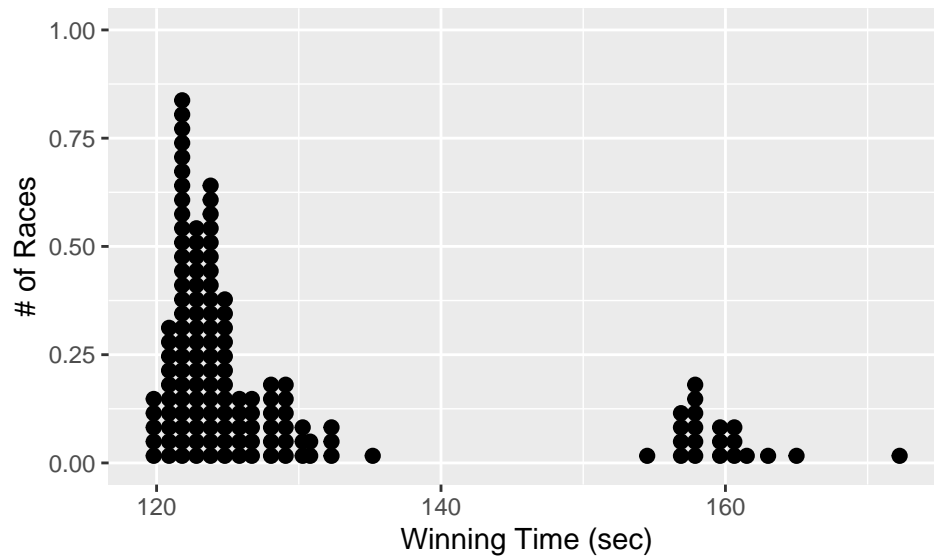
```
##
## The decimal point is 1 digit(s) to the right of the |
##
## 5 | 7
## 6 | 13444
## 6 | 556668888899
## 7 | 0012223333444
## 7 | 5557777888889
## 8 | 0112233
## 8 | 6
```

Figure 2.9, page 27

```
Derby <- read_csv("http://nhorton.people.amherst.edu/is5/data/Kentucky_Derby_2016.csv")
```

Dotplot

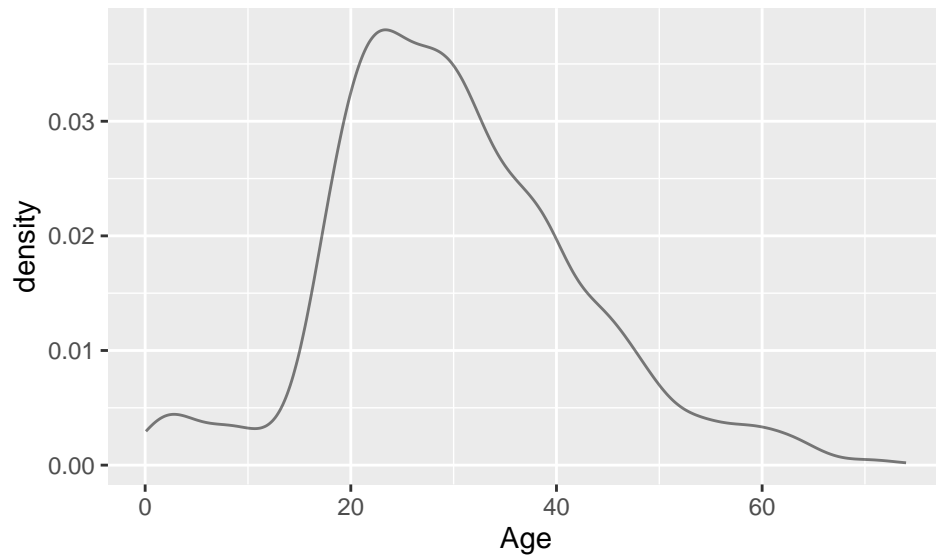
```
##
## -- Column specification -----
## cols(
##   Year = col_double(),
##   Year_no = col_double(),
##   Date = col_character(),
##   Winner = col_character(),
##   Mins = col_double(),
##   Secs = col_double(),
##   Time_Sec = col_double(),
##   Distance = col_double(),
##   Speed_mph = col_double()
## )
gf_dotplot(~Time_Sec, data = Derby, binwidth = 1) %>%
  gf_labs(x = "Winning Time (sec)", y = "# of Races")
```



```
# Figure 2.10, page 27
gf_dens(~Age, data = Titanic)
```

Density Plots

```
## Warning: Removed 3 rows containing non-finite values (stat_density).
```



Section 2.3: Shape

See displays on pages 28-29.

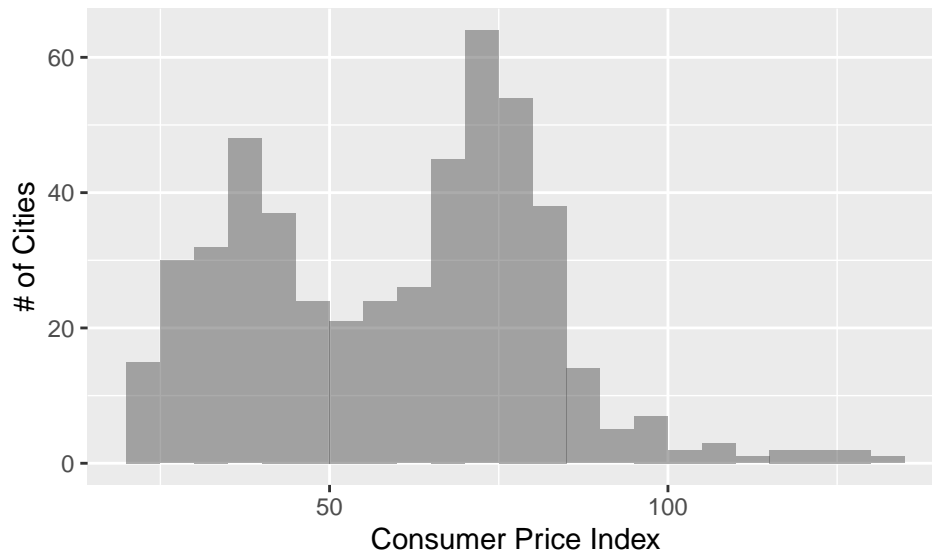
```
CPI <- read_csv("http://nhorton.people.amherst.edu/is5/data/CPI_Worldwide.csv") %>%
  janitor::clean_names()
names(CPI)
```

Consumer Price Index

```
## [1] "city" "consumer_price_index"
## [3] "rent_index" "consumer_price_plus_rent_index"
## [5] "groceries_index" "restaurant_price_index"
## [7] "local_purchasing_power_index"
```

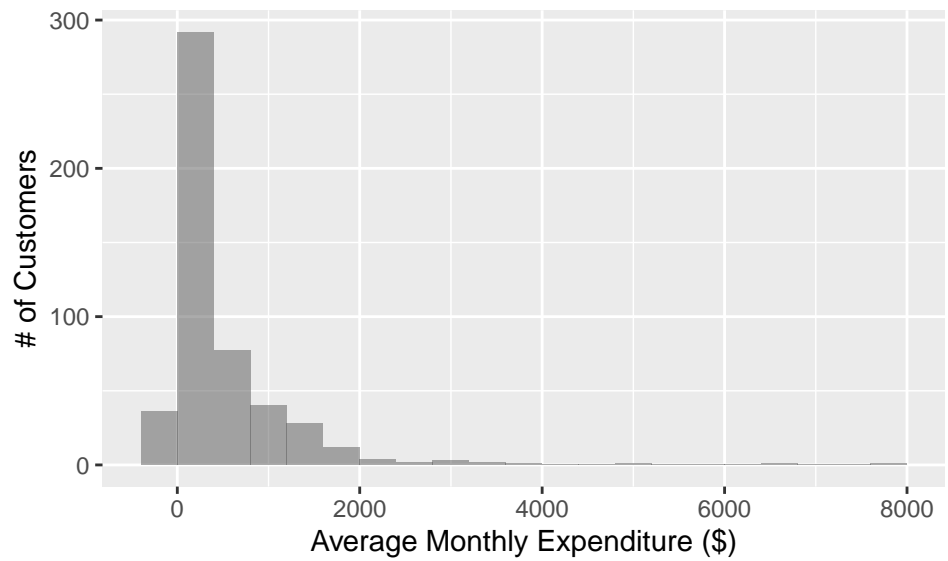
Example 2.5, page 30

```
gf_histogram(~consumer_price_index,
  data = CPI, ylab = "# of Cities",
  xlab = "Consumer Price Index", binwidth = 5, center = 5 / 2
)
```



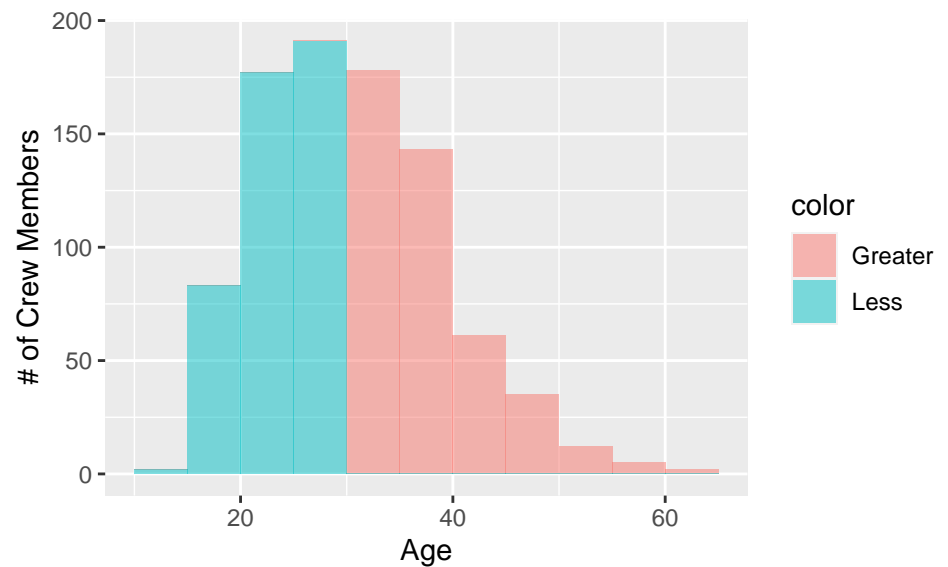
We can use `clean_names()` from the `janitor` package to format the names of the columns when necessary. You can use the `names()` function to check the reformatted names. The pipe operator (`%>%`) takes the output of the line of code and uses it in the next.

```
CreditCardEx <- read_csv("http://nhorton.people.amherst.edu/is5/data/Credit_card_charges.csv") %>%
  janitor::clean_names()
# Figure 2.6, page 30
gf_histogram(~charges,
  data = CreditCardEx, ylab = "# of Customers",
  xlab = "Average Monthly Expenditure ($)", binwidth = 400, center = 200
)
```



Section 2.4: Center

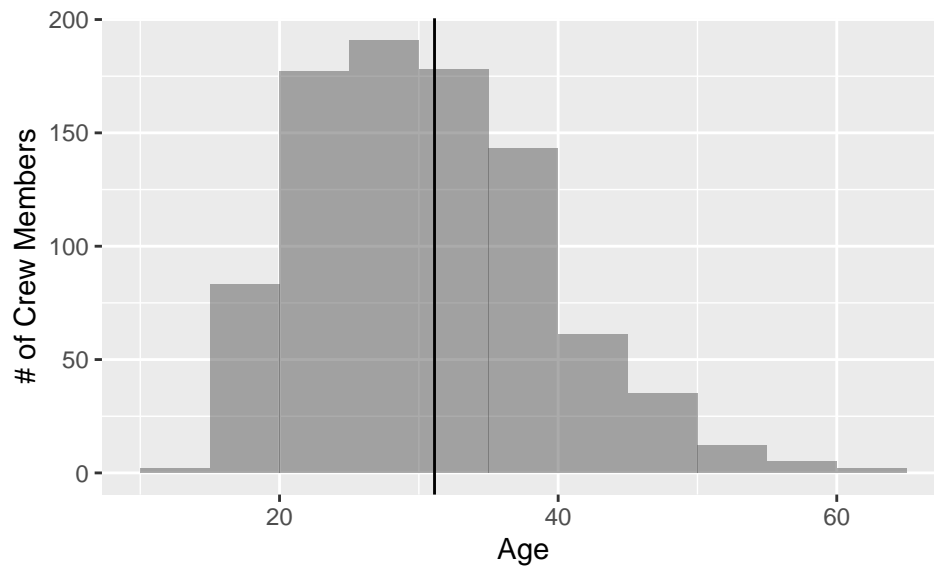
```
TitanicCrew <- filter(Titanic, Class == "Crew")
# Figure 2.15, page 32
TitanicCrew %>%
  mutate(color = ifelse(Age <= median(Age), "Less", "Greater")) %>%
  gf_histogram(~Age, fill = ~color, binwidth = 5, center = 5 / 2, ylab = "# of Crew Members")
```



Finding Median and Mean

```
# Figure 2.16
gf_histogram(~Age, data = TitanicCrew, ylab = "# of Crew Members", binwidth = 5, center = 5 / 2) %>%
  gf_vline(xintercept = mean(~Age, data = TitanicCrew))
```

```
## Warning: geom_vline(): Ignoring `mapping` because `xintercept` was provided.
```



```
df_stats(~Age, data = TitanicCrew)
```

```
##   response min Q1 median Q3 max mean   sd   n missing
## 1      Age  14 24    30 37  62 31.1 8.55 889      0
```

Another way to generate summary statistics is the `favstats()` command (we will stick to `df_stats()` because it is more flexible).

```
favstats(~Age, data = TitanicCrew)
```

```
##   min Q1 median Q3 max mean   sd   n missing
##   14 24    30 37  62 31.1 8.55 889      0
```

Section 2.5: Spread

```
range(~Age, data = TitanicCrew)
```

The Range

```
## [1] 14 62
```

```
diff(range(~Age, data = TitanicCrew))
```

```
## [1] 48
```

The `range()` function returns the maximum and minimum values, so we can use the `diff()` function to find the difference between the two values.

```
df_stats(~Age, data = TitanicCrew)
```

The Interquartile Range

```
##   response min Q1 median Q3 max mean   sd   n missing
## 1      Age  14 24    30 37  62 31.1 8.55 889      0
```

```
IQR(~Age, data = TitanicCrew)
```

```
## [1] 13
```


Using the `IQR()` function allows us to avoid having to manually find the IQR by subtracting Q1 from Q3 from the `df_stats()` output.

```
sd(~Age, data = TitanicCrew)
```

Standard Deviation

```
## [1] 8.55
```

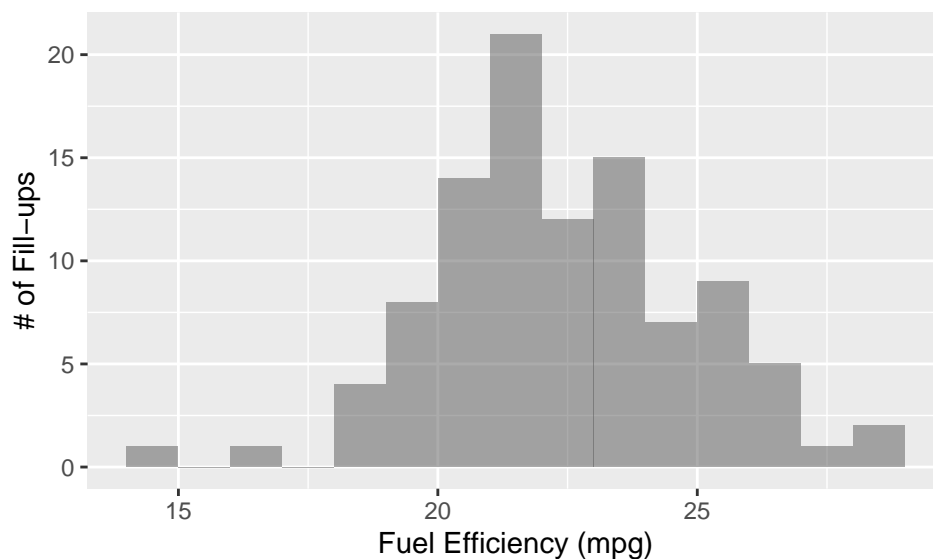
```
var(~Age, data = TitanicCrew)
```

```
## [1] 73.1
```

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

Summarizing a Distribution

```
##  
## -- Column specification -----  
## cols(  
##   mpg = col_double()  
## )  
  
# Step-by-Step Example, page 39  
gf_histogram(~mpg,  
  data = Nissan, binwidth = 1, xlab = "Fuel Efficiency (mpg)",  
  ylab = "# of Fill-ups", center = 5 / 2  
)
```

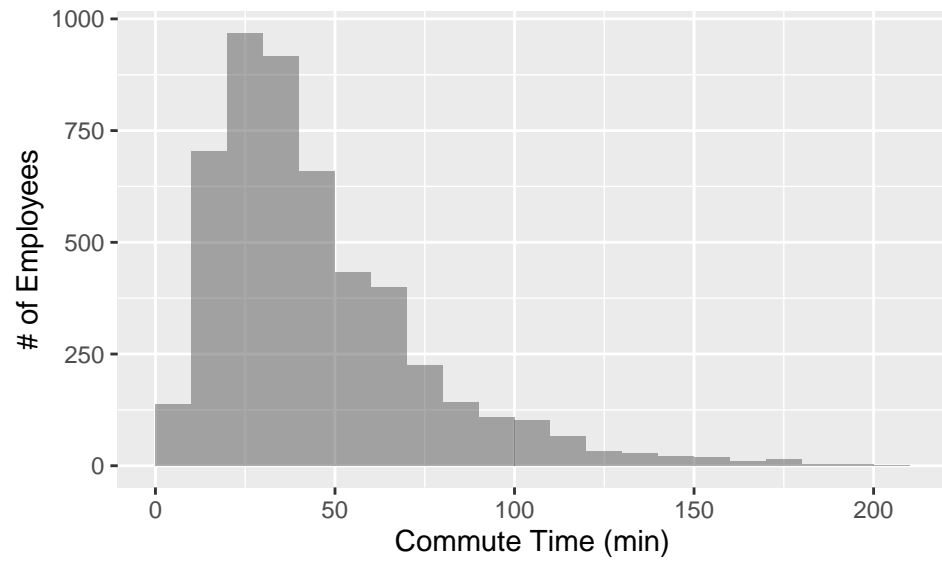


```
df_stats(~mpg, data = Nissan)
```

```
##   response min   Q1 median Q3   max mean   sd   n missing  
## 1      mpg 14.7 20.8   22.1 24  28.2 22.4 2.45 100      0
```

```
Commute <- read_csv("http://nhorton.people.amherst.edu/is5/data/Population_Commute_Times.csv") %>%  
  janitor::clean_names()
```

```
# Figure 2.19, page 40
gf_histogram(~commute_time,
  data = Commute, binwidth = 10, xlab = "Commute Time (min)",
  ylab = "# of Employees", center = 5
)
```



Random Matters