# Histograms

## Introduction

One of the most critical steps in data science – one I'll confess I often overlook in my rush to look at treatment differences – is to inspect the distribution of data. This allows us to confirm whether the distribution follows a bell shape, and also to identify extreme values that should be further investigated. There are a multiple ways to do this. This week, we will focus on the histogram.

## Case Study

The first example we will use is a cotton uniformity trial conducted in Greece in 1938. Yes, it is a bit dated, but the data are open-source and the concept is one we deal with every day in agricultural research: how consistent are our plots? In other words, what is the variance among plots?

In a uniformity trial, a several plots are managed identically so that the distribution of their yields can be measured. These data can be used to determine whether that field is a good site for a research trial, or how data might vary n a similar field.

Our first step is to load the data, which is in a .csv file. To do this, we will use the `read.csv` function.

```
cotton = read.csv("data/cotton_uniformity.csv")
```

## Basic Histogram

Drawing a simple histogram in R requires two steps:

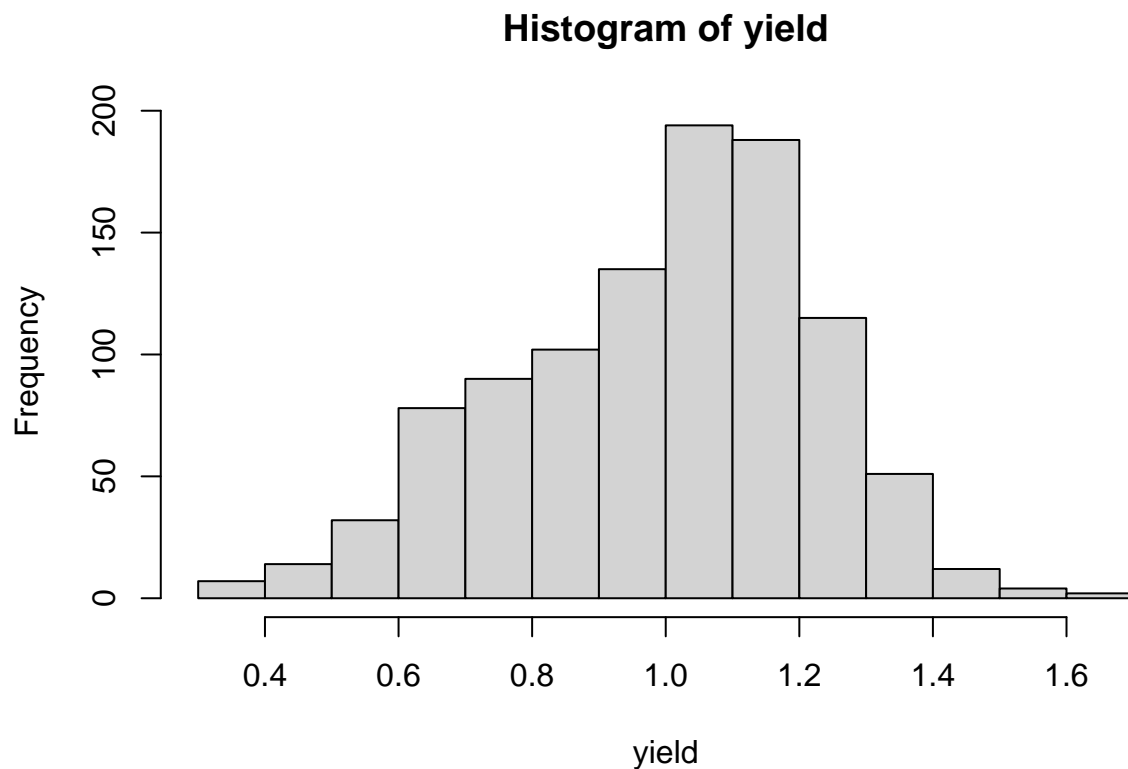1. Define the column in which we are interested.

2. Create the histogram.

Since we are using a data.frame, we need to tell R which column to plot. We will first extract the data from the data.frame We will call this column "yield". To extract the yield column from our data.frame, we reference the column of the data.frame as `cotton$yield`. To the left of the dollar sign ($) is the name of the data.frame, *cotton*. To the right of the dollar sign is the name of the column, *yield*.
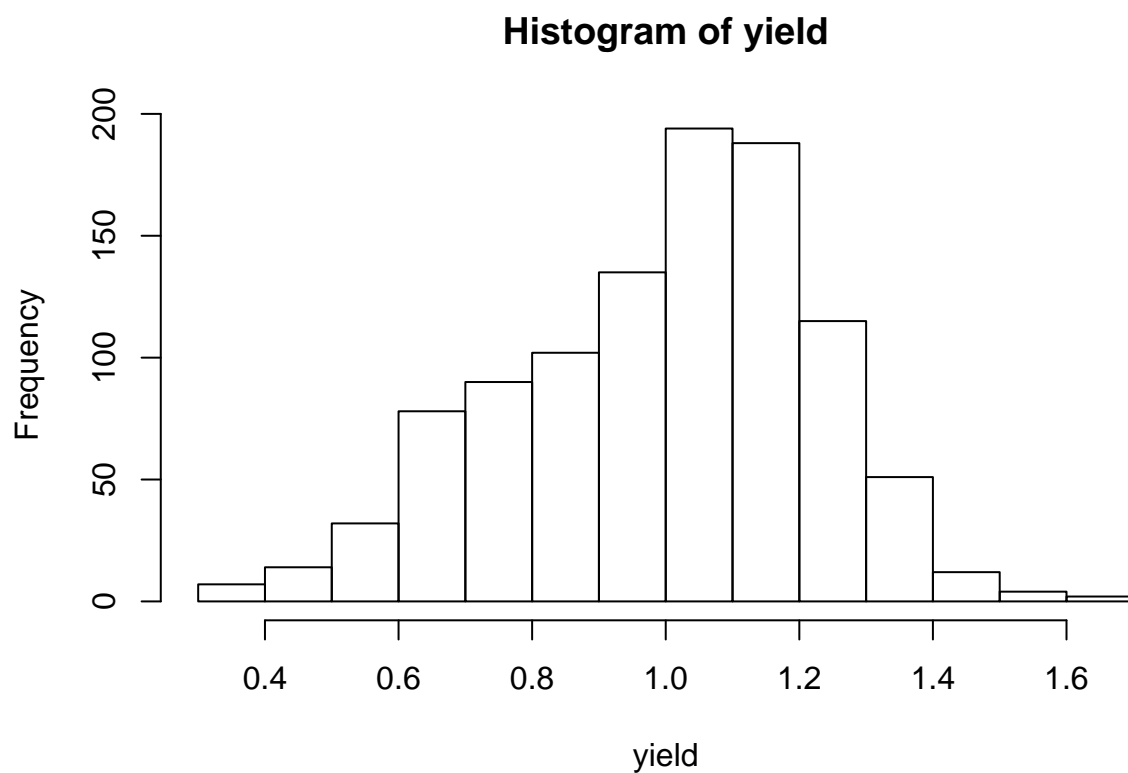
```
yield = cotton$yield
```

We then can create the histogram using the `hist()` function. We will assign the output to another R object, *histogram*. As we will see, that object contains much more information than just our histogram plot.

In the following line, we tell R to plot the histogram, using the `plot()` function.

```
histogram = hist(yield)
```

**Histogram of yield**



```
plot(histogram)
```

**Histogram of yield**



And, *voila*!, we have our histogram. Each bar represents a range in values called a **bin**.

If we view the histogram object itself, however, we will see it is composed of many different statistics. histogram is a kind of R object called a **list**. It is a collection of many data frames.

```
histogram
```

```
## $breaks
##  [1] 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0 1.1 1.2 1.3 1.4 1.5 1.6 1.7
##
## $counts
##  [1]   7  14  32  78  90 102 135 194 188 115  51  12   4   2
##
## $density
##  [1] 0.06835938 0.13671875 0.31250000 0.76171875 0.87890625 0.99609375
##  [7] 1.31835938 1.89453125 1.83593750 1.12304688 0.49804687 0.11718750
## [13] 0.03906250 0.01953125
##
## $mids
##  [1] 0.35 0.45 0.55 0.65 0.75 0.85 0.95 1.05 1.15 1.25 1.35 1.45 1.55 1.65
##
## $xname
## [1] "yield"
##
## $equidist
## [1] TRUE
##
## attr(,"class")
## [1] "histogram"
```

Each of these objects can be viewed separately, using the dollar sign method we used above.

For example, we can see the upper and lower yield limits of each bins by running `histogram$breaks`.

```
histogram$breaks
```

```
##  [1] 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0 1.1 1.2 1.3 1.4 1.5 1.6 1.7
```

Alternatively, we could see the midpoints, or yield values that define the middle of each bin by running `histogram$mids` as a command.
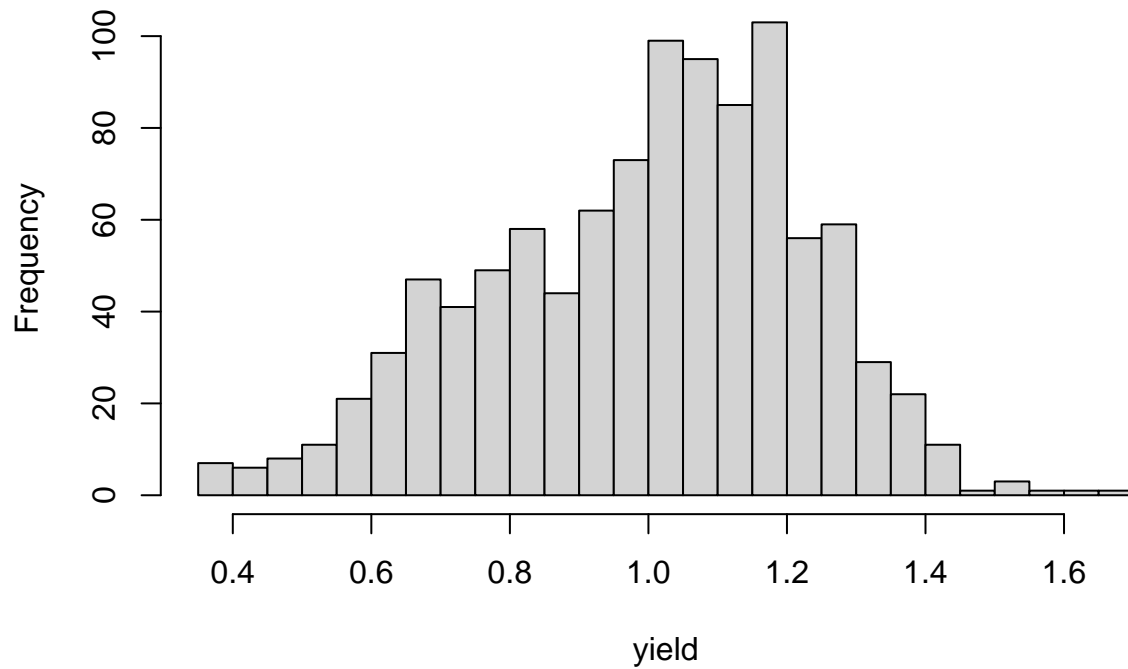
```
histogram$mids
```

```
##  [1] 0.35 0.45 0.55 0.65 0.75 0.85 0.95 1.05 1.15 1.25 1.35 1.45 1.55 1.65
```

As we saw in the lesson, varying the number of columns can affect how we see patterns in the data. In the plot above, we have 14 bars. Each bar represents a bin width of 0.1. What if we tell R to use 28 bins to draw the histogram? We can do that by adding the argument `breaks=28` to our `hist()` function.
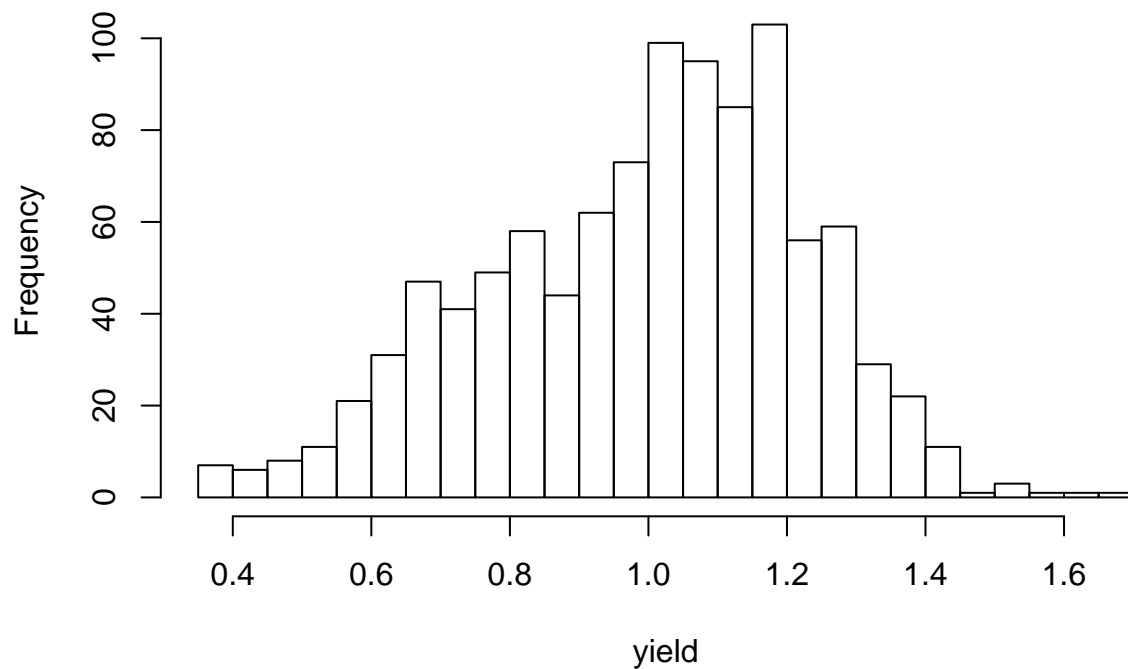
```
histogram = hist(yield, breaks = 28)
```

**Histogram of yield**



```
plot(histogram)
```

**Histogram of yield**



Note that we ended up with 27, not 28 bins, This is because it takes 28 breaks to define 27 bins. You can count the breaks, or sides of the bars in the histogram, to prove this to yourself. Each bar now has a bin width that is 0.05 kg. We can verify this by again inspecting the data by running `histogram$breaks`.

```
histogram$breaks
```

```
##  [1] 0.35 0.40 0.45 0.50 0.55 0.60 0.65 0.70 0.75 0.80 0.85 0.90 0.95 1.00 1.05
## [16] 1.10 1.15 1.20 1.25 1.30 1.35 1.40 1.45 1.50 1.55 1.60 1.65 1.70
```

# Histograms with ggplot

One of the great features of R is that there are hundreds of **packages** that can be added to it. These are like add-ins that you make have used with *Word* or *Excel.* Packages, like R, are open source, and are composed of specialized functions for many, many areas of research: agronomy, spatial data, ecology, medicine, machine learning, working with online databases like *SSURGO,* and so forth.

`ggplot2` is a one very powerful and population package in R, used to creating highly-customized plots. I use `ggplot2` to create a most of the plots I use in industry. if you can envision a plot, you almost certainly can draw it using `ggplot2`, .

There are entire books written on `ggplot2` (here is one I have consult regularly: https://r-graphics.org/index .html). I don't have time to teach it to you in depth, but I will show you many basics in this course.

The first time we want to use a package in R, we need to install it using the `install.packages()` function. While I try to install all packages you need for this course in *RStudio Cloud*, if you download *RStudio* to your desktop, you will need to install packages yourself. To install `ggplot2`, we run the following command:

```
install.packages("ggplot2")
```

```
## Installing package into '/cloud/lib/x86_64-pc-linux-gnu-library/4.1'
## (as 'lib' is unspecified)
```
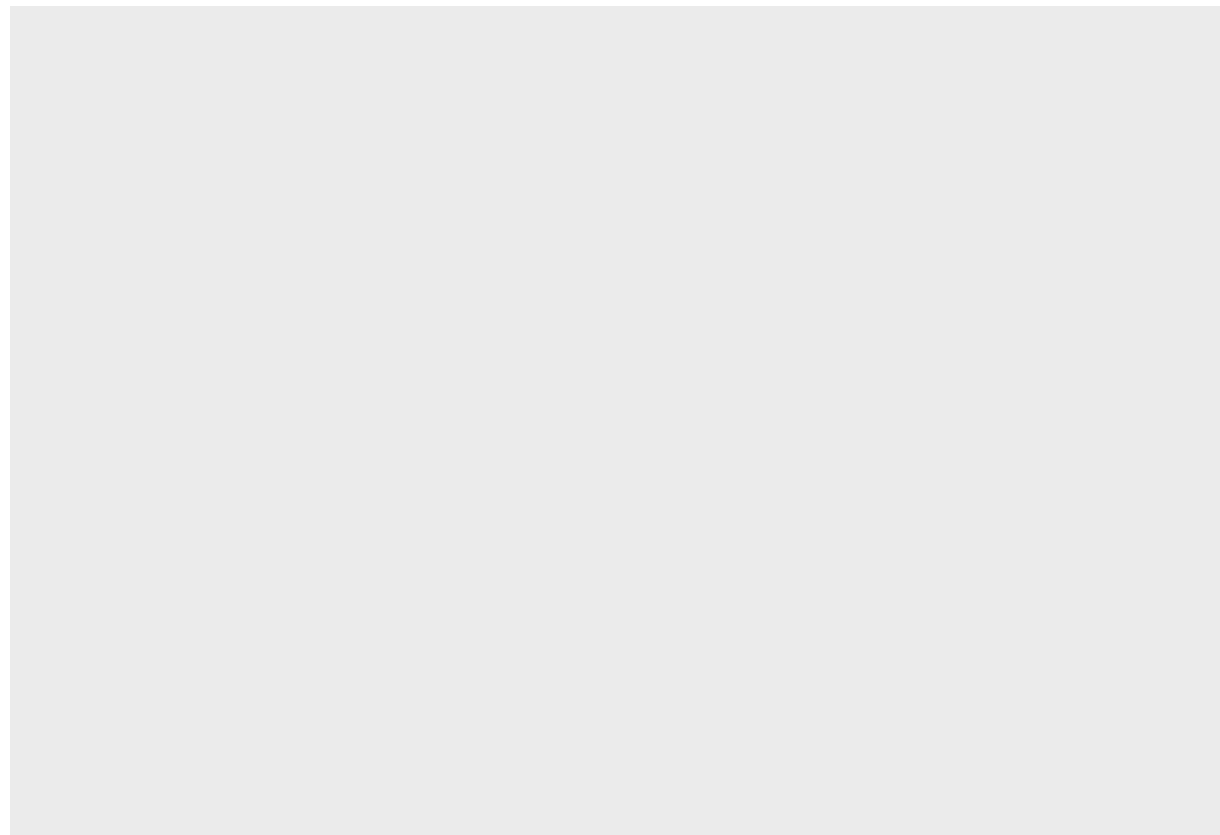
Note: you only need to install a package the very first time you use it. We can now start the package using the `library()` function.

```
library(ggplot2)
```

You only need to load a package once per R session. But if you restart R on your desktop, or log out and back into *RStudio Cloud*, you may need to reload the packages you intend to use. I make it a point to load the package in each notebook, rather than assuming it is loaded from another notebook I was working on at the same time.

The first line of code calls `ggplot()`. Although the package is called `ggplot2`, the function we use is `ggplot.` Confusing, I know.We use `data=cotton` to tell R we are drawing data from the "cotton" data frame.
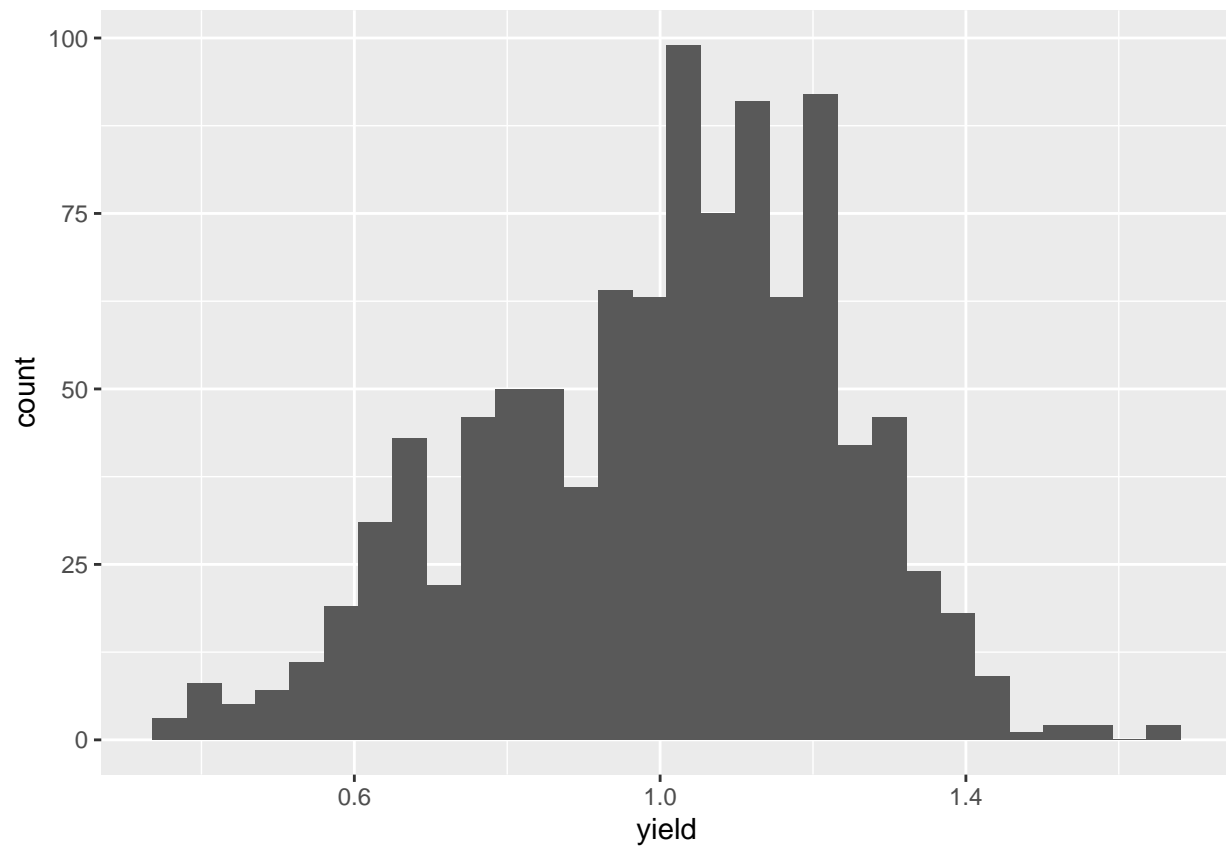
```
ggplot(data=cotton)
```

Running that first like alone will give us a blank plot. This is because we need to specify one or more aesthetics for each plot. An **aesthetic** is any property of the plot that relates to a variable, that is, a column in our data frame. In this case, we use aes(x=yield) to tell R that our aesthetic is positioned horizontally according to the value of yield.

Notice the *x-axis* now contains the range of yield values in our dataset. But there are no bars yet. We need to tell R what shape to use to represent the values of yield.

In the second line, we use `geom_histogram` to tell R we are drawing a histogram. Knowing this, R will automatically assign bins and count the number of observations in each bin.

```
ggplot(data=cotton, aes(x=yield)) +
  geom_histogram()
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

By default, R chooses to create 30 bins. We can easily specify a different number of bins adding the `bins=15` argument to the `geom_histogram()` function.

```
ggplot(data=cotton, aes(x=yield)) +
  geom_histogram(bins=15)
```

Alternatively, we can set a particular binwidth using the `binwidth =0.05` argument.

```
ggplot(data=cotton, aes(x=yield)) +
  geom_histogram(binwidth=0.05)
```

If we want to make our plot a little less bland, we can tell ggplot to use a different color to fill the bars by adding the `fill = darkolivegreen` argument.

```
ggplot(data=cotton, aes(x=yield)) +
  geom_histogram(binwidth=0.05, fill="darkolivegreen")
```

Finally, we can outline the bars using the `color ="black"` argument.

```
ggplot(data=cotton, aes(x=yield)) +
  geom_histogram(binwidth=0.05, fill="darkolivegreen", color="black")
```

We will introduce more ways to fine tune our plots as the course goes on. Here is a table with named colors you can use for fill or lines. (Note: If you right-click on it and select *Open in new tab*, you can enlarge it for much better legibility

# R Colors By Name

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| white | coral4 | deepskyblue | gray28 | gray88 | grey40 | grey100 | lightpink2 | mistyrose2 | plum | slategray2 |
| aliceblue | cornflowerblue | deepskyblue1 | gray29 | gray89 | grey41 | honeydew | lightpink3 | mistyrose3 | plum1 | slategray3 |
| antiquewhite | cornsilk | deepskyblue2 | gray30 | gray90 | grey42 | honeydew1 | lightpink4 | mistyrose4 | plum2 | slategray4 |
| antiquewhite1 | cornsilk1 | deepskyblue3 | gray31 | gray91 | grey43 | honeydew2 | lightsalmon | moccasin | plum3 | slategrey |
| antiquewhite2 | cornsilk2 | deepskyblue4 | gray32 | gray92 | grey44 | honeydew3 | lightsalmon1 | navajowhite | plum4 | snow |
| antiquewhite3 | cornsilk3 | dimgray | gray33 | gray93 | grey45 | honeydew4 | lightsalmon2 | navajowhite1 | powderblue | snow1 |
| antiquewhite4 | cornsilk4 | dimgrey | gray34 | gray94 | grey46 | hotpink | lightsalmon3 | navajowhite2 | purple | snow2 |
| aquamarine | cyan | dodgerblue | gray35 | gray95 | grey47 | hotpink1 | lightsalmon4 | navajowhite3 | purple1 | snow3 |
| aquamarine1 | cyan1 | dodgerblue1 | gray36 | gray96 | grey48 | hotpink2 | lightseagreen | navajowhite4 | purple2 | snow4 |
| aquamarine2 | cyan2 | dodgerblue2 | gray37 | gray97 | grey49 | hotpink3 | lightskyblue | navy | purple3 | springgreen |
| aquamarine3 | cyan3 | dodgerblue3 | gray38 | gray98 | grey50 | hotpink4 | lightskyblue1 | navyblue | purple4 | springgreen1 |
| aquamarine4 | cyan4 | dodgerblue4 | gray39 | gray99 | grey51 | indianred | lightskyblue2 | oldlace | red | springgreen2 |
| azure | darkblue | firebrick | gray40 | gray100 | grey52 | indianred1 | lightskyblue3 | olivedrab | red1 | springgreen3 |
| azure1 | darkcyan | firebrick1 | gray41 | green | grey53 | indianred2 | lightskyblue4 | olivedrab1 | red2 | springgreen4 |
| azure2 | darkgoldenrod | firebrick2 | gray42 | green1 | grey54 | indianred3 | lightslateblue | olivedrab2 | red3 | steelblue |
| azure3 | darkgoldenrod1 | firebrick3 | gray43 | green2 | grey55 | indianred4 | lightslategray | olivedrab3 | red4 | steelblue1 |
| azure4 | darkgoldenrod2 | firebrick4 | gray44 | green3 | grey56 | ivory | lightslategrey | olivedrab4 | rosybrown | steelblue2 |
| beige | darkgoldenrod3 | floralwhite | gray45 | green4 | grey57 | ivory1 | lightsteelblue | orange | rosybrown1 | steelblue3 |
| bisque | darkgoldenrod4 | forestgreen | gray46 | greenyellow | grey58 | ivory2 | lightsteelblue1 | orange1 | rosybrown2 | steelblue4 |
| bisque1 | darkgray | gainsboro | gray47 | grey | grey59 | ivory3 | lightsteelblue2 | orange2 | rosybrown3 | tan |
| bisque2 | darkgreen | ghostwhite | gray48 | grey0 | grey60 | ivory4 | lightsteelblue3 | orange3 | rosybrown4 | tan1 |
| bisque3 | darkgrey | gold | gray49 | grey1 | grey61 | khaki | lightsteelblue4 | orange4 | royalblue | tan2 |
| bisque4 | darkkhaki | gold1 | gray50 | grey2 | grey62 | khaki1 | lightyellow | orangered | royalblue1 | tan3 |
| black | darkmagenta | gold2 | gray51 | grey3 | grey63 | khaki2 | lightyellow1 | orangered1 | royalblue2 | tan4 |
| blanchedalmond | darkolivegreen | gold3 | gray52 | grey4 | grey64 | khaki3 | lightyellow2 | orangered2 | royalblue3 | thistle |
| blue | darkolivegreen1 | gold4 | gray53 | grey5 | grey65 | khaki4 | lightyellow3 | orangered3 | royalblue4 | thistle1 |
| blue1 | darkolivegreen2 | goldenrod | gray54 | grey6 | grey66 | lavender | lightyellow4 | orangered4 | saddlebrown | thistle2 |
| blue2 | darkolivegreen3 | goldenrod1 | gray55 | grey7 | grey67 | lavenderblush | limegreen | orchid | salmon | thistle3 |
| blue3 | darkolivegreen4 | goldenrod2 | gray56 | grey8 | grey68 | lavenderblush1 | linen | orchid1 | salmon1 | thistle4 |
| blue4 | darkorange | goldenrod3 | gray57 | grey9 | grey69 | lavenderblush2 | magenta | orchid2 | salmon2 | tomato |
| blueviolet | darkorange1 | goldenrod4 | gray58 | grey10 | grey70 | lavenderblush3 | magenta1 | orchid3 | salmon3 | tomato1 |
| brown | darkorange2 | gray | gray59 | grey11 | grey71 | lavenderblush4 | magenta2 | orchid4 | salmon4 | tomato2 |
| brown1 | darkorange3 | gray0 | gray60 | grey12 | grey72 | lawngreen | magenta3 | palegoldenrod | sandybrown | tomato3 |
| brown2 | darkorange4 | gray1 | gray61 | grey13 | grey73 | lemonchiffon | magenta4 | palegreen | seagreen | tomato4 |
| brown3 | darkorchid | gray2 | gray62 | grey14 | grey74 | lemonchiffon1 | maroon | palegreen1 | seagreen1 | turquoise |
| brown4 | darkorchid1 | gray3 | gray63 | grey15 | grey75 | lemonchiffon2 | maroon1 | palegreen2 | seagreen2 | turquoise1 |
| burlywood | darkorchid2 | gray4 | gray64 | grey16 | grey76 | lemonchiffon3 | maroon2 | palegreen3 | seagreen3 | turquoise2 |
| burlywood1 | darkorchid3 | gray5 | gray65 | grey17 | grey77 | lemonchiffon4 | maroon3 | palegreen4 | seagreen4 | turquoise3 |
| burlywood2 | darkorchid4 | gray6 | gray66 | grey18 | grey78 | lightblue | maroon4 | paleturquoise | seashell | turquoise4 |
| burlywood3 | darkred | gray7 | gray67 | grey19 | grey79 | lightblue1 | mediumaquamarine | paleturquoise1 | seashell1 | violet |
| burlywood4 | darksalmon | gray8 | gray68 | grey20 | grey80 | lightblue2 | mediumblue | paleturquoise2 | seashell2 | violetred |
| cadetblue | darkseagreen | gray9 | gray69 | grey21 | grey81 | lightblue3 | mediumorchid | paleturquoise3 | seashell3 | violetred1 |
| cadetblue1 | darkseagreen1 | gray10 | gray70 | grey22 | grey82 | lightblue4 | mediumorchid1 | paleturquoise4 | seashell4 | violetred2 |
| cadetblue2 | darkseagreen2 | gray11 | gray71 | grey23 | grey83 | lightcoral | mediumorchid2 | palevioletred | sienna | violetred3 |
| cadetblue3 | darkseagreen3 | gray12 | gray72 | grey24 | grey84 | lightcyan | mediumorchid3 | palevioletred1 | sienna1 | violetred4 |
| cadetblue4 | darkseagreen4 | gray13 | gray73 | grey25 | grey85 | lightcyan1 | mediumorchid4 | palevioletred2 | sienna2 | wheat |
| chartreuse | darkslateblue | gray14 | gray74 | grey26 | grey86 | lightcyan2 | mediumpurple | palevioletred3 | sienna3 | wheat1 |
| chartreuse1 | darkslategray | gray15 | gray75 | grey27 | grey87 | lightcyan3 | mediumpurple1 | palevioletred4 | sienna4 | wheat2 |
| chartreuse2 | darkslategray1 | gray16 | gray76 | grey28 | grey88 | lightcyan4 | mediumpurple2 | papayawhip | skyblue | wheat3 |
| chartreuse3 | darkslategray2 | gray17 | gray77 | grey29 | grey89 | lightgoldenrod | mediumpurple3 | peachpuff | skyblue1 | wheat4 |
| chartreuse4 | darkslategray3 | gray18 | gray78 | grey30 | grey90 | lightgoldenrod1 | mediumpurple4 | peachpuff1 | skyblue2 | whitesmoke |
| chocolate | darkslategray4 | gray19 | gray79 | grey31 | grey91 | lightgoldenrod2 | mediumseagreen | peachpuff2 | skyblue3 | yellow |
| chocolate1 | darkslategrey | gray20 | gray80 | grey32 | grey92 | lightgoldenrod3 | mediumslateblue | peachpuff3 | skyblue4 | yellow1 |
| chocolate2 | darkturquoise | gray21 | gray81 | grey33 | grey93 | lightgoldenrod4 | mediumspringgreen | peachpuff4 | slateblue | yellow2 |
| chocolate3 | darkviolet | gray22 | gray82 | grey34 | grey94 | lightgoldenrodyellow | mediumturquoise | peru | slateblue1 | yellow3 |
| chocolate4 | deeppink | gray23 | gray83 | grey35 | grey95 | lightgray | mediumvioletred | pink | slateblue2 | yellow4 |
| coral | deeppink1 | gray24 | gray84 | grey36 | grey96 | lightgreen | midnightblue | pink1 | slateblue3 | yellowgreen |
| coral1 | deeppink2 | gray25 | gray85 | grey37 | grey97 | lightgrey | mintcream | pink2 | slateblue4 | |
| coral2 | deeppink3 | gray26 | gray86 | grey38 | grey98 | lightpink | mistyrose | pink3 | slategray | |
| coral3 | deeppink4 | gray27 | gray87 | grey39 | grey99 | lightpink1 | mistyrose1 | pink4 | slategray1 | |

## Practice

In the data folder there are three practice files: barley_uniformity.csv, peanut_uniformity.csv, and tomato_uniformity.csv. Practice creating histograms using both the basic and ggplot methods. To start you off, here is the command to load the barley file.

```r
barley = read.csv("data/barley_uniformity.csv")
```

Remember to use CTRL+ALT+I to insert new code chunks below.