

Post02

Ming Chen

December 2, 2017

Make correlation plot by using `corrplot` and `ggplot2`

Introduction and Motivation

We have learned the `cor` function in lecture to calculate the correlation coefficient between variables in data set. However, we will get a messy correlation matrix when the data set contains lots of variables. It is hard to read every numbers to analyze the correlations. Instead, it is helpful if we can visualize the correlation matrix by making a correlation plot.

There are several ways to make correlation plot:

- using `ggplot2` and `reshape2`
- using `corrplot`

I would like to introduce a useful package `corrplot` to help with making correlation plot in this post. Since we learned `ggplot2` in lecture, I will also discuss about how to make correlation plot by using package `ggplot2` and `reshape2`.

1. Data preparation

I choose the `mtcars` dataset which R already provided as an example. You can also try any other dataset by yourself.

```
#inspect the data
dim(mtcars)
```

```
## [1] 32 11
```

```
#get the correlations of the first 6 variables in columns of the dataset
data <- cor(mtcars[1:6])
data
```

```
##           mpg       cyl       disp       hp       drat       wt
## mpg    1.0000000 -0.8521620 -0.8475514 -0.7761684  0.6811719 -0.8676594
## cyl   -0.8521620  1.0000000  0.9020329  0.8324475 -0.6999381  0.7824958
## disp  -0.8475514  0.9020329  1.0000000  0.7909486 -0.7102139  0.8879799
## hp    -0.7761684  0.8324475  0.7909486  1.0000000 -0.4487591  0.6587479
## drat   0.6811719 -0.6999381 -0.7102139 -0.4487591  1.0000000 -0.7124406
## wt    -0.8676594  0.7824958  0.8879799  0.6587479 -0.7124406  1.0000000
```

2. Using `corrplot` package

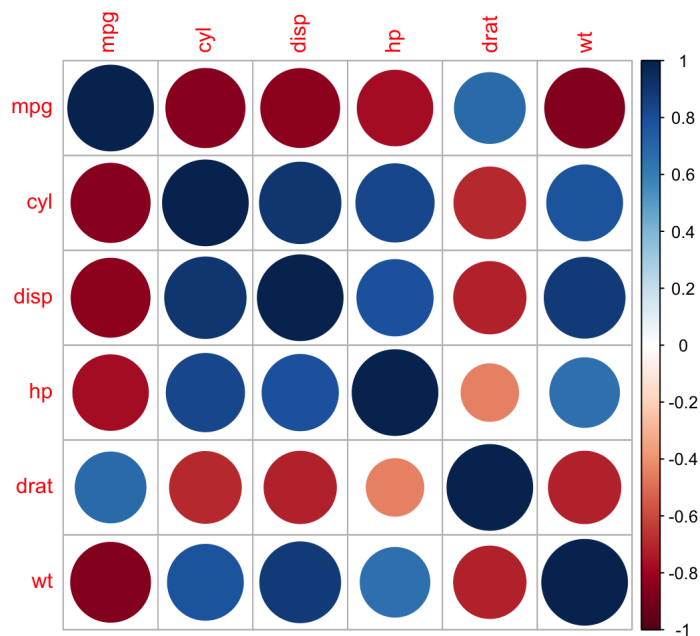
At first, install the package and load it. Then, let's try `corrplot` function to make a simple correlation plot.

```
#install.packages('corrplot')
library(corrplot)
```

```
## Warning: package 'corrplot' was built under R version 3.4.2
```

```
## corrplot 0.84 loaded
```

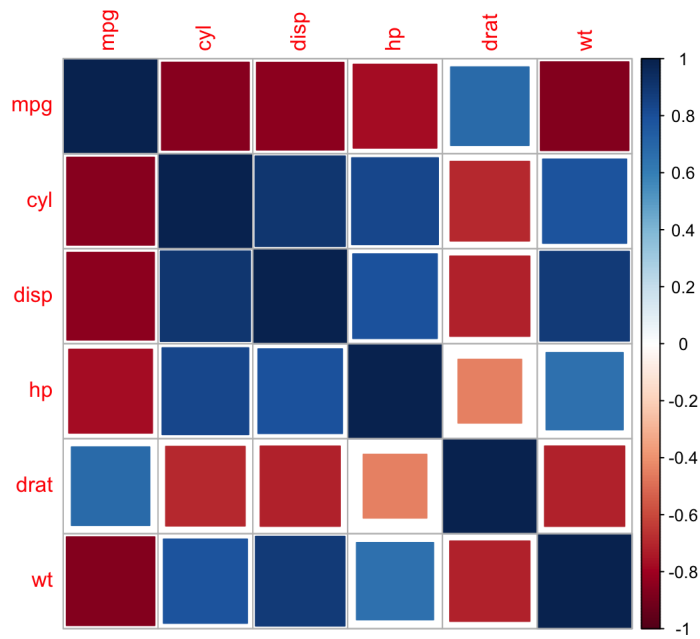
```
corrplot(corr = data)
```



It is easy to make a simple plot. Moreover, we have different options to make the plot more beautiful.

- change argument `method` to plot different shapes or just numbers.

```
corrplot(corr = data, method = 'square')
```

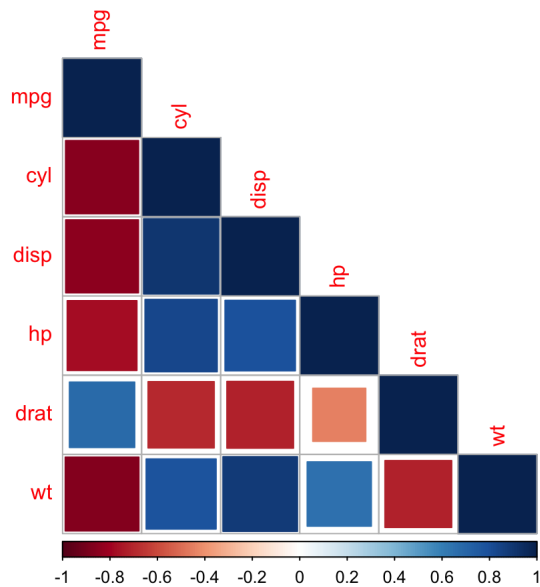


```
corrplot(corr = data, method = 'number')
```

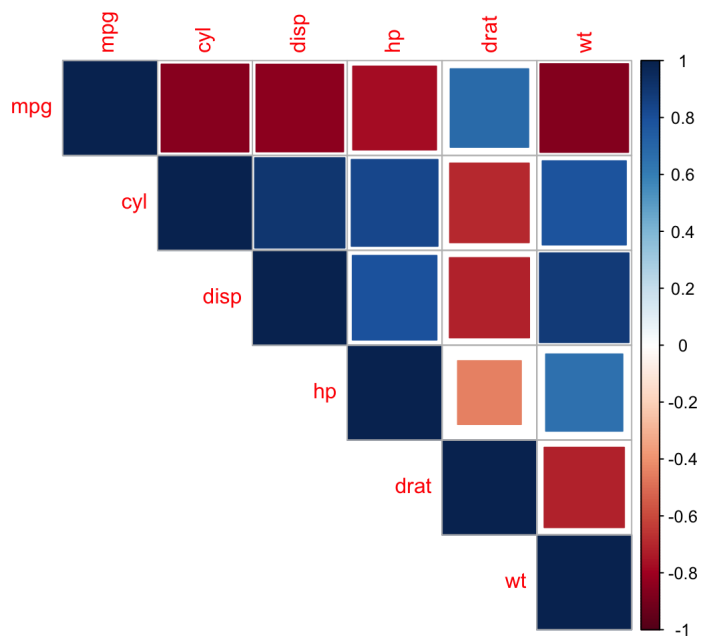


• change argument `type` to represent the plot in different ways: `full`, `lower`, `upper`

```
corrplot(corr = data, method = 'square', type = 'lower')
```

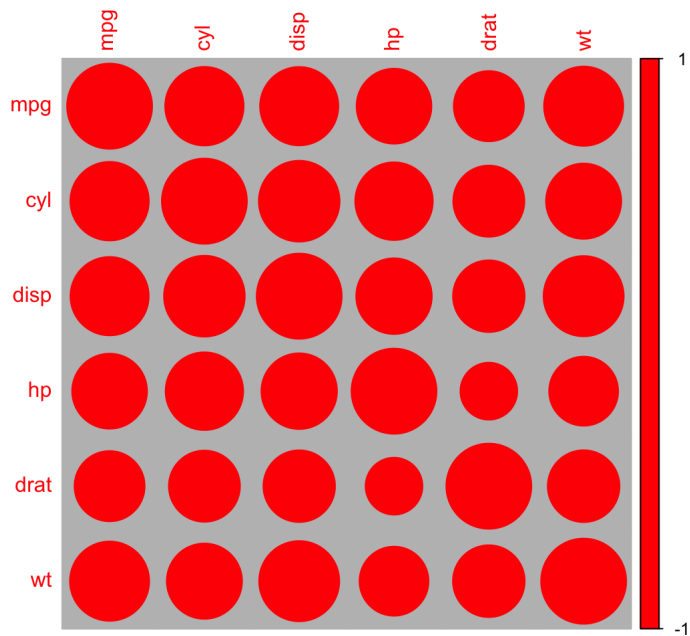


```
corrplot(corr = data, method = 'square', type = 'upper')
```



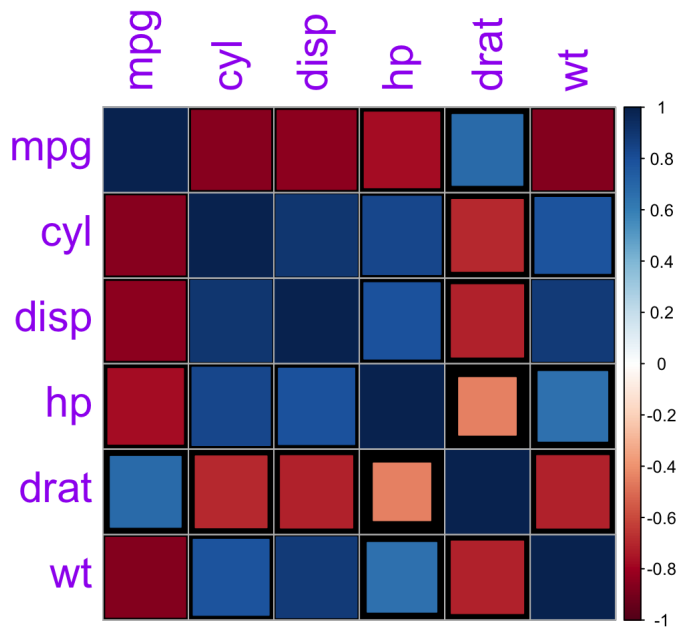
- change argument `col` to decide the color of shapes and using argument `bg` to decide the color of the background of the plot.

```
corrplot(corr = data, method = 'circle', type = 'full', col = 'red', bg = 'grey')
```



- using argument `tl.cex` to change the size of the text label
- using argument `tl.col` to change the color of the text label

```
corrplot(corr = data, method = 'square', type = 'full', bg = 'black',
          tl.cex = 2, tl.col = 'purple')
```



3. Using ggplot2 package

At first, install the package and load it. Here, we need another package `reshape2` to help with making a simple correlation plot. Therefore, I also install package `reshape2` and load it.

```
#install.packages('ggplot2')
#install.packages('reshape2')
library(ggplot2)
library(reshape2)
```

Data reorganization:

```
#use `melt` function to break up the correlation table
melted_data <- melt(data)

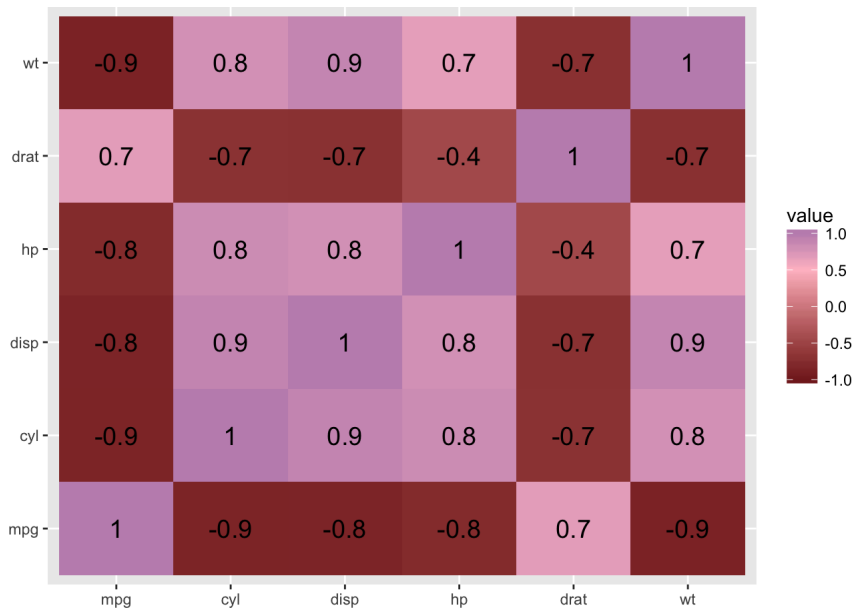
#inspect the melted data
head(melted_data, 10)
```

```
##   Var1 Var2      value
## 1  mpg  mpg  1.0000000
## 2  cyl  mpg -0.8521620
## 3 disp  mpg -0.8475514
## 4   hp  mpg -0.7761684
## 5 drat  mpg  0.6811719
## 6   wt  mpg -0.8676594
## 7  mpg  cyl -0.8521620
## 8  cyl  cyl  1.0000000
## 9 disp  cyl  0.9020329
## 10  hp  cyl  0.8324475
```

For ggplot2, we have to use `melt` function to break up the correlation table into a long format table. So we get Var1 as x-axis, Var2 as y-axis and value to plot.

Let's make a full correlation plot by using `ggplot` function:

```
ggplot(data = melted_data, aes(x = Var1, y = Var2, fill = value)) +
  geom_tile() +
  scale_fill_gradient2(midpoint = 0.5, mid = 'pink', limits = c(-1, 1)) +
  labs(x = '', y = '') +
  geom_text(aes(x = Var1, y = Var2, label = round(value, 1)), size = 5)
```



Some functions mentioned when I make correlation plot by using `ggplot2` :

- `scale_fill_gradient2` function is used with the argument `limit = c(-1,1)` as correlation coefficients range from -1 to 1.
- `geom_tile` function is used to tile Plane With Rectangles
- `geom_text` function is used to adds text directly to the plot.
- `labs` function is used to change axis labels and legend titles.

You can find more functions here:

Data Visualization with ggplot2 : : CHEAT SHEET



Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same components: a **data set**, a **coordinate system**, and **geoms**—visual marks that represent data points.

To display values, map variables in the data to visual properties of the geom (**aesthetics**) like **size**, **color**, and **x** and **y** locations.

Complete the template below to build a graph.

```
ggplot(data = <DATA>) + <GEOM_FUNCTION>
(mapping = aes(<MAPPINGS>),
stat = <STAT>, position = <POSITION>) +
<COORDINATE_FUNCTION> +
<FACET_FUNCTION> +
<SCALE_FUNCTION> +
<THEME_FUNCTION>
```

ggplot(data = mpg, aes(x = cty, y = hwy)) Begins a plot that you finish by adding layers to. Add one geom function per layer.

geom (geom = "point") Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

last_plot() Returns the last plot

ggsave("plot.png", width = 5, height = 5) Saves last plot as 5" x 5" file named "plot.png" in working directory. Matches file type to file extension.

Geoms

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

GRAPHICAL PRIMITIVES

```
a <- ggplot(economics, aes(date, unemployment))
b <- ggplot(seals, aes(x = long, y = lat))
```

a + geom_blank() (Useful for expanding limits)

b + geom_curve(aes(yend = lat + 1, xend = long + 1, curvature = z)) - x, yend, y, end, alpha, angle, color, curvature, linetype, size

a + geom_path(linetype = "butt", linejoin = "round", linewidth = 1) x, y, alpha, color, group, linetype, size

a + geom_polygon(aes(group = group)) x, y, alpha, color, fill, group, linetype, size

b + geom_rect(aes(xmin = long, ymin = lat, xmax = long + 1, ymax = lat + 1)) - xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size

a + geom_ribbon(aes(ymin = unemployment - 900, ymax = unemployment + 900)) - x, ymax, ymin, alpha, color, fill, group, linetype, size

LINE SEGMENTS

```
b + geom_abline(aes(intercept = 0, slope = 1))
b + geom_hline(aes(yintercept = lat))
b + geom_vline(aes(xintercept = long))
```

b + geom_segment(aes(yend = lat + 1, xend = long + 1))

b + geom_spoke(aes(angle = 1:1155, radius = 1))

ONE VARIABLE continuous

```
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)
```

c + geom_area(stat = "bin") x, y, alpha, color, fill, linetype, size

c + geom_density(kernel = "gaussian") x, y, alpha, color, fill, group, linetype, size, weight

c + geom_dotplot() x, y, alpha, color, fill

c + geom_freqpoly() x, y, alpha, color, group, linetype, size

c + geom_histogram(binwidth = 5) x, y, alpha, color, fill, linetype, size, weight

c2 + geom_qq(aes(sample = hwy)) x, y, alpha, color, fill, linetype, size, weight

discrete

```
d <- ggplot(mpg, aes(fill))
d + geom_bar()
```

x, alpha, color, fill, linetype, size, weight

TWO VARIABLES

continuous x, continuous y

```
e <- ggplot(mpg, aes(cty, hwy))
```

e + geom_label(aes(label = cty), nudge_x = 1, nudge_y = 1, check_overlap = TRUE) x, y, label, alpha, angle, color, family, fontface, hjust, linewidth, size, vjust

e + geom_jitter(height = 2, width = 2) x, y, alpha, color, fill, shape, size

e + geom_point() x, y, alpha, color, fill, shape, size, stroke

e + geom_quantile() x, y, alpha, color, group, linetype, size, weight

e + geom_rug(sides = "bl") x, y, alpha, color, linetype, size

e + geom_smooth(method = lm) x, y, alpha, color, fill, group, linetype, size, weight

e + geom_smooth(method = lm, nudge_x = 1, nudge_y = 1, check_overlap = TRUE) x, y, label, alpha, angle, color, family, fontface, hjust, linewidth, size, vjust

discrete x, continuous y

```
f <- ggplot(mpg, aes(cty, hwy))
```

f + geom_col() x, y, alpha, color, fill, group, linetype, size

f + geom_boxplot() x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight

f + geom_dotplot(binaxis = "y", stackdir = "center") x, y, alpha, color, fill, group

f + geom_violin(scale = "area") x, y, alpha, color, fill, group, linetype, size, weight

continuous bivariate distribution

```
h <- ggplot(diamonds, aes(carat, price))
```

h + geom_bin2d(binwidth = c(0.25, 500)) x, y, alpha, color, fill, linetype, size, weight

h + geom_density2d() x, y, alpha, color, group, linetype, size

h + geom_hex() x, y, alpha, color, fill, size

continuous function

```
i <- ggplot(economics, aes(date, unemployment))
```

i + geom_area() x, y, alpha, color, fill, linetype, size

i + geom_line() x, y, alpha, color, group, linetype, size

i + geom_step(direction = "hv") x, y, alpha, color, group, linetype, size

visualizing error

```
df <- data.frame(murder = c("A", "B"), fit = 4.5, se = 1.2)
j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

j + geom_crossbar(latten = 2) x, y, ymax, ymin, alpha, color, fill, group, linetype, size

j + geom_errorbar() x, ymax, ymin, alpha, color, group, linetype, size, width (also **geom_errorbarh()**)

j + geom_linerange() x, ymin, ymax, alpha, color, group, linetype, size

j + geom_pointrange() x, y, ymin, ymax, alpha, color, fill, group, linetype, shape, size

visualizing error

```
data <- data.frame(murder = USArrests$Murder, state = tolower(rownames(USArrests)))
map <- map_data("state")
k <- ggplot(data, aes(fill = murder))
```

k + geom_map(aes(map_id = state), map = map)

+ expand_limits(x = map\$long, y = map\$lat)

map_id, alpha, color, fill, linetype, size

THREE VARIABLES

```
seals2 <- with(seals, sqrt(delta_long^2 + delta_lat^2))
l <- ggplot(seals, aes(long, lat))
```

l + geom_contour(aes(z = z)) x, y, z, alpha, color, group, linetype, size, weight

l + geom_raster(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE) x, y, alpha, fill

l + geom_tile(aes(fill = z), x, y, alpha, color, fill, linetype, size, width



RStudio® is a trademark of RStudio, Inc. • CC BY RStudio • info@rstudio.com • 844-448-1212 • rstudio.com • Learn more with browseVignettes(package = "dplyr", "tibble") • dplyr 0.5.0 • tibble 1.2.0 • Updated: 2017-01

4. Conclusion

From the class, we learned to use ggplot2 to make histogram, scatterplot, barchart an so on. Visualization is an improtant tool to analyze data. Making correlation plot is also ver helpful when you deal with many correlation coefficients. It is quite easy to plot it by using `corrplot` function. You can also use `ggplot2`, which is a little bit more complicated. `ggplot2` is also fun and have lots of functions to explore.

Reference

An Introduction to `corrplot` Package

Plot variable Correlations

Rdocument

ggcorrplot

ggcorr

Plot correlation matrix with R in specific data range

Seven Easy Graphs to Visualize Correlation Matrices in R

Plot a correlation matrix with ggplot2

Correlations