

# Data Visualization with ggplot2 (Part 3)

Plots for Specific Data Types (Chapter 2)

*Seun Odeyemi*

2019-04-15

## Contents

Load Libraries . . . . .	1
The Graphics of Large Data . . . . .	1
Many Observations . . . . .	2
Many variables . . . . .	9
Create a Correlation Matrix in ggplot2 . . . . .	13
Ternary Plots . . . . .	15
Proportional / Stacked Bar Plots . . . . .	20
Networks Plots . . . . .	24
Build the Network (1) . . . . .	24
Build the Network (2) . . . . .	25
Adjusting the Network . . . . .	26
Diagnostic Plots . . . . .	28

## Load Libraries

```
library(readr)
library(dplyr)
library(ggplot2)
# library(ggplot2movies)
library(tidyr)
library(skimr)
library(knitr)
library(kableExtra)
library(RColorBrewer)
library(grid)
library(ggthemes)
library(forcats)
library(GGally)
library(here)
library(hexbin)
```

## The Graphics of Large Data

In this chapter we will continue our discussion of specialty plot types turning our attention to those suited for specific data types. My goal here is to familiarize you with unique plot types so that you may call upon them when you have suitable data even when that may not be very often. Remember to think of the purpose of a plot: *the more plot types you have in your data viz bag of tricks, the creative and fitting your data visualizations will be.*

To lead us into this topic we review and round out knowledge of working with large data sets. So far we've dealt with relatively small data sets, but once we start working on large data sets we'll run into a number of issues. How we define a large dataset will dictate the problems we'll encounter for visualizations. For example, large can refer to the:

1. Many observations (rows, records, etc.)
  - a. Very high resolution time series
  - b. Large surveys
  - c. Website analytics
2. Many variables (features, characteristics, columns, etc.)
  - Multidimensional data
3. A combination

Each situation demands its own solution. There is an obvious overlap between this topic and data handling and storage for large data sets. However, data munging details are beyond this course. We'll focus on the visualization aspects.

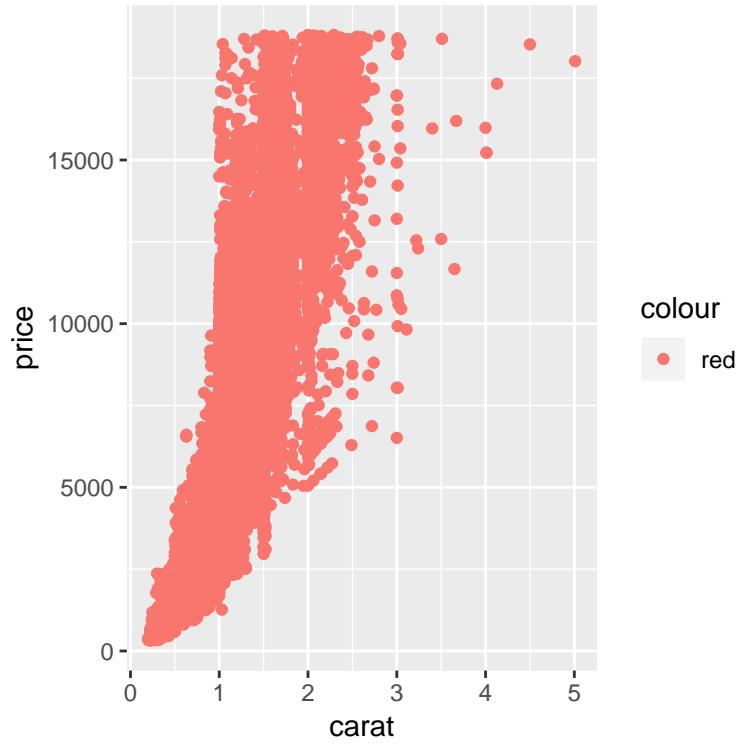
## Many Observations

Let's consider many observations. The most straightforward thing to do is to adjust the `geom` we are using. Using the `diamonds` data set as an example—it contains about 50,000 observations relating to diamonds such as price, carat, color, and clarity. For example in this case of high density points such as the scatter plot shown below, we can adjust things like the plotting symbol, point size, and alpha blending to make trends in our data set visible. In the last chapter, we learned that we can also use a 2D density plot as shown below with contour lines. However, we should realize that this is not a one-stop solution.

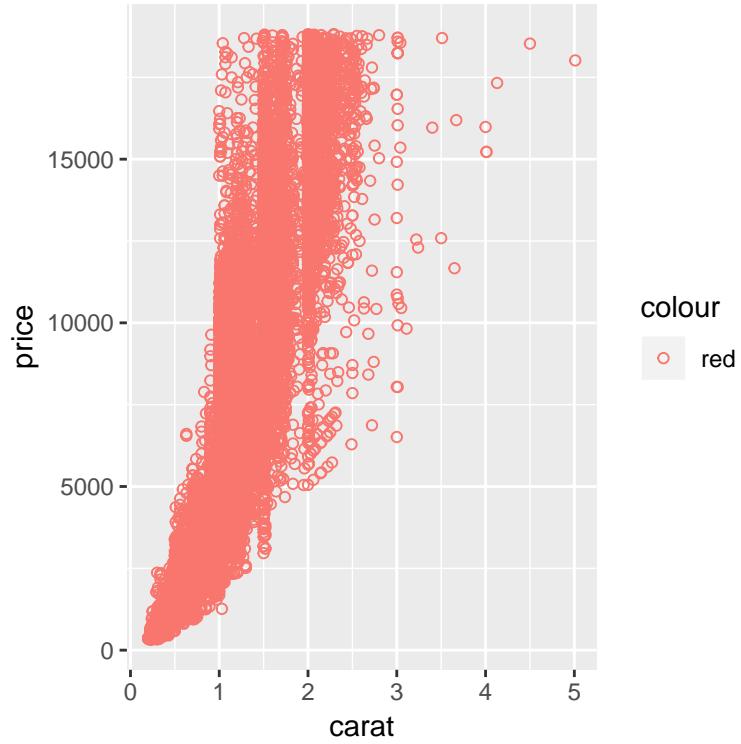
In this case, there is such a high density at the lower end of the scales that we lose too much detail in the rest of the plot. Recall that we can map the actual density, shown by each contour line, to a continuous color scale. However, comparing these plots to the scatter plot with alpha blending shows that neither of 2D density plots really shows an accurate representation of the distribution. An alternative to 2D density plots is to simply bin the values into a grid. This is simply a 2D version of a histogram, which means we can change the `bin` number to increase or decrease the resolution of our plot just like a 1D histogram.

An extension of this, which is popular in infographics is the use of hex binning. This looks just like the binning we've already seen only using hexagons instead of squares. We can adjust the bin size here also. Methods for dealing with many observations are basically concerned with reducing over-plotting or reducing the amount of information that is plotted. This is done by aggregating the data into two-dimensions further removing us from the large amount of raw observations. This is in the hopes of seeing some interesting revelations in the data set. The choice will depend on the data set at hand. So, it's worth experimenting with a wide variety of geoms.

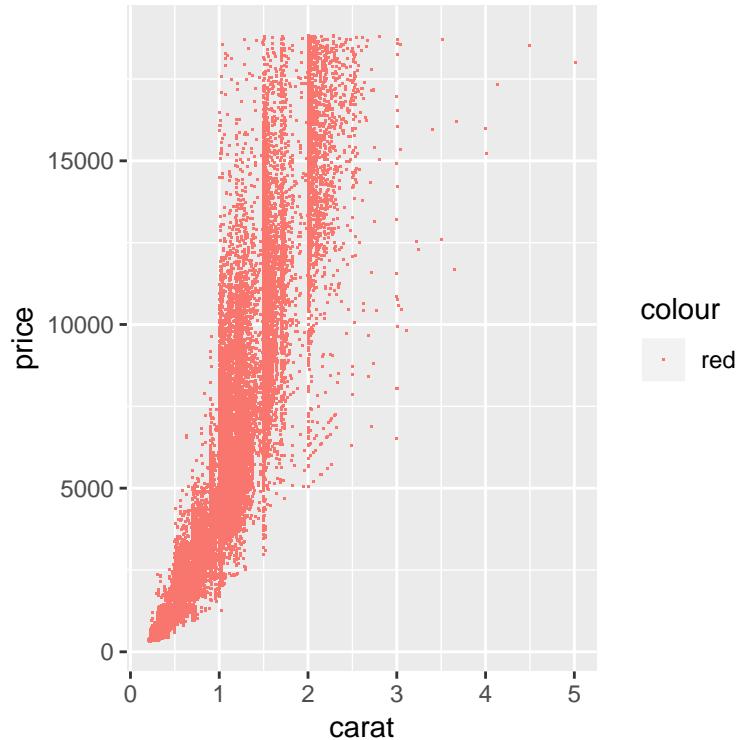
```
ggplot(diamonds, aes(x = carat, y = price, color = "red")) +
  geom_point() #simplified
```



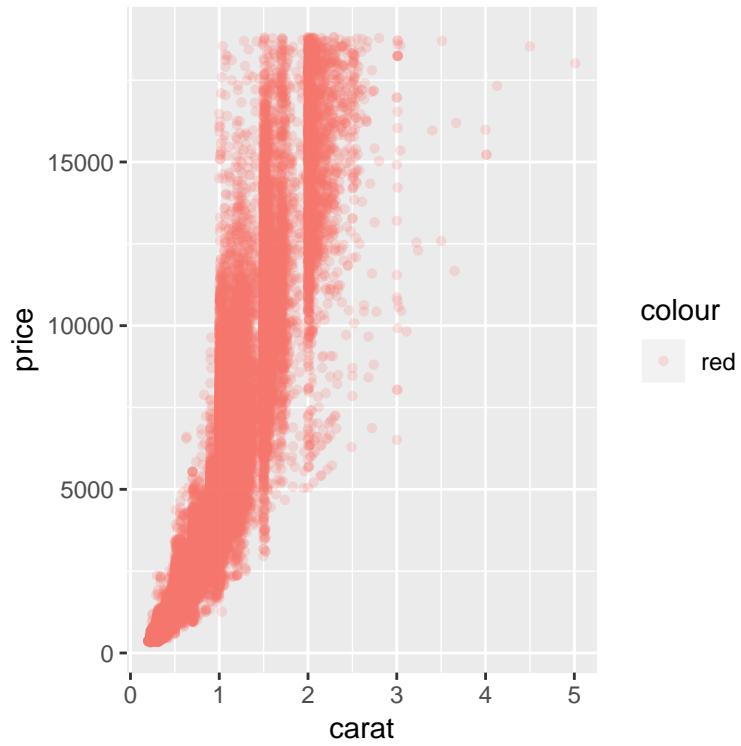
```
ggplot(diamonds, aes(x = carat, y = price, color = "red")) +  
  geom_point(shape = 1) #plotting symbol adjusted using shape in the geom layer
```



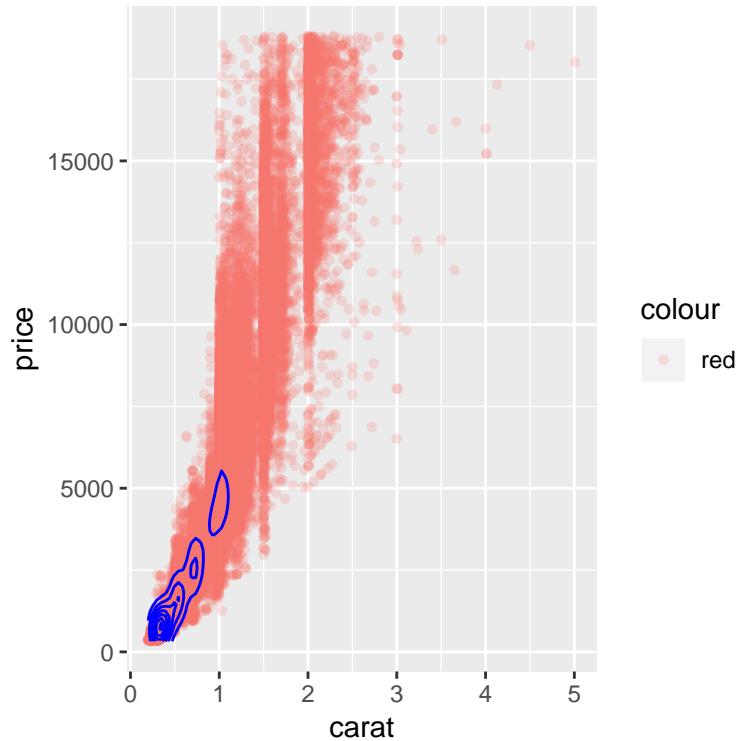
```
ggplot(diamonds, aes(x = carat, y = price, color = "red")) +  
  geom_point(shape = ".") #pointsize adjusted
```



```
ggplot(diamonds, aes(x = carat, y = price, color = "red")) +  
  geom_point(shape = 16, alpha = 0.2) #alpha blending
```



```
ggplot(diamonds, aes(x = carat, y = price, color = "red")) +  
  geom_point(shape = 16, alpha = 0.2) + #alpha blending  
  stat_density_2d(color = "blue")
```

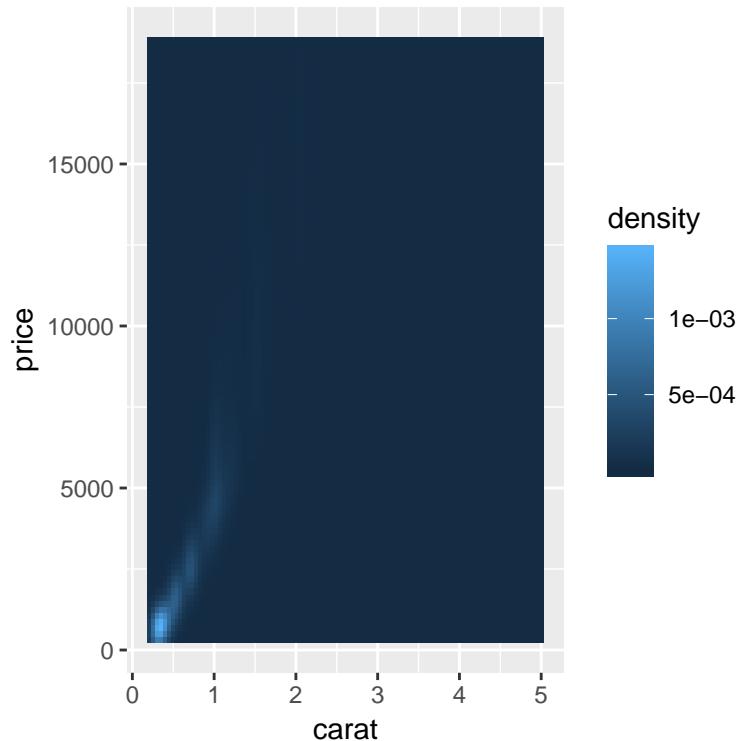


```

# diamonds <- diamonds %>% mutate(category = cut(carat, breaks = 3, labels = c("low", "medium", "high")))
#
# ggplot(diamonds, aes(x = carat, y = price, color = "red")) +
#   geom_point(shape = 16, alpha = 0.2) + #alpha blending
#   stat_density2d(aes(fill = ..density..))

ggplot(diamonds, aes(x = carat, y = price)) +
  stat_density2d(geom = "tile",
                 aes(fill = ..density..),
                 contour = FALSE)

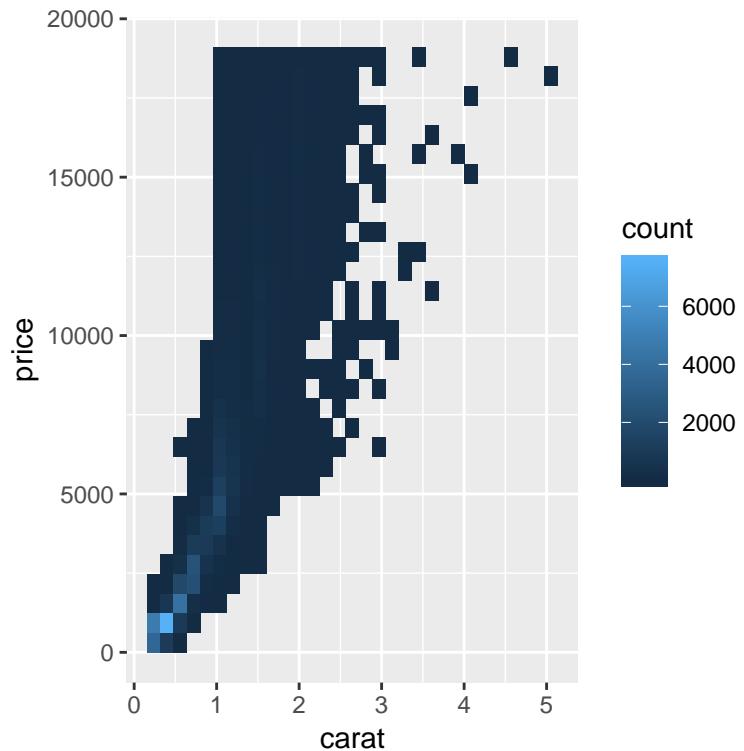
```



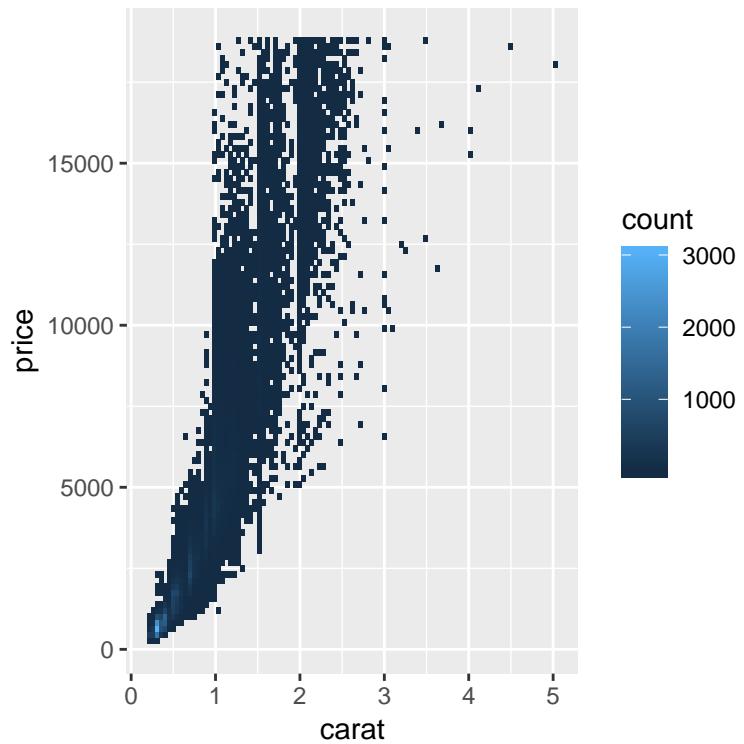
```

# 2D histogram
ggplot(diamonds, aes(x = carat, y = price)) +
  geom_bin2d()

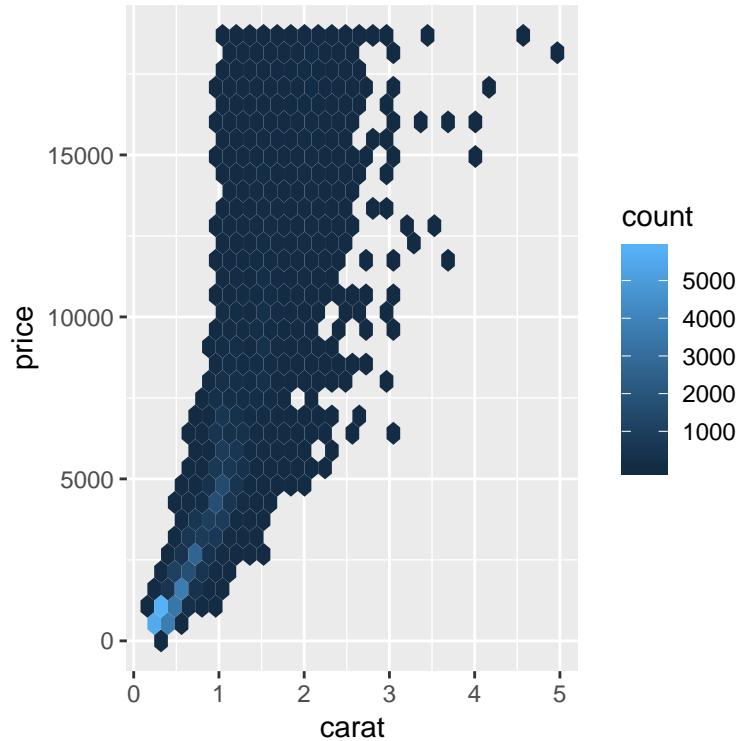
```



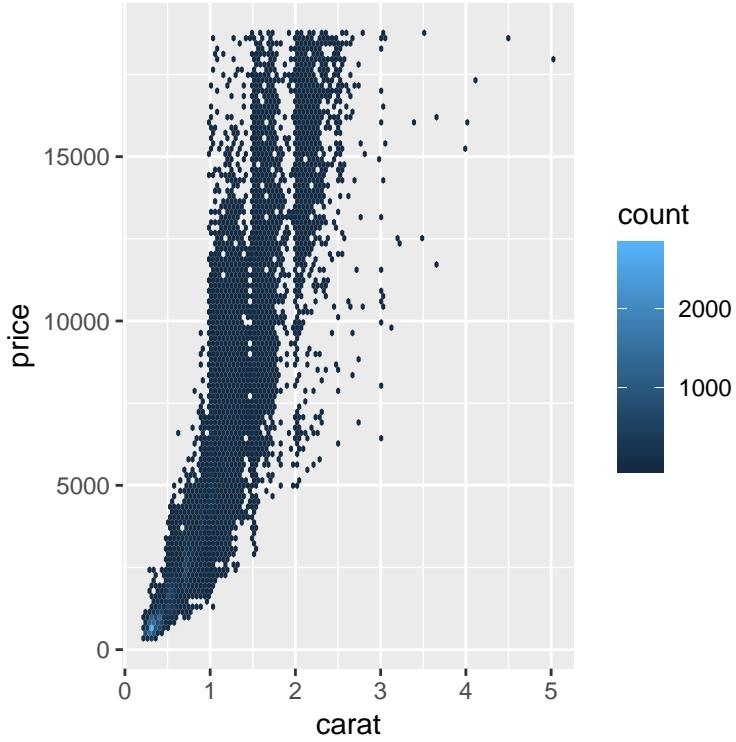
```
ggplot(diamonds, aes(x = carat, y = price)) +  
  geom_bin2d(bins = 100)
```



```
# hex binning
ggplot(diamonds, aes(x = carat, y = price)) +
  geom_hex()
```



```
ggplot(diamonds, aes(x = carat, y = price)) +
  geom_hex(bins = 100)
```

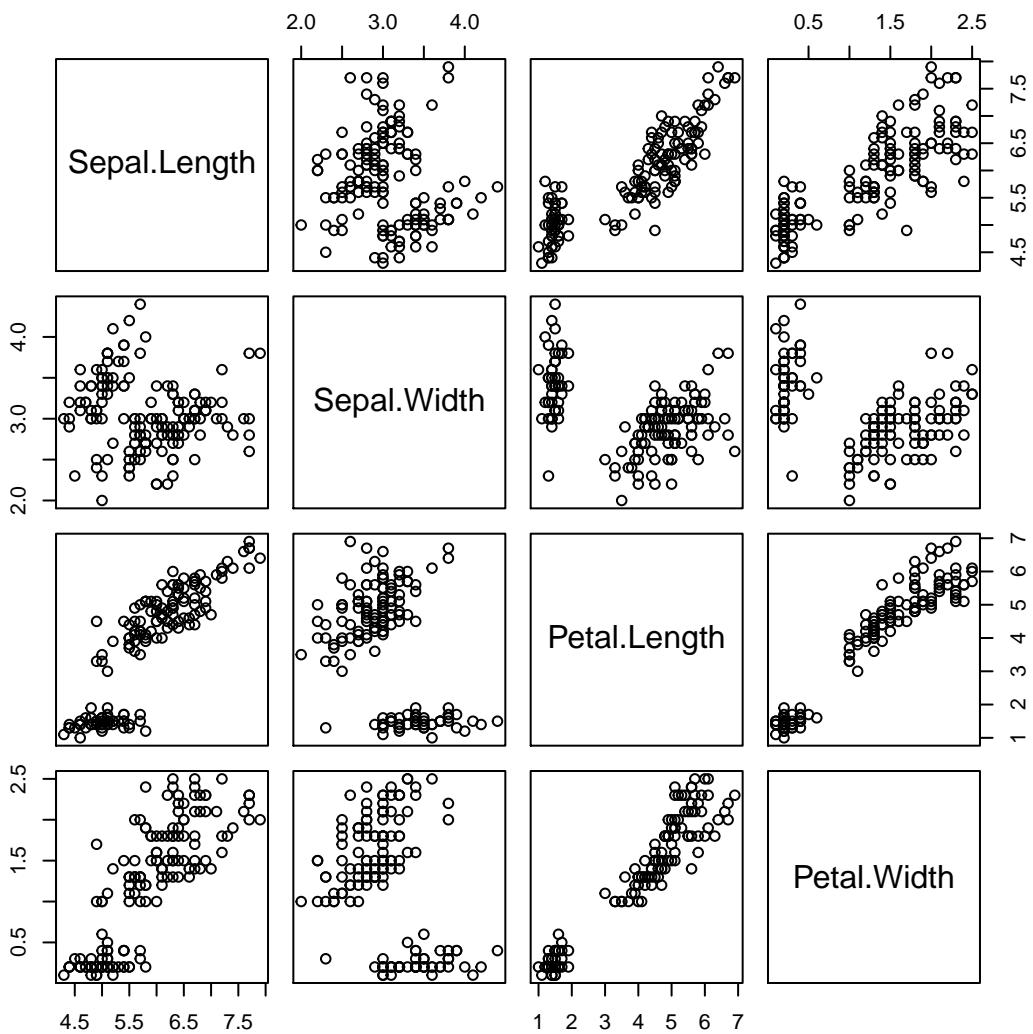


## Many variables

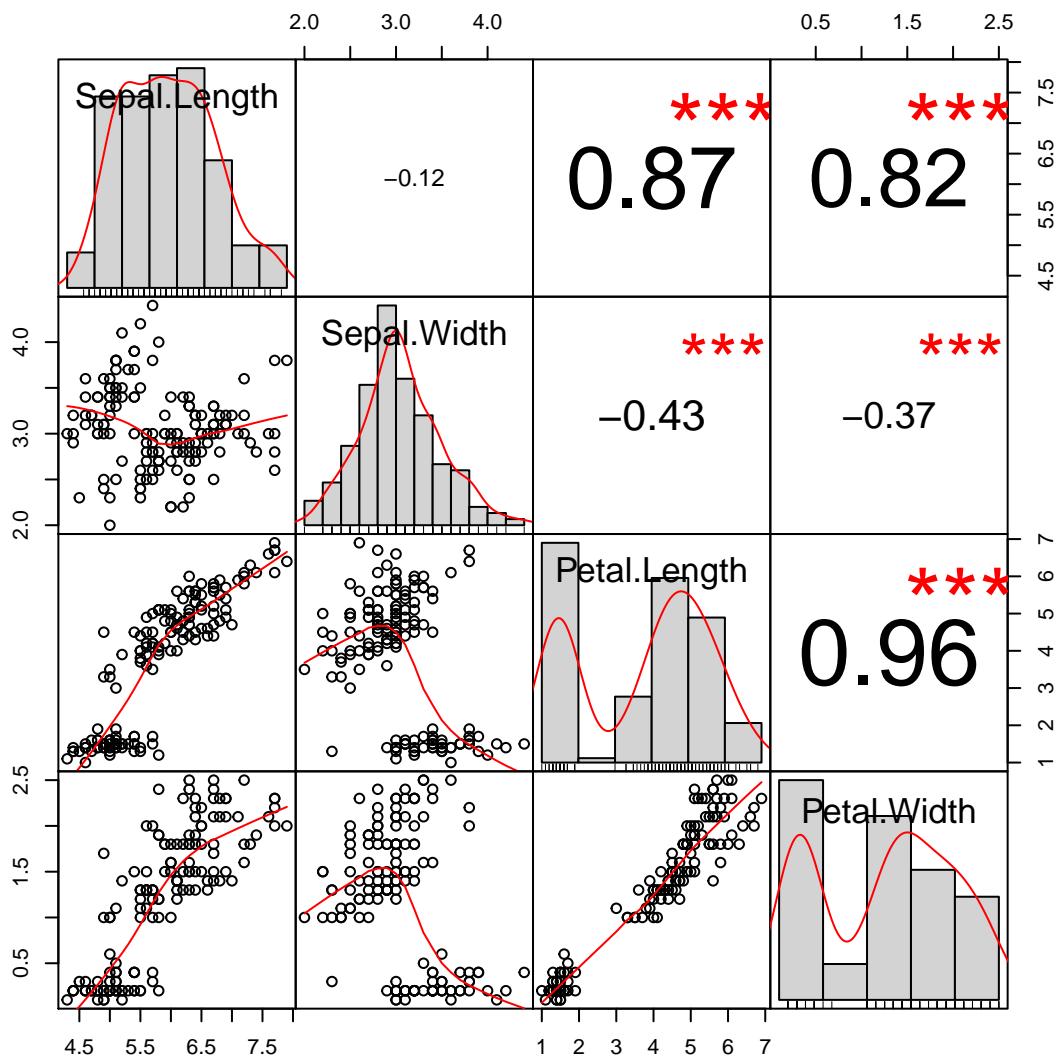
The second type of large data is when we have many variables referred to as multi-variate or high-dimensional data. Before, we begin plotting there are many steps we can take in dealing with the understanding of relationships between the variables like data reduction methods such as principal component analysis (PCA). We do consider data visualization methods for PCA results later on in our section on diagnostic plots. We've already seen examples of how to deal with many dimensions using `facets`, which allow us to treat levels within a factor variable as components in rows or columns of plots. What happens when we have many levels? Facetting then becomes cumbersome and computationally ineffective. We'll see a nice solution to this problem when we deal we discuss animations in the next chapter.

Here, we'll take a look at two special plot types that are particularly useful. The first type of plot is called *SPLOM*, which stands for *Scatter Plot Matrix*. Here, I have made a SPLOM using the four continuous variables in the `iris` dataset using the base package `pairs` function. There are many variations of this theme such as this correlation matrix in the `PerformanceAnalytics` package or this method in the `GGally` package, which uses the `mtcars` data set. This function can handle different variable types not just continuous data. Another popular plot type for dealing with many variables is the parallel coordinate plot, which we encountered in the second course. This example takes the four continuous variables in the `iris` data set and places them on parallel vertical axes. Typically, it will be completely taboo to draw lines between individual values in different nominal variables, but in the case as an exception to the rule, we do want to compare many different variables including categorical and continuous variables together. Variables that are on completely different axes can be lumped into one large visualization. The goal here is to look for trends in how particular variables are related. This should give you an idea of the variety of options we have available to visualize large data sets.

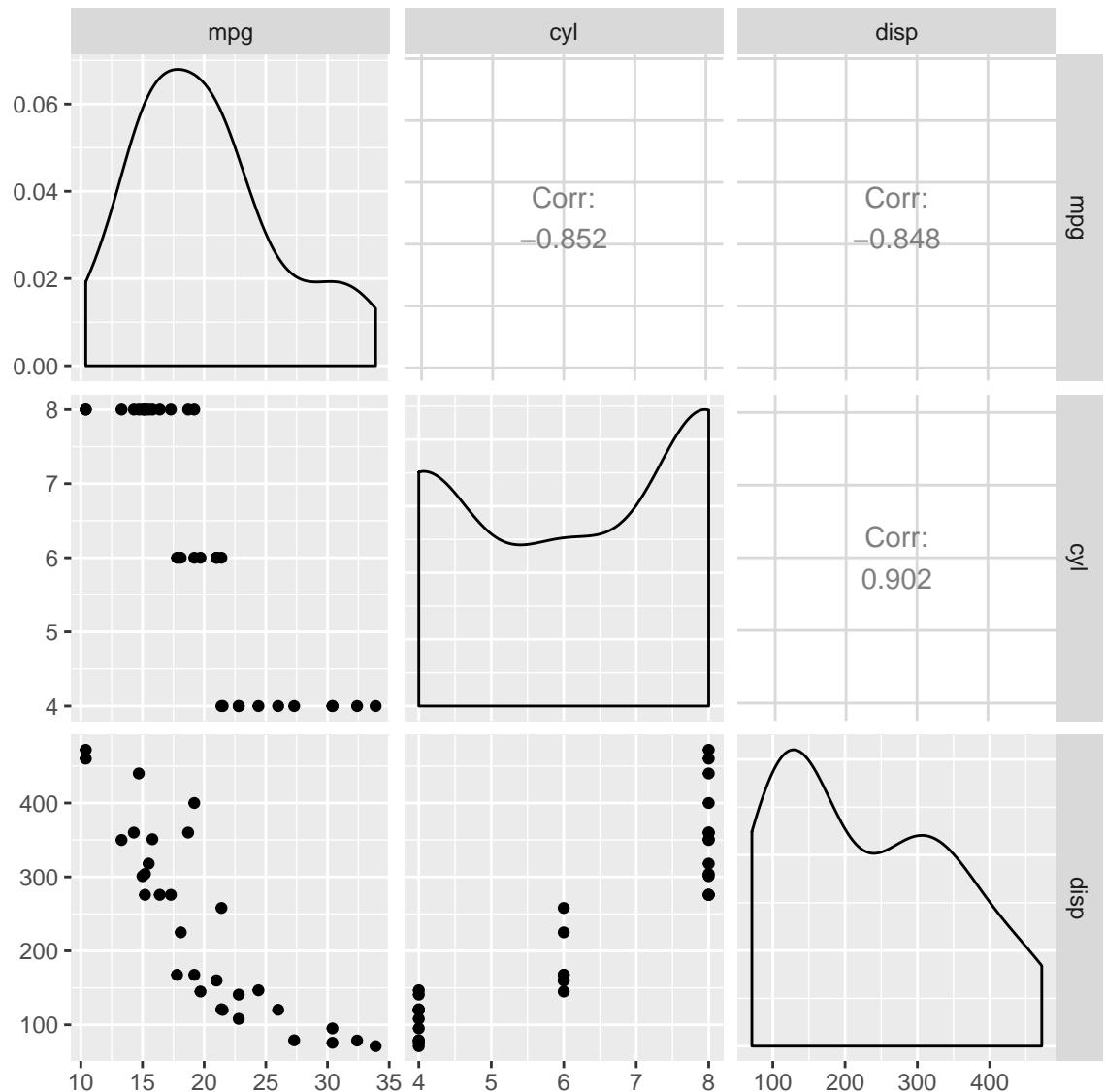
```
pairs(iris[-5])
```



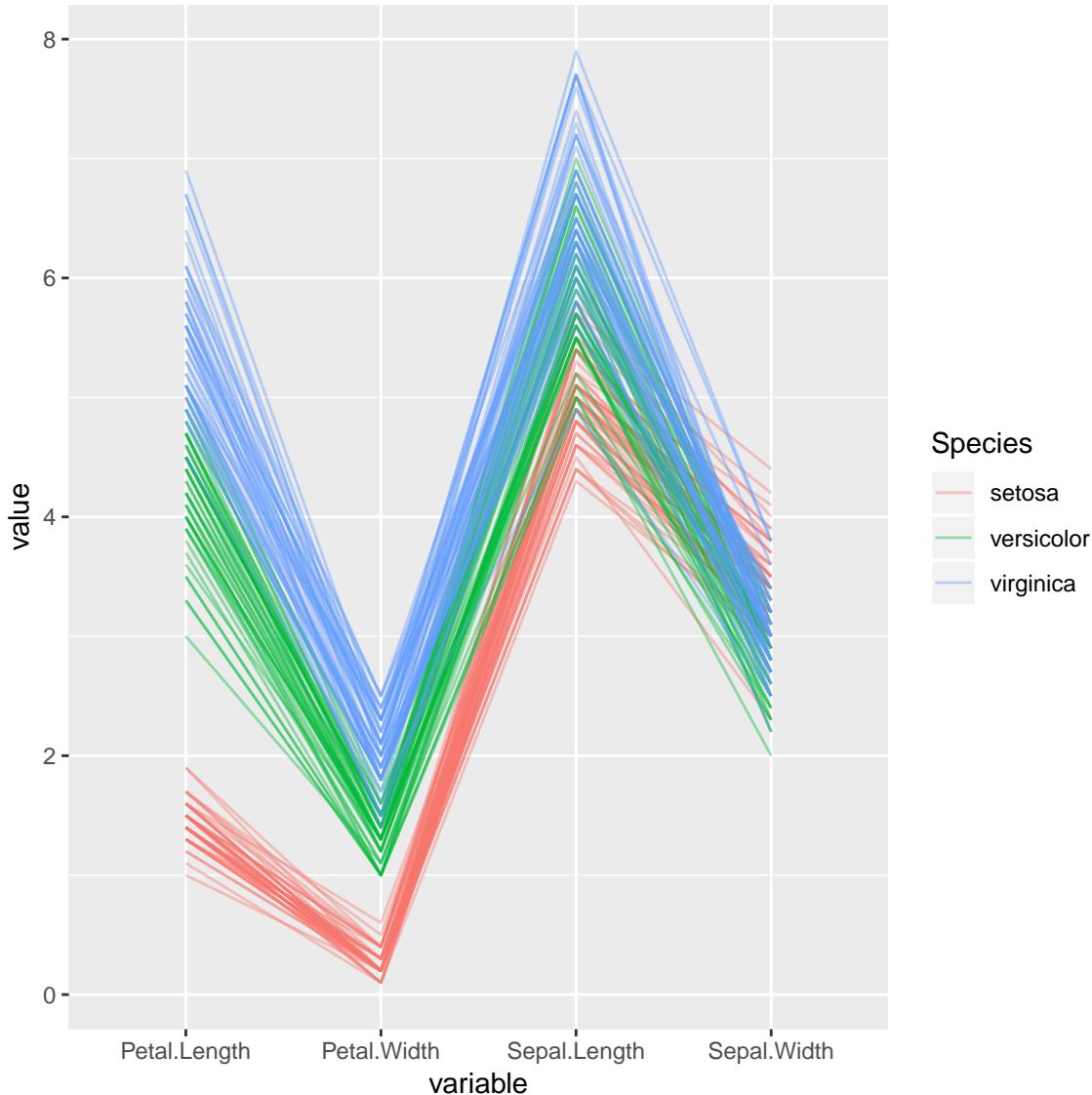
```
library(PerformanceAnalytics)
chart.Correlation(iris[-5],
                  pch = 0:18,
                  cex.labels = 4,
                  labels = c("Sepal Length", "Sepal Width", "Petal Length", "Petal Width"))
```



```
library(GGally)
ggpairs(mtcars[1:3])
```



```
ggparcoord(iris, columns = 1:4,
            groupColumn = 5,
            scale = "globalminmax",
            order = "anyClass", alphaLines = 0.4)
```



SPLOM! These are some great-looking *Scatter PLOT Matrices*.

### Create a Correlation Matrix in ggplot2

Instead of using an off-the-shelf correlation matrix function, you can of course create your own plot. Just for fun, in this exercise, you'll re-create the scatterplot you see on the right. The strength of the correlation is depicted by the size and color of the points and labels.

For starters, a correlation matrix can be calculated using, for example, `cor(dataframe)` (if all variables are numerical). Before you can use your data frame to create your own correlation matrix plot, you'll need to get it in the right format.

In the editor, you can see the definition of `cor_list()`, a function that re-formats the data frame `x`. Here, `L` is used to add the points to the lower triangle of the matrix, and `M` is used to add the numerical values as text to the upper triangle of the matrix. With `reshape2::melt()`, the correlation matrices `L` and `M` are each converted into a three-column data frame: the `x` and `y` axes of the correlation matrix make up the first two columns and the corresponding correlation coefficient makes up the third column. These become the new variables "`points`" and "`labels`", which can be mapped onto the `size` aesthetic for the points in the lower

triangle and onto the `label` aesthetic for the text in the upper triangle, respectively. Their values will be the same, but their positions on the plot will be symmetrical about the diagonal! Merging L and M, you have everything you need.

If you're not familiar with `reshape2` - don't worry, the only reason we use that instead of `tidyR` is that `reshape2::melt()` can handle a matrix, whereas `tidyR::gather()` requires a data frame. At this point you just need to understand how to use the output from `cor_list()`.

You'll first use `dplyr` to execute this function on the continuous variables in the `iris` data frame (the first four columns), but separately for each species. Please refer to the course on `dplyr` if you are not familiar with these functions.

Next, you'll actually plot the resulting data frame with `ggplot2` functions.

```
library(reshape)

cor_list <- function(x) {
  L <- M <- cor(x)

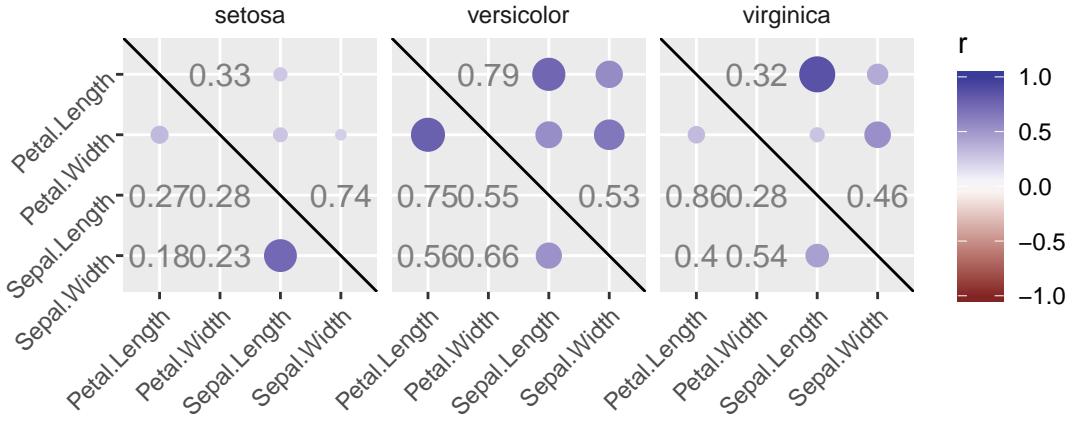
  M[lower.tri(M, diag = TRUE)] <- NA
  M <- melt(M)
  names(M)[3] <- "points"

  L[upper.tri(L, diag = TRUE)] <- NA
  L <- melt(L)
  names(L)[3] <- "labels"

  merge(M, L)
}

# Calculate xx with cor_list
library(dplyr)
xx <- iris %>%
  group_by(Species) %>%
  do(cor_list(.[1:4])) %>%
  select( "Var1" = X1, "Var2" = X2, everything())

# Finish the plot
# We use abs() to get the absolute value here since correlations can be positive or negative,
# but size can only be positive. Don't set
ggplot(xx, aes(x = Var1, y = Var2)) +
  geom_point(aes(col = points, size = abs(points)), shape = 16) +
  geom_text(aes(col = points, size = abs(points), label = round(labels, 2))) +
  scale_size(range = c(0, 6)) +
  scale_color_gradient2("r", limits = c(-1, 1)) +
  scale_y_discrete("") +
  scale_x_discrete("") +
  guides(size = FALSE) +
  geom_abline(slope = -1, intercept = nlevels(xx$Var1) + 1) +
  coord_fixed() +
  facet_grid(. ~ Species) +
  theme(axis.text.y = element_text(angle = 45, hjust = 1),
        axis.text.x = element_text(angle = 45, hjust = 1),
        strip.background = element_blank())
#> Warning: Removed 30 rows containing missing values (geom_point).
#> Warning: Removed 30 rows containing missing values (geom_text).
```



Massive Matrices! Another great custom visualisation!

## Ternary Plots

The first type of specialized plots we'll consider is the **Ternary Plot** also known as **Triangle Plot**. These plots are a ideal way to depict compositional trivariate data. This means we have a parts-of-a-whole problem, where our parameters add up to 100%. A data type that is typically represented with a Ternary plot is soil composition. One way to describe soil is to understand the proportional composition of three measures: Sand, Silt, and Clay content. Let's take a look at the `africa` dataset containing over 40000 soil samples from all over the continent of Africa.

The sum of each row add up to 100. At this point in your dataviz career, you will probably make something like a **proportional stacked bar plot**.

Ternary plots have three axes: one for each each variable of interest rearranged to form a triangle. (Recall that in geometry an *edge* is a surface like a side of a triangle and a *vertex* is the corner point). We typically expect that the edges of a plot are its axes like in a scatterplot. That's kind of what is happening here, but the confusing part is that the edges are not the whole story. The axes traverse the middle of the triangle like this for the x axis. The lower left vertex is 100% and the edge opposite is 0%. The tick marks of the x axis are found on the left side of the triangle. For the y axis, the top vertex is 100% and the base edge is 0%, and the tick marks are shown on the right of the triangle. The z axis works in the same way. The lower right vertex is 100% and the edge opposite is 0%, and the tick marks are shown on the bottom of the triangle. If we consider the interplay between a three axes, it becomes clear that this is a parts-of-a-whole problem. If we find a value for any point in the triangle, they will all add up to 100%. The point directly in the middle is where all the axes have equivalent input so 33.3% of each.

Let's see this in action with some of the soil data. Here the three axes are the Sand, Silt, and Clay content. In the first location, the composition is overwhelmingly made up of Sand. So, the point in our Ternary plot is very close to the vertex of this axis, indicating a high percentage. It must therefore be close to the edges of the remaining axes, which reflects a low abundance there. In the seventh site, Clay is predominant. So, our point is close to the Clay vertex. Of the remaining two axes there is a uneven split. Silt having a high

proportion than Sand. So there the point is closer to the Sand axis edge and further away from the Silt edge. As another example, Silt makes the largest proportion at site 8, which is why our third point leans toward the Silt vertex. The positioning of the remaining axes tells us that there is slightly more Sand and Clay content at this site.

The advantages of this plot type are numerous and akin to a scatter plot. Our complete data set has over 40,000 locations. It will not be helpful to plot them all as a bar chart. But we can keep adding points on a Ternary plot (like the plot below). So, now we can get an overview of the distribution of the dataset and identify any single value within that distribution. Of course, the advantage here is that we can use `geoms` other than points. We can look at the density as a series of contour lines or use a fill to explore a region of high and lower density.

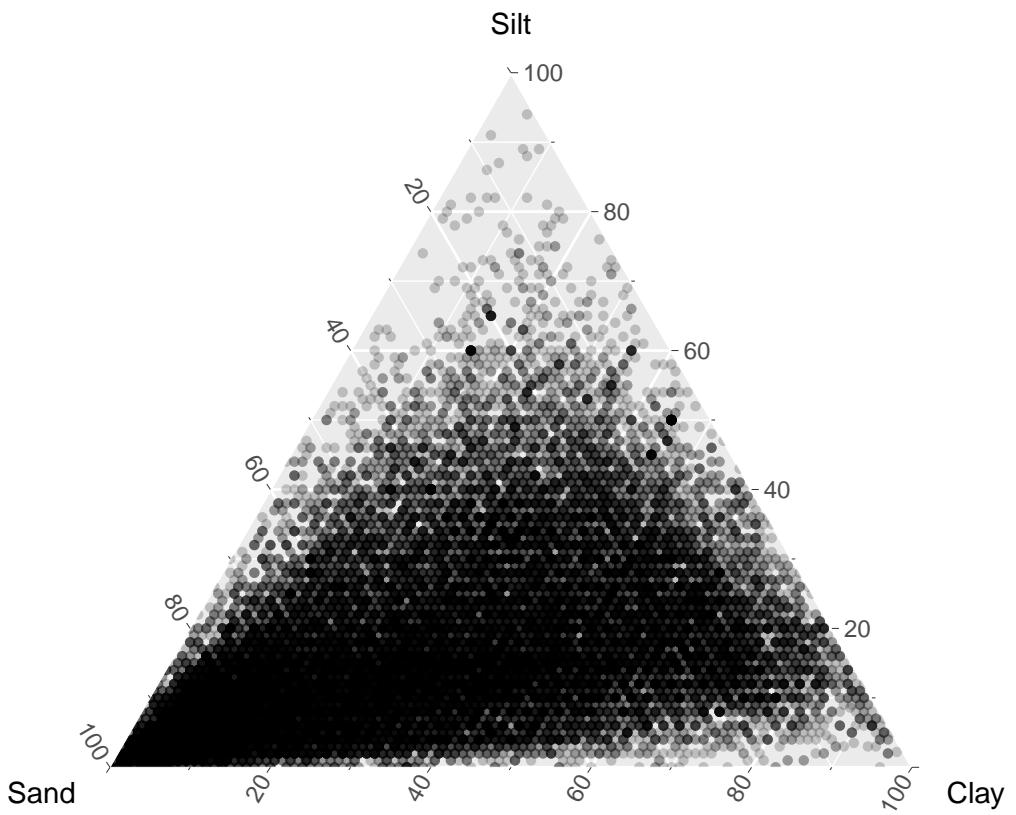
```
load("datasets/africa.RData")
```

```
dim(africa)
#> [1] 40093      3
```

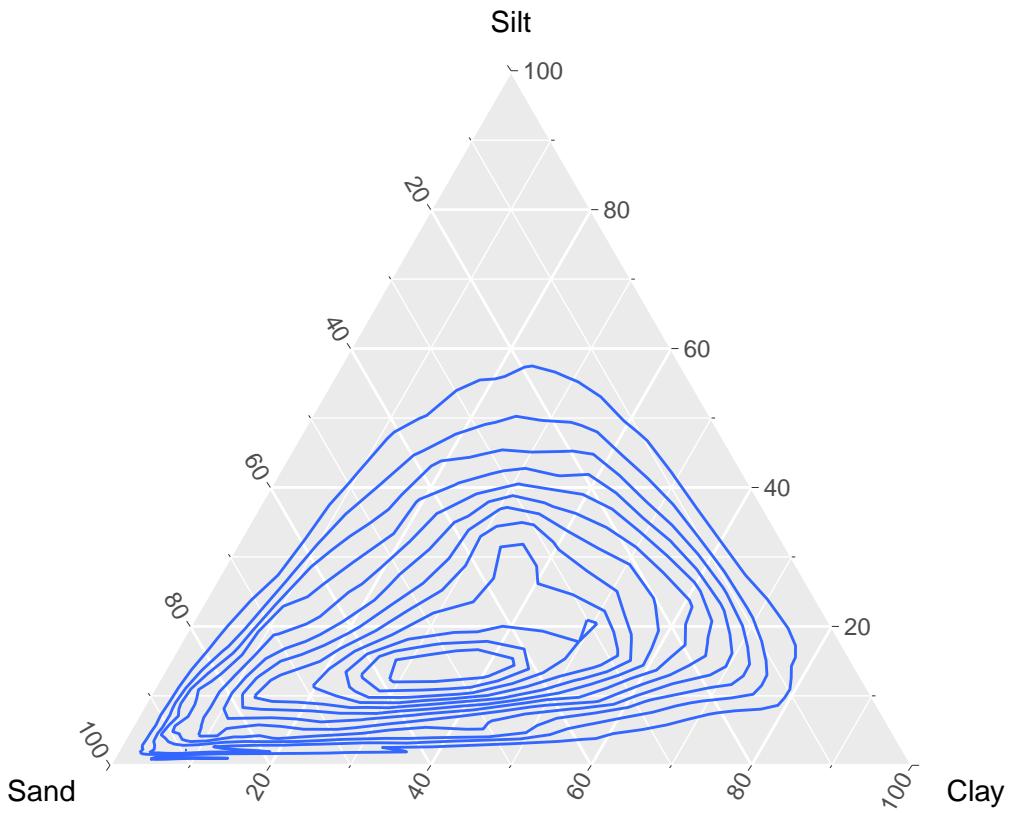
```
head(africa)
```

	Sand	Silt	Clay
18100	24	12	64
18103	36	14	50
18113	56	18	26
18115	52	21	27
18142	65	3	32
18158	43	14	43

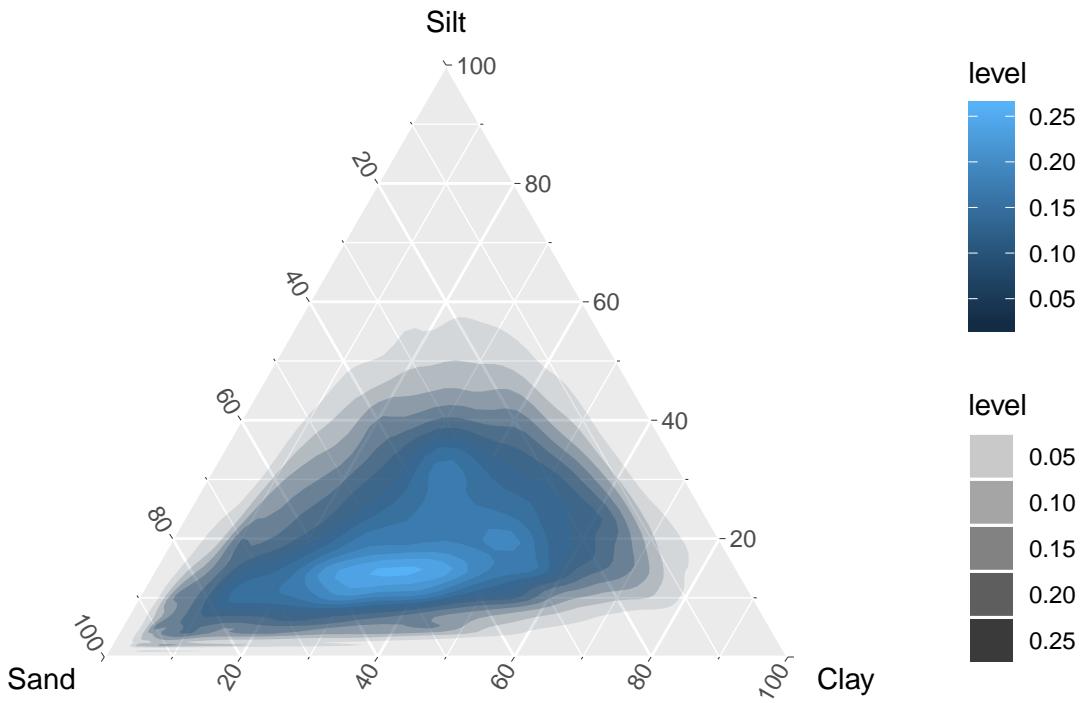
```
library(ggtern)
#> --
#> Remember to cite, run citation(package = 'ggtern') for further info.
#> --
#>
#> Attaching package: 'ggtern'
#> The following objects are masked from 'package:ggplot2':
#>
#>     %+%, aes, annotate, calc_element, ggplot, ggplot_build,
#>     ggplot_gtable, ggplotGrob, ggsave, layer_data, theme,
#>     theme_bw, theme_classic, theme_dark, theme_gray, theme_light,
#>     theme_linedraw, theme_minimal, theme_void
ggtern(africa, aes(Sand, Silt, Clay)) +
  geom_point(alpha = 0.2, shape = 16)
```



```
ggtern(africa, aes(Sand, Silt, Clay)) +
  geom_density_tern()
#> Warning: Removed 420 rows containing non-finite values (StatDensityTern).
```



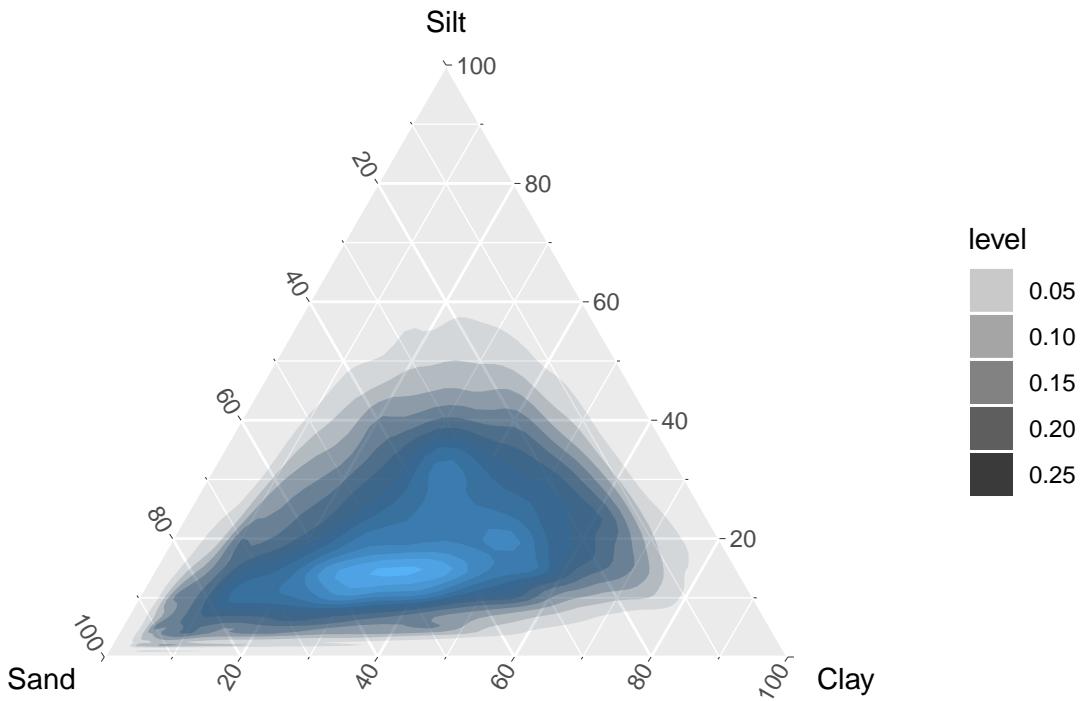
```
ggtern(africa, aes(Sand, Silt, Clay)) +
  stat_density_tern(geom = "polygon",
    n = 200,
    aes(fill = ..level..,
        alpha = ..level..))
#> Warning: Removed 420 rows containing non-finite values (StatDensityTern).
```



```

ggtern(africa, aes(Sand, Silt, Clay)) +
  stat_density_tern(geom = "polygon",
    n = 200,
    aes(fill = ..level..,
        alpha = ..level..)) +
  guides(fill = FALSE)
#> Warning: Removed 420 rows containing non-finite values (StatDensityTern).

```



Terrific ternary! Those are some awesome-looking plots!

### Proportional / Stacked Bar Plots

Before you head over to ternary plots, let's try to make a classical proportional/stacked bar plot of a subset of the data. We'll use a stacked bar plot and the `coord_flip()` function to flip the x and y axes.

The data frame for the African Soil Profiles Database is available in your workspace as `africa` and can be found in the `GSIF` package. It contains three columns: `Sand`, `Silt` and `Clay`. A smaller version, containing only 50 observations is stored in `africa_sample`.

In the first course we mentioned that in the data layer, the structure of the data should reflect how you wish to plot it. For a ternary plot, you need to have three separate variables, for example, Sand, Silt and Clay in `africa`. However, for a proportional/stacked bar plot, you just need two. The type should be defined as three levels within a single factor variable. That is, you want tidy data.

It's also useful to maintain the site IDs as a variable within the data frame, currently, they are stored at row names, which is poor style and not useful.

```
(africa.sample <- africa  
  %>% mutate(ID = row_number())  
  %>% select (ID, everything())  
  %>% sample_n(50))
```

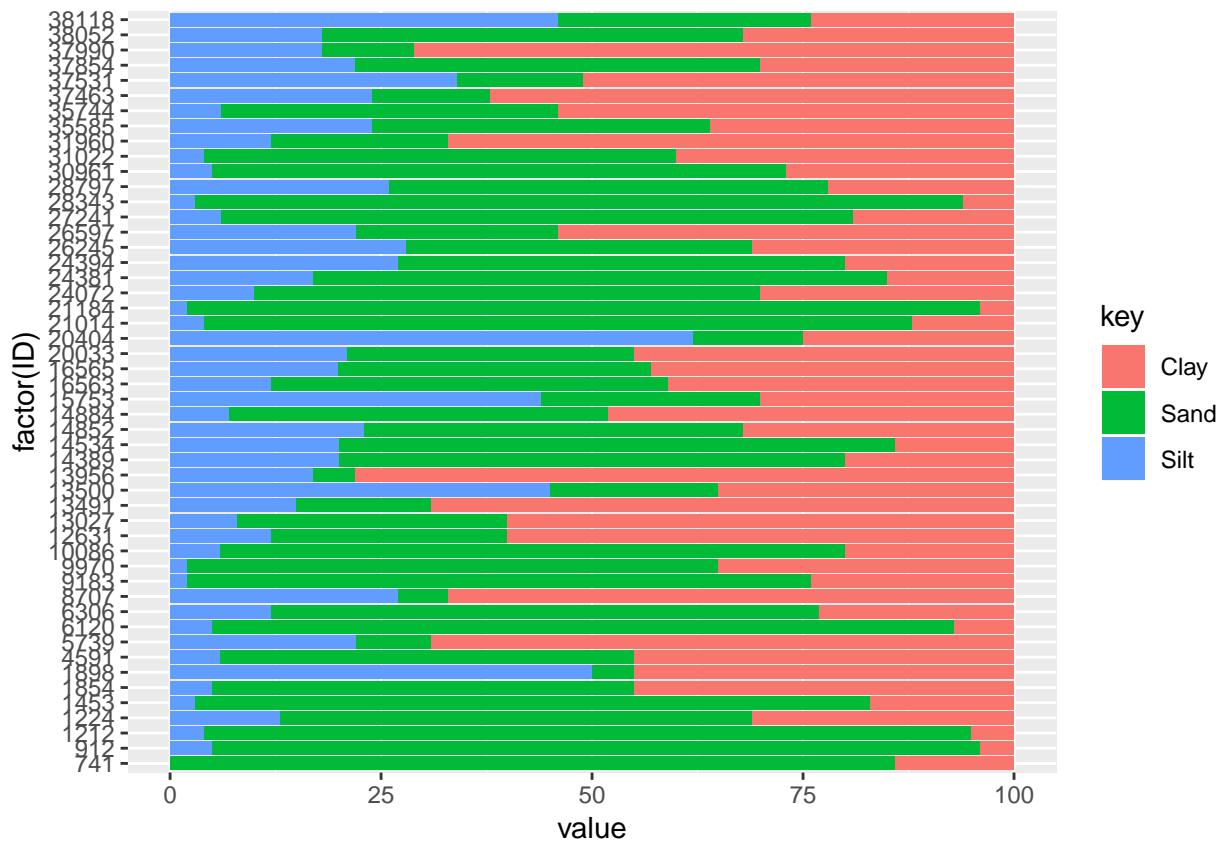
ID	Sand	Silt	Clay
10086	74	6	20
8707	6	27	67
37531	15	34	51
30961	68	5	27
14534	66	20	14
14852	45	23	32
13500	20	45	35
24381	68	17	15
26597	24	22	54
1212	91	4	5
15753	26	44	30
6120	88	5	7
4591	49	6	45
9970	63	2	35
35585	40	24	36
13027	32	8	60
741	86	0	14
38052	50	18	32
14884	45	7	48
21184	94	2	4
27241	75	6	19
24394	53	27	20
6306	65	12	23
37463	14	24	62
1224	56	13	31
35744	40	6	54
16563	47	12	41
37854	48	22	30
20404	13	62	25
13491	16	15	69
1898	5	50	45
1453	80	3	17
16565	37	20	43
21014	84	4	12
20033	34	21	45
37990	11	18	71
31960	21	12	67
13956	5	17	78
912	91	5	4
24072	60	10	30
9183	74	2	24
12631	28	12	60
28343	91	3	6
5739	9	22	69
14389	60	20	20
1854	50	5	45
28797	52	26	22
31022	56	4	40
26245	41	28	31
38118	30	46	24

```
africa.sample.tidy <- gather(africa.sample, key, value, -ID)
```

```
head(africa.sample.tidy)
```

ID	key	value
10086	Sand	74
8707	Sand	6
37531	Sand	15
30961	Sand	68
14534	Sand	66
14852	Sand	45

```
ggplot(africa.sample.tidy, aes(x = factor(ID), y = value, fill = key)) +  
  geom_col() +  
  coord_flip()
```



```
# africa.sample.2 <- africa  
#   %>% mutate(ID = row_number())  
#   %>% gather(`Sand`, `Silt`, `Clay`, key = "Component", value = "Quantity", -ID)  
#   %>% sample_n(50))  
#  
# ggplot(africa.sample.2, aes(x = factor(ID), y = Quantity, fill = Component)) +  
#   geom_col() # using geom_bar I got this Error: stat_count() must not be used with a y aesthetic.  
#   coord_flip()
```

Stoked for stacked! Remember stacked bar charts can be used as an alternative to pie charts.

## Networks Plots

Network analysis is a field of study onto itself, but even though it is not your focus, it is not unusual to encounter relationship data at some point. There are several R packages which can be used to visualize relationship data. Integrating network data with ggplot2 allows you to take advantage of all the tools we've learned so far in these series. Let's begin with a simple example – a network of blood type donors. There are four different blood types: A, B, AB, O. Each present as  $\rho$  negative or positive. In total we've eight different blood types. Not all blood types are compatible i.e. in this case, the question we want to address with our network is the directionality of blood type donors: Which blood type can be used as a donor for the other blood types? There are 27 relationships which we want to visualize which are stored in the dataframe `blood_donors` shown below.

At this point in our data visualization understanding, we may be tempted to just make a heat map. Here we are just looking at where we have a relationship or not. This is already a good starting point, and maybe this visualization already answers your questions since we can patterns in donors and recipient profiles. For example, we can already see that blood type O- is a universal donor. An alternative will be a dot plot, which would allow us to further use the size of the dots to represent how prevalent a donor is in our data set. However, at this point we just want to look at the relationships so we turn to a network.

In a network each observation is called a `vertex`, and in the `blood_donor` dataframe, we've eight vertices arranged in a circle. The arrangement of vertices can alter our perception of the relationship. So, a circle is a good starting point since it will allow us to see all possible relationships clearly. The connections between each vertex is an `edge`. The shape of the network are the nature of the relationships and can be quantified in various ways. Here, each edge has a arrow head since we are talking about a directional relationship, for which blood type is each vertex a donor – only the relationship present are visualized. Since, this is a `ggplot2` object, we can apply other layers such as facetting. For example, if we have had information about the predominance of each blood type in different ethnicities, we can scale this as of each dot to the variable and facet the network according to ethnicity.

### Build the Network (1)

Network data may be stored in a variety of ways.

For this example, you'll use an undirected network of romantic relationships in the TV show Mad Men: `geomnet::madmen`.

```
# Load geomnet & examine structure of madmen
# Load the geomnet package and examine the structure of the madmen dataset.
# Notice that it is a list of two data frames called edges and vertices.
# This is enough information to build the network.

library(geomnet)
str(madmen)
#> List of 2
#> $ edges  :'data.frame': 39 obs. of  2 variables:
#>   ..$ Name1: Factor w/ 9 levels "Betty Draper",...: 1 1 2 2 2 2 2 2 2 ...
#>   ..$ Name2: Factor w/ 39 levels "Abe Drexler",...: 15 31 2 4 5 6 8 9 11 21 ...
#> $ vertices:'data.frame': 45 obs. of  2 variables:
#>   ..$ label : Factor w/ 45 levels "Abe Drexler",...: 5 9 16 23 26 32 33 38 39 17 ...
#>   ..$ Gender: Factor w/ 2 levels "female","male": 1 2 2 1 2 1 2 2 2 2 ...

# Merge edges and vertices
# Use merge() to merge the madmen$edges and madmen$vertices data frames.
# You can connect observations: by.x is the Name1 variable in the edges data frame,
# which matches by.y, the label variable in the vertices data frame.
# Store the result inside a new data frame mmnet.
mmnet <- merge(madmen$edges, madmen$vertices,
```

```

    by.x = "Name1", by.y = "label",
    all = TRUE)

# Examine structure of mmnet
str(mmnet)
#> 'data.frame': 75 obs. of 3 variables:
#> $ Name1 : Factor w/ 45 levels "Betty Draper",...: 1 1 2 2 2 2 2 2 2 ...
#> $ Name2 : Factor w/ 39 levels "Abe Drexler",...: 15 31 2 4 5 6 8 9 11 21 ...
#> $ Gender: Factor w/ 2 levels "female","male": 1 1 2 2 2 2 2 2 2 ...

```

Remember, the first step is to get the data structure right.

## Build the Network (2)

Now that your data is in the correct format, you can build the actual network plot.

You'll use the `geom_net()` function, a `ggplot` layer that's in the geomnet package. The `ggnetwork` package is a popular alternative, but we will not discuss that here.

Can you finish the `ggplot()` command?

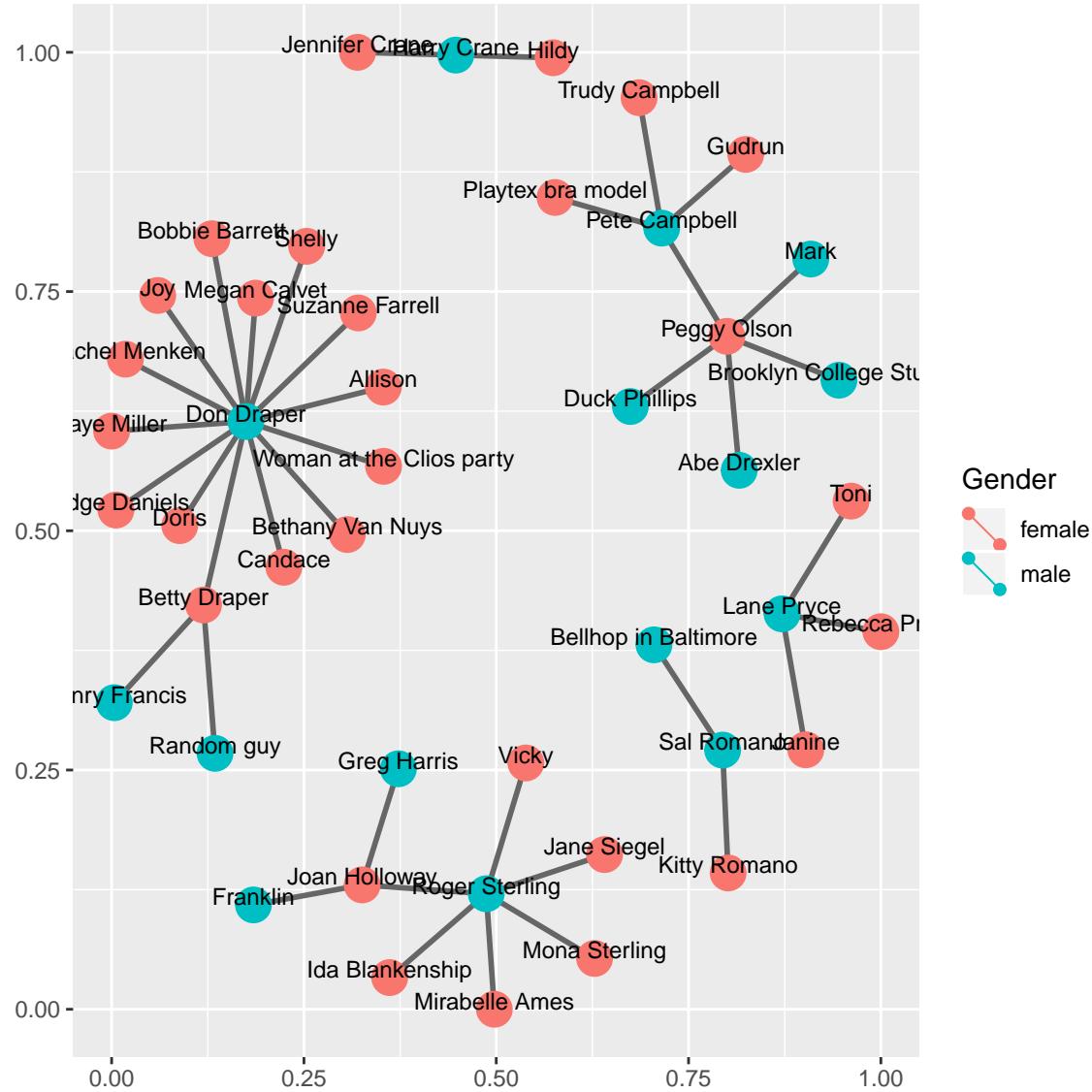
```

# geomnet is pre-loaded

# Merge edges and vertices
mmnet <- merge(madmen$edges, madmen$vertices,
               by.x = "Name1", by.y = "label",
               all = TRUE)

# Finish the ggplot command
ggplot(data = mmnet, aes(from_id = Name1, to_id = Name2)) +
  geom_net(aes(col = Gender),
           size = 6,
           linewidth = 1,
           labelon = TRUE,
           fontsize = 3,
           labelcolour = "black")

```



Mad about Networks! Can you spot the show's womanizer? Draper, is by far, the show's womanizer.

### Adjusting the Network

Let's clean up the network a bit. As you can see, since this is in the `ggplot2` framework, you can manually adjust the scales like you have always done.

Here you're going to use the `void` theme, `theme_void()`, to remove all background elements like axes and grid lines to make a clean network plot.

```
# geomnet is pre-loaded and mmnet is defined
head(mmnet)
```

Name1	Name2	Gender
Betty Draper	Henry Francis	female
Betty Draper	Random guy	female
Don Draper	Allison	male
Don Draper	Bethany Van Nuys	male
Don Draper	Betty Draper	male
Don Draper	Bobbie Barrett	male

```
# Node colors
pink_and_blue <- c(female = "#FF69B4", male = "#0099ff")

# Tweak the network plot
ggplot(data = mmnet, aes(from_id = Name1, to_id = Name2)) +
  geom_net(aes(col = Gender),
           size = 6,
           linewidth = 1,
           labelon = TRUE,
           fontsize = 3,
           labelcolour = "black",
           # Make the graph directed
           directed = TRUE) +
  # Add manual color scale
  scale_color_manual(values = pink_and_blue) +
  # Set x-axis limits
  xlim(-0.05, 1.05) +
  # Set void theme
  theme_void()
```



Notable Networks! Although blue and pink are typical colors from male and female, feel free to break gender stereotypes!

## Diagnostic Plots

Please visit the development page of the [prettydoc](#) package for latest updates and news. Comments, bug reports and pull requests are always welcome.