

# Basic Course on **R**: Using ggplot2

*David Nieuwenhuijse*

*May 18<sup>th</sup> - 24<sup>th</sup>, 2017*

## Contents

<b>1</b>	<b>Using ggplot2: Lecture</b>	<b>2</b>
1.1	Grammar of graphics . . . . .	2
1.2	Data . . . . .	2
1.3	Aesthetics . . . . .	5
1.4	Geometries . . . . .	6
1.4.1	Scatter plots . . . . .	6
1.4.2	Line graphs . . . . .	8
1.4.3	Bar charts . . . . .	9
1.4.4	Histograms . . . . .	11
1.4.5	Box plots . . . . .	13
1.4.6	Heatmaps . . . . .	17
1.5	Facets . . . . .	18
1.6	Scales . . . . .	21
1.7	Themes . . . . .	22
1.8	Saving ggplots . . . . .	24
<b>2</b>	<b>Using ggplot: Practicals</b>	<b>25</b>
<b>3</b>	<b>Using ggplot: Answers</b>	<b>26</b>

# 1 Using ggplot2: Lecture

The basic plotting functions of R are nice to quickly visualize data, however, the strength of R is that it has many packages with extra functionality created by others that you can use. ggplot2 is one of those packages, created by Hadley Wickham. ggplot2 allows you to build more sophisticated plots in R using a “Grammar of Grappics” (therefore GGplot).

If ggplot2 is not yet available for you we need to install and load the ggplot2 library like so.

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

```
## Warning: package 'ggplot2' was built under R version 3.3.3
```

ggplot2 is a very extensive package, so it is impossible to show all the possibilities during this course. If you would like to experiment further, the ggplot manual at <http://ggplot2.tidyverse.org/reference/index.html> is the place to go. A large variety of functionallities are listed there, including example plots.

## 1.1 Grammar of graphics

The grammar of graphics consists of a number of “layers”, which generate a ggplot when correctly put together. These five basic layers will be discussed here:

- Data
- Aesthetics
- Geometries
- Facets
- Scales
- Themes

## 1.2 Data

The basis of each plot is the data. To plot data with ggplot the data needs to be of the class “data.frame”. The structure of your data.frame is important when you want to plot it with ggplot. Each variable you want to use has to be in a separate column in your data.frame. Take for example the “Indometh” dataset which comes with R.

```
?Indometh
```

The “Indometh” dataset describes the pharmacokinetics of indomethacin after intravenous injection of indomethacin in six subjects. The data.frame has three columns one containing the subject number, one containing the time of sampling in hours, and one containing the indomethacin concentration ug/ml. This dataset is easily plotted with ggplot because it is a data.frame and there is a good separation of variables in the columns.

```
##   Subject time conc
## 1      1  0.25 1.50
## 2      1  0.50 0.94
## 3      1  0.75 0.78
## 4      1  1.00 0.48
## 5      1  1.25 0.37
## 6      1  2.00 0.19
```

An example of a badly formatted dataset for visualization with ggplot is a dataset of monthly deaths from lung diseases in the UK:

## ?ldeaths

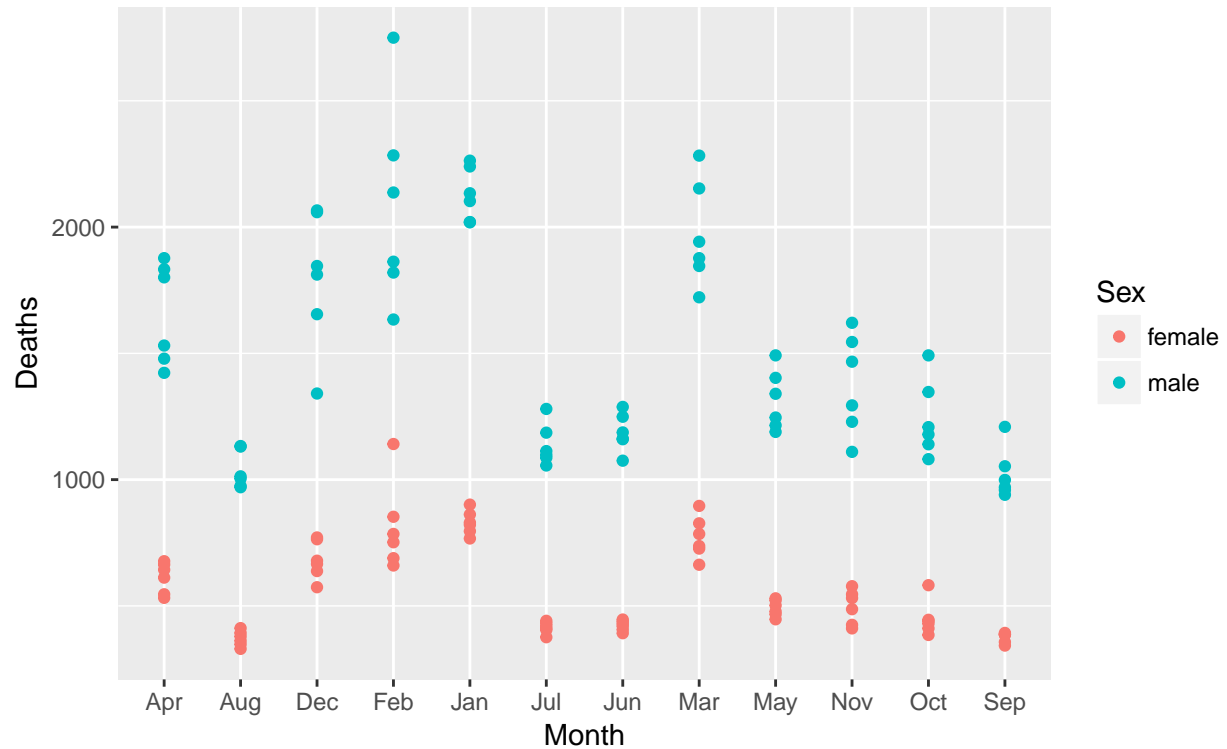
This dataset contains three separate datasets for female, male and both male and female deaths (ldeaths, fdeaths, and mdeaths). The columns contain values per month and the rows indicate the year. (also they are not data.frames)

```
##      Jan  Feb  Mar  Apr  May  Jun  Jul  Aug  Sep  Oct  Nov  Dec
## 1974 3035 2552 2704 2554 2014 1655 1721 1524 1596 2074 2199 2512
## 1975 2933 2889 2938 2497 1870 1726 1607 1545 1396 1787 2076 2837
## 1976 2787 3891 3179 2011 1636 1580 1489 1300 1356 1653 2013 2823
## 1977 3102 2294 2385 2444 1748 1554 1498 1361 1346 1564 1640 2293
## 1978 2815 3137 2679 1969 1870 1633 1529 1366 1357 1570 1535 2491
## 1979 3084 2605 2573 2143 1693 1504 1461 1354 1333 1492 1781 1915
```

Since there are two separate datasets, there is no column indicating the sex. There is also no column indicating the year and no column indicating the month. This means that we cannot use these variables in ggplot to visualize, for instance, the difference between the number of deaths in men and women in a certain year. We can use some more advanced R data reshaping to get these datasets in the right shape.

```
##   Deaths Year Month   Sex
## 1     901 1974   Jan female
## 2     689 1974   Feb female
## 3     827 1974   Mar female
## 4     677 1974   Apr female
## 5     522 1974   May female
## 6     406 1974   Jun female
```

```
ggplot(data, aes(x=Month, y=Deaths, color=Sex)) + geom_point()
```



However, it is much more easy to keep in mind that any variable that you want to use as a part of your ggplot should get its own column when you are designing your dataset (in Excel). With the data in the right class and the right format we can easily call the ggplot function, and add a scatter plot layer to it.

### 1.3 Aesthetics

After figuring out what data to plot it is necessary to indicate how to plot it. Aesthetics are used to indicate how to plot what. As we saw in the plot of the `ldeaths` data we start with indicating what `data.frame` to plot and use `aes()` to indicate what to plot on the x axis, the y axis and what to use as color. We do not have to indicate what to plot where and which color to give to which category. We simply map a column of our `data.frame` to an aesthetic. Depending on the geom we use there are a number of aesthetics that can be mapped to a variable:

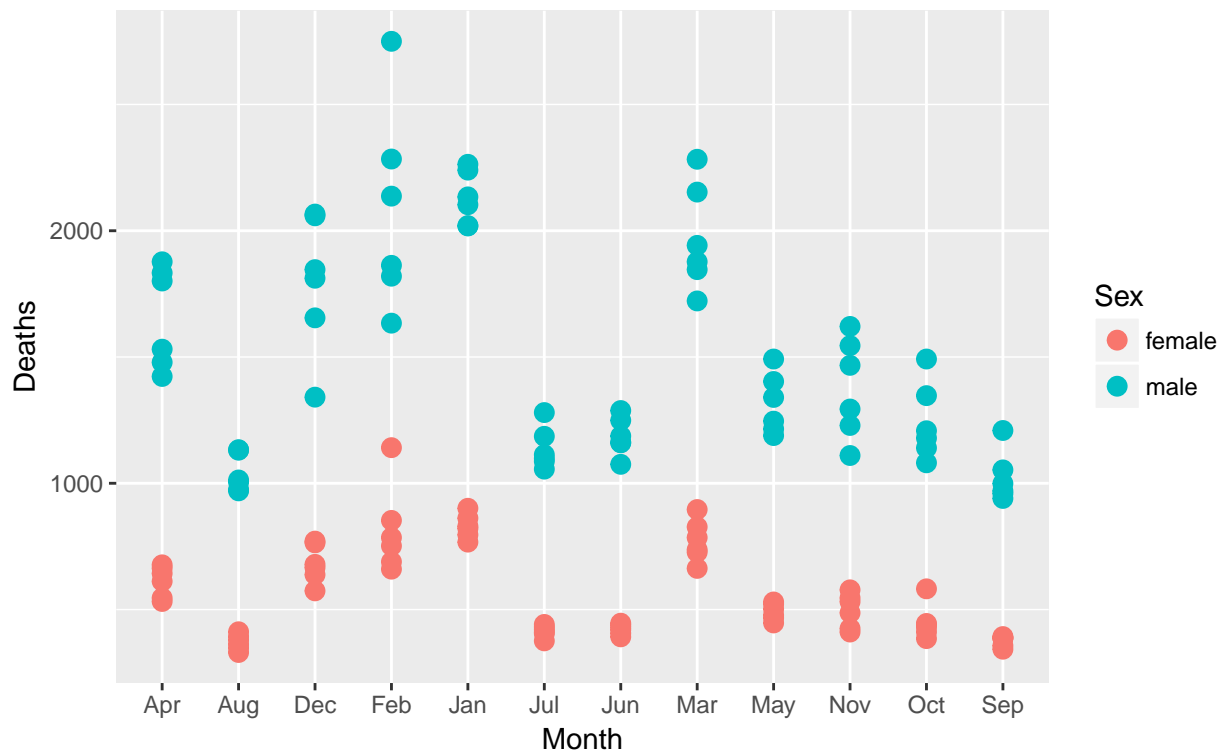
- x
- y
- color
- fill
- size
- alpha
- linetype
- labels
- shape
- group

Aesthetics can also be specified per geom layer. This allows plotting multiple geom layers with different aesthetics. However, when using one data source and one geom layer this will result in exactly the same plot.

```
ggplot(data) + geom_point(aes(x=Month, y=Deaths, color=Sex))
```

We can also use plotting parameters without mapping them to a variable to change all dots in the `geom_point`. This can only be done in the geom itself.

```
ggplot(data, aes(x=Month, y=Deaths, color=Sex)) + geom_point(size=3)
```



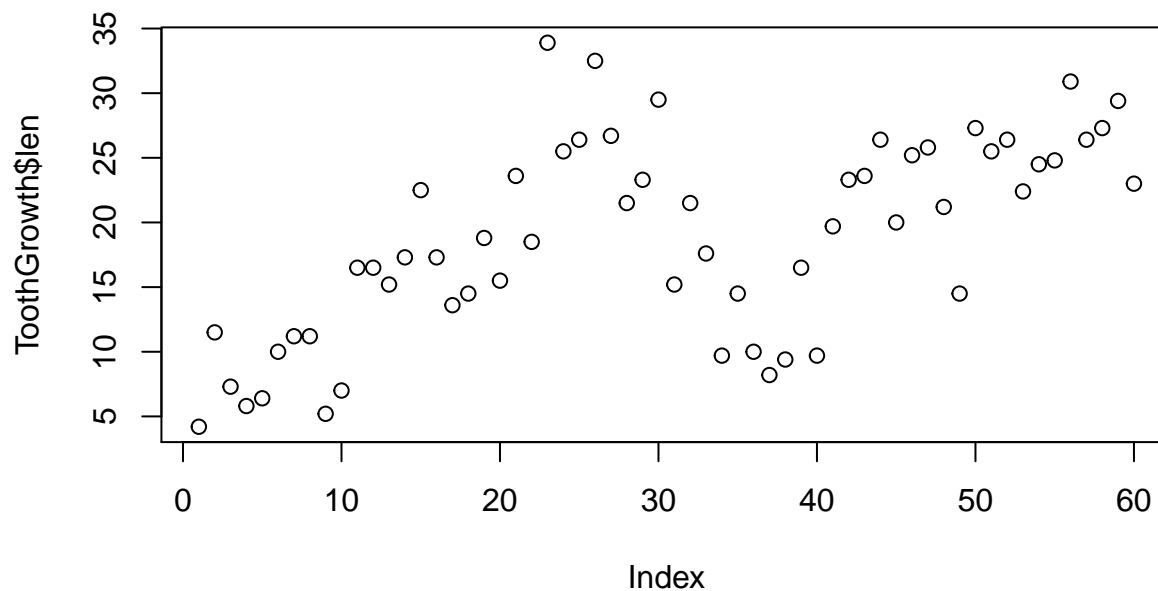
## 1.4 Geometries

For the previous ggplots we have used the `geom_point` function which generates a scatter plot layer that can be added to the ggplot. “point” is one of the many geometry layers available in the ggplot library. There are over thirty different geoms in ggplot which can be found at <http://docs.ggplot2.org/>. We will have a look at the plot types that we have already encountered using the basic R plotting functions and replace them with their ggplot counterparts.

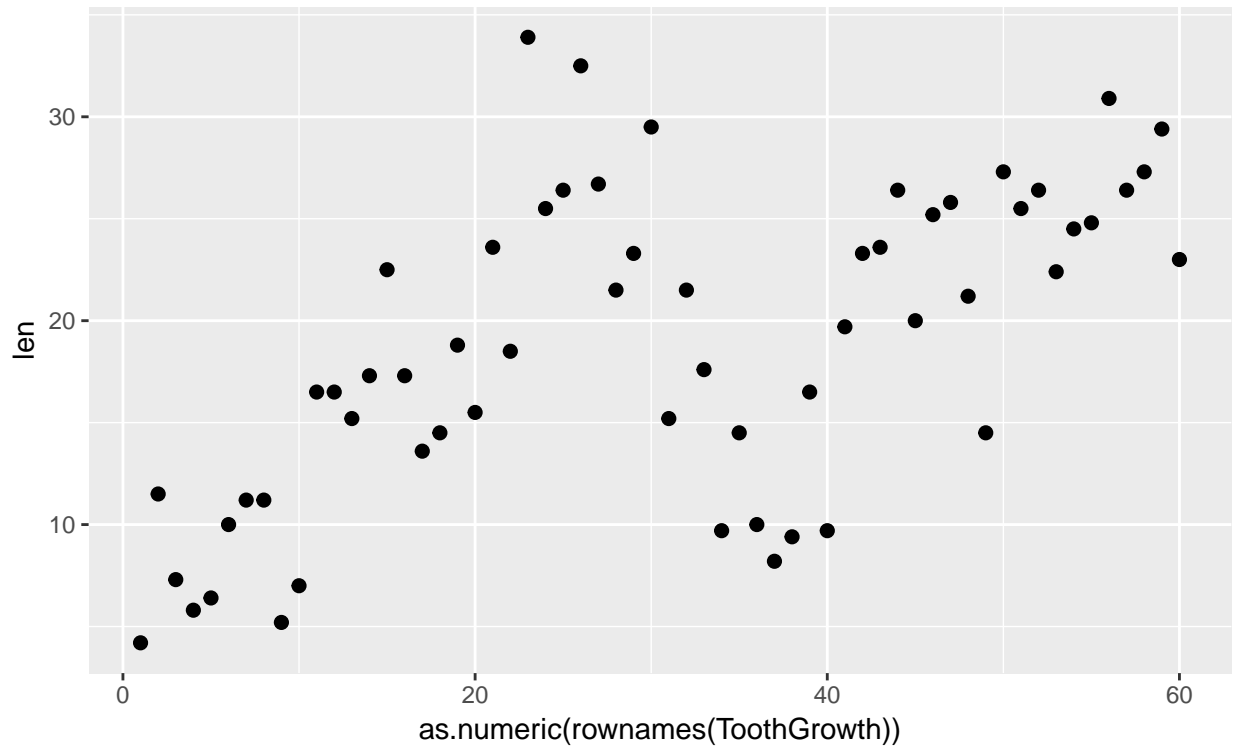
### 1.4.1 Scatter plots

Just like with other objects in R a ggplot object can be saved to a variable using “<-”. This makes it possible to store the basic plot in a variable and add different geoms to it. Geometries are added to a plot using “+”.

```
#Basic graphics:  
plot(ToothGrowth$len)
```



```
p <- ggplot(ToothGrowth)
p + geom_point(aes(x=as.numeric(rownames(ToothGrowth))), y=len), size=2)
```



We now map the numeric rownames of ToothGrowth to the x aesthetic, however, we could also create a new column called index in the data.frame.

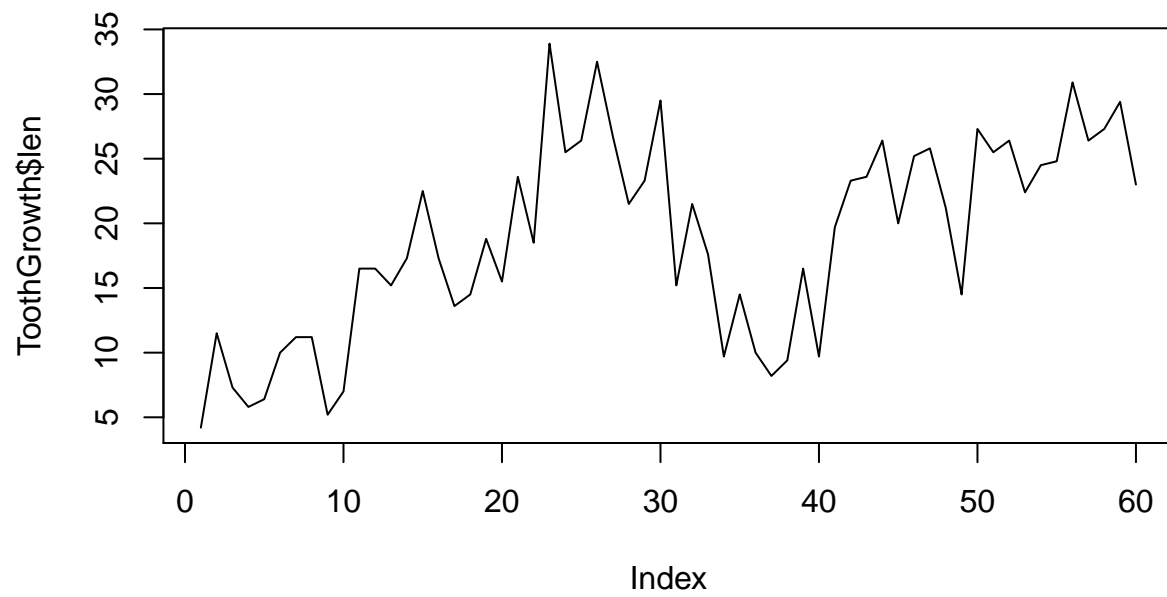
```
ToothGrowth$index <- as.numeric(rownames(ToothGrowth))
```

When we change the data of the ggplot object we saved in variable p we need to recreate the ggplot object with the new dataset.

```
p <- ggplot(ToothGrowth)
```

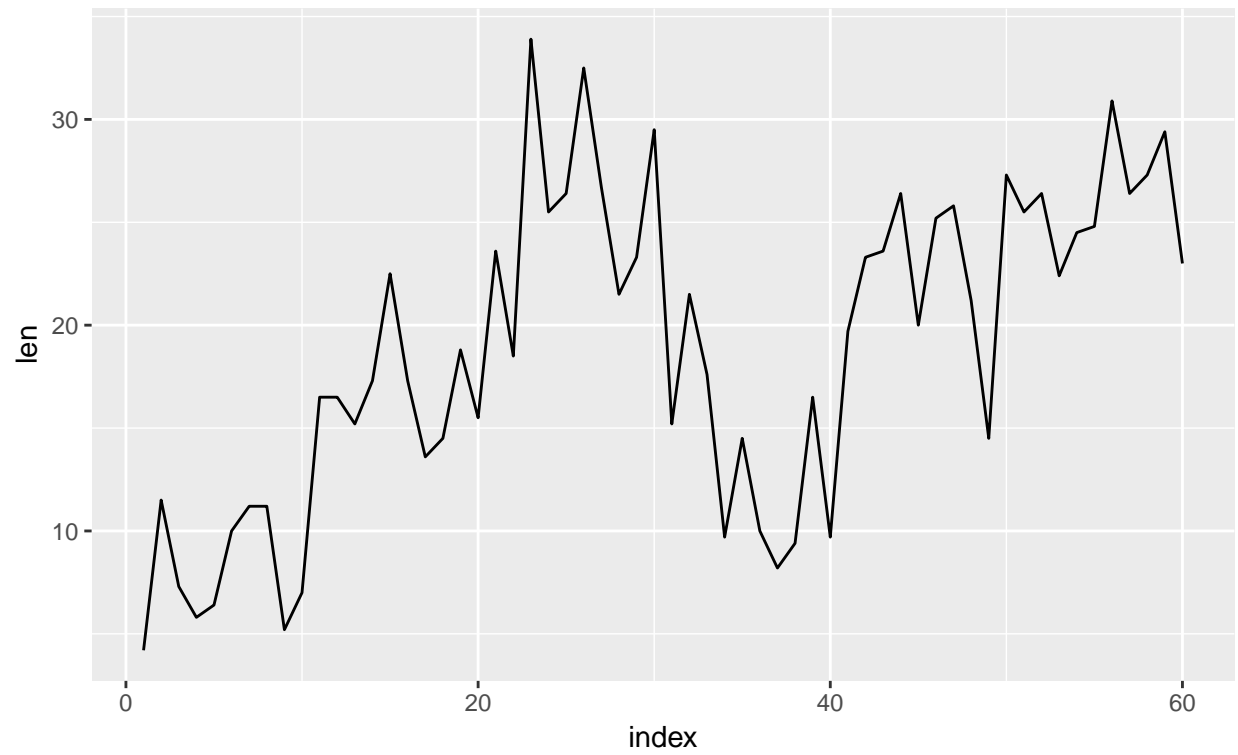
### 1.4.2 Line graphs

```
#Basic graphics:  
plot(ToothGrowth$len, type = "l")
```



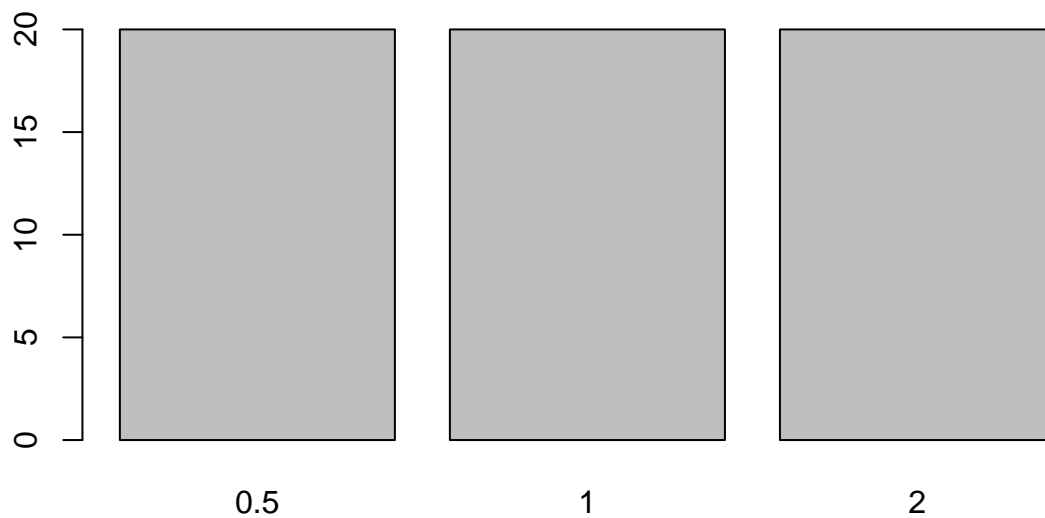
```
p + geom_line(aes(x=index, y=len))
```



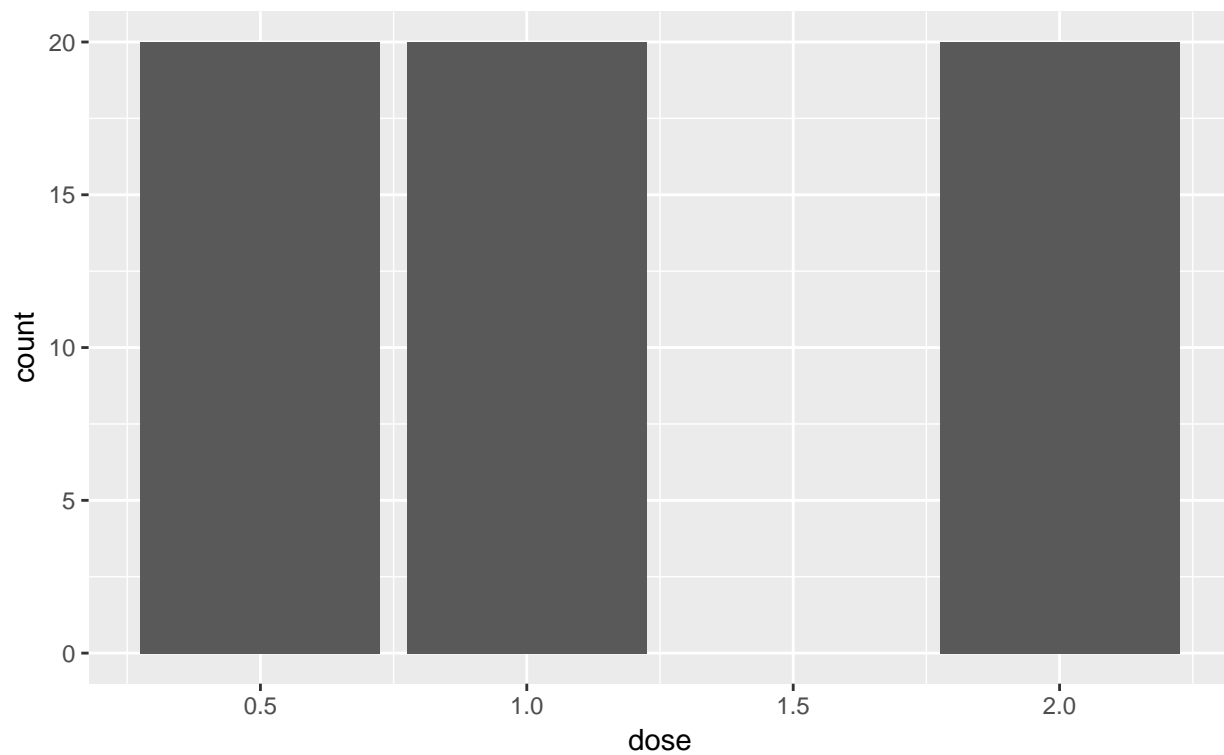


### 1.4.3 Bar charts

```
#Basic graphics:  
barplot(table(ToothGrowth$dose))
```

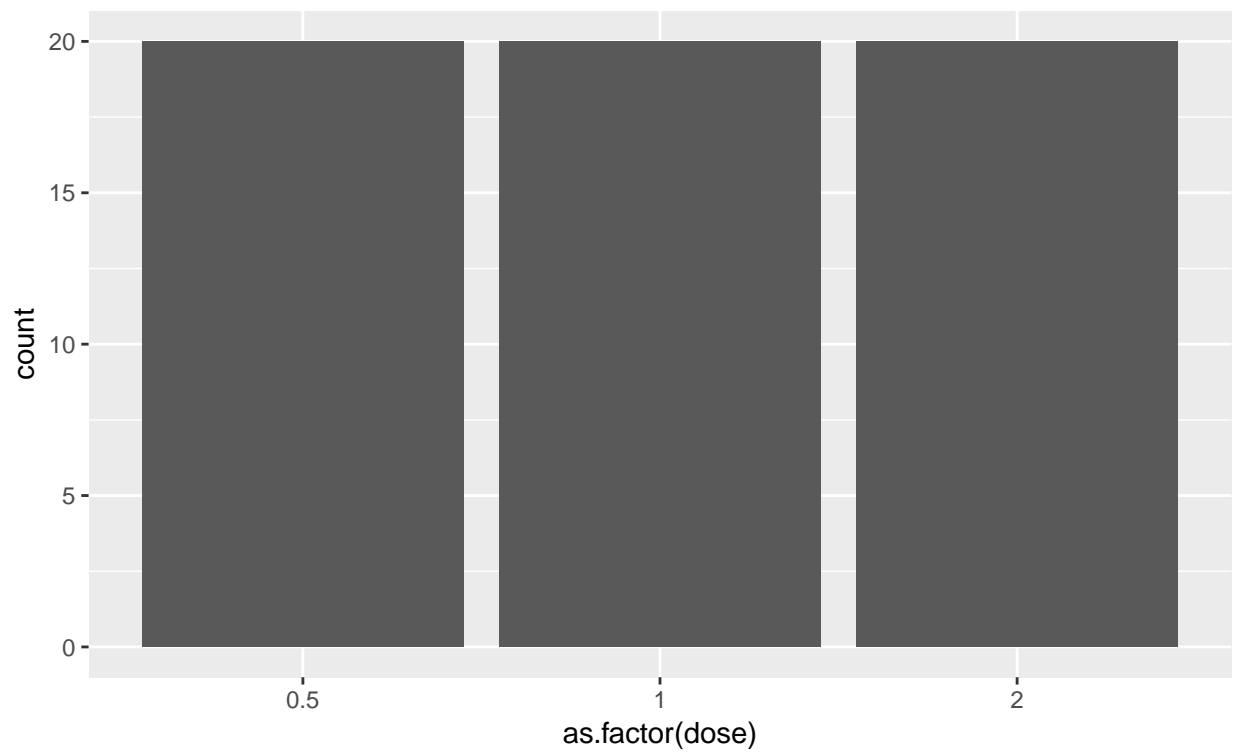


```
p + geom_bar(aes(x=dose))
```



Because dose is a continuous variable, ggplot creates a x axis with a continuous scale, which is why there is an empty spot for the 1.5 dose. If we want to use a continuous value as a categorical value we can change it into a factor on the fly in the geom function:

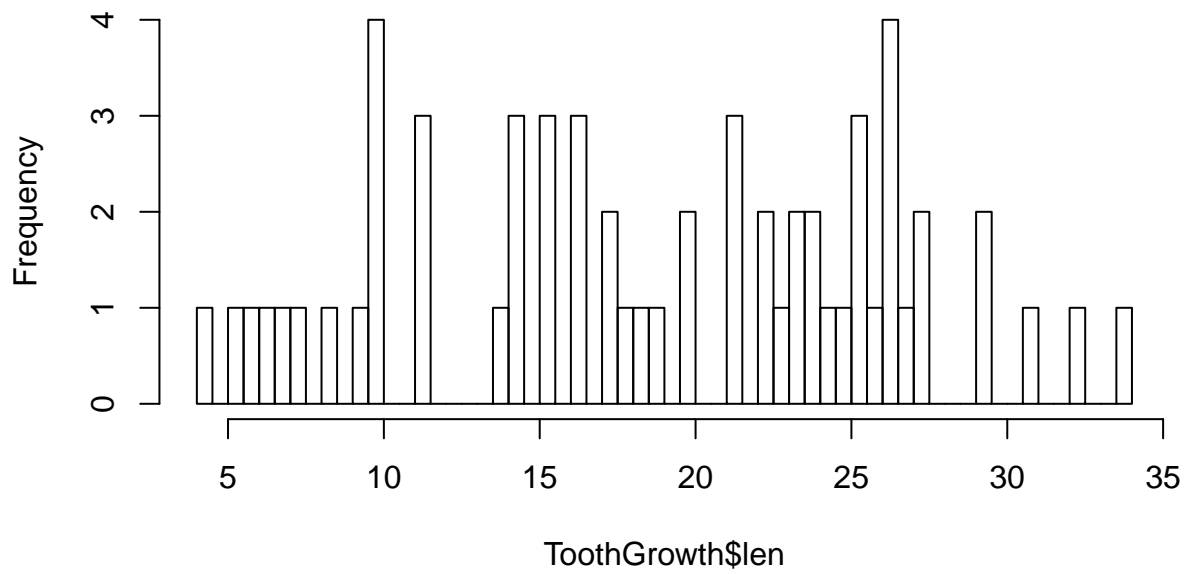
```
p + geom_bar(aes(x=as.factor(dose)))
```



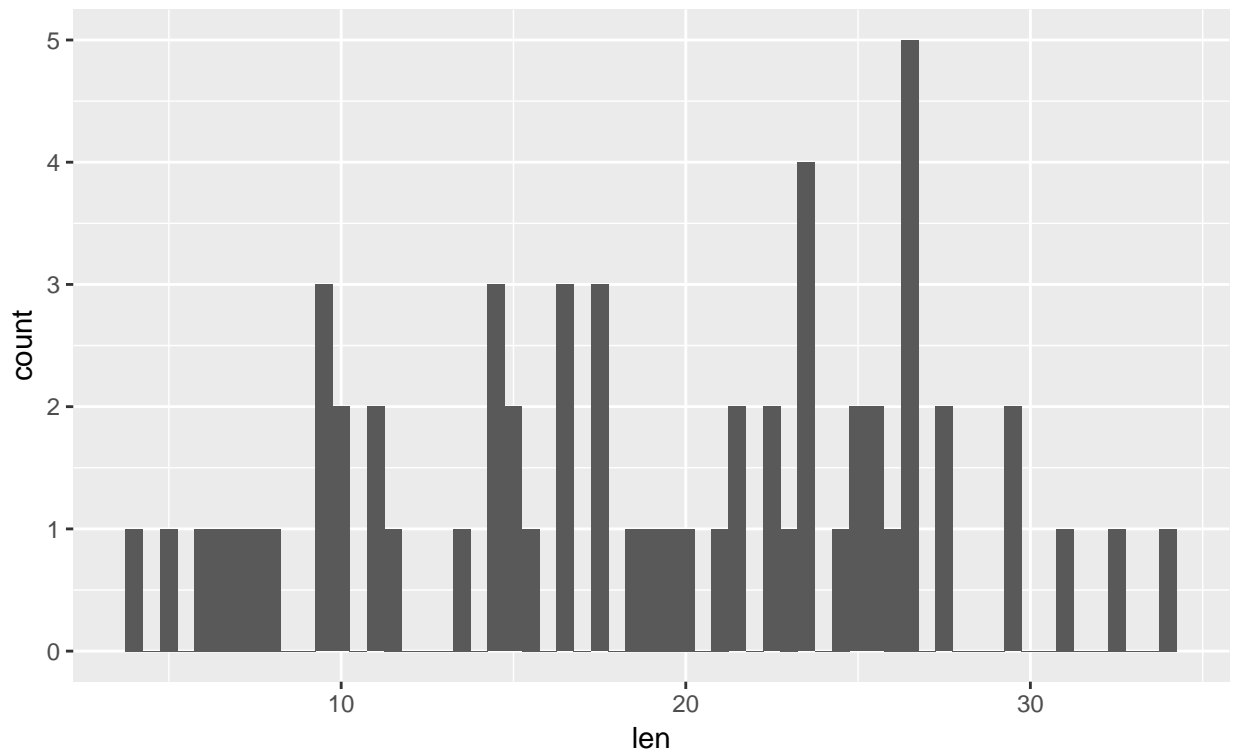
#### 1.4.4 Histograms

```
#Basic graphics:  
hist(ToothGrowth$len, breaks = 50)
```

**Histogram of ToothGrowth\$len**

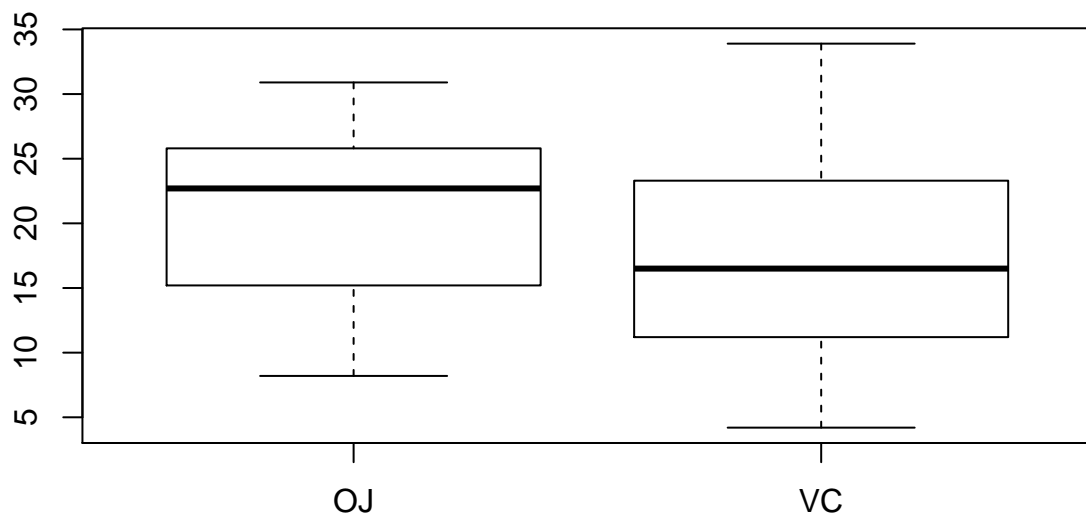


```
p + geom_histogram(aes(x=len), binwidth = 0.5)
```

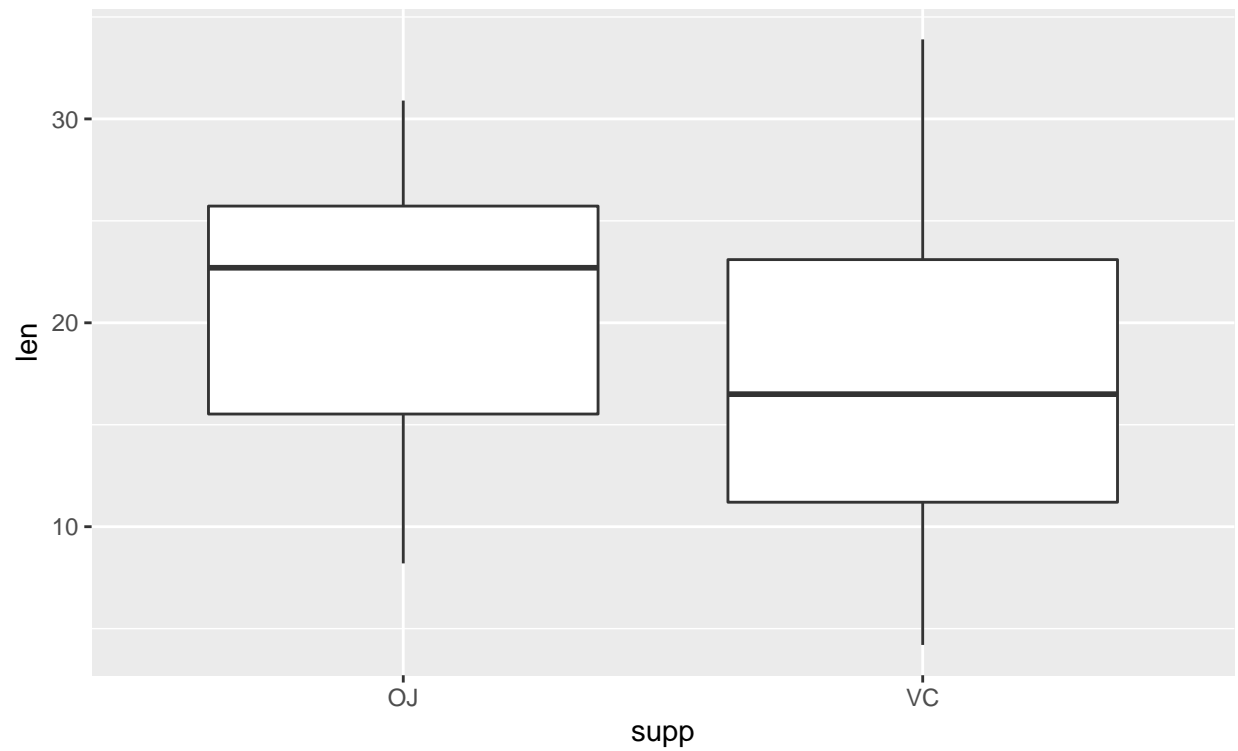


### 1.4.5 Box plots

```
#Basic graphics:  
boxplot(ToothGrowth$len~ToothGrowth$supp)
```

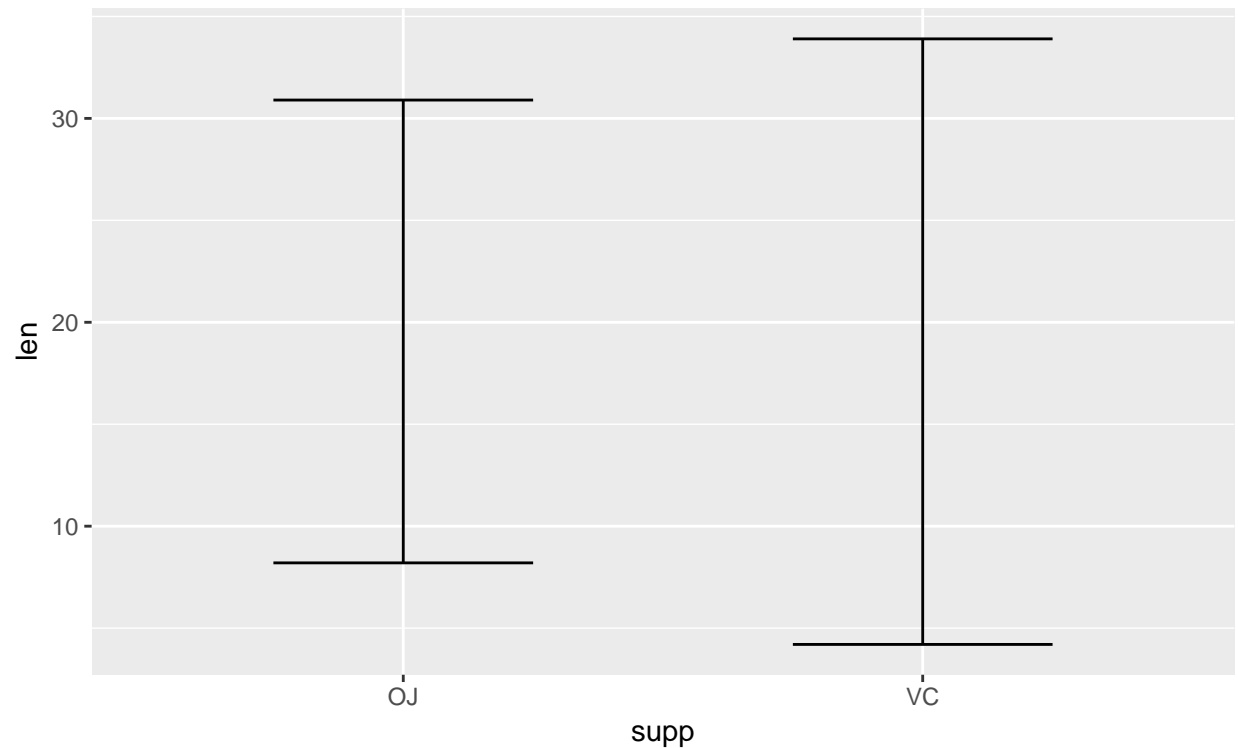


```
p + geom_boxplot(aes(x=supp,y=len))
```



If we want whiskers in our plot we can add an errorbar geom

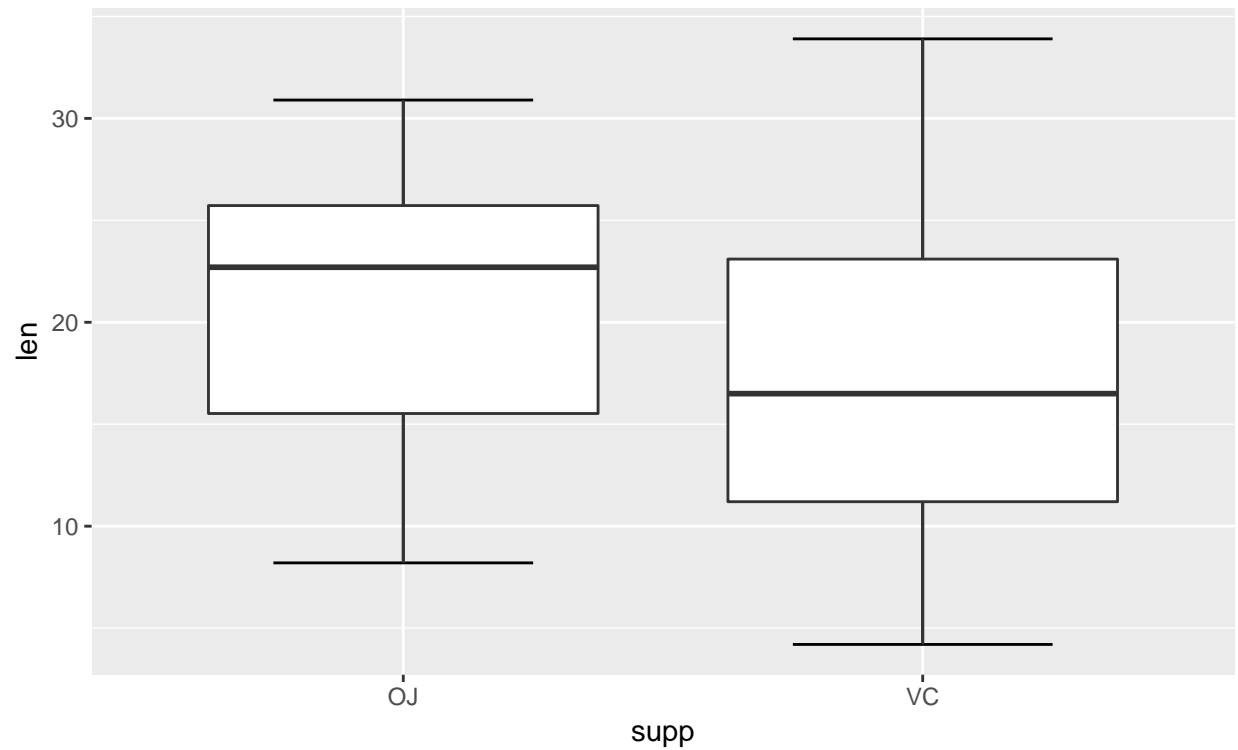
```
p + geom_errorbar(aes(x=supp, ymin=..., ymax=...),width=0.5)
```



However, as you can see, this geom requires `ymin` and `ymax` as aesthetics and does not calculate the values by itself. We could calculate the min and max values for toothlength grouped by supplement, but we will use a trick using ggplot's boxplot statistics function. *(if you want to know more, please read more about these "stat" functions online)*

First, we add the errorbars with the `"stat_boxplot"` function, and then add the boxplots on top of them.

```
p + stat_boxplot(aes(x=supp,y=len), geom="errorbar", width=0.5) +  
  geom_boxplot(aes(x=supp,y=len))
```



As we can see the aesthetics of both geoms have to be defined separately in this plot. To reduce typing we could also write the ggplot as follows.

```
p <- ggplot(ToothGrowth, aes(x=supp,y=len))  
p + stat_boxplot(geom="errorbar", width=0.5) + geom_boxplot()
```

When the aesthetics are specified in the base ggplot object it is not necessary to specify them in the added layers.

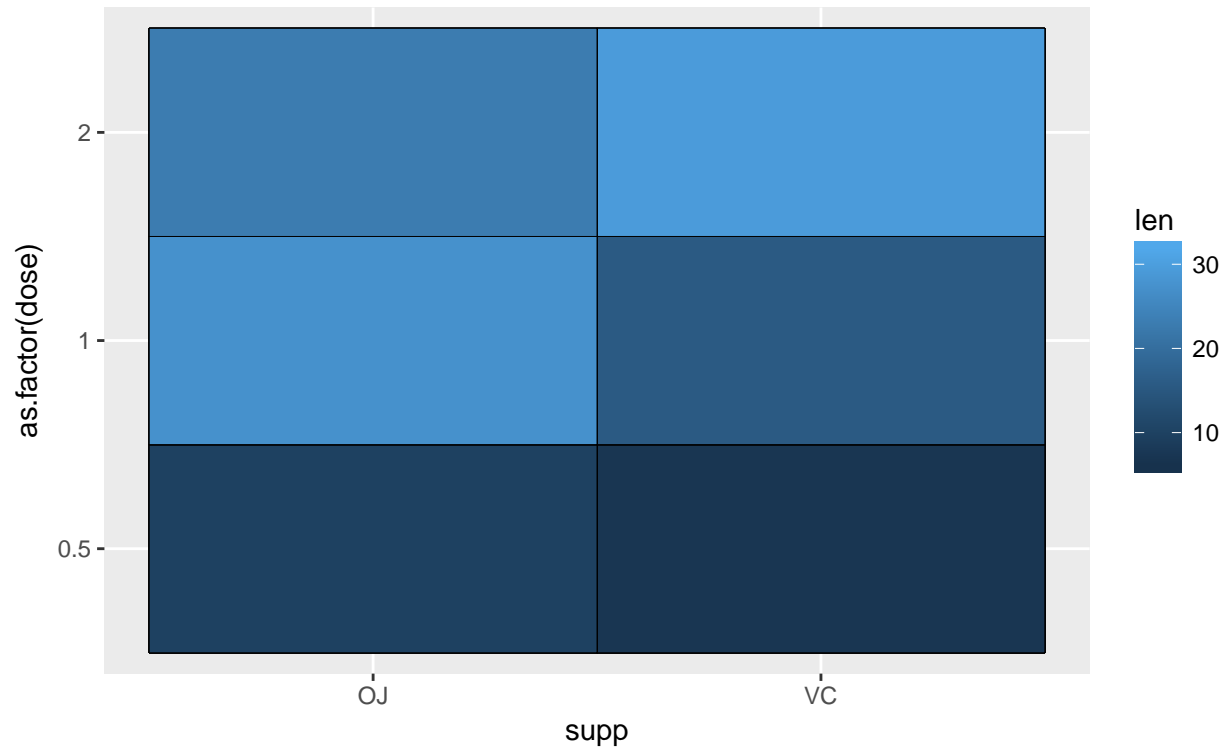


### 1.4.6 Heatmaps

One popular way of displaying large grids of data is using a heatmap. The tile geom generates a tile at each provided x and y value and uses the fill aesthetic to fill the tiles according to another variable. The border color of the tiles can be changed by using the color aesthetic.

```
p <- ggplot(ToothGrowth)

p + geom_tile(aes(x = supp, y = as.factor(dose), fill = len), color="black")
```



There are dedicated packages for plotting heatmaps, which also perform horizontal and vertical clustering of the tiles. These are, however, out of the scope of this course, but it should be easy to find them online.

## 1.5 Facets

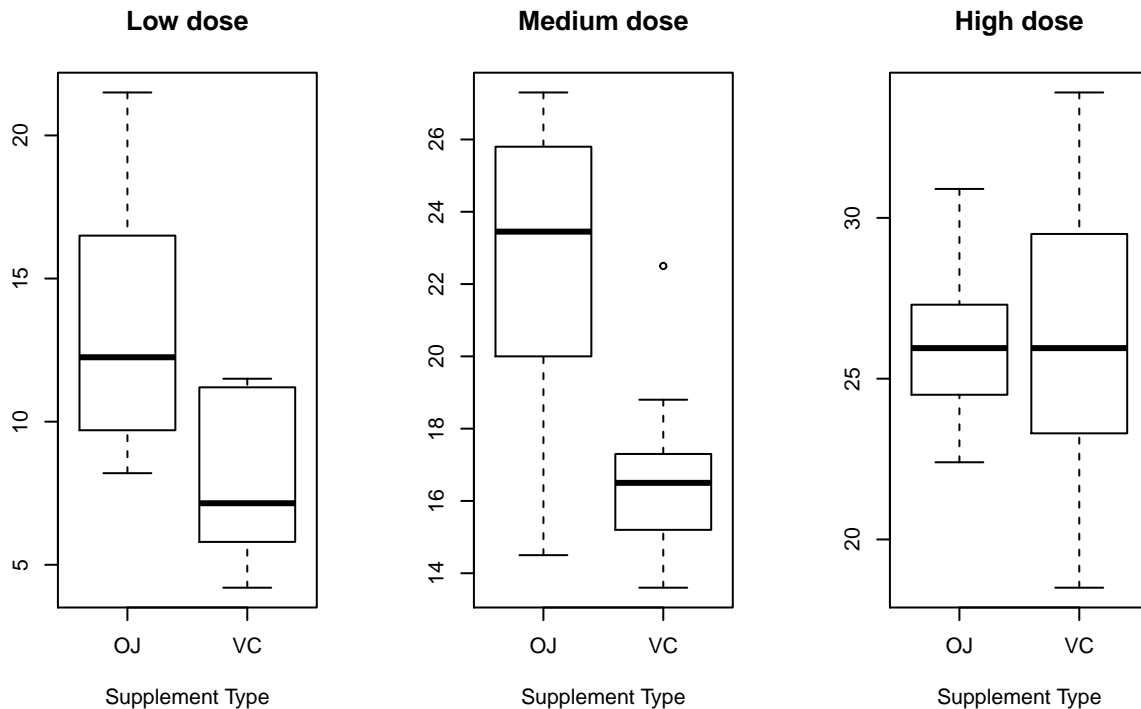
Apart from using aesthetics to map variables to, ggplot has another feature which can help visualizing data called “facets”. We have added multiple plots to the viewing window before, using “`par(mfrow=c(...,...))`”. However, after specifying how many plots should be printed in the plot viewer we have to call the plot function multiple times to populate the viewer.

If we want to compare the effect of orange juice versus vitamin C depending on the dose we can plot three boxplots next to each other.

```
par(mfrow = c(1, 3))

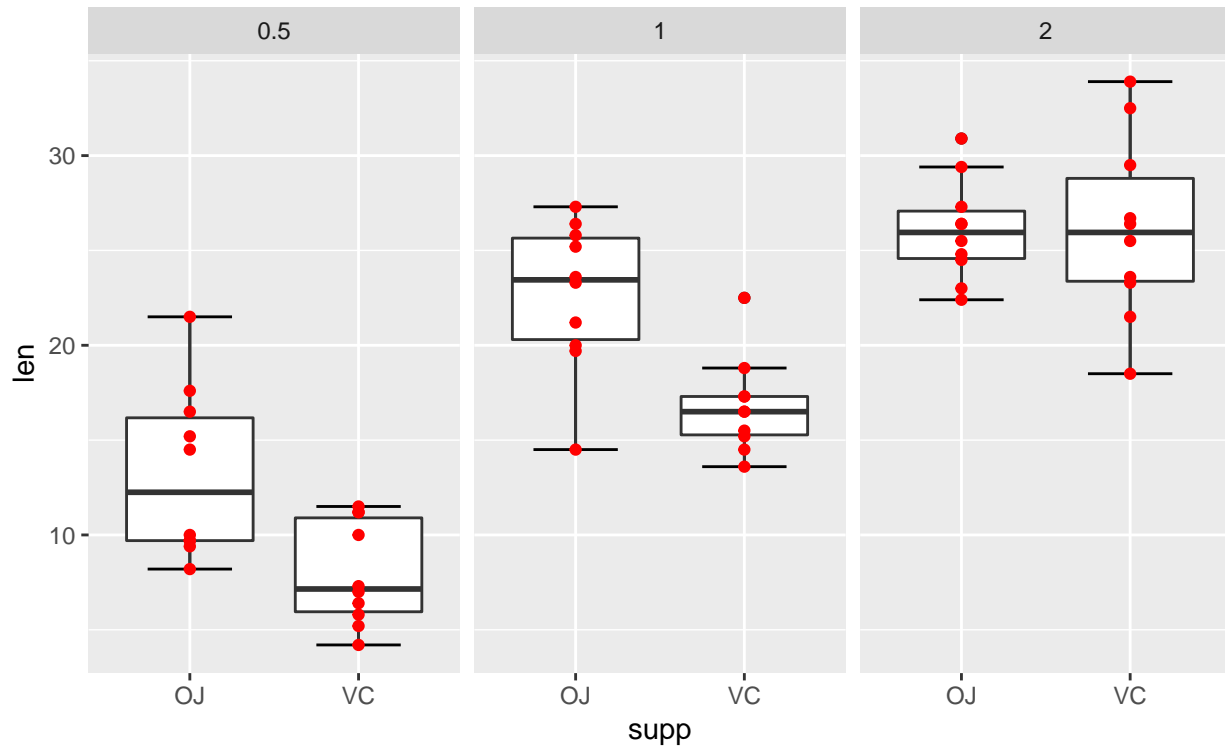
low <- ToothGrowth[which(ToothGrowth$dose==0.5),]
med <- ToothGrowth[which(ToothGrowth$dose==1),]
high <- ToothGrowth[which(ToothGrowth$dose==2),]

boxplot(low$len~low$supp, main = "Low dose", xlab = "Supplement Type")
boxplot(med$len~med$supp, main = "Medium dose", xlab = "Supplement Type")
boxplot(high$len~high$supp, main = "High dose", xlab = "Supplement Type")
```



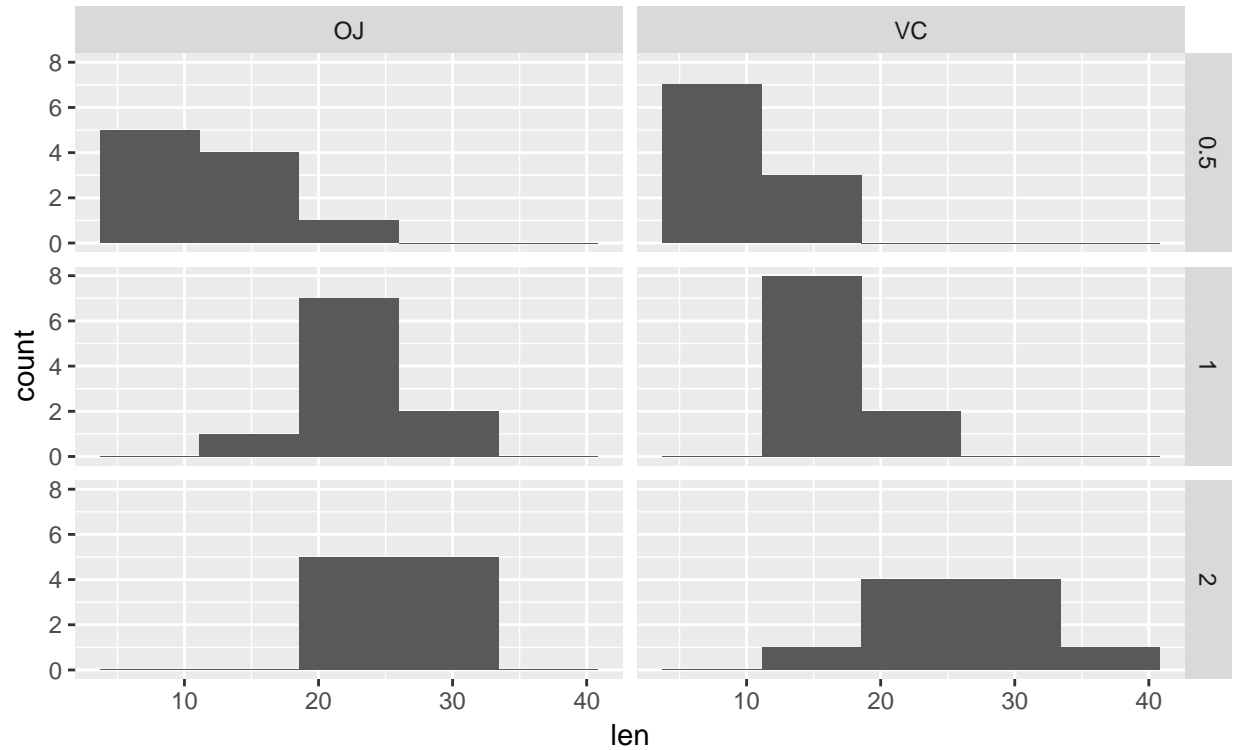
However, using ggplot we can simply map the dose column to a ggplot facet grid. The mapping of variables is not done via a `aes()`, but in formula form. Variables on the right indicate the columns in the grid and variables on the left indicate the rows of the grid. A “.” can be used if we want to use only one variable for facetting. This will also preserve the scales, making it more easy to compare the three plots. For fun, let’s add the actual datapoints in an additional geom on top of the boxplot.

```
ggplot(ToothGrowth, aes(x=supp,y=len)) +
  stat_boxplot(geom="errorbar", width=0.5) +
  geom_boxplot() + geom_point(color="red") +
  facet_grid(. ~ as.factor(dose))
```



We can do the same with two variables in a grid.

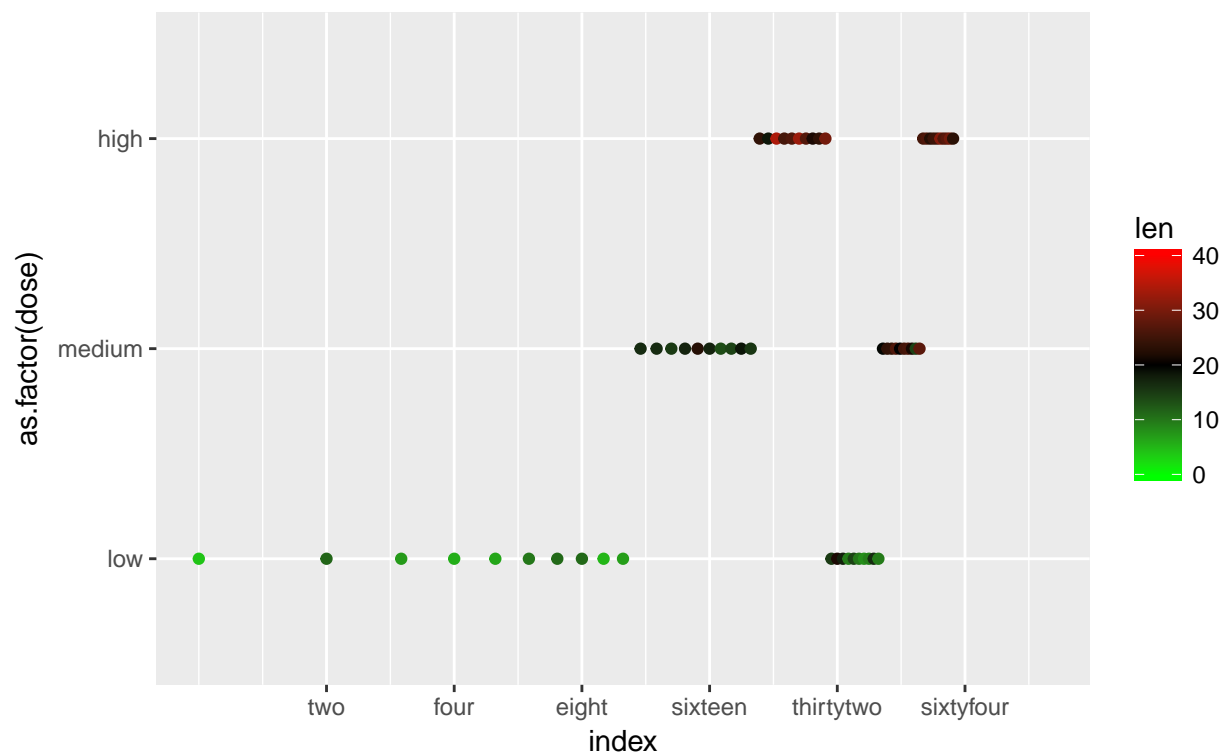
```
ggplot(ToothGrowth, aes(x=len)) +  
  geom_histogram(bins = 5) +  
  facet_grid(dose ~ as.factor(supp))
```



## 1.6 Scales

When mapping a variable to one of the aesthetics of a ggplot, ggplot uses the default scales. It is however possible to adjust the scales of the plot and their color using the scales functions. The scales that are generally used the most are those for the x, y, color and fill of the plot. The type of function to adjust the scale of one of the aesthetics depends on the aesthetic and the type of variable (continuous or discrete). Some of the possible scale adjustments, such as limits, transformation, breaks, labels, and color, are exemplified by the plot below.

```
ggplot(ToothGrowth, aes(x=index, y=as.factor(dose), color=len)) +  
  geom_point() +  
  scale_x_continuous(limits=c(1,100), trans = "log2",  
    breaks=c(2,4,8,16,32,64),  
    label=c("two","four","eight","sixteen","thirtytwo","sixtyfour")) +  
  scale_y_discrete(label=c("low","medium","high")) +  
  scale_color_gradient2(limits=c(0,40), low = "green", mid = "black", high = "red",  
    midpoint = 20)
```



Have a look at the scales chapter at <http://ggplot2.tidyverse.org/reference/> to explore the other possibilities!

## 1.7 Themes

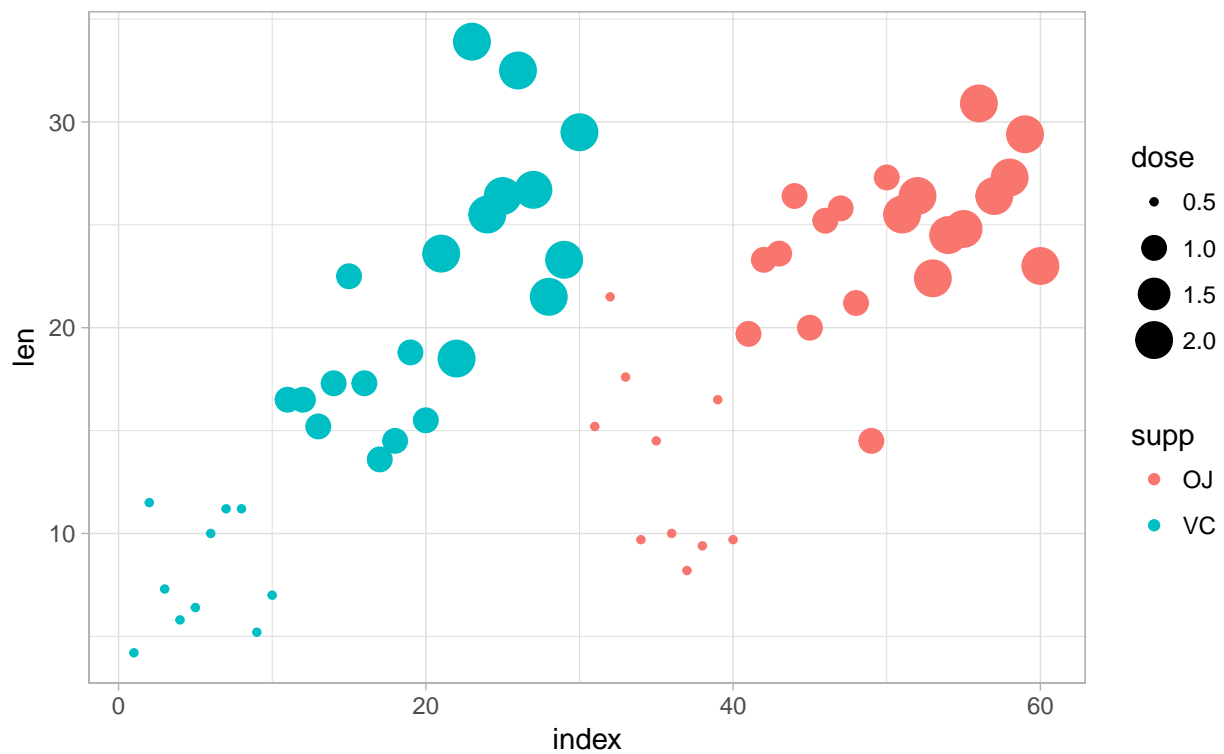
Until now we have simply mapped a variable from our dataset to an aesthetics parameter, potentially change the scales, and let ggplot handle the styling of the graph. However, ggplot has many options to customize the style of your plot using themes. There are several standard themes that we can use, however it is possible to create our own style from scratch.

The standard themes available in the ggplot package are:

- `theme_gray`
- `theme_bw`
- `theme_linedraw`
- `theme_light`
- `theme_dark`
- `theme_minimal`
- `theme_classic`
- `theme_void`

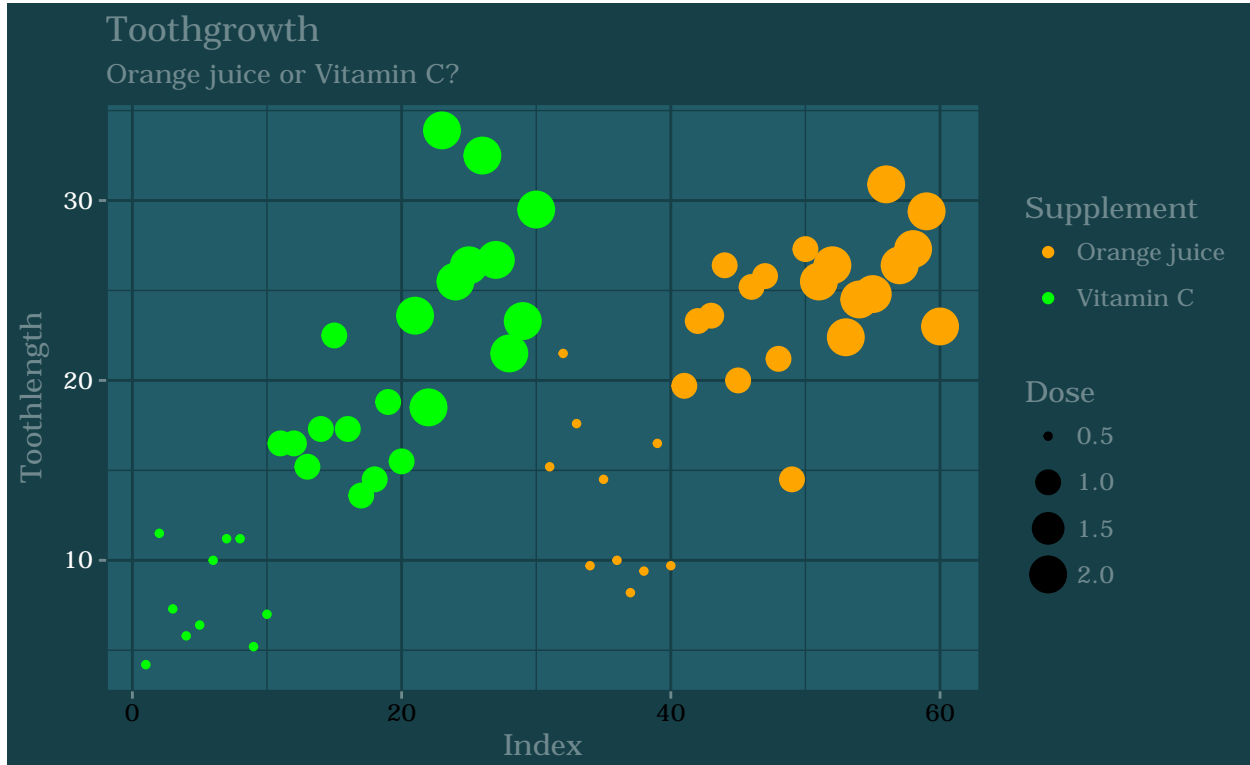
To add a standard theme we simply add the theme to the ggplot like we would with a geom.

```
p <- ggplot(ToothGrowth)
p + geom_point(aes(x = index, y = len, color = supp, size = dose)) + theme_light()
```



Custom themes can be made from scratch using the “theme()” function. To have a look at all customizable parameters visit <http://ggplot2.tidyverse.org/reference/theme.html>. Labels can be edited by adding “labs()” to the plot.

```
p <- ggplot(ToothGrowth)
p + geom_point(aes(x = index, y = len, color = supp, size = dose)) +
  theme(text = element_text(family = "Bookman", colour = "#6f898e"),
        line = element_line(color = "#163f47"),
        rect = element_rect(fill = "#163f47", color = "#163f47"),
        axis.text.x = element_text(color="black"),
        axis.text.y = element_text(color="white"),
        axis.ticks = element_line(color = "#6f898e"),
        axis.line = element_line(color = "#163f47", linetype = 1),
        legend.background = element_blank(),
        legend.key = element_blank(),
        panel.background = element_rect(fill = "#215c68", colour = "#163f47"),
        panel.border = element_blank(),
        panel.grid = element_line(color = "#163f47"),
        panel.grid.major = element_line(color = "#163f47"),
        panel.grid.minor = element_line(color = "#163f47"),
        plot.background = element_rect(fill = NULL, colour = NA, linetype = 0)
  ) +
  labs(title="Toothgrowth",
        subtitle = "Orange juice or Vitamin C?", x="Index", y="Toothlength",
        size="Dose", color="Supplement") +
  scale_color_manual(label=c("Orange juice", "Vitamin C"),
                     values = c("VC"="green", "OJ"="orange"))
```



## 1.8 Saving ggplots

We can save out plots by using the `ggsave()` function. `ggsave` takes two argument, the first argument being the ggplot and the second being the location where the file should be saved. If we do not specify the complete path, it will save the plot in our current working directory (`getwd()`). The type of the file depends on the file extension we use. Possible file extensions are “eps”, “ps”, “tex” (pictex), “pdf”, “jpeg”, “tiff”, “png”, “bmp”, “svg” or “wmf” (windows only).

```
p <- ggplot(ToothGrowth)

myplot <- p + geom_point(aes(x = index, y = len, color = supp, size = dose))

ggsave("my_plot.pdf",myplot)
```



## 2 Using ggplot: Practicals

1. Load the ggplot2 library (install it if you have to) and the diamonds dataset using `data()`
2. Explore the dataset using `dim()`, `str()` and `help()`, which variables are continuous, which variables are discrete? Is this dataset ready for plotting with ggplot?
3. Use ggplot to plot a scatterplot of the relationship between the diamonds' carat and their price
4. Make all dots darkblue and set the alpha value to 0.1
5. Visualize the influence of the color of a diamond on its price by mapping the diamond color to the color aesthetic
6. Use a ggplot barplot to visualize diamond clarity depending on color, map diamond color to x and diamond clarity to fill
7. Create a boxplot of the carat of a diamond based on its clarity and add whiskers using `stat_boxplot`
8. Add a `geom_point` layer to the previous plot mapping the diamonds price to the color
9. Create a histogram of the price of the diamonds and separate the histograms into facets using diamond color, choose a good binwidth or number of bins
10. Create a grid of facets of the same histogram by comparing both color and cut
11. Use `'aggregate(diamonds, by = list(cut = diamonds$cut, color = diamonds$color), mean)'` to calculate the mean of all variables by cut and color. Create a heatmap of the mean prices by cut and color using `geom_tile`
12. Change the title of the heatmap to "Average prices"
13. Change the gradient of the fill scale using `'scale_fill_gradient2'`. Have it go from darkblue to white to darkred, set the midpoint to 4500
14. Choose and add a theme to the heatmap, or create a theme yourself using the options listed at <http://ggplot2.tidyverse.org/reference/theme.html>

### 3 Using ggplot: Answers

1. Load the ggplot2 library (install it if you have to) and the diamonds dataset using data()

```
# if needed install.packages("ggplot2")
library(ggplot2)
data("diamonds")
```

2. Explore the dataset using dim(), str() and help(), which variables are continuous, which variables are discrete? Is this dataset ready for plotting with ggplot?

```
dim(diamonds)
str(diamonds)
help(diamonds)
```

Cut, color, and clarity are factors, and therefore discrete. The others are numeric continuous variables.

3. Use ggplot to plot a scatterplot of the relationship between the diamonds' carat and their price

```
ggplot() + geom_point(data=diamonds, aes(x=carat, y=price))

#or

ggplot(diamonds) + geom_point(aes(x=carat, y=price))

#or

ggplot(diamonds, aes(x=carat, y=price)) + geom_point()
```

4. Make all dots darkblue and set the alpha value to 0.1

```
ggplot(diamonds, aes(x=carat, y=price)) + geom_point(color="darkblue", alpha=0.1)
```

5. Visualize the influence of the color of a diamond on its price by mapping the diamond color to the color aesthetic

```
#The color of the dots will be overwritten if we specify it statically
#in the geom_point function itself
ggplot(diamonds, aes(x=carat, y=price, color=color)) + geom_point(alpha=0.1)
```

6. Use a ggplot barplot to visualize diamond clarity depending on color, map diamond color to x and diamond clarity to fill

```
ggplot(diamonds, aes(x=color, fill=clarity)) + geom_bar()
```

7. Create a boxplot of the carat of a diamond based on its clarity and add whiskers using stat\_boxplot

```
ggplot(diamonds, aes(x=clarity, y=carat)) +
  stat_boxplot(geom="errorbar", width=0.5) +
  geom_boxplot()
```

8. Add a geom\_point layer to the previous plot mapping the diamonds price to the color

```
ggplot(diamonds, aes(x=clarity, y=carat)) +
  stat_boxplot(geom="errorbar", width=0.5) +
  geom_boxplot() +
  geom_point(aes(color=price))
```

9. Create a histogram of the price of the diamonds and separate the histograms into facets using diamond color, choose a good binwidth or number of bins

```
ggplot(diamonds, aes(x=price)) +  
  geom_histogram(binwidth = 100) +  
  facet_grid(color ~ .)
```

10. Create a grid of facets of the same histogram by comparing both color and cut

```
ggplot(diamonds, aes(x=price)) +  
  geom_histogram(binwidth = 100) +  
  facet_grid(color ~ cut)
```

11. Use 'aggregate(diamonds, by = list(cut = diamonds\$cut, color = diamonds\$color), mean)' to calculate the mean of all variables by cut and color. Create a heatmap of the mean prices by cut and color using geom\_tile

```
#Aggregate uses a function (in this case mean) to aggregate all variables  
#in a given data.frame by a list of variables given in "by"  
mean.price <- aggregate(diamonds,  
  by = list(cut = diamonds$cut, color = diamonds$color),  
  mean)  
ggplot(mean.price, aes(x=cut, y=color, fill=price)) +  
  geom_tile()
```

12. Change the title of the heatmap to "Average prices"

```
ggplot(mean.price, aes(x=cut, y=color, fill=price)) +  
  geom_tile() +  
  labs(title="Average prices")
```

13. Change the gradient of the fill scale using 'scale\_fill\_gradient2'. Have it go from darkblue to white to darkred, set the midpoint to 4500

```
ggplot(mean.price, aes(x=cut, y=color, fill=price)) +  
  geom_tile() +  
  labs(title="Average prices") +  
  scale_fill_gradient2(low="darkblue", mid="white", high="darkred", midpoint = 4500)
```

14. Choose and add a theme to the heatmap, or create a theme yourself using the options listed at <http://ggplot2.tidyverse.org/reference/theme.html>

```
ggplot(mean.price, aes(x=cut, y=color, fill=price)) +  
  geom_tile() +  
  labs(title="Average prices") +  
  scale_fill_gradient2(low="darkblue", mid="white", high="darkred", midpoint = 4500) +  
  theme_minimal()
```