# Intro to R - Part I

## Installing RStudio

First, install the R framework. Go to the CRAN website and in the section *Download and Install R* choose a distribution for your operating system.

In order to install the RStudio, go to the download page of the RStudio website. In the section *Installers for Supported Platforms*, choose an installer for your operating system. Run the downloaded installer and complete the installation process.

## Data Types

### Vectors

A datum occurring by itself in an expression is taken as a vector of length one. Following are most used classes of vectors.

#### Numeric

```
x <- 41.5
x
```

```
## [1] 41.5
```

```
print(class(x))
```

```
## [1] "numeric"
```

#### Integer

```
x <- 5L
x
```

```
## [1] 5
```

```
print(class(x))
```

```
## [1] "integer"
```

#### Character

```
x <- "Hello"
x
```

```
## [1] "Hello"
```

```
print(class(x))
```

```
## [1] "character"
```

### Logical (true/false)

```
x <- TRUE
x
```

```
## [1] TRUE
```

```
print(class(x))
```

```
## [1] "logical"
```

When you want to create a vector with more than one element, use the *c()* function:

```
x <- c(1.4, 5.6, 3.1)
x
```

```
## [1] 1.4 5.6 3.1
```

## Vector arithmetic

An arithmetic operation on a vector results in a vector containing values of that operation applied to each element of the vector.

```
x + 1
```

```
## [1] 2.4 6.6 4.1
```

```
-x
```

```
## [1] -1.4 -5.6 -3.1
```

```
2*x + 2
```

```
## [1]  4.8 13.2  8.2
```

```
x/2
```

```
## [1] 0.70 2.80 1.55
```

```
x^2
```

```
## [1]  1.96 31.36  9.61
```

Common arithmetic functions are available: *log, exp, sin, cos, tan, sqrt, …*

*max* and *min* select the largest and smallest elements of a vector respectively. *range* is a function that returns a vector containing the minimum and maximum of all the given arguments, namely *c(min(x), max(x))*. *length(x)* is the number of elements in x, *sum(x)* gives the total of the elements in *x*, and *prod(x)* their product [2].

```
x1 <- c(1, 2, 3, 4, 5)
```

```r
# get the max value
max(x1)
```

```
## [1] 5
```

```r
# get the min value
min(x1)
```

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

```r
# get the range of values
range(x1)
```

```
## [1] 1 5
```

```r
# calculate the sum of all elements
sum(x1)
```

```
## [1] 15
```

```r
# calculate the product of all elements
prod(x1)
```

```
## [1] 120
```

## Generating regular sequences

```r
# generate a sequence from 1 to 10
1:10
```

```
##  [1]  1  2  3  4  5  6  7  8  9 10
```

```r
# generate a sequence from 10 to 1
10:1
```

```
##  [1] 10  9  8  7  6  5  4  3  2  1
```

```r
# generate a sequence from 3.2 to 4.7, with a step 0.2
seq(3.2, 4.7, by = 0.2)
```

```
## [1] 3.2 3.4 3.6 3.8 4.0 4.2 4.4 4.6
```

## Factor values

Factors are used to represent categorical variables.

```r
x2 <- c("cold", "mild", "mild", "hot", "cold")
x2
```

```
## [1] "cold" "mild" "mild" "hot"  "cold"
```

```r
# convert x2 to a factor variable
temp <- factor(x2)
```

```r
# print all levels
levels(temp)
```

```
## [1] "cold" "hot"  "mild"
```

```r
# print the summary
summary(temp)
```

```
## cold  hot mild
##    2    1    2
```

## Missing values

When an element or value is "not available" or a "missing value" in the statistical sense, a place within a vector may be reserved for it by assigning it the special value *NA*. The function is.na(x) gives a logical vector of the same size as x with value TRUE if and only if the corresponding element in x is NA. [2].

```r
x3 <- c(1, NA, 3, 4, NA)
```

```r
# check which values are NAs
is.na(x3)
```

```
## [1] FALSE  TRUE FALSE FALSE  TRUE
```

A second kind of "missing" values which are produced by numerical computation, the so-called Not a Number, *NaN*, values [2].

```r
0/0
```

```
## [1] NaN
```

*is.na(xx)* is TRUE both for NA and NaN values. *is.nan(xx)* is only TRUE for NaNs.


## Data frame

Data frames are matrix-like structures, where the columns can be of different types. It consists of a list of vectors of equal length. It is used for storing data tables.

```r
# create a data frame from vectors
weather <- data.frame(x1, x2, x3)
weather
```

```
##   x1   x2 x3
## 1  1 cold  1
## 2  2 mild NA
## 3  3 mild  3
## 4  4  hot  4
## 5  5 cold NA
```

## Loading existing data

Data can be loaded from files. The most used file formats are CSV (Comma Separated Values) and TSV (Tab Separated Values). In order to read a CSV file *The Beatles songs dataset, v0.csv*, we can use the *read.csv* function. By setting the argument *stringsAsFactors* to FALSE, strings will not be converted to factors (which is the default setting).

### read.csv

```
# reading a data frame from the CSV file
the.beatles.songs <- read.csv("The Beatles songs dataset, v0.csv", stringsAsF
actors = FALSE)
the.beatles.songs
```

```
##                       Title Year Duration
## 1            12-Bar Original 1965      174
## 2           A Day in the Life 1967      335
## 3           A Hard Day's Night 1964      152
## 4 A Shot of Rhythm and Blues 1963      104
## 5            A Taste of Honey 1963      163
## 6          Across the Universe 1968      230
## 7              Act Naturally 1965      139
## 8             Ain't She Sweet 1961      150
## 9           All I've Got to Do 1963      124
```

## Inspecting data frame

Print the number of rows.

```
nrow(the.beatles.songs)
```

```
## [1] 9
```

Print the number of columns

```
ncol(the.beatles.songs)
```

```
## [1] 3
```

### summary

Used to produce the summary of data.

```
summary(the.beatles.songs)
```

```
##     Title              Year          Duration
## Length:9          Min.   :1961   Min.   :104.0
## Class :character  1st Qu.:1963   1st Qu.:139.0
## Mode  :character  Median :1964   Median :152.0
##                   Mean   :1964   Mean   :174.6
##                   3rd Qu.:1965   3rd Qu.:174.0
##                   Max.   :1968   Max.   :335.0
```

## str

An alternative to *summary* (to an extent), the function *str* can also be used to inspect the data and compactly display the internal structure of a dataframe.

```
str(the.beatles.songs)

## 'data.frame':    9 obs. of  3 variables:
##  $ Title   : chr  "12-Bar Original" "A Day in the Life" "A Hard Day's Nigh
t" "A Shot of Rhythm and Blues" ...
##  $ Year    : int  1965 1967 1964 1963 1963 1968 1965 1961 1963
##  $ Duration: int  174 335 152 104 163 230 139 150 124
```

## head and tail

Prints an example of data inside a data frame, first several items.

```
# get first several rows
head(the.beatles.songs)

##                          Title Year Duration
## 1            12-Bar Original 1965      174
## 2          A Day in the Life 1967      335
## 3          A Hard Day's Night 1964     152
## 4 A Shot of Rhythm and Blues 1963      104
## 5            A Taste of Honey 1963     163
## 6          Across the Universe 1968    230

# get last several rows
tail(the.beatles.songs)

##                          Title Year Duration
## 4 A Shot of Rhythm and Blues 1963      104
## 5            A Taste of Honey 1963     163
## 6          Across the Universe 1968    230
## 7              Act Naturally 1965      139
## 8              Ain't She Sweet 1961    150
## 9            All I've Got to Do 1963   124
```

## names (colnames)

*names()* function prints column names. *colnames* can also be used as it works the same on data frames (but not on other data types!).

```
# get column names
names(the.beatles.songs)

## [1] "Title"     "Year"      "Duration"
```

We can similarly change column names.

```r
# make a copy of the original data frame
the.beatles.songs1 <- the.beatles.songs

# update column names
names(the.beatles.songs1) <- c("song_name", "release_year", "duration")
the.beatles.songs1

##                   song_name release_year duration
## 1          12-Bar Original         1965      174
## 2          A Day in the Life        1967      335
## 3          A Hard Day's Night       1964      152
## 4 A Shot of Rhythm and Blues       1963      104
## 5          A Taste of Honey         1963      163
## 6         Across the Universe       1968      230
## 7             Act Naturally         1965      139
## 8             Ain't She Sweet        1961      150
## 9          All I've Got to Do       1963      124
```

### Removing a column

```r
# remove column duration
the.beatles.songs1$duration <- NULL
the.beatles.songs1

##                   song_name release_year
## 1          12-Bar Original         1965
## 2          A Day in the Life        1967
## 3          A Hard Day's Night       1964
## 4 A Shot of Rhythm and Blues       1963
## 5          A Taste of Honey         1963
## 6         Across the Universe       1968
## 7             Act Naturally         1965
## 8             Ain't She Sweet        1961
## 9          All I've Got to Do       1963
```

## Subsetting

Print a specific element from a data frame. Indexes start from 1.

```r
# get an element from the row 3, column 1
song <- the.beatles.songs[3, 1]
song

## [1] "A Hard Day's Night"
```

Retrieve a specific row or several rows by index.

```r
# get the third row
beatles.subset <- the.beatles.songs[3,]
beatles.subset
```

```
##                Title Year Duration
## 3 A Hard Day's Night 1964      152
```

```r
# get rows at positions from 3 to 6
beatles.subset1 <- the.beatles.songs[3:5,]
beatles.subset1
```

```
##                          Title Year Duration
## 3          A Hard Day's Night 1964      152
## 4 A Shot of Rhythm and Blues 1963      104
## 5             A Taste of Honey 1963      163
```

```r
# get rows at positions 3 and 6
beatles.subset2 <- the.beatles.songs[c(3,6),]
beatles.subset2
```

```
##                   Title Year Duration
## 3  A Hard Day's Night 1964      152
## 6 Across the Universe 1968      230
```

Retrieve a specific column by index or by name.

```r
# get the second column
years <- the.beatles.songs[,2]
years
```

```
## [1] 1965 1967 1964 1963 1963 1968 1965 1961 1963
```

```r
# get the Year column
years <- the.beatles.songs$Year
years
```

```
## [1] 1965 1967 1964 1963 1963 1968 1965 1961 1963
```

Lastly, we can retrieve rows with a logical index vector.

```r
# retrieve all songs released before year 1965
songsBefore1965 <- the.beatles.songs[the.beatles.songs$Year < 1965,]
songsBefore1965
```

```
##                          Title Year Duration
## 3          A Hard Day's Night 1964      152
## 4 A Shot of Rhythm and Blues 1963      104
## 5             A Taste of Honey 1963      163
## 8              Ain't She Sweet 1961      150
## 9          All I've Got to Do 1963      124
```

```r
# retrieve all songs released before year 1965 with duration lower than 150 seconds
shortSongsBefore1965 <- the.beatles.songs[the.beatles.songs$Year < 1965 & the.beatles.songs$Duration < 150,]
shortSongsBefore1965
```

```
##                        Title Year Duration
## 4 A Shot of Rhythm and Blues 1963      104
## 9          All I've Got to Do 1963      124
```

## Task 1

Print a number of songs from 1963 that last more than 2 minutes (120 seconds).
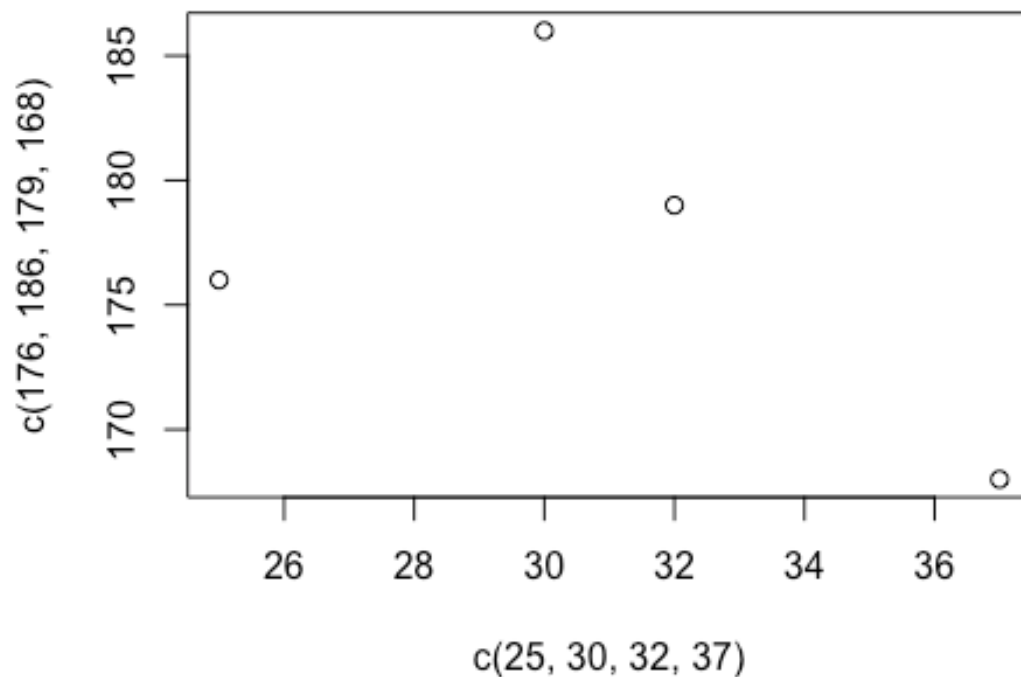
Answer:

```
nrow(the.beatles.songs[the.beatles.songs$Year == 1965 & the.beatles.songs$Dur
ation > 120,])
```
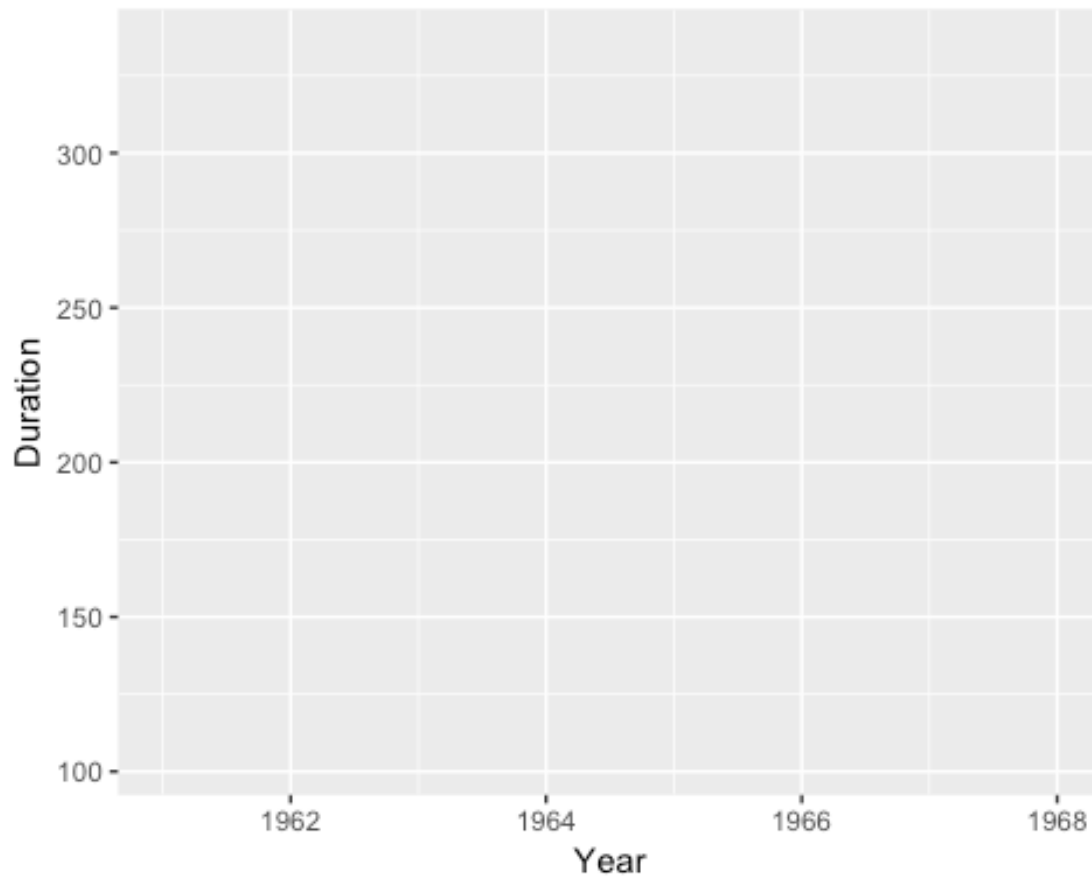
```
## [1] 2
```

## Plotting

The most basic function to create a plot is plot(x, y), where x and y are numeric vectors denoting the x and y axes.

```
# create a plot for the given vectors
plot(c(25,30,32,37), c(176,186,179,168))
```

For more advanced plots, *ggplot2* library is recommended. *ggplot* function works with data frames and not individual vectors. All information that is part of the dataframe has to be specified inside the *aes()* function that specifies x and y axes.

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

# render a plot for the given data frame
ggplot(the.beatles.songs, aes(x=Year, y=Duration))
```
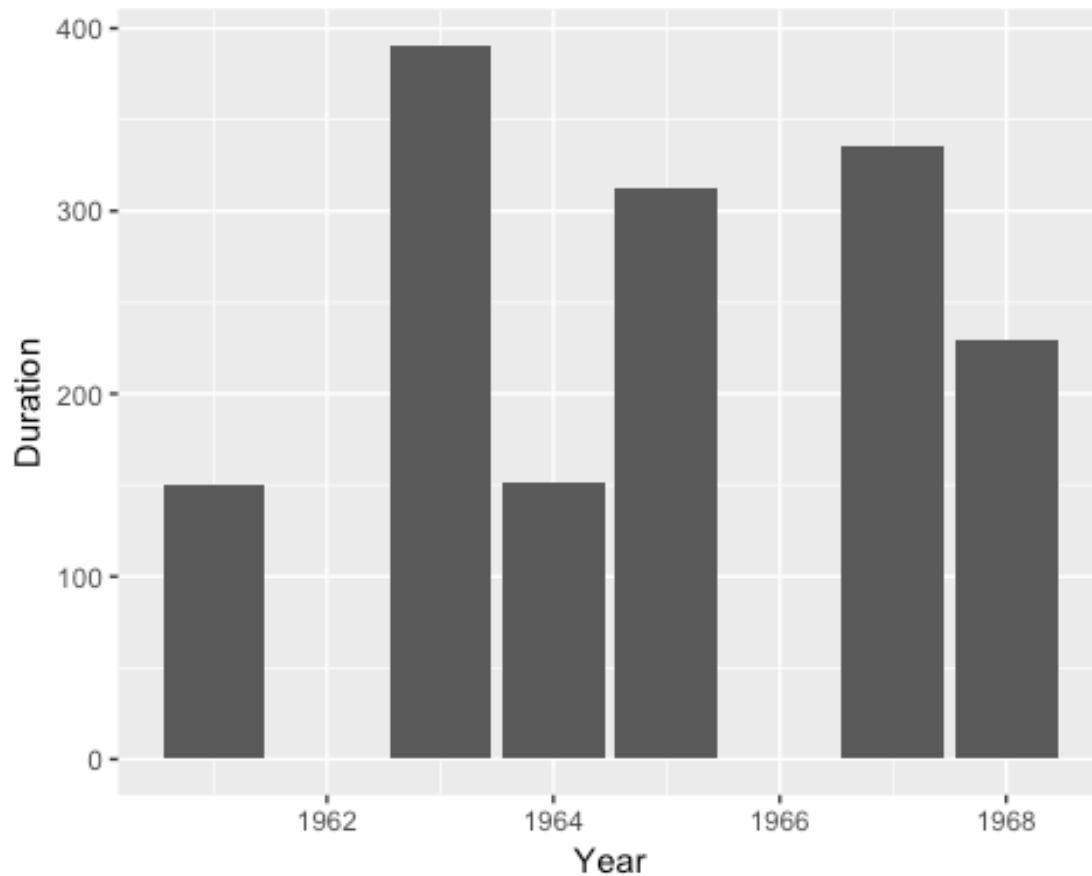


A blank plot is drawn. Even though the x and y are specified, there are no points or lines in it. This is because, ggplot doesn't assume which plot we want, a scatterplot or a line chart. We have only told ggplot what dataset to use and what columns should be used for x- and y-axis. We haven't explicitly asked it to draw any points. Plots are drawn by adding layers to the basic plot generated by the *ggplot* function. More specificly, in order to generate a scatter plot, we us the *geom_point()* layer.

```
# render a plot for the given data frame with points
ggplot(the.beatles.songs, aes(x=Year, y=Duration)) + geom_point()
```

Similarly, we can plot a bar chart by adding a *geom_col()* layer.

```
# render a bar chart
ggplot(the.beatles.songs, aes(x=Year, y=Duration)) +
  geom_col()
```

In the previous bar chart, we can observe that the x-axis has all values from 1961 to 1968, which is the value interval of the *Year* attribute. These include years 1962 and 1966 that with no songs. If we want to omit these two values, we can convert the *Year* variable to be a factor, and now the *Year* variable will have the possible values: 1961, 1963, 1964, 1965, 1967, and 1968. Only these values will be plotted.

```
# convert the Year variable to factor
the.beatles.songs$Year <- factor(the.beatles.songs$Year)

# render a bar chart
ggplot(the.beatles.songs, aes(x=Year, y=Duration)) +
  geom_col()
```
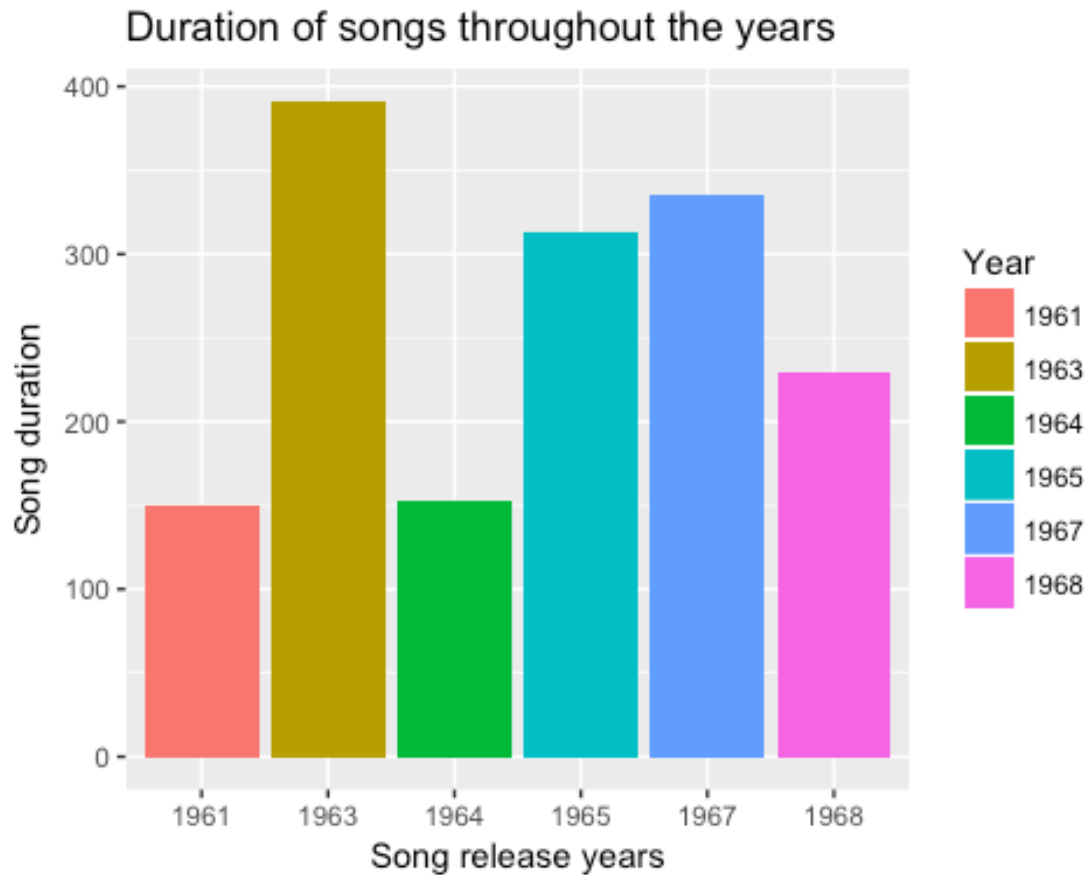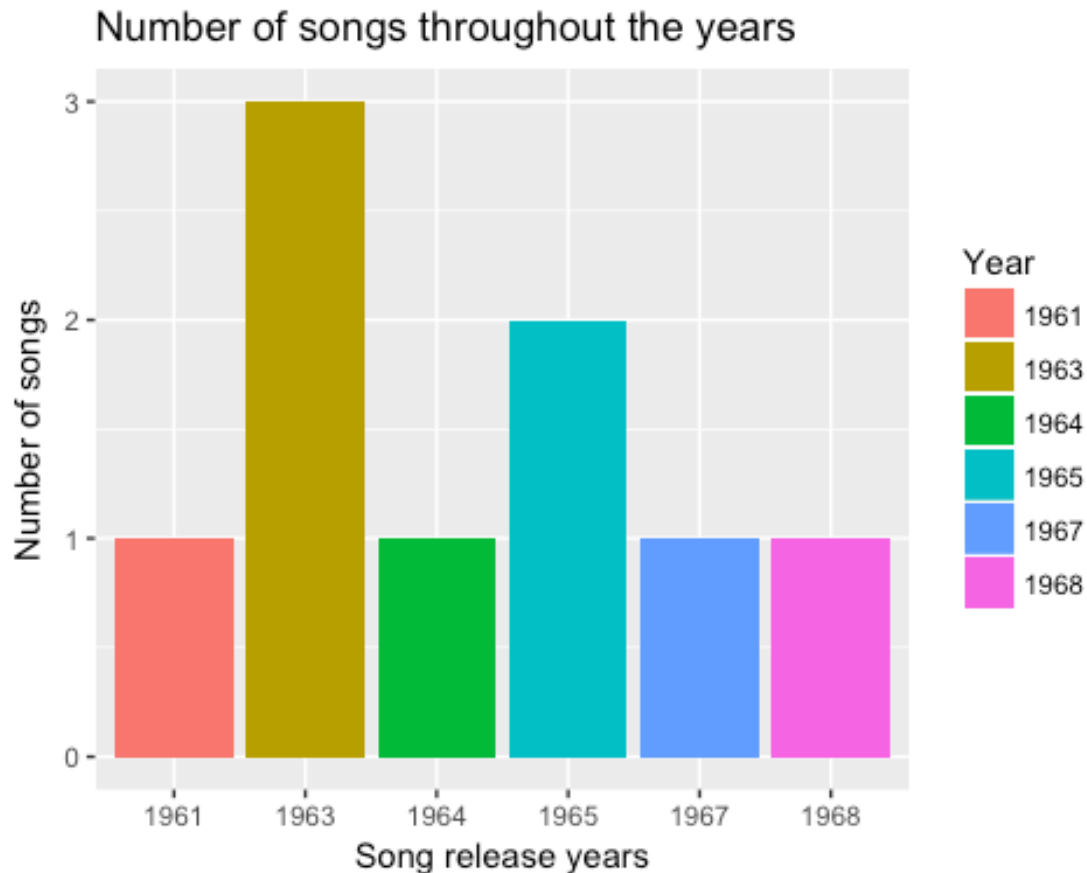
Chart can be further enhanced by adding a custom x- and y-axis labels, and a chart title. The aesthetic fill will take different colouring scales by setting the *fill* to be equal to a factor variable.

```
# render a bar chart with custom title and axes labels
ggplot(the.beatles.songs, aes(x=Year, y=Duration, fill = Year)) +
  geom_col() +
  xlab("Song release years") + ylab("Song duration") +
  ggtitle("Duration of songs throughout the years")
```

Duration of songs throughout the years

*geom_bar()* makes the height of the bar proportional to the number of cases in each group. In this case, we provide only one variable for x-axis.
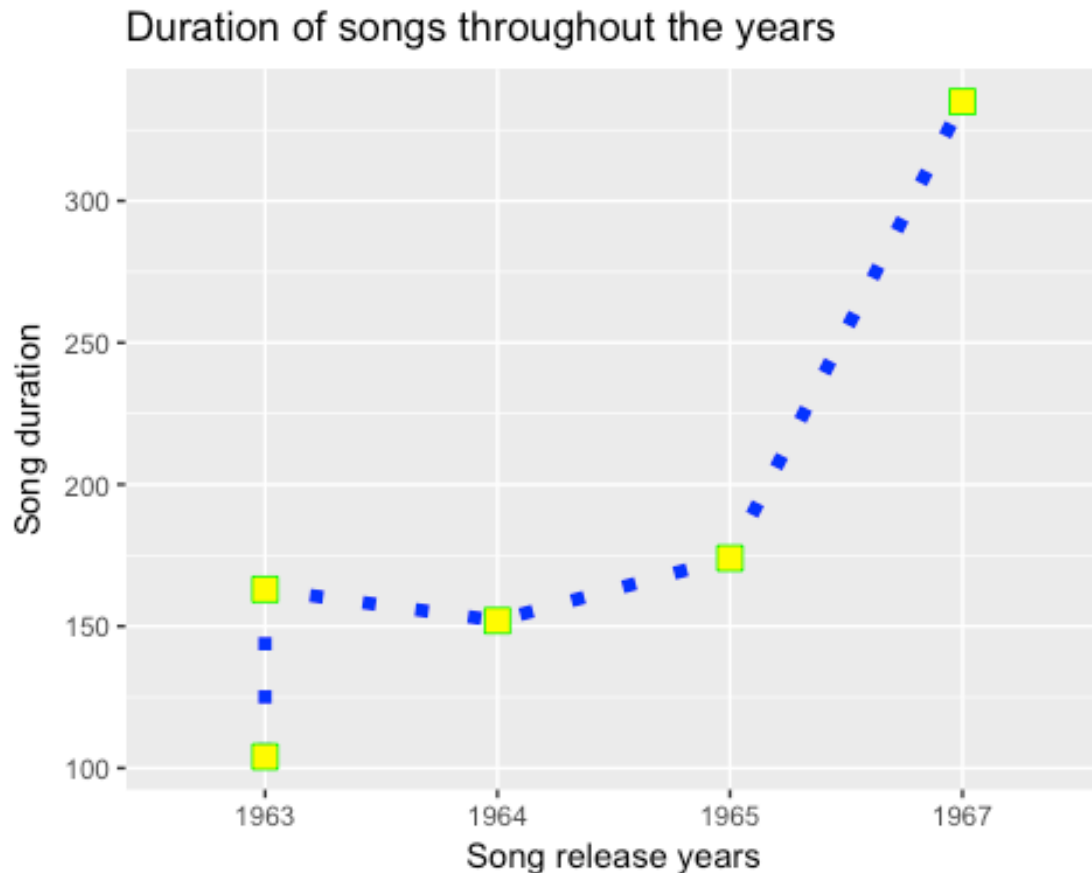
```
# render a bar chart where the y-axis displays number of cases for each value
on the x-axis
ggplot(the.beatles.songs, aes(x=Year, fill=Year)) +
  geom_bar() +
  xlab("Song release years") + ylab("Number of songs") +
  ggtitle("Number of songs throughout the years")
```

## Number of songs throughout the years



Line chart can be added by adding the *geom_line* layer. *geom_line()* tries to connect data points that belong to the same group. Different levels of a factor variable belong to different groups. By specifying *group=1* we indicate we want a single line connecting all the points.

```
# render a line chart for the first five songs with specific line and ponts p
roperties
ggplot(the.beatles.songs[1:5,], aes(x=Year, y=Duration, group = 1)) +
  geom_line(colour = "blue", linetype = "dotted", size = 2) +
  geom_point(colour="green", size = 4, shape = 22, fill = "yellow") +
  xlab("Song release years") + ylab("Song duration") +
  ggtitle("Duration of songs throughout the years")
```

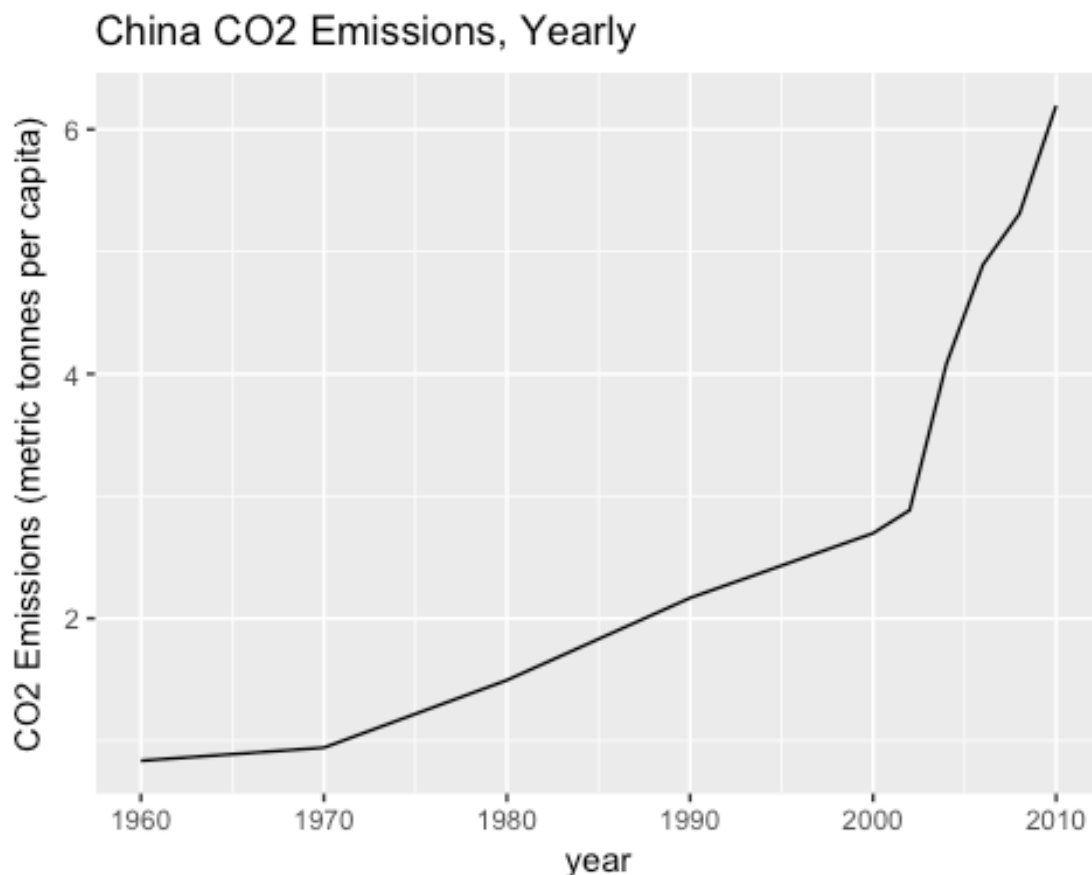Duration of songs throughout the years

## Task 2

Create a line chart with x-axis containing the following years: 1960, 1970, 1980, 1990, 2000, 2002, 2004, 2006, 2008, 2010. The y-axis should have the following values of CO2 emissions: 0.836046900792028, 0.942934534989582, 1.49525074931082, 2.16770307659104, 2.69686243322549, 2.88522504139331, 4.08013890554173, 4.89272709798477, 5.31115185538876, 6.19485757472686. Chart title should say "China CO2 Emissions, Yearly" and y-axis should have a label "CO2 Emissions (metric tonnes per capita)".

Answer:

```
co2.emissions <- data.frame(
  year = c(1960, 1970, 1980, 1990, 2000, 2002, 2004, 2006, 2008, 2010),
  emission = c(0.836046900792028, 0.942934534989582, 1.49525074931082, 2.1677
0307659104, 2.69686243322549, 2.88522504139331, 4.08013890554173,  4.89272709
798477, 5.31115185538876, 6.19485757472686)
)

ggplot(co2.emissions, aes(x = year, y = emission, group = 1)) +
  geom_line() +
  ylab("CO2 Emissions (metric tonnes per capita)") +
  ggtitle("China CO2 Emissions, Yearly")
```

China CO2 Emissions, Yearly

## Homework - Complete interactive R tutorials with Swirl

Swirl is an interactive R tutorial that teaches you R from the R console. All you need to do is install Swirl package for R and issue a *swirl()* command which will start the tutorial.

```
#install.packages("swirl")
library("swirl")
swirl()
```

Once a Swirl session is started, you will be prompted with an option to install a Swirl course. For our course, you need to install the *R Programming: The basics of programming in R* Swirl course and go through the following tutorials:

1. Basic Building Blocks
2. Workspace and Files
3. Sequences of Numbers
4. Vectors
5. Missing Values
6. Subsetting Vectors
7. Matrices and Data Frames

8. Logic
9. Functions
10. Looking at Data

## References

[1] https://www.tutorialspoint.com/r/r_data_types.htm [2] An Introduction to R, https://cran.r-project.org/doc/manuals/R-intro.pdf