

## What is Analytics?

The science of using data to build models that lead to better decisions that add value to individuals, companies and institutions.

(Source: Prof. Dimitris Bertsimas, MIT)  
Prof. Thomas Davenport, Babson College

## The Analytics Edge

- 1) Analytics provides a competitive edge to individuals, companies and institutions.
- 2) Analytics is often critical to the success of a company and this is increasingly becoming the norm.

Analytics is often classified into:

- 1) Descriptive analytics
- 2) Predictive analytics
- 3) Prescriptive analytics

Source:

Davenport & Harris.  
Competing on  
Analytics

Companies such as IBM traditionally a PC company has now invested over \$20 billion since 2005 in growing its analytics business.

One of the main reasons for the advanced interest in analytics is that data is often unstructured and comes in various forms.

However such data is becoming increasingly available and one tries to develop models that can use such data.

Data might be incomplete and does not come with labels and not directly given as regression or optimization problems.

Analytics is often used in conjunction with Big Data. While more data is better, we can also do analytics with small datasets. In this course, we will work with both small and large datasets.

Predictive analytics is also often used to make predictions at a personal (individual) level.

This brings powerful insights but also has privacy issues that need to be dealt with for any organization. While we will not deal with privacy issues in this class, it is important in practice.

## Watson & Jeopardy

Jeopardy is a popular TV quiz show launched in 1964. The challenge is that the quiz needs Watson to have Natural Language Processing that deals with understanding the human language. The quiz is specifically designed so that the clues are meant to be hard to obtain meaning from. Furthermore the answer must be given in around 5 seconds if the contestant wants to successfully buzz in.

Watson used a fundamental measure accuracy that is the fraction of clues for which it was successfully confident of buzzing in and precision the fraction of times the response was correct.

Watson generated a set of candidate answers using a wide range of encyclopedias, dictionaries, news articles. From identifying relevant articles. For each answer it generated a confidence level to determine the most likely correct ones. By choosing when to buzz using a threshold for the confidence level, it could play more defensively or aggressively. All this needed huge computational memory & power. The key to success was to exploit the strength of a computer along with massive parallelization to play a game against the top human. Ability of Watson to work with unstructured ~~information~~ information & answer questions gives it lots of new possibilities.

Ensemble methods are an important part of making predictions. Instead of using a single model to make predictions, one can use multiple models to make predictions (ensemble) and then use a majority vote or an average to make better predictions. One has to tradeoff between accuracy and interpretability often in such models.

Watson made extensive use of ensemble methods in predicting the accuracy of answers.

# An Introduction to R

[www.r-project.org](http://www.r-project.org)

## Why R?

- 1) R is free and open.
- 2) R is an integrated suite of software facilities for data manipulation, calculation and graphical display.
- 3) R provides an environment within which many statistical techniques are implemented and this can be extended easily via packages.
- 4) R has an extensive online support and discussion forum.

## Origins of R (First appeared: 1993)

R is an implementation of the S language developed at Bell Laboratories (formerly AT&T, now Lucent) by John Chambers and colleagues. R was created by Ross Ihaka and Robert Gentleman at the University of Auckland and currently developed by the R Development Core Team of which Chambers is a member.

# An Introduction To R Commands

## Preliminaries

getwd( )

Returns the current working directory

help(getwd)

Provides help on a specific function

?getwd

Alternative to get help on a function

setwd("C:/Users/earhnik/Desktop/Analytics")      Sets the working directory

dir( )

Lists all files in directory

ls( )

Lists (Display) objects stored within R

$x \leftarrow 50$

Assign the number 50 to x

ls( )

Display objects stored within R

$50 \rightarrow x$

Alternate way of assigning 50 to x

$x = 50$

Alternate way of assigning 50 to x

$50 * x$

Multiply x with 50

$x^2$

Square the value in x

$\exp(x)$

Compute  $e^x$

$1/x$

Divide 1 by x

$y \leftarrow 1/x$

Assigns y the value of  $1/x$   
which is 0.02 here

rm(y)

Removes the variable y

## Numbers & vectors

$$x \leftarrow c(1, -3, 7, 9.4, \pi)$$

Concatenates (combines) numbers to form a vector

$$y \leftarrow c(4) \text{ is equivalent to } y \leftarrow 4$$

$$1/x$$

Reciprocal of the numbers are computed in the vector

$$y \leftarrow c(x, 0, x)$$

Combines the vector  $x$ , 0 and  $x$  to create a vector of length 11

$$\exp(y)$$

Applies the exponentiation operation to each element of  $y$

$$y + x$$

Creates a new vector of length 11 by recycling the shorter vector (even fractionally) till it matches longer vector

$$\max(x)$$

Find maximum and minimum elements in  $x$

$$\min(x)$$

$$\text{which.max}(x)$$

Determines location (index) of first max

$$\text{which}(x == \max(x))$$

Determines all locations where max is obtained

$$\text{sum}(x)$$

Sum and product of entries in  $x$

$$\text{prod}(x)$$

`mean(x)`

Mean of elements in  $x$

`var(x)`

Variance of elements in  $x$

Equivalent to :

$$\text{sum}((x - \text{mean}(x))^2) / (\text{length}(x) - 1)$$

`max(x, y)`

Returns maximum entry in the vector with entries from  $x$  and  $y$

`pmax(x, y)`

Parallel maximum returns a vector (of length equal to their longest argument) that contains in each element the largest element in that position of any vectors in input

`summary(x)`

For a vector, `summary` provides a six number summary including the min, max, mean, 1st quartile, median and 3rd quartile.

This function can also be used with other objects as we will see later

`rm(list = ls())`

Removes all variables from the workspace



In general it is preferred to use  $\leftarrow$  for assignment instead of  $=$ .

The difference for example is clear when you use it inside a function

```
> exp(x = 1:10)
```

```
> x
```

Object x is not found

```
> exp(x ← 1:10)
```

```
> x
```

1 2 3 ... 10  
⏟  
Here x is defined in  
the user workspace

Note  $x \leftarrow 5$  is equivalent to  $5 \rightarrow x$

However  $x = 5$  is not equivalent to  $5 = x$

Note:

$x \leftarrow 3$

(Assigns 3 to x)

$x < -3$

(Checks if x is less than -3)

$-4:4$  Generates the vector  $c(-4, -3, -2, \dots, 4)$   
 $\text{seq}(-4, 4, \text{by}=0.2)$  Generates the vector  $c(-4, -3.8, \dots, 4)$   
 $\text{rep}(x, \text{times}=2)$  Puts two copies of  $x$  end to end  
 $\text{rep}(x, \text{each}=2)$  Repeats each element of  $x$  twice before moving on

$x > 1$  Returns a logical vector with values TRUE & FALSE by comparing each element of  $x$  with 1.

$\text{is.na}(x)$  Returns a logical vector with values TRUE & FALSE where TRUE is given if the element is not available or a missing value (NA).

$\text{is.na}(c(1:3, \text{NA}))$  Returns the vector FALSE FALSE FALSE TRUE

$\text{mean}(c(1:3, \text{NA}), \text{na.rm}=\text{TRUE})$  Returns mean after removing NA value

$c("a", "b")$  Returns a character vector

$x[4]$  Returns the fourth element of  $x$

$x[c(1, 5)]$  Returns the first and fifth elements of  $x$

$x[! \text{is.na}(x)]$  Returns a vector containing the non-missing values of  $x$  in the same order

$x[4] \leftarrow 10$  Sets the fourth element of  $x$  to 10

$Z \leftarrow 0:9$

`as.integer(as.character(z))` Returns back the integer vector 0 1 2 ... 9

`class(x)`

Name of the class of the object  $x$  such as numeric, integer, character, logical, list, matrix, data.frame

`as.character(x)`

Returns  $x$  as a vector of characters

`as.logical(x)`

Returns  $x$  as a vector of TRUE & FALSE terms where it is FALSE if entry is 0.

`C(1, "a")`

Returns the vector "1" "a" by overriding the class of the vector to character

`C(T, F, T)`

Returns a logical vector with TRUE, FALSE, TRUE entries

`T`

Global variable (logical) whose initial value is set to TRUE

## Factors

```
x <- c("yes", "no", "yes", "maybe", "maybe",  
       "no", "no", "no")
```

class(x)

Character

factor(x) → y

Transforms a vector into  
a factor (categorical variable)

levels(y)

Returns levels of a factor  
Here yes, no, maybe

summary(y)

Provides a summary of  
the factor vector

table(x)

Builds a table from vector

```
income <- c(500, 1200, 4000, 2300, 2300,  
            1234, 5439, 432, 4555)
```

tapply(income, x, mean)

Computes mean of the  
subvectors in income  
from the factors in x.

maybe	No	Yes
3346.3	1855.25	2250

## Matrices, arrays

$m \leftarrow \text{matrix}(c(3, 4, 5, 6, 7, 8), \text{ncol} = 2)$

Creates a matrix of size  $3 \times 2$

by filling entries by column

by default

$\text{dim}(m)$

Returns dimension of matrix  $m$

$m[4]$

Returns fourth entry in matrix counting by columns

$m[1, 2]$

Returns first row, second column entry in matrix

$m[1, ]$

Returns first row in matrix

$\text{class}(m)$

Returns matrix as class here

$z \leftarrow \text{array}(c(3, 4, 5, 6, 7, 8), c(3, 2))$

Creates an array of size  $3 \times 2$  as the matrix

Arrays can have more than 2 dimensions unlike a matrix

$z \leftarrow 1:50$

Treats  $z$  as an array of

$\text{dim}(z) \leftarrow c(5, 2, 5)$

dimension  $5 \times 2 \times 5$

$z[1, 1]$

Returns value 1

$z[5, 2, 5]$

Returns value 50

$z[5, 2, 1:5]$

Returns vector 10 20 30 40 50

$\text{diag}(10)$

Returns a diagonal matrix of size  $10 \times 10$  with 1 on diagonal

$\text{cbind}(c(1, 2, 3), c(4, 5, 6))$

Returns matrix

$\begin{matrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{matrix}$

$\text{rbind}(c(1, 2, 3), c(4, 5, 6))$

Returns matrix

$\begin{matrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{matrix}$

$Z \leftarrow \text{matrix}(c(5, 7, 9, 6, 3, 4), \text{nrow} = 3, \text{ncol} = 2)$

$Y \leftarrow \text{matrix}(c(1, 3, 0, 9, 5, -1), \text{nrow} = 3, \text{ncol} = 2)$

$Z * Y$

Returns a matrix with  
componentwise multiplication

$t(Y)$

Returns transpose of matrix  $Y$

$Z \% * \% t(Y)$

Returns a  $3 \times 3$  matrix  
using matrix multiplication  $ZY'$

$\text{help}(\"% * \%")$

$a \leftarrow \text{array}(c(2, 1, -1, 2), c(2, 2))$

Solves linear equations

$b \leftarrow c(4, 4)$

$$2x - y = 4$$

$$x + 2y = 4$$

$\text{solve}(a, b)$

to give  $x = 2.4, y = 0.8$

$\text{solve}(a)$

Returns inverse of  
matrix  $a$

$t\text{eigen}(a)$

Returns eigenvectors and  
eigenvalues of  $a$  with  
 $t$  & values and  $t$  & vectors

## Lists & data frames

A list is an object consisting of an ordered collection of objects that can be of different or same type

```
karthik <- list(age = 37, sex = "M", child.ages = c(1, 1))
```

```
kim <- list(age = 24, sex = "F", child.ages = NA)
```

Creates two lists for karthik & kim respectively

```
karthik[[1]]
```

Returns number 37

```
class(karthik)
```

list

```
karthik[[2]]
```

Returns character "M"

```
karthik$age
```

Returns number 37

```
karthik$child.ages
```

Returns vector 1 1

```
karthik$child.ages[1]
```

Returns number 1

```
karthik[1]
```

Returns sublist with age and 37

```
class(karthik[1])
```

Returns "list"

```
class(karthik[[1]])
```

Returns "numeric"

```
t <- c(karthik, kim)
```

Returns a list with all the objects concatenated into one large vector & loses the dimensions

Data frames are a tightly coupled collection of variables that share many of the properties of matrices & lists, & is often the fundamental data structure used by R modeling software.

A data frame is a list of variables/vectors of the same length

```
t <- data.frame(names = c("Karthik", "Sam", "Jim"),  
               ages = c(36, 34, 40),  
               sex = c("M", "F", "M"),  
               children = c(2, 0, 1))
```

Creates a data frame where columns  
can be viewed as attributes and rows  
can be viewed as observations

Often the data will be read from files  
using the `read.csv` command or `read.table` command

```
t$names                Lists out "Karthik", "Sam", "Jim"
```

Any character vector will by default be stored  
in a data frame as a factor variable.

```
t$spouse <- c("Kit", "Spenser", "Blake")
```



## Reading data

`data()`

Lists available datasets in R

`data(faithful)`

Loads the Old Faithful geyser data set in R

`str(faithful)`

Compactly displays internal structure of R object

Here it is a dataframe with 272 observations and 2 variables

`edit(faithful) → new`

Allows you to edit the data and assign it to new

`edit(data.frame()) → new`

Enter new data via spreadsheet interface

? `read.csv`

Read data from a .csv file

## Old Faithful Geyser

A geyser is a hot spring that occasionally becomes unstable and erupts hot water and steam in the air. The Old Faithful Geyser is at Yellowstone Park, Wyoming.

## Simple plots

`plot (faithful)`

Plots a scatter plot of the eruptions & waiting time observations

`hist (faithful $ eruptions)`

Plots histogram of the eruptions

`hist (faithful $ eruptions, seq(1.6, 5.2, 0.2))` More detailed histogram

`plot (faithful $ eruptions)`

Plots all the observations

`plot (faithful $ eruptions, type="l")`

Plots observation with lines

`plot.ecdf (faithful $ eruptions)`

Plots the empirical cumulative distribution function

`qqnorm (faithful $ eruptions)`

Provides Q-Q plot to create normal quantile plots

`qqnorm (faithful $ eruptions [faithful $ eruptions > 3])`

Q-Q plot with subset of entries

`qqline (faithful $ eruptions [faithful $ eruptions > 3])`

Adds a line to data

`boxplot (faithful $ eruptions)`

Provides boxplot of data

Old Faithful Geyser dataset observations:

- 1) Eruption times and waiting time between successive eruptions exhibit highly oscillatory behavior: low followed by high and high followed by low
- 2) Eruption times have a highly bimodal distribution
- 3) Lower eruption times                      Lower waiting times  
Higher eruption times                      Higher waiting times

This can be used to predict when the next geyser eruption will occur.

During a short eruption, less water and heat are used, so both are restored in shorter time.

During longer eruptions, more time is needed for it to rebuild.

`faithful1 <- subset(faithful, faithful$eruptions < 3)`

`faithful <- subset(faithful, faithful$eruptions >= 3)`

Creates subset of dataframes  
with all rows corresponding to  
cases where the condition is satisfied

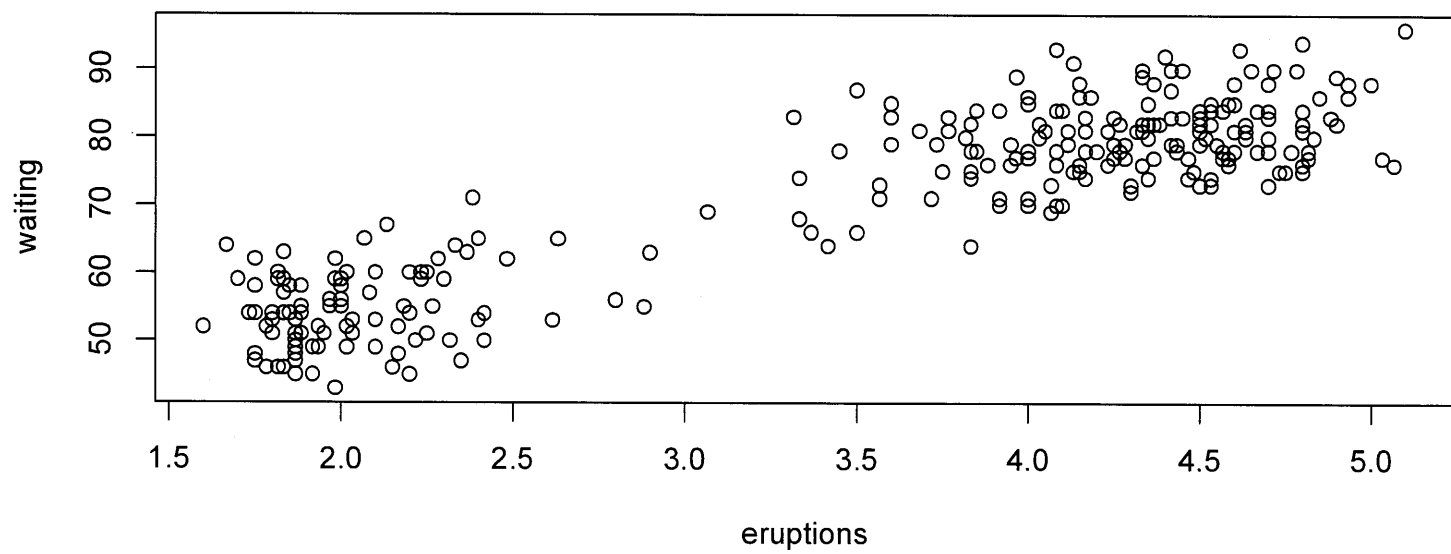
`t.test(faithful1)`  
\$waiting

`t.test(faithful2)`  
\$waiting

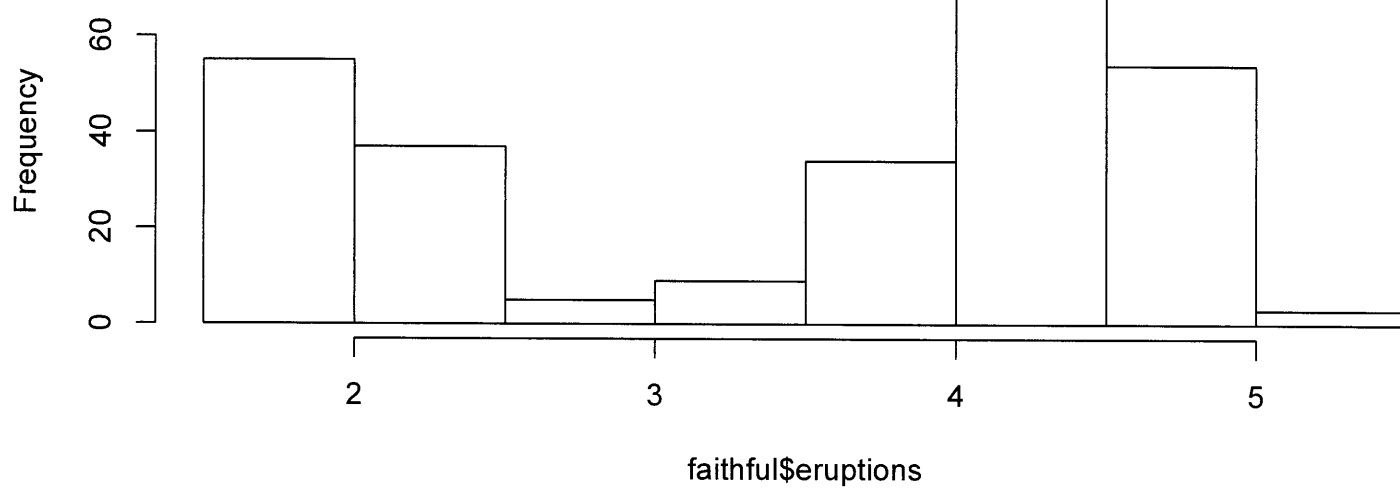
Performs a statistical one sample  
t-test to test the hypothesis  
if mean value = 0 and  
derive 95% confidence interval  
for mean parameter

In this case with 95% confidence, if the  
eruption time < 3, the average waiting time  
is between 53.3 and 55.67.

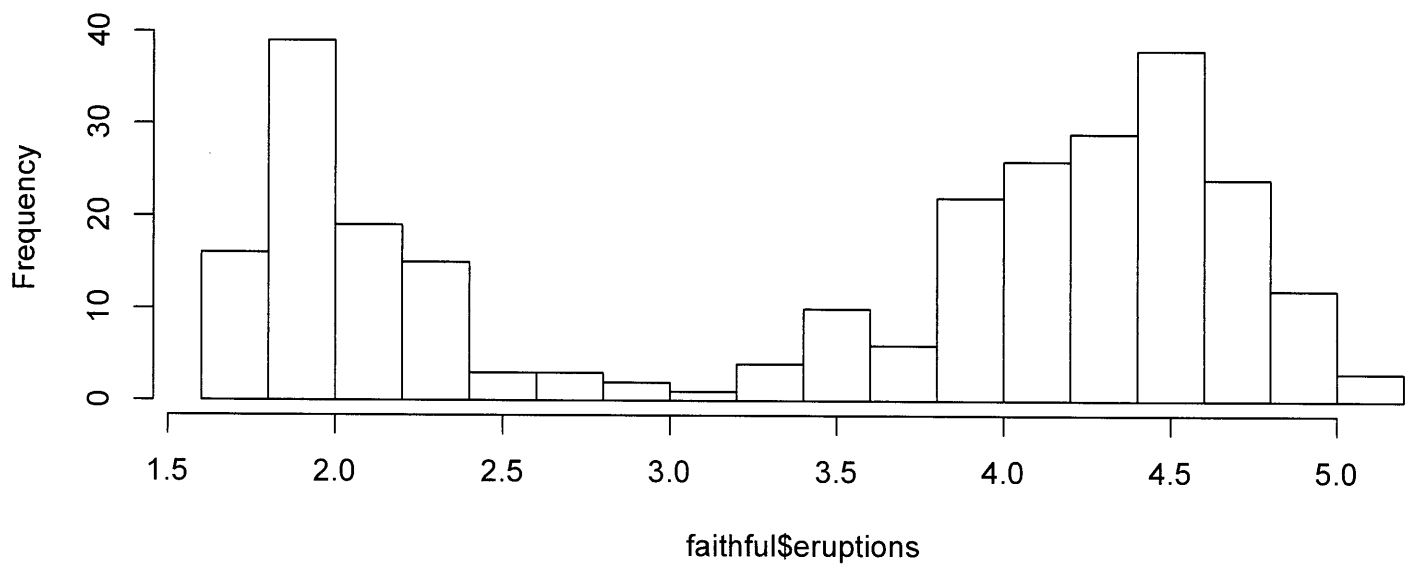
If the eruption time >= 3, the average waiting time  
with 95% confidence is between 79.1 and 80.89

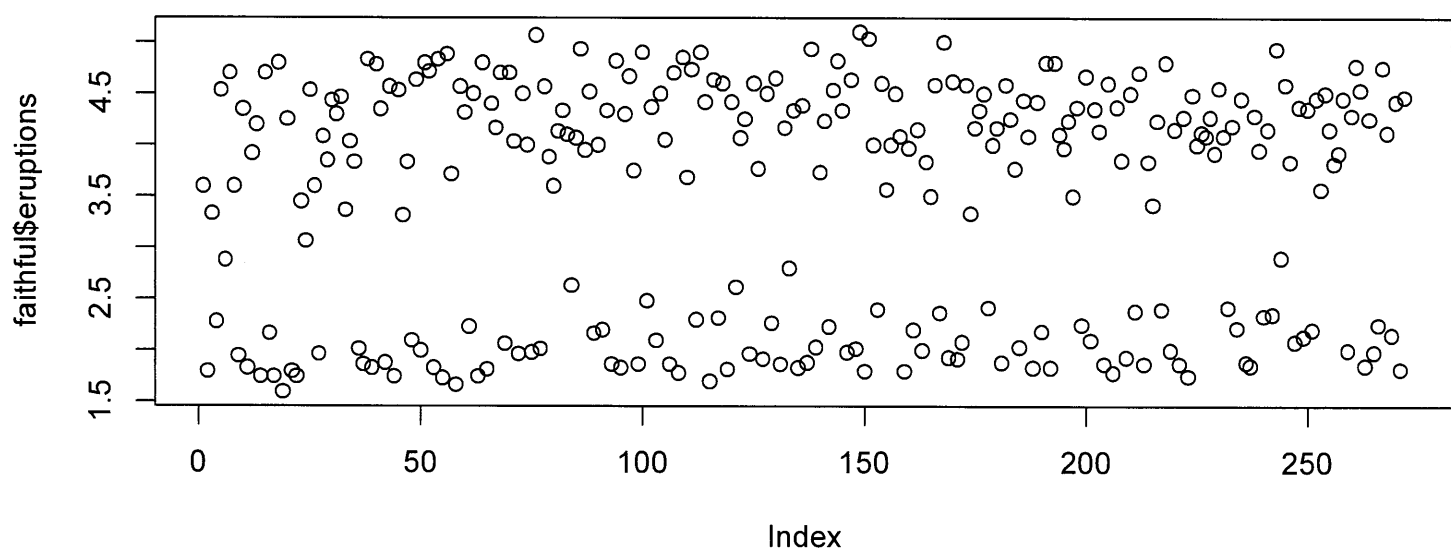


**Histogram of faithful\$eruptions**

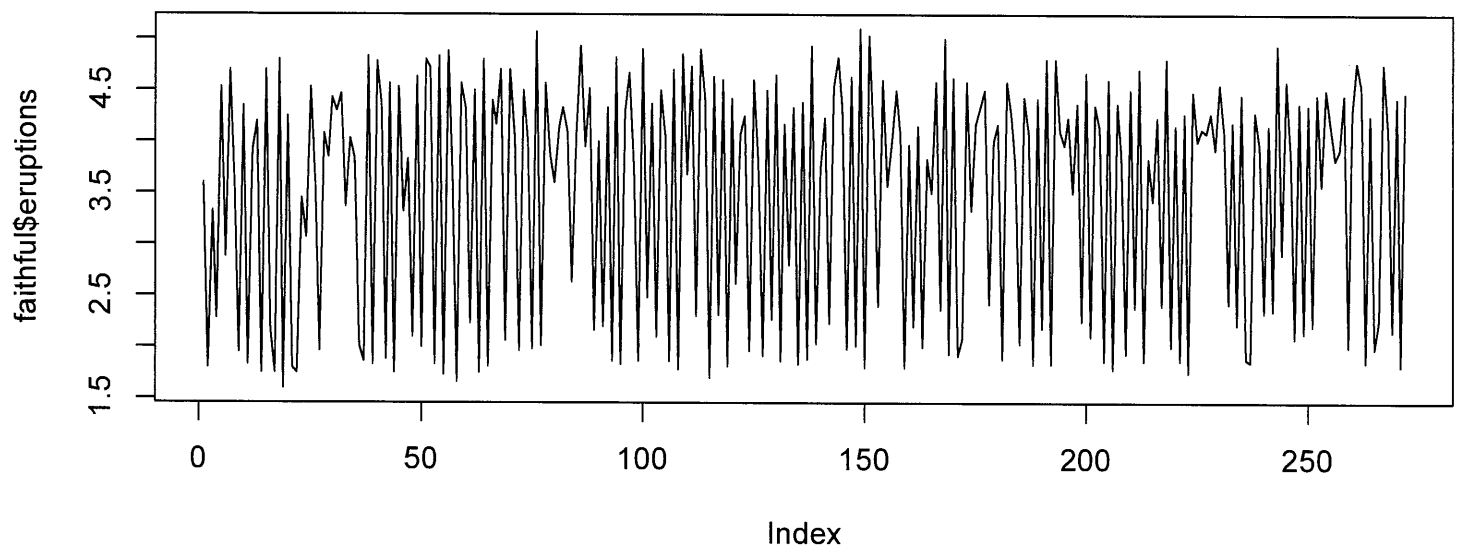


**Histogram of faithful\$eruptions**

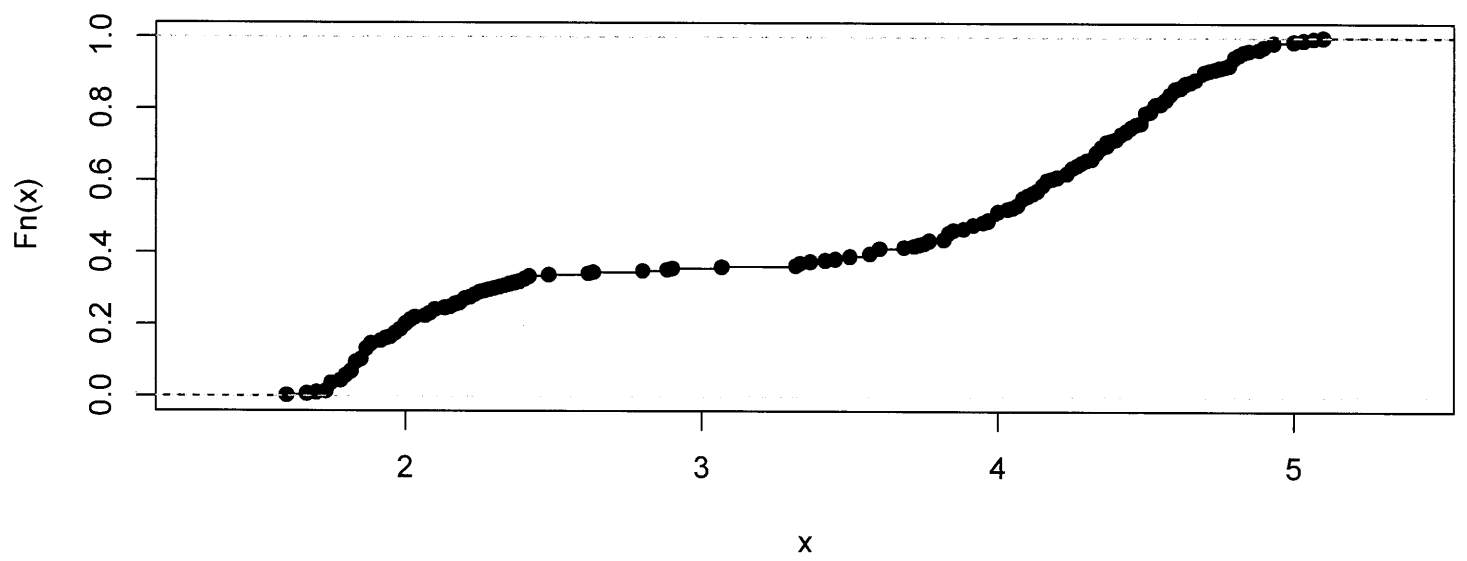




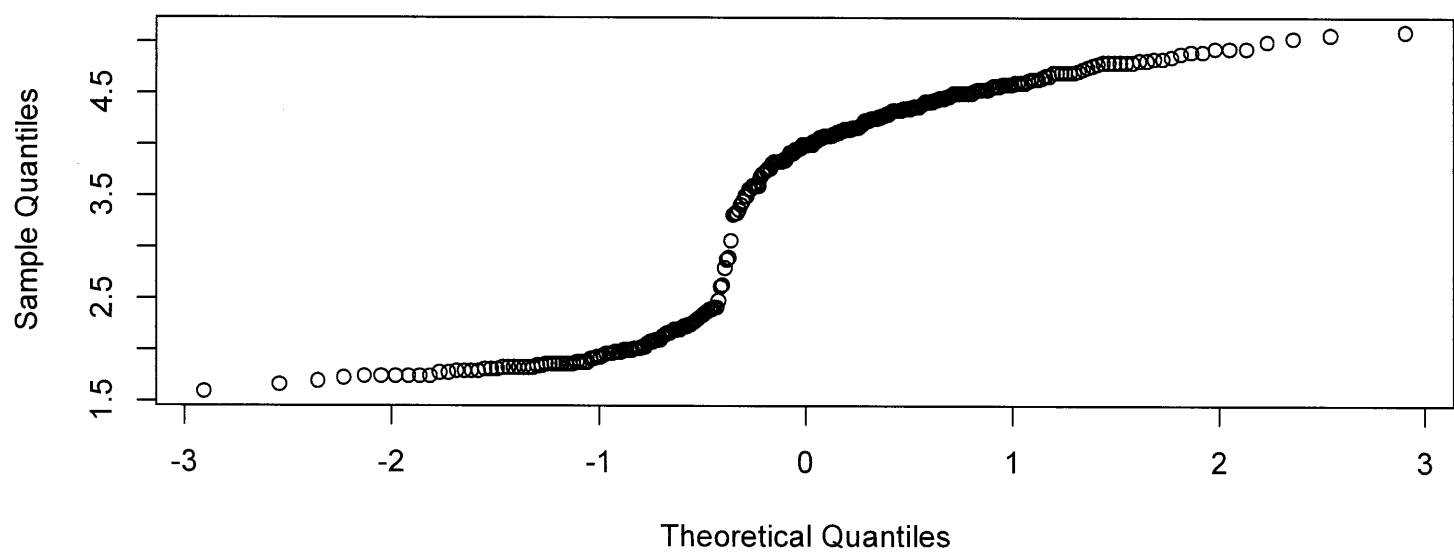




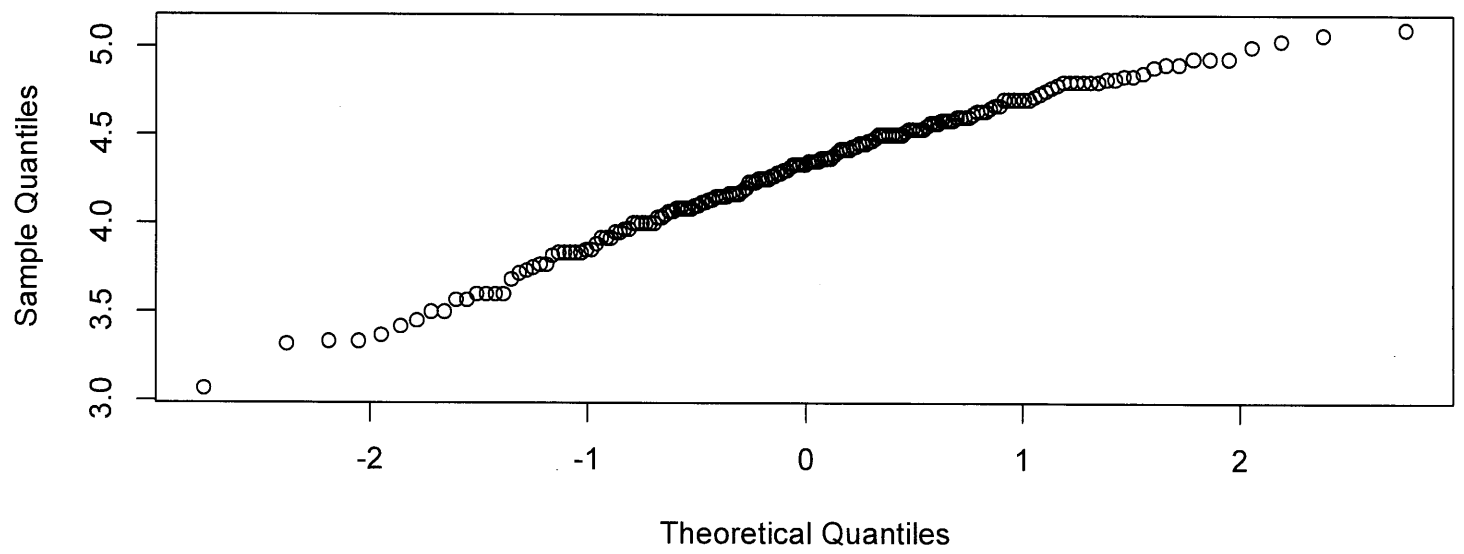
ecdf(x)



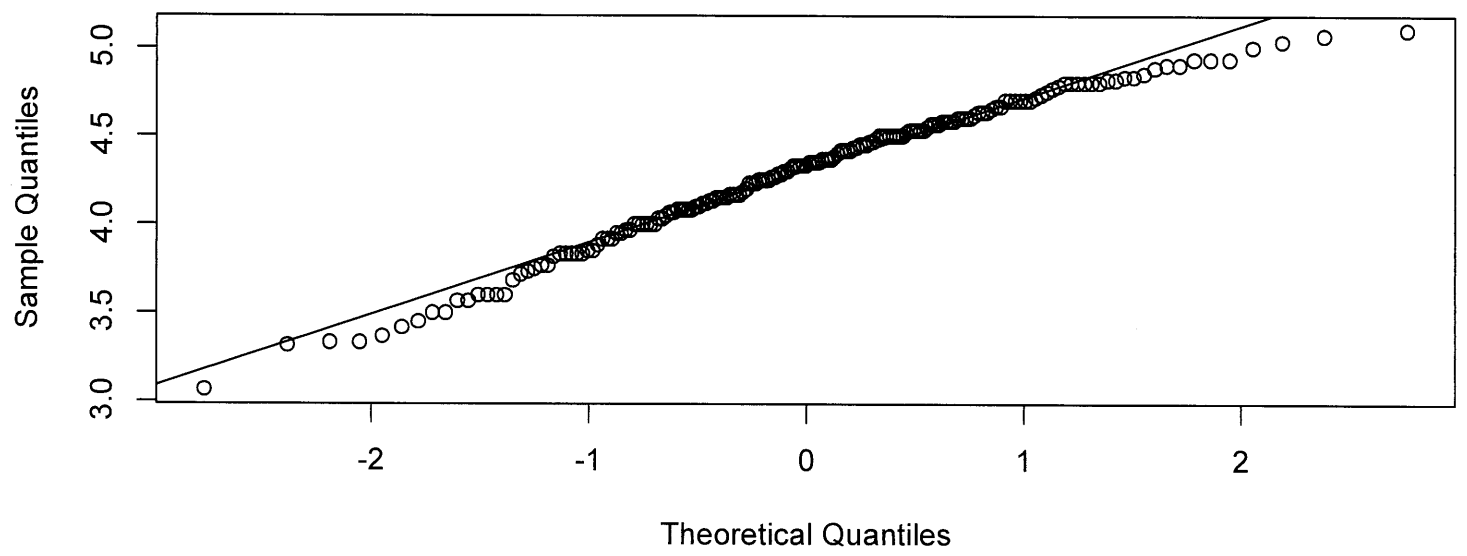
Normal Q-Q Plot

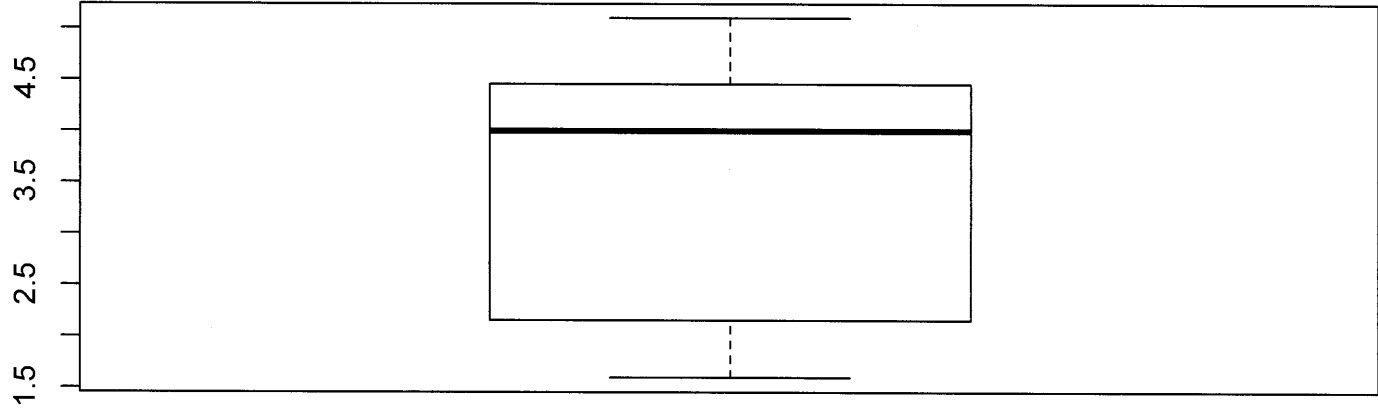


Normal Q-Q Plot



Normal Q-Q Plot





# Data from the World Health Organization

[www.who.int/gho/en](http://www.who.int/gho/en)

(Global Health Observatory data)

```
w ← read.csv("WHO.csv")
```

str(w)

This data set contains data on 194 countries from

WHO. The variables are:

Country - name of country

Region - region the country lies in

Population - population in thousands

Under 15 - % of population under 15 years of age

Over 60 - % of population over 60 years of age

FertilityRate - Average number of children per woman

Life Expectancy - Life expectancy in years

Literacy Rate - Literacy rate among adults at least 15 years of age

GNI - gross national income per capita

```
which(w$Country == "Singapore")
```

This identifies the row of the dataframe with Singapore

```
Summary(w$Over60)
```

```
Summary(w$Under15)
```

```
w[155,]
```

As this indicates, Singapore has a high ratio of population Over 60 & low ratio under 15 compared to many countries

Plot (w \$ GNI, w \$ Fertility Rate)

The figure indicates an inverse mediation between income & fertility. This suggests that reproductive restraint arises often as a consequence of economic progress or lower fertility leads to more resources being available per child, making them more productive.

We recreate this with ggplot2 package.

```
ggplot(w, aes(x = GNI, y = Fertility Rate)) + geom_point()
```

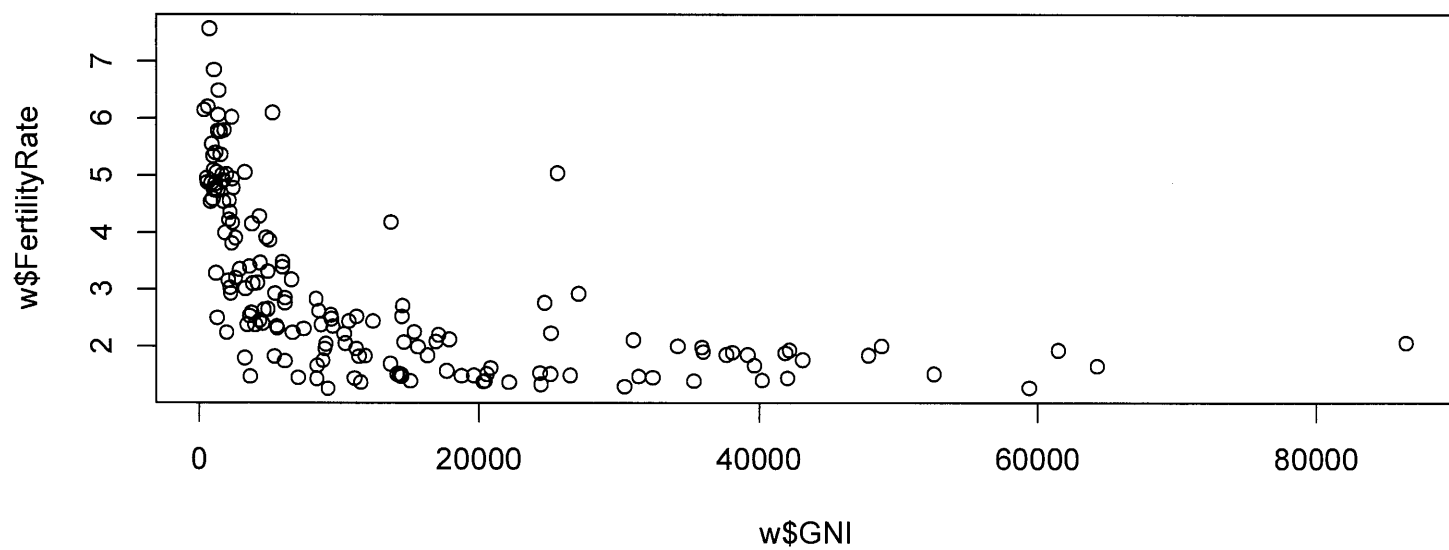
This provides a visual improvement where we add the data and the aesthetic mapping and then the layer with the points.

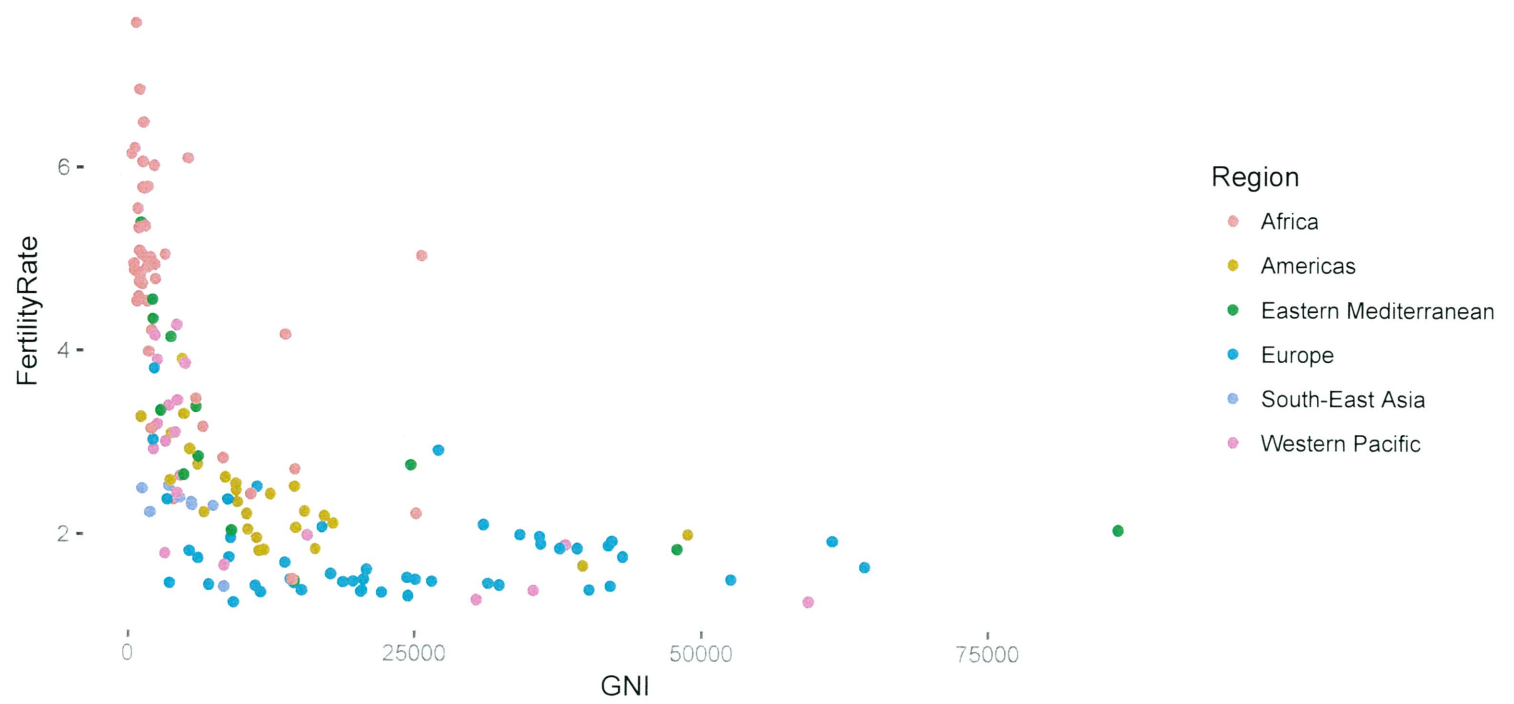
2 ggplot

```
ggplot(w, aes(x = GNI, y = Fertility Rate, color = Region))  
+ geom_point()
```

This colours the points as per the regions by adding the colors argument to aesthetic option. This helps see that there are trends by regions such as Europe has high GNI, low fertility rates while Africa has low GNI, high fertility rates.







## More advanced visualization in R

One of the advantages of software such as R is that you can access and install packages that provide advanced techniques. The ggplot2 package in R is one such package that improves the basic plotting functionality in R.

In ggplot2, you build a graphic by building layer upon layer where the layers are:

- Data and aesthetic mapping

Data turns an abstract graphic into a concrete graphic

- Statistical transformation

Transforms the data by summarizing it such as smoother which calculates the mean of  $y$  given  $x$  while ensuring smoothness

- Geometric object

Controls the type of plot you can create

`install.packages("ggplot2")`

`library(ggplot2)`

Install and load the ggplot2 package before using. The next time you want to run it, you only need to use the library command.

## Date from crime location and time in Chicago

Increasingly government organizations are providing access to data to citizens on statistics such as occurrence of crimes. One such website is <https://data.cityofchicago.org>.

```
Crime <- read.csv("Crime.csv")
```

```
Str(Crime)
```

This dataset consists of 191641 observations with date and time, latitude and longitude of location of motor theft crimes.

```
Crime$date <- strptime(Crime$date,  
  format = "%m/%d/%y %H%M")
```

```
Crime$Weekdays <- weekdays(Crime$date)
```

```
WeekdayCounts <- as.data.frame(table(Crime$Weekdays))
```

The first line converts the date variable to a format that R can work with. The second line extracts the weekday from the date while the third line creates a dataframe that counts the number of crimes on each day of the week. Here `WeekdayCounts$Var1` gives the day & `WeekdayCounts$Freq` gives the number of crimes on the day.

We now create a heatmap to plot the crime data.

```
Crime$Hour <- Crime$Date$hour
```

This creates a new variable that captures the hour which is easy to retrieve from the R object now.

```
WeekdayHourCounts <- as.data.frame(  
  table(Crime$Weekdays, Crime$Hour))
```

This creates a data frame with 168 observations where Var 1 is the weekday, Var 2 is the hour of the crime and Freq is the number of crimes.  $(168 = \underbrace{24}_{\text{Hours}} \times \underbrace{7}_{\text{Days}})$

```
WeekdayHourCounts$Var1 <-
```

```
factor(WeekdayHourCounts$Var1, Ordered = TRUE,  
  levels = c("Monday", "Tuesday", "Wednesday",  
    "Thursday", "Friday", "Saturday", "Sunday"))
```

This helps convert it to chronological order rather than alphabetical order.

```
WeekdayHourCounts$Var2 <-
```

```
as.numeric(as.character(WeekdayHourCounts$Var2))
```

This converts Var 2 to a numeric vector from factor. We need to use the intermediate as.character() to deal with how R stores factors.

```
ggplot (Weekday Hour Counts, aes (x = Var2, y = Var1))  
+ geom_tile (aes (fill = Freq)) +  
xlab ("Hour of day") + ylab ("")
```

The figure is a heat map where the lighter color shows more motor vehicle thefts on that day and that hour.

```
ggplot (Weekday Hour Counts, aes (x = Var2, y = Var1))  
+ geom_tile (aes (fill = Freq)) + xlab ("Hour of day")  
+ ylab ("") + scale_fill_gradient (low = "white",  
high = "red")
```

This shows higher frequency values by red and lower by white helping policemen identify hotspots more easily in practice.

As the figure shows, there are more crimes for example Friday nights around 22:00.

We can also overlay such heat maps on geographical maps.

```
install.packages("maps")
```

```
library(maps)
```

```
install.packages("ggmap")
```

```
library(ggmap)
```

} Use packages maps and ggmap for this display of maps and to combine with Google Maps and other such information

```
Chicago <- get_map(location = "Chicago")
```

```
ggmap(Chicago)
```

This gets a map of Chicago and plots it

```
LatLongCounts <- as.data.frame(table(round(Crime$Longitude, 2), round(Crime$Latitude, 2)))
```

This creates a new data frame that captures latitude & longitude coordinates but rounds it so that the number of points reduce. We have 1638 observations where Var1 is longitude, Var2 is latitude & Freq is number of thefts in that area.

```
LatLongCounts$Var1 <- as.numeric(as.character(LatLongCounts$Var1))
```

```
LatLongCounts$Var2 <- as.numeric(as.character(LatLongCounts$Var2))
```

This helps convert factor to numeric.

```
ggmap(Chicago) + geom_tile(data = LatLongCounts, aes(x = Var1, y = Var2, alpha = Freq), fill = "red")
```

