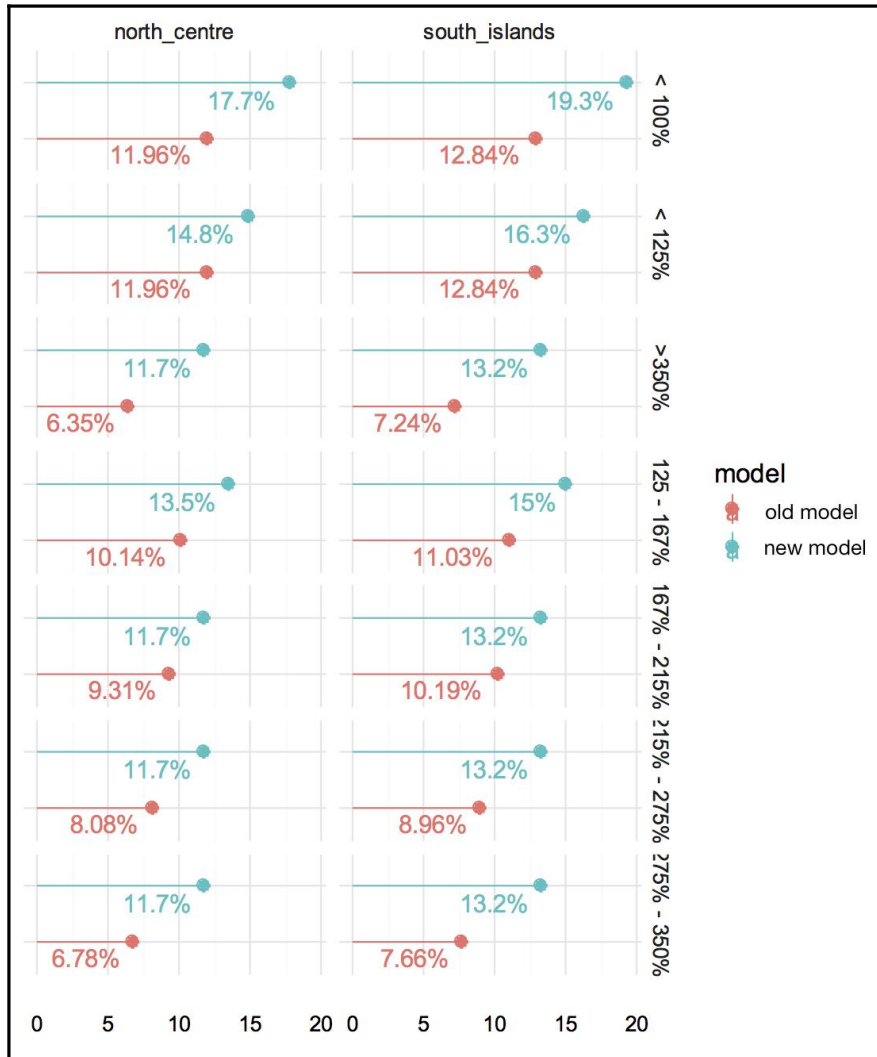


Chapter 1: Why to Choose R for Your Data Mining and Where to Start





The Comprehensive R Archive Network

Download and Install R

Precompiled binary distributions of the base system and contributed packages, **Windows and Mac** users most likely want one of these versions of R:

- [Download R for Linux](#)
- [Download R for \(Mac\) OS X](#)
- [Download R for Windows](#)

R is part of many Linux distributions, you should check with your Linux package management system in addition to the link above.

Source Code for all Platforms

Windows and Mac users most likely want to download the precompiled binaries listed in the upper box, not the source code. The sources have to be compiled before you can use them. If you do not know what this means, you probably do not want to do it!

- The latest release (Monday 2016-10-31, Sincere Pumpkin Patch) [R-3.3.2.tar.gz](#), read [what's new](#) in the latest version.
- Sources of [R alpha and beta releases](#) (daily snapshots, created only in time periods before a planned release).
- Daily snapshots of current patched and development versions are [available here](#). Please read about [new features and bug fixes](#) before filing corresponding feature requests or bug reports.
- Source code of older versions of R is [available here](#).
- Contributed extension [packages](#)

Questions About R

- If you have questions about R like how to download and install the software, or what the license terms are, please read our [answers to frequently asked questions](#) before you send an email.

CRAN
[Mirrors](#)
[What's new?](#)
[Task Views](#)
[Search](#)

About R
[R Homepage](#)
[The R Journal](#)

Software
[R Sources](#)
[R Binaries](#)
[Packages](#)
[Other](#)

Documentation
[Manuals](#)
[FAQs](#)
[Contributed](#)

The screenshot displays the R environment interface. On the left, a script editor window titled 'first_script.R' contains the code `print("hello world")`. The central R Console window shows the R version 3.3.0 (2016-05-03) and platform information (x86_64-apple-darwin13.4.0). It also displays the R license notice and instructions for using help and package management functions. The right-hand side shows a 'Search in History' panel with a list of commands entered in the console, including `Rprof()`, `rw2d4 = function(n) { steps = sample(c(-1, 1), n); Rprof(); library(shiny); install.packages('shiny'); }`, and `print("hello world")`.

~/Dropbox/R_projects/r_datamining - RStudio

```

1 types <- c("source", "win.binary",
2           "mac.binary", "mac.binary.mavericks")
3
4 CRANmirror <- "http://cran.revolutionanalytics.com"
5
6 pdb <- lapply(types, function(x){
7   cran <- contrib.url(repos = CRANmirror,
8                       type = x)
9   available.packages(contriburl = cran, type = x)
10 })
11 names(pdb) <- types
12 str(pdb, max.level = 1)
13
14
15 # Number of available packages
16 data <- sapply(pdb, nrow)
17
17.1 (Top Level)

```

Environment History

Global Environment

Values

CRANmirror	"http://cran.revolutionanalytic...
data	Named int [1:4] 9691 9613 0 9589
pdb	Large list (4 elements, 9.9 Mb)
types	chr [1:4] "source" "win.binary"...

Files Plots Packages Help Viewer

Console

```

.. attr(*, "dimnames")=List of 2
$ mac.binary.mavericks: chr [1:9589, 1:17] "A3" "abbyrR" "abc.data" "ABC.RAP" ...
.. attr(*, "dimnames")=List of 2
>
> # Number of available packages
> sapply(pdb, nrow)
      source      win.binary      mac.binary mac.binary.mavericks
      9691           9613              0             9589
> data <- sapply(pdb, nrow)
> plot(data)
> barplot(data)
>

```

rproject13 - Microsoft Visual Studio

```

10 mtcars$am <- factor(mtcars$am, levels = c(0, 1),
11 labels = c("Automatic", "Manual"))
12 mtcars$cyl <- factor(mtcars$cyl, levels = c(4, 6, 8),
13 labels = c("4cyl", "6cyl", "8cyl"))
14
15 # Kernel density plots for mpg
16 # grouped by number of gears (indicated by color)
17 # qplot(mpg, data = mtcars, geom = "density", fill = gear, alpha = I(.5),
18 main = "Distribution of Gas Mileage", xlab = "Miles Per Gallon",
19 ylab = "Density")
20
21 s <- sample()
22 # Scatterplot
23 # in each fa
24 # sample takes a sample of the specified size from the elements of x using either with
25 # or without replacement.
26 # xlab = "H x: Either a vector of one or more elements from which to choose, or a positive integer.
27 # See Details.

```

Variable Explorer

Name	Value	Class	Type
mtcars	32 obs. of 11 variables	data.frame	list
@.Data	List of 11	list	list
@.names	chr [1:11] "mpg" "cyl" "disp" "hp" "drat" "wt"	character	character
@.rownames	chr [1:32] "Mazda RX4" "Mazda RX4 Wag" "Datsun 710" "Hornet 4 Drive" "Hornet Sportabout" "Valiant" "Duster 360" "Merc 240D" "Merc 230" "Merc 280" "Merc 280C" "Merc 450SE"	character	character

R Data: base::GlobalEnv\$mtcars

	mpg	cyl	disp	hp	drat	wt
Mazda RX4	21.0	6cyl	160.0	110	3.90	2.620
Mazda RX4 Wag	21.0	6cyl	160.0	110	3.90	2.875
Datsun 710	22.8	4cyl	108.0	93	3.85	2.320
Hornet 4 Drive	21.4	6cyl	258.0	110	3.08	3.215
Hornet Sportabout	18.7	8cyl	360.0	175	3.15	3.440
Valiant	18.1	6cyl	225.0	105	2.76	3.460
Duster 360	14.3	8cyl	360.0	245	3.21	3.570
Merc 240D	24.4	4cyl	146.7	62	3.69	3.190
Merc 230	22.8	4cyl	140.8	95	3.92	3.150
Merc 280	19.2	6cyl	167.6	123	3.92	3.440
Merc 280C	17.8	6cyl	167.6	123	3.92	3.440
Merc 450SE	16.4	8cyl	275.8	180	3.07	4.070

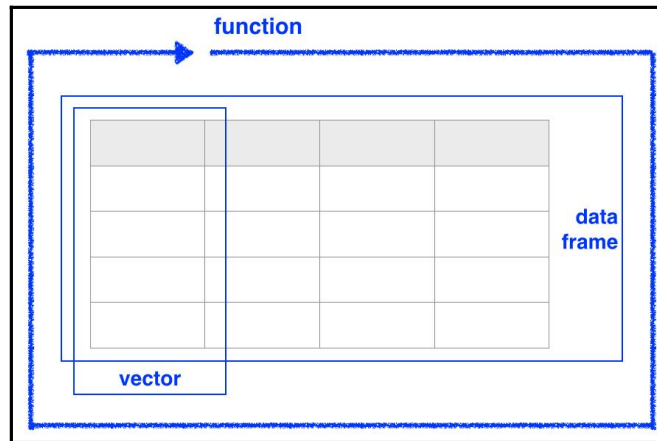
R Interactive

```

> library(ggplot2)
> # create factors with value labels
> mtcars$gear <- factor(mtcars$gear, levels = c(3, 4, 5),
+ labels = c("3gears", "4gears", "5gears"))
> mtcars$am <- factor(mtcars$am, levels = c(0, 1),
+ labels = c("Automatic", "Manual"))
> mtcars$cyl <- factor(mtcars$cyl, levels = c(4, 6, 8),
+ labels = c("4cyl", "6cyl", "8cyl"))
> # Kernel density plots for mpg
> # grouped by number of gears (indicated by color)
> # qplot(mpg, data = mtcars, geom = "density", fill = gear, alpha = I(.5),
+ main = "Distribution of Gas Mileage", xlab = "Miles Per Gallon",
+ ylab = "Density")
> # Scatterplot of mpg vs. hp for each combination of gears and cylinders
>

```

R Plot



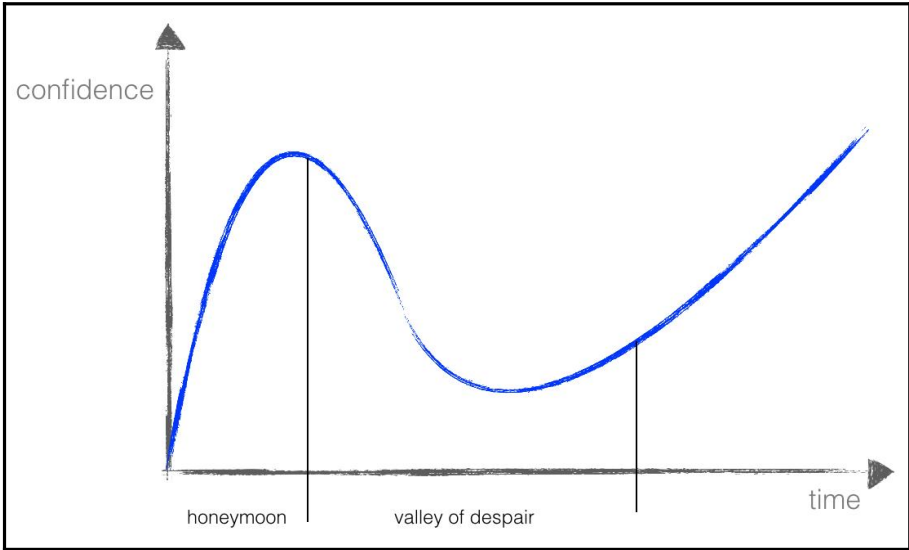
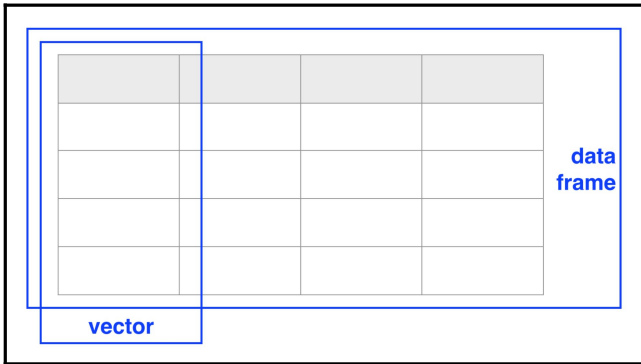
```
Console ~/ | ↻  
  
> demo()  
> demo(scoping)  
  
      demo(scoping)  
      ----  
  
Type <Return> to start :  
  
> ## Here is a little example which shows a fundamental difference between  
> ## R and S. It is a little example from Abelson and Sussman which models  
> ## the way in which bank accounts work.      It shows how R functions can  
> ## encapsulate state information.  
> ##  
> ## When invoked, "open.account" defines and returns three functions
```

The screenshot shows the RStudio Source Editor window titled "RStudio Source Editor" with a tab for "first_r_script.R". The editor contains the following code:

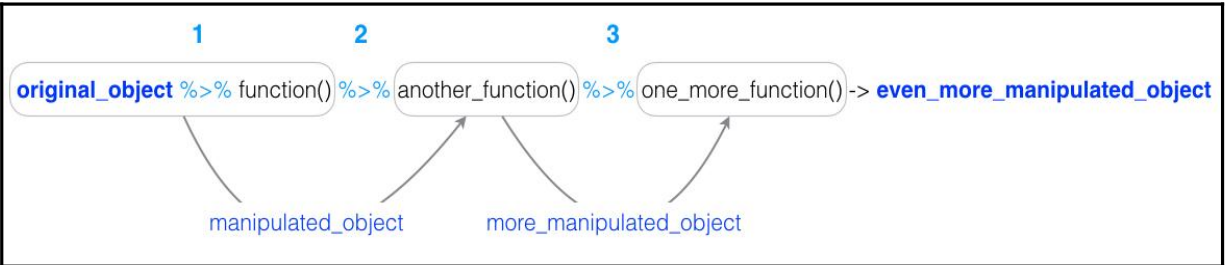
```
1 print("hello world")  
2 |
```

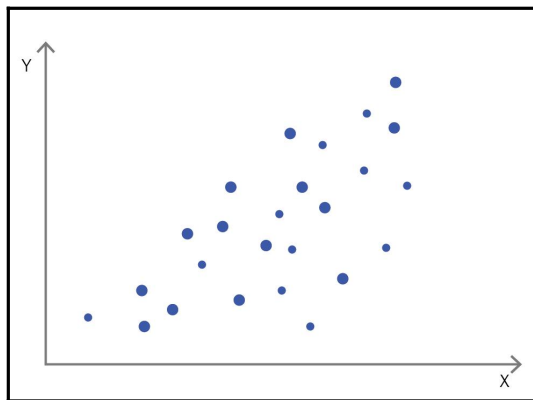
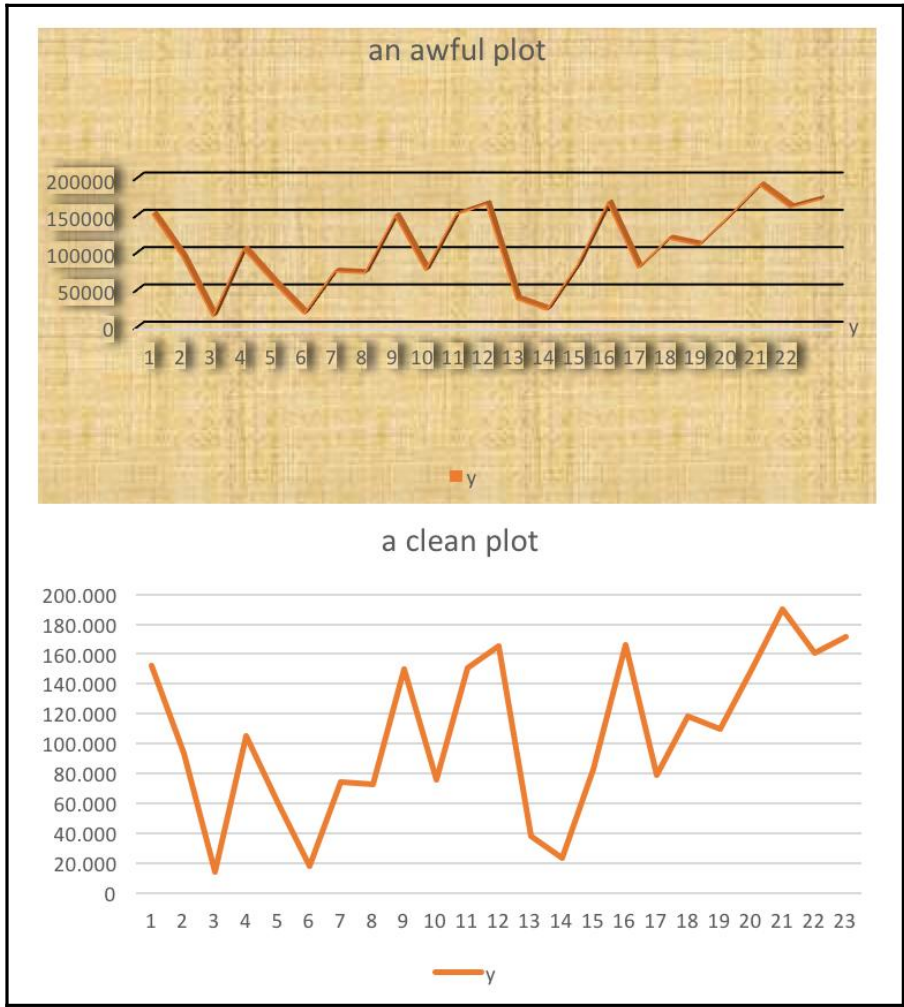
The status bar at the bottom indicates "2:1 (Top Level) ↕" and "R Script ↕".

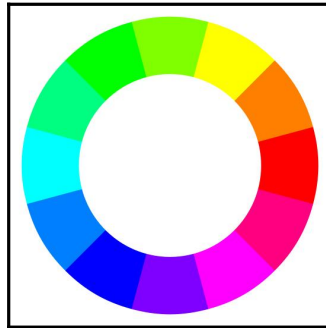
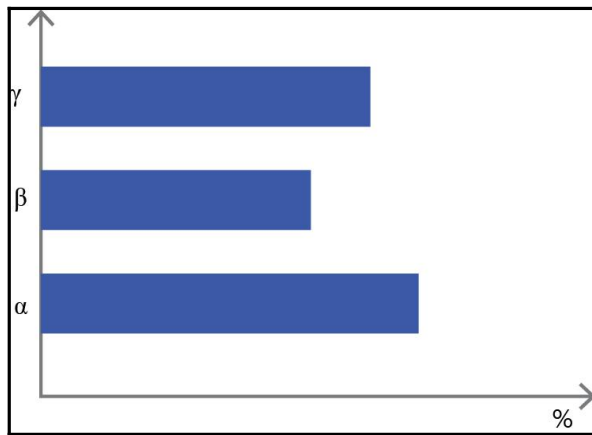
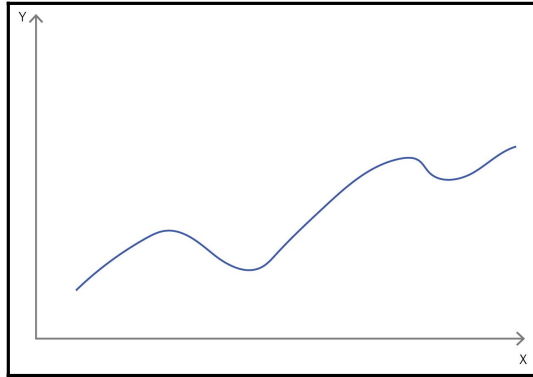
100	20
20	100
40	90
15	40
90	15

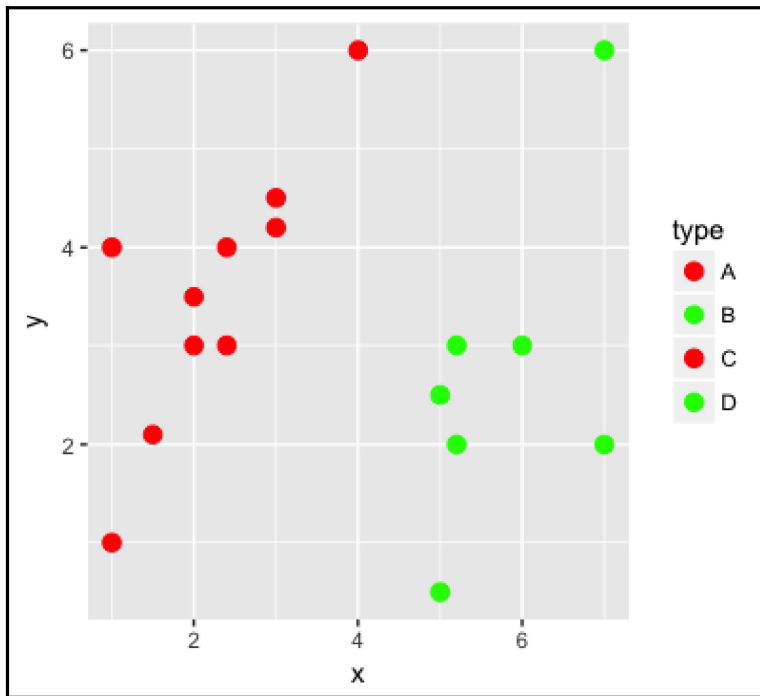


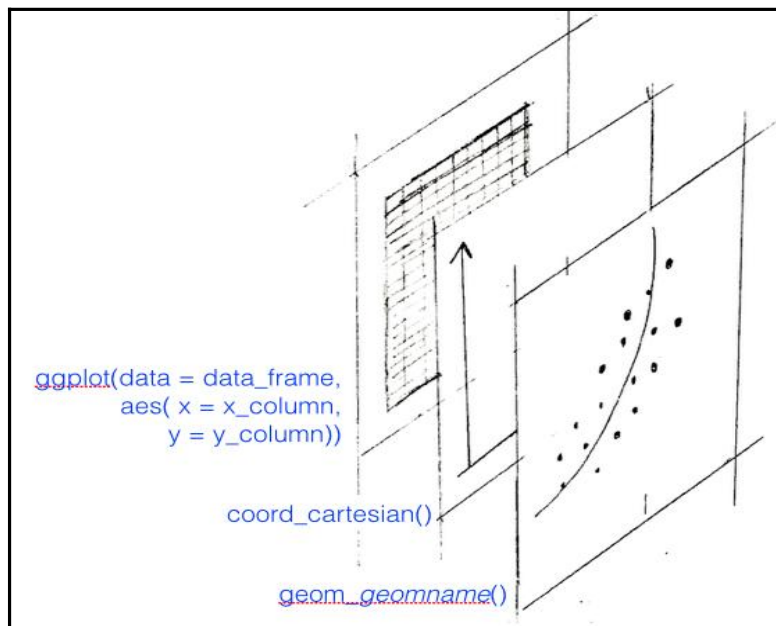
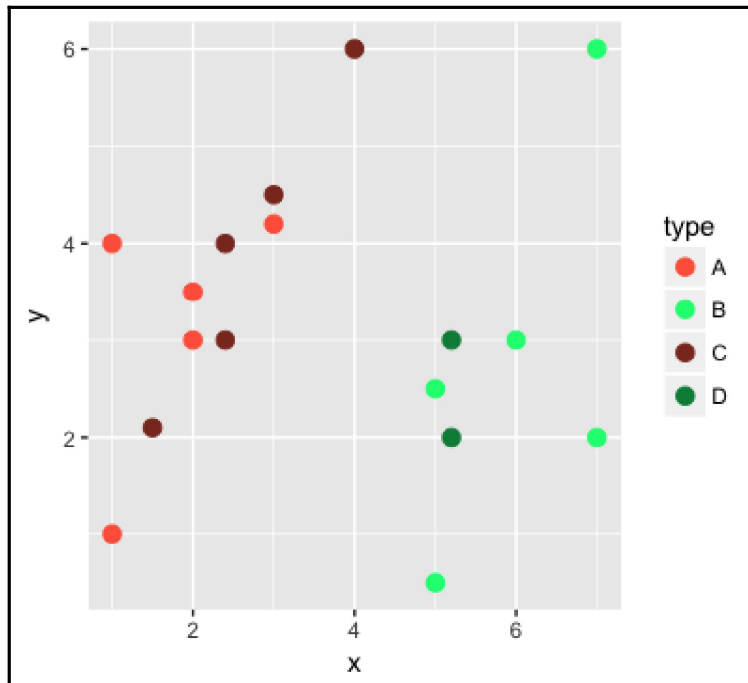
Chapter 2: A First Primer on Data Mining Analysing Your Bank Account Data

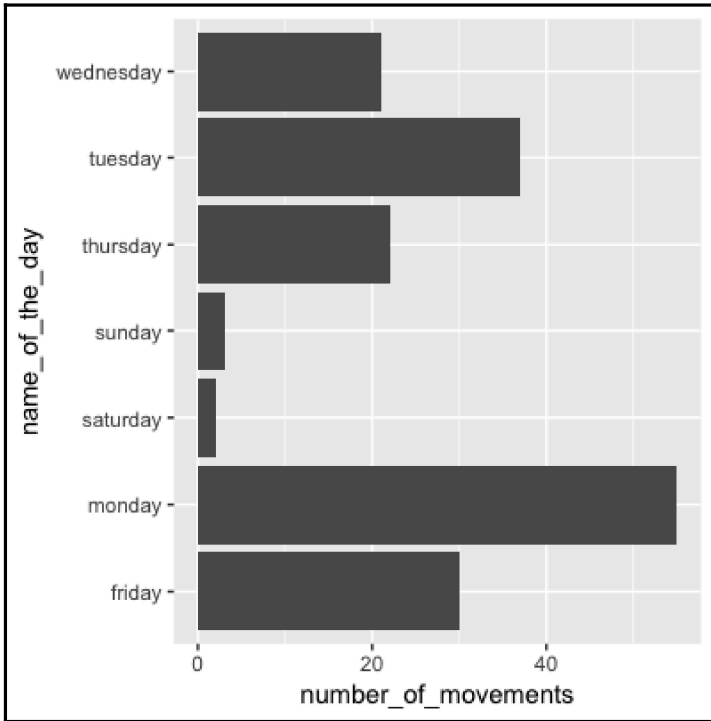
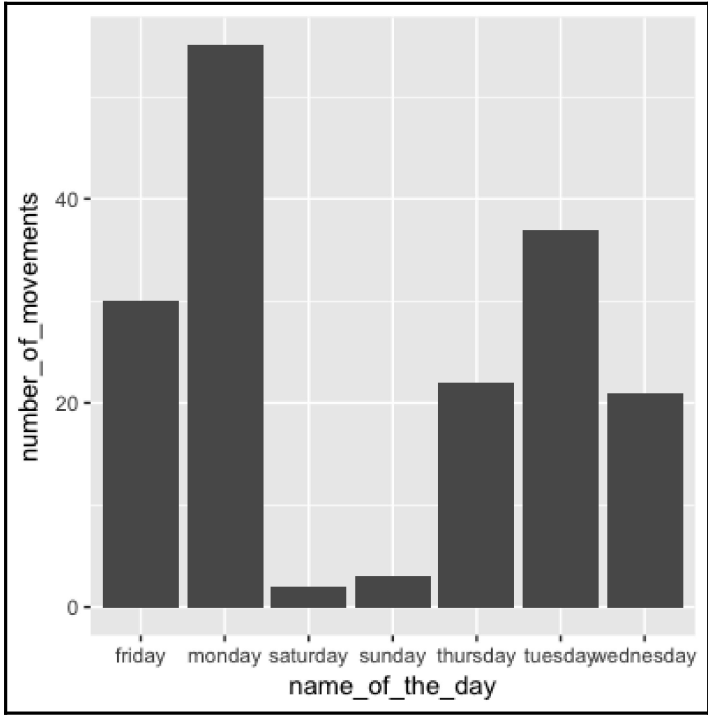


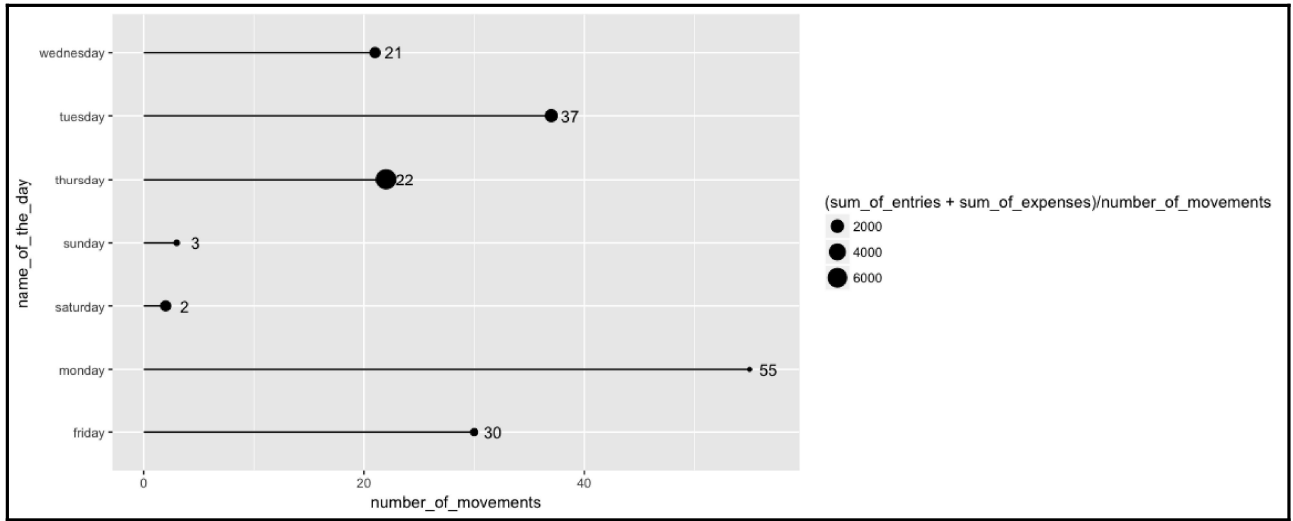




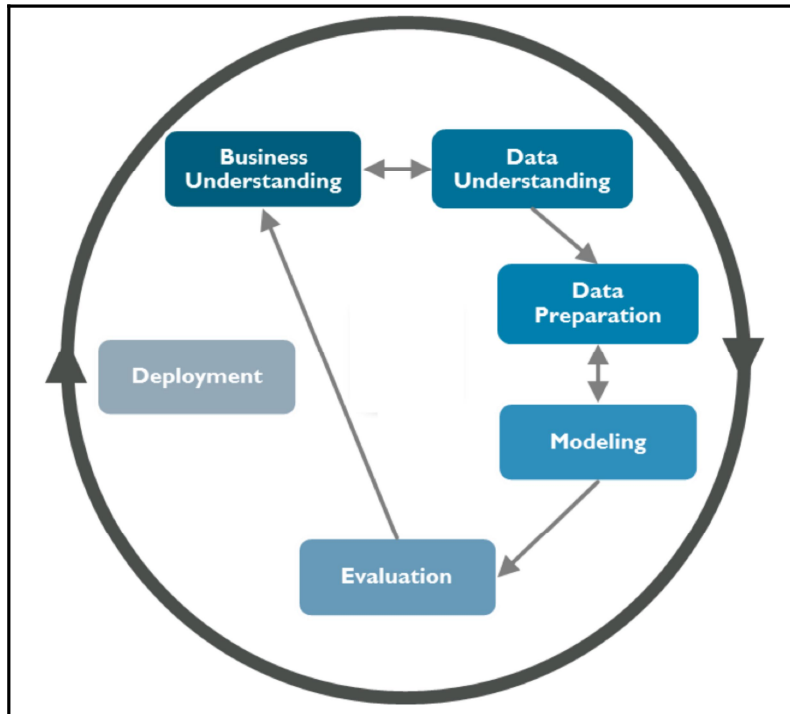


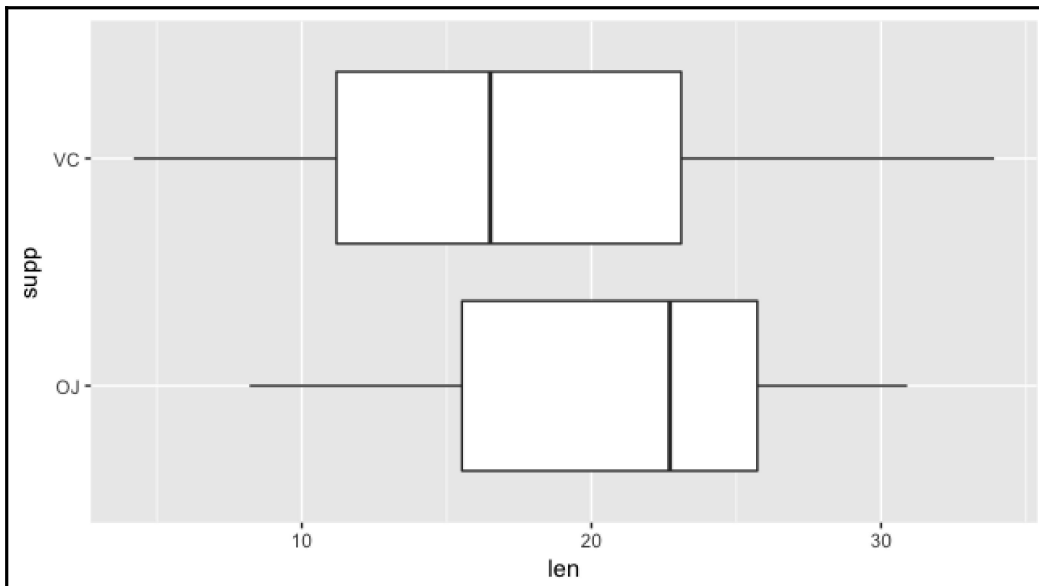
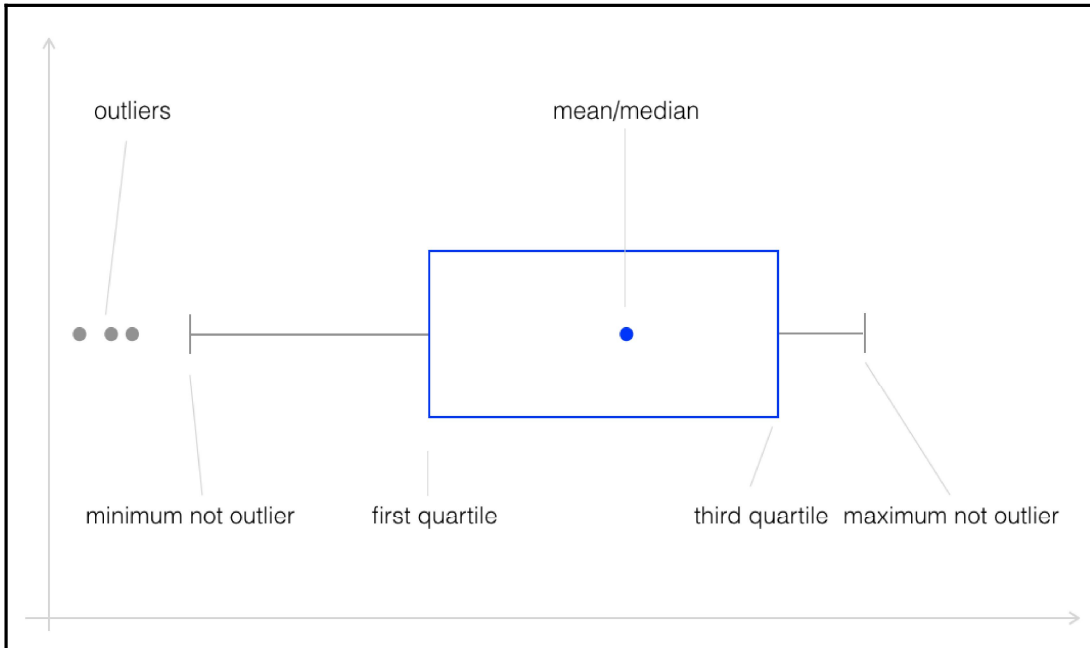


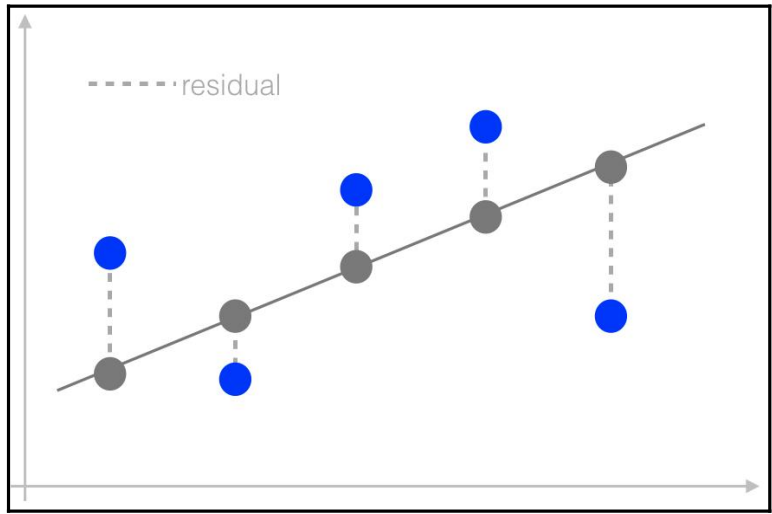
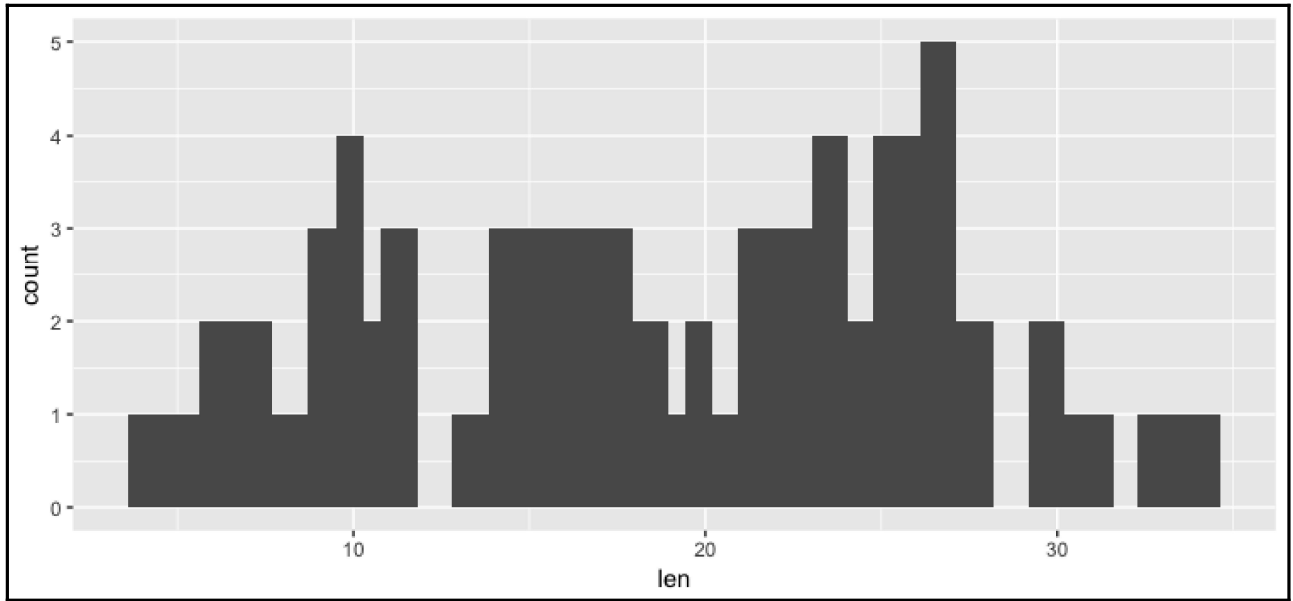


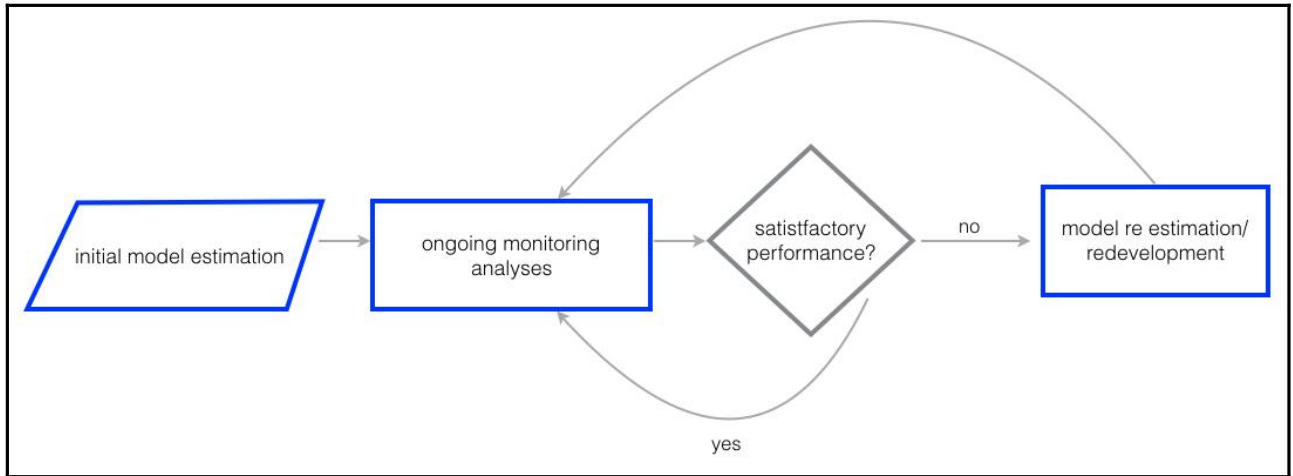


Chapter 3: The Data Mining Process - CRISP-DM Methodology

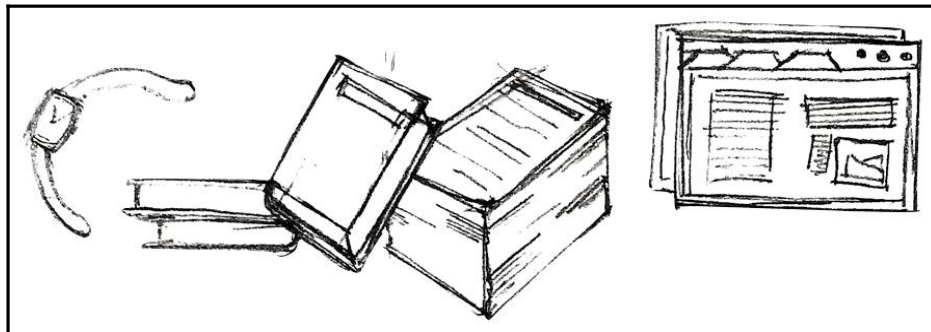
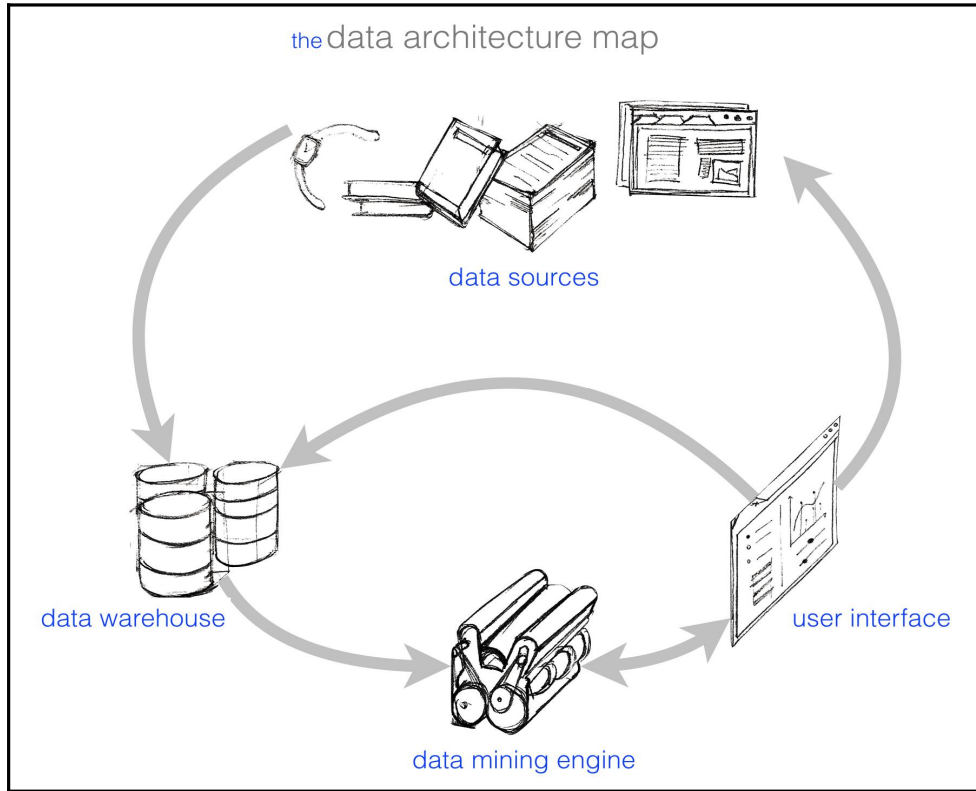


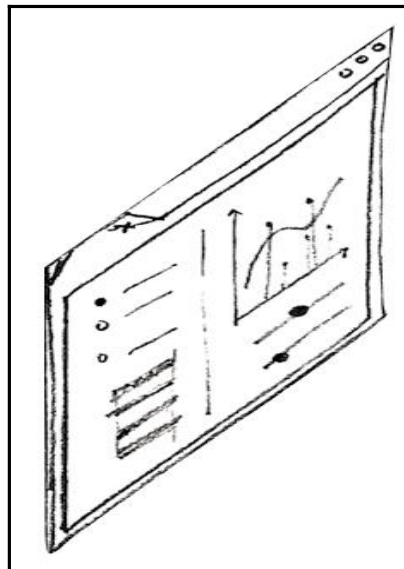
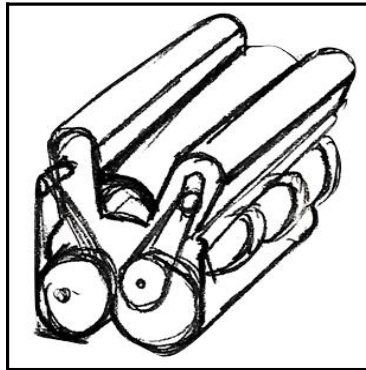
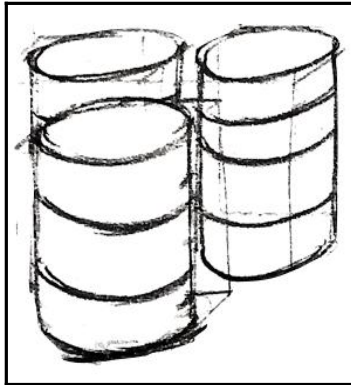






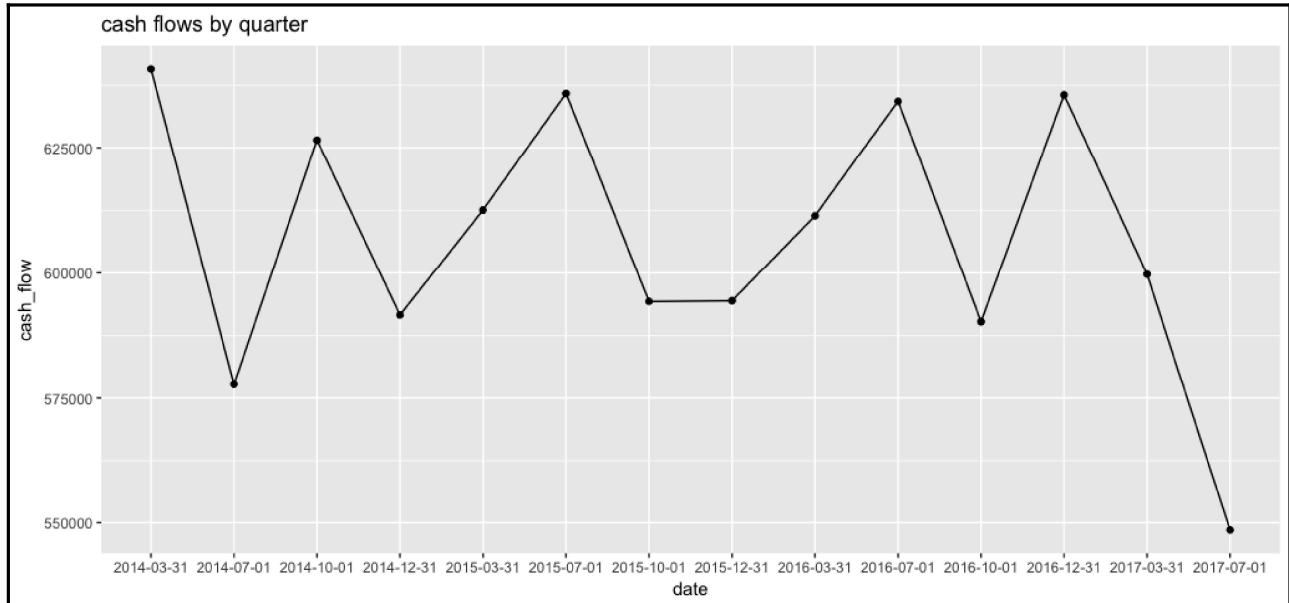
Chapter 4: Keeping the House Clean – The Data Mining Architecture



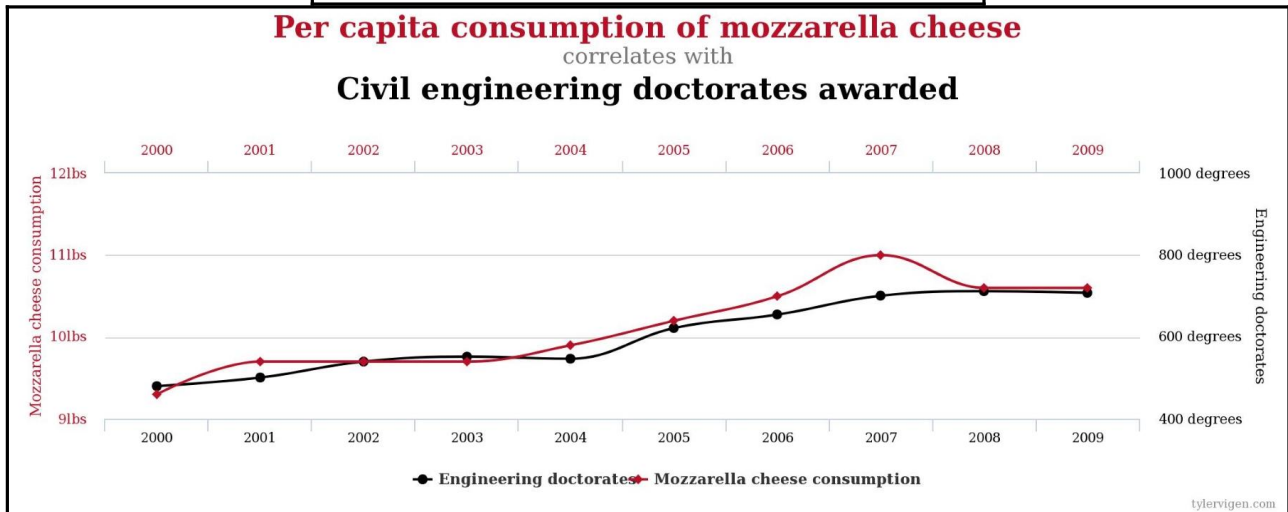
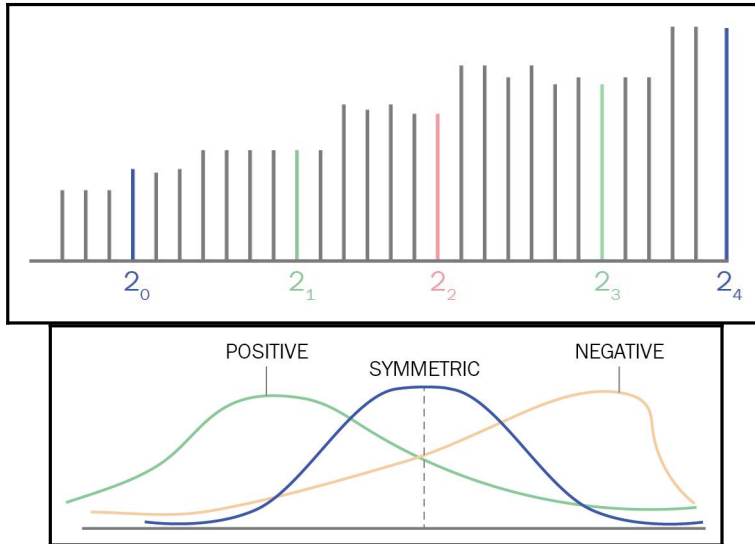


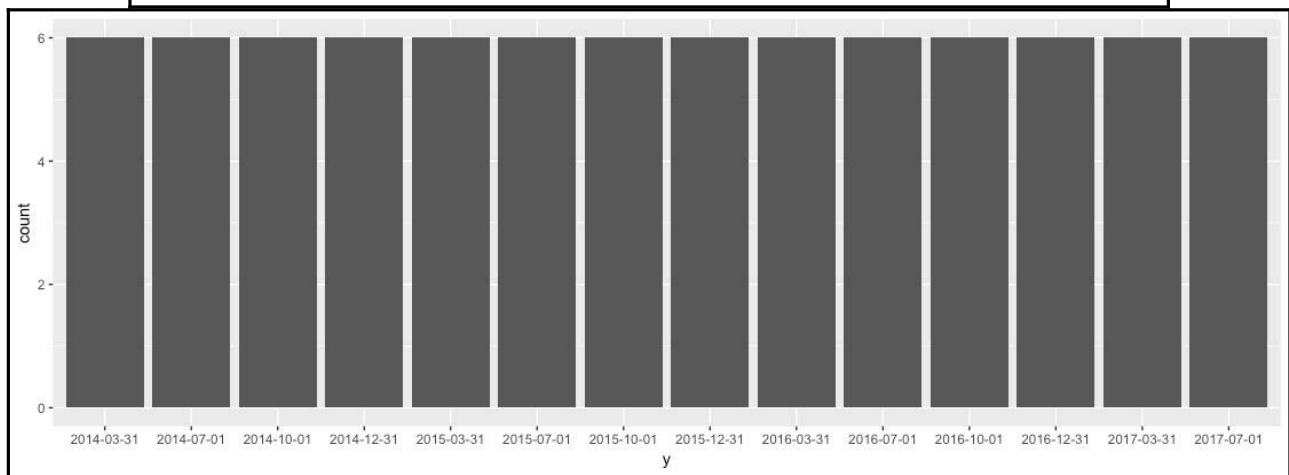
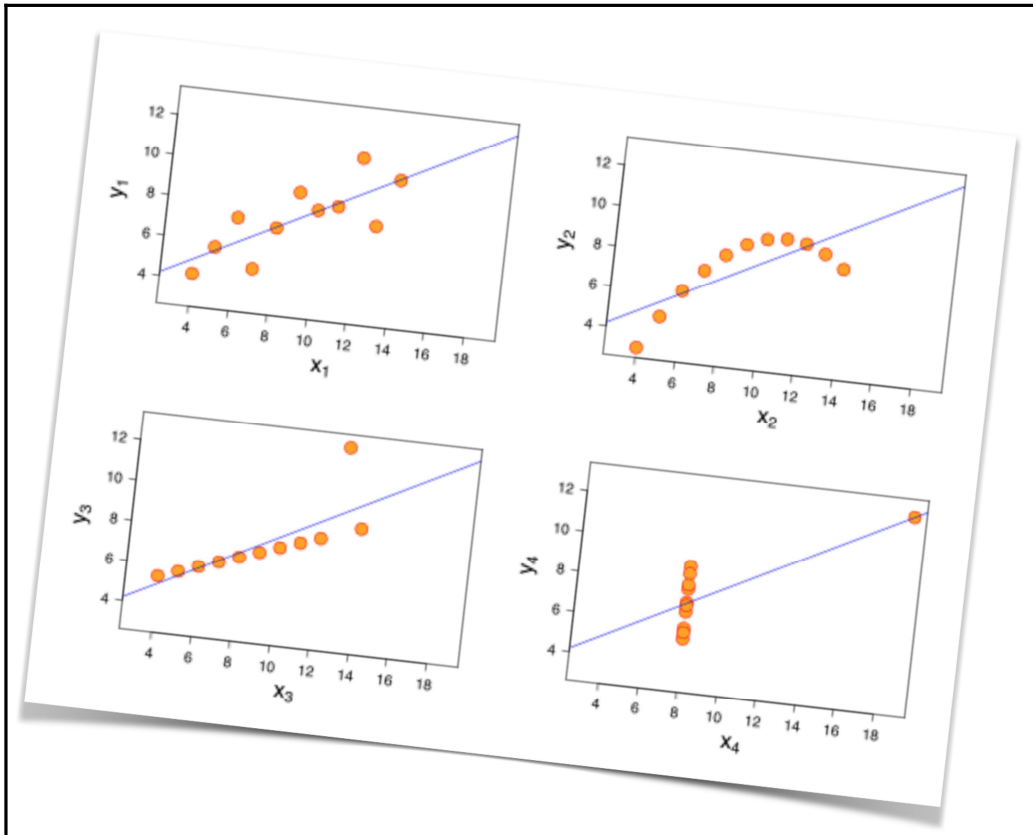


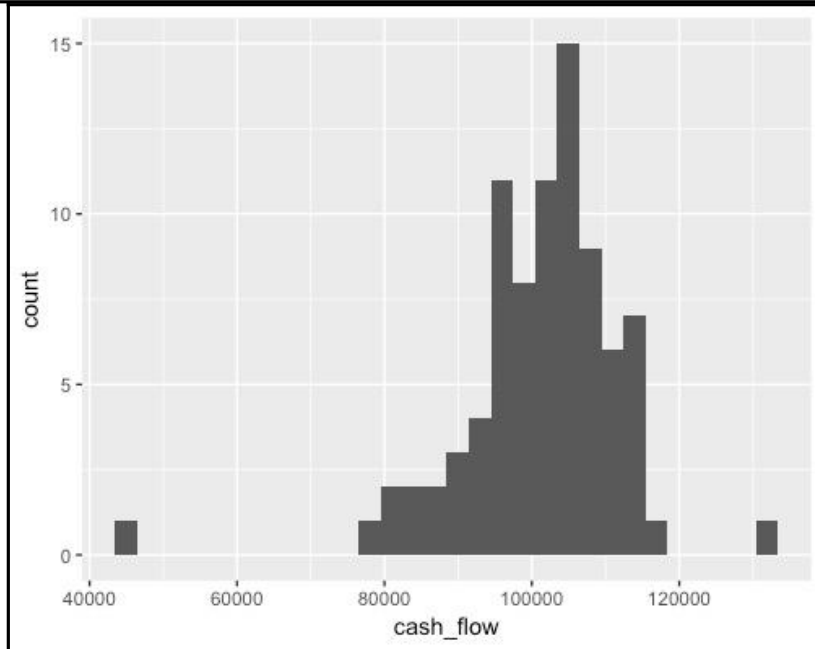
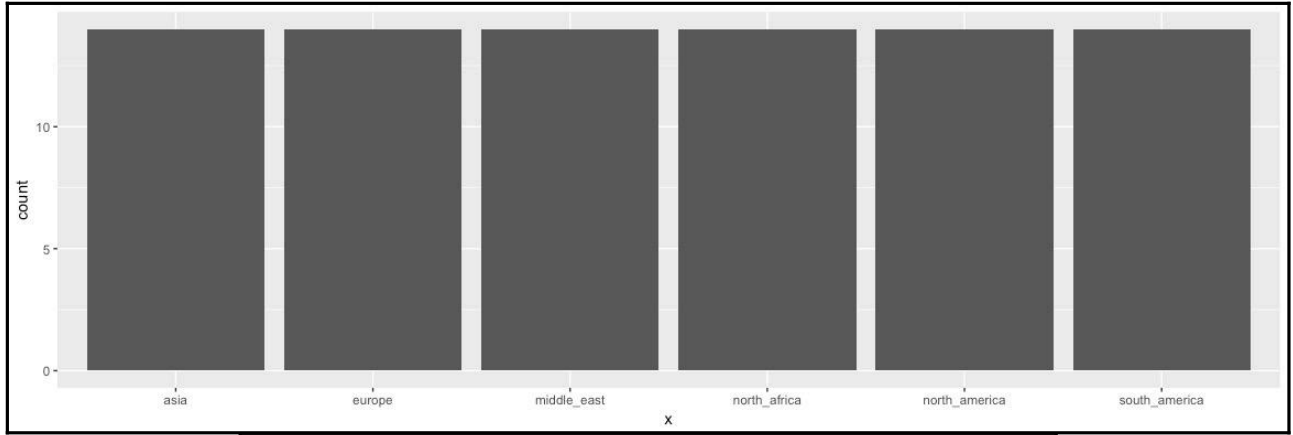
Chapter 5: How to Address a Data Mining Problem – Data Cleaning and Validation

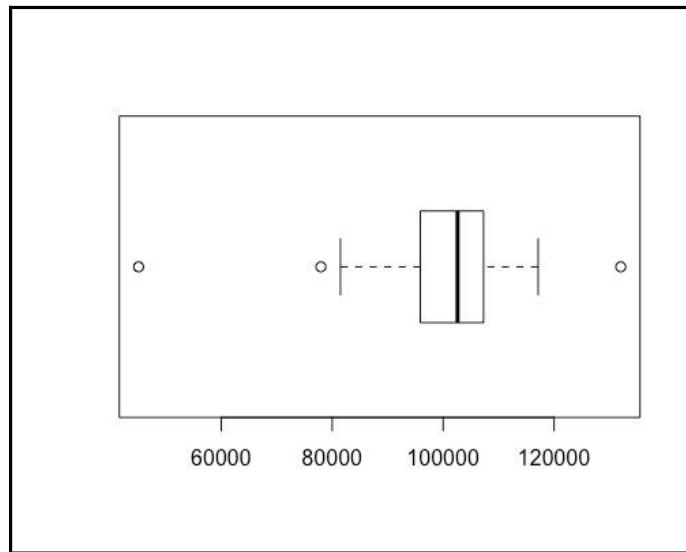
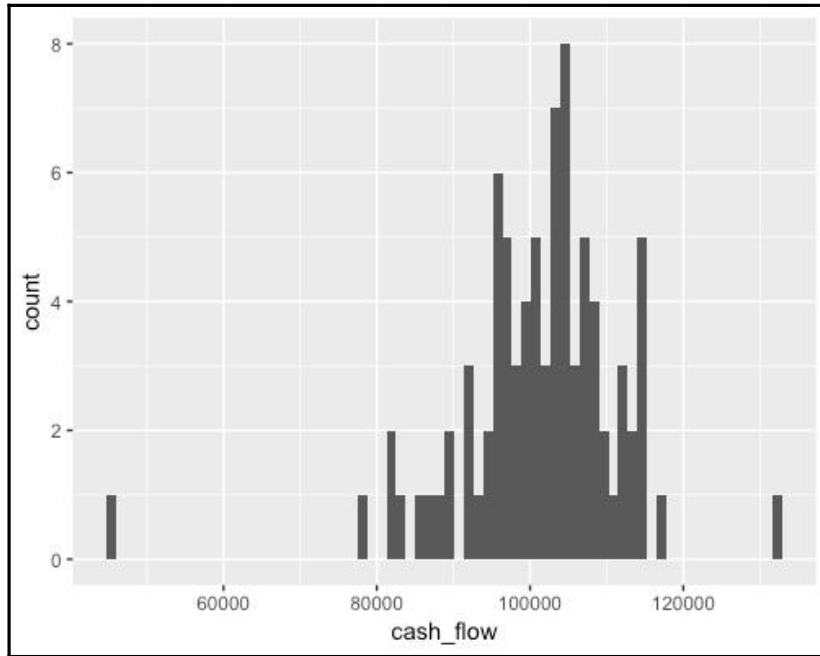


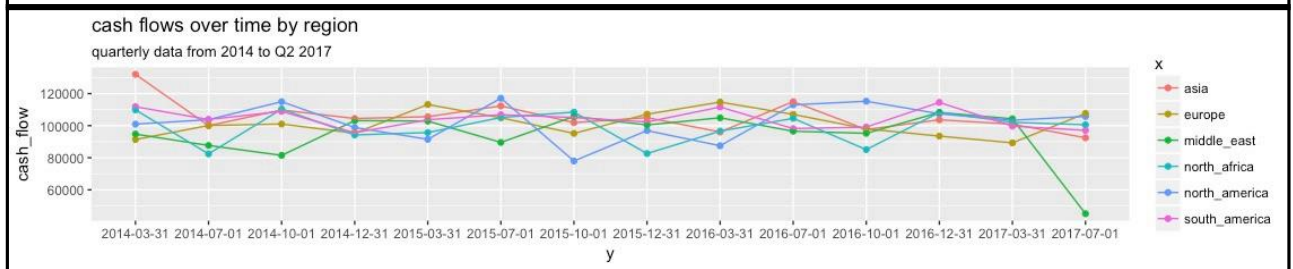
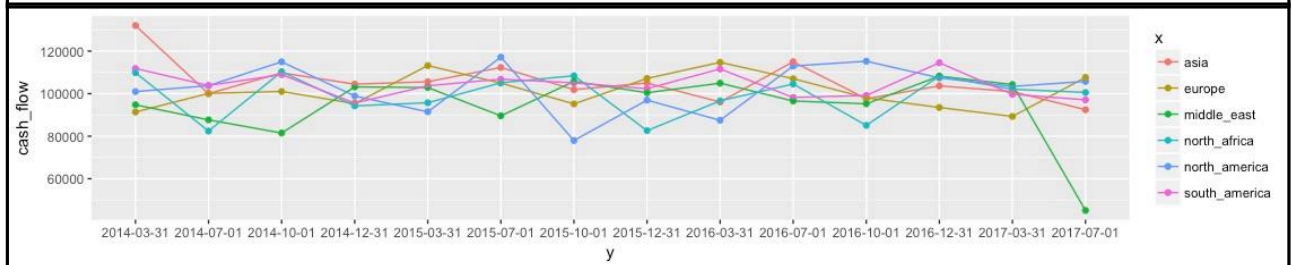
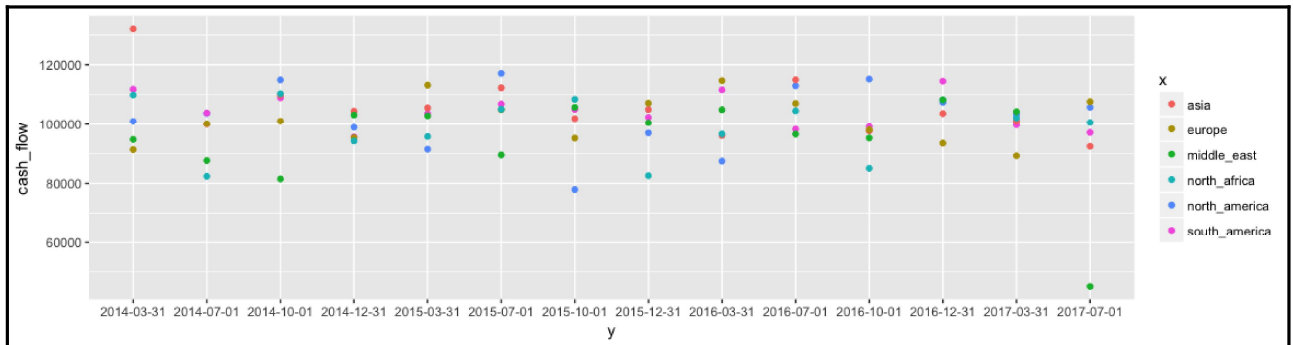
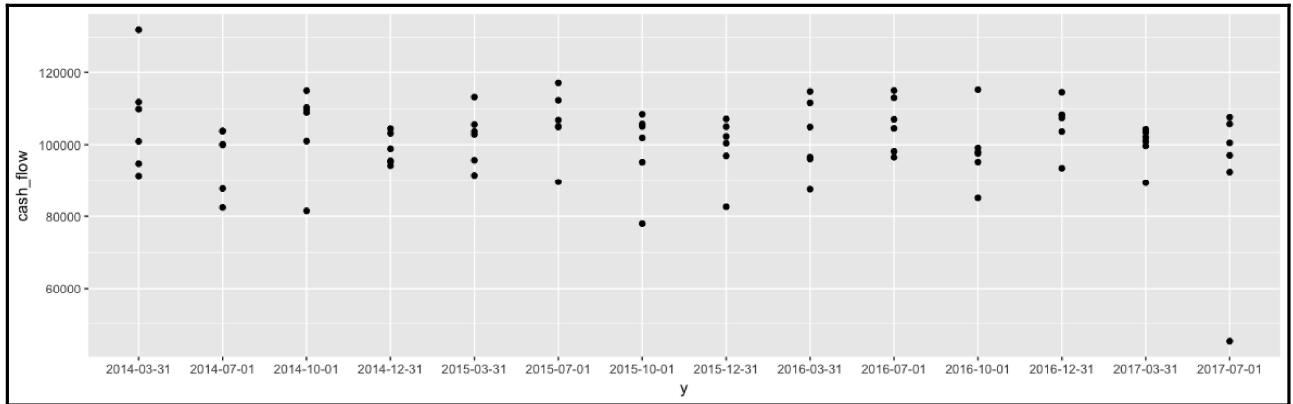
Chapter 6: Looking into Your Data Eyes – Exploratory Data Analysis

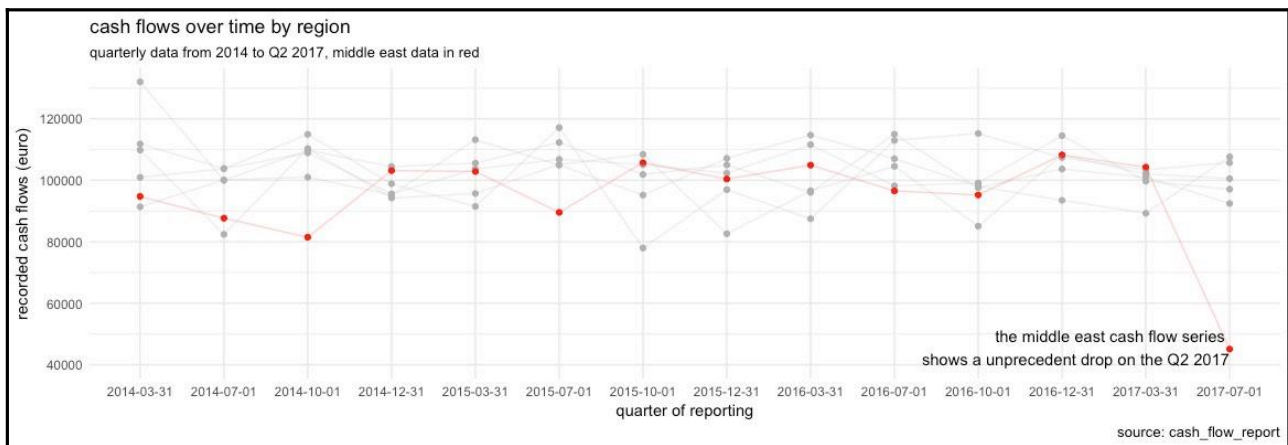
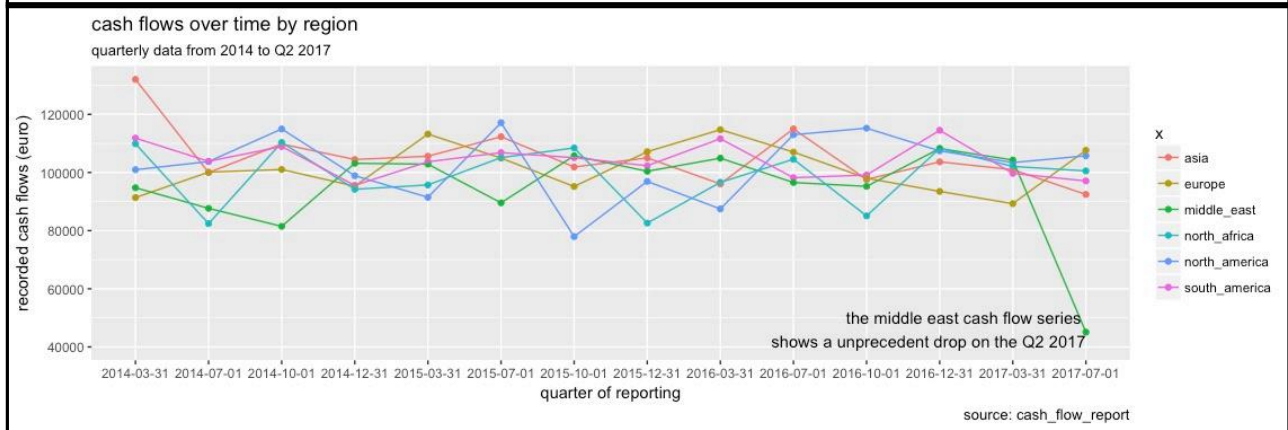
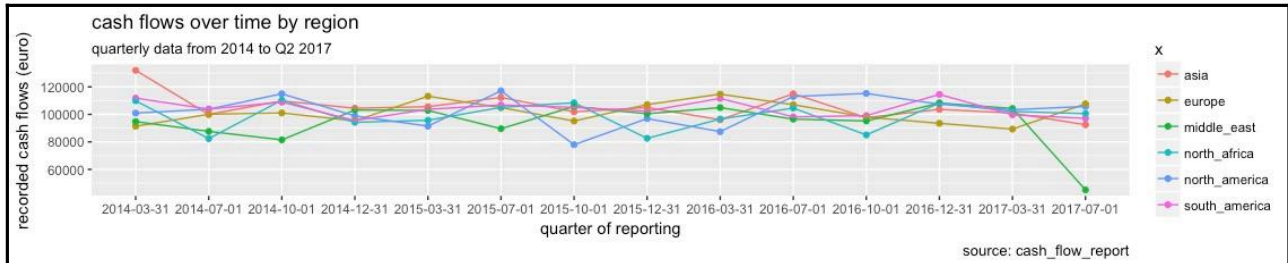




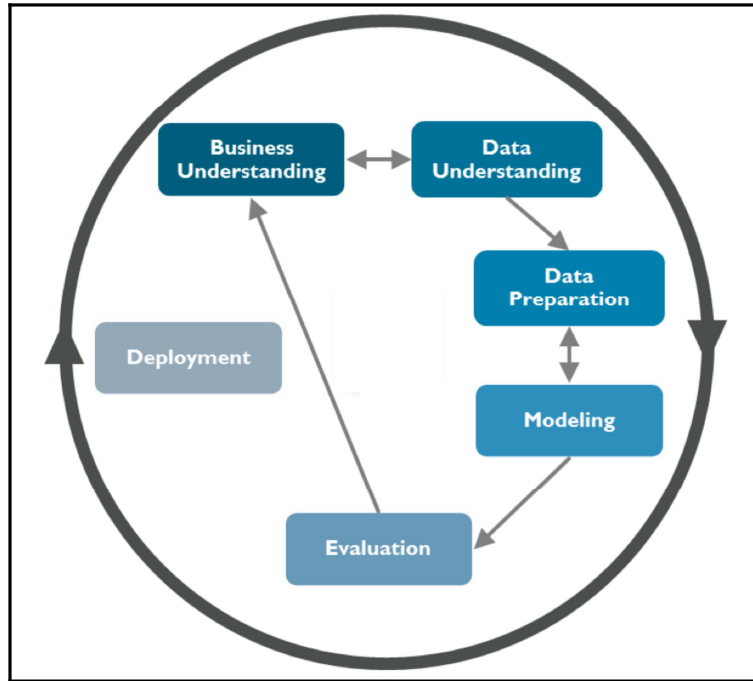


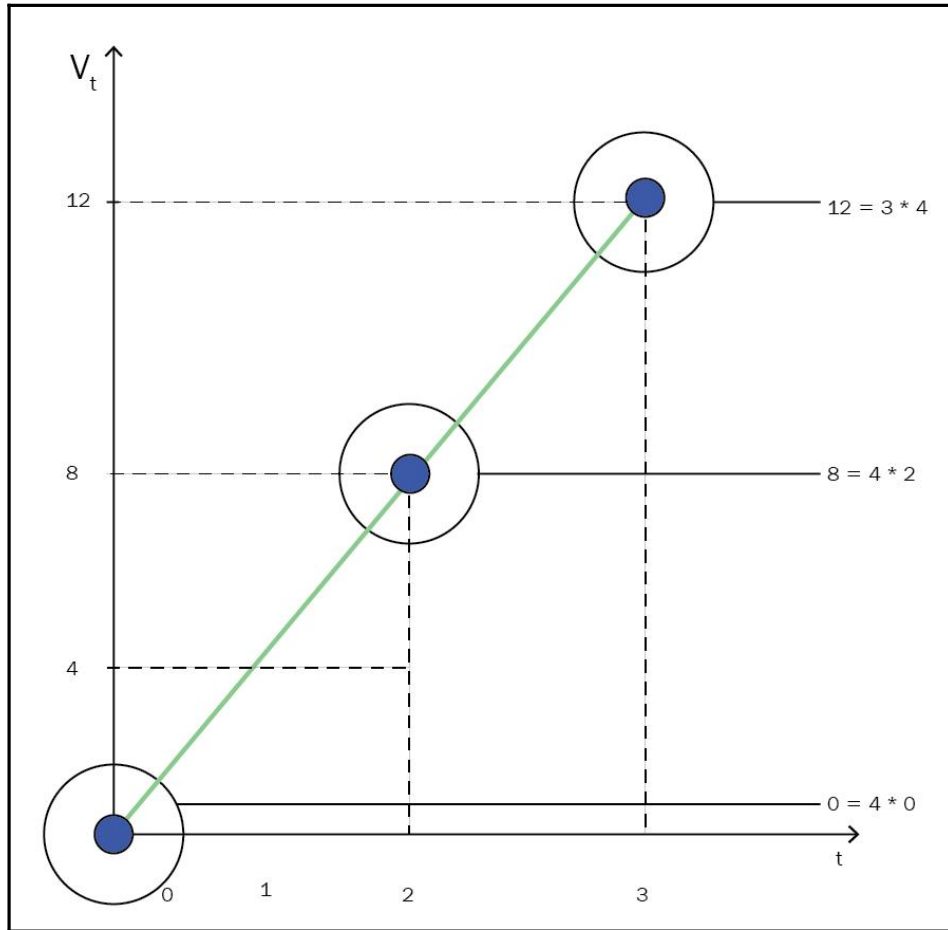


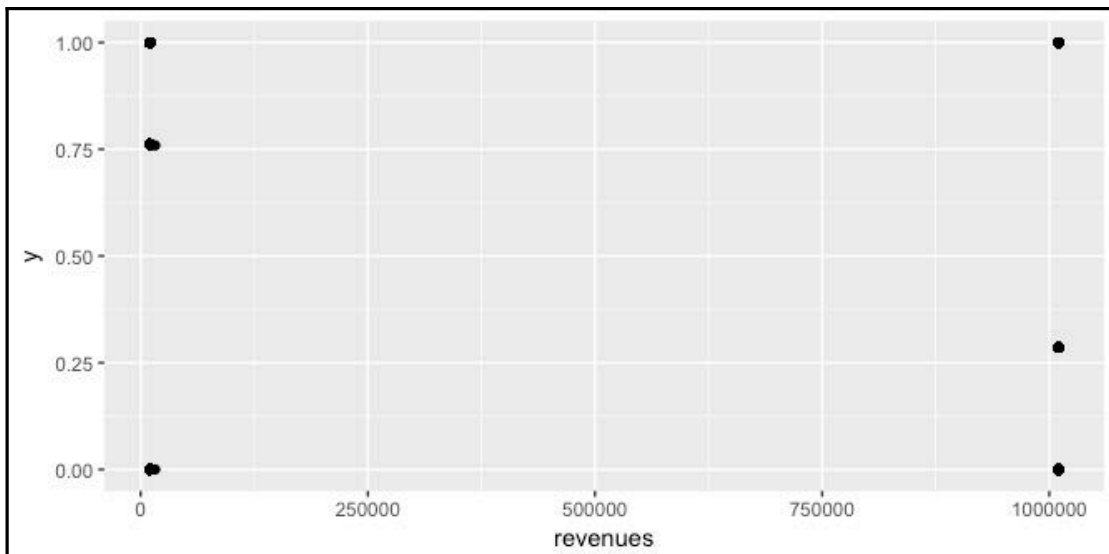
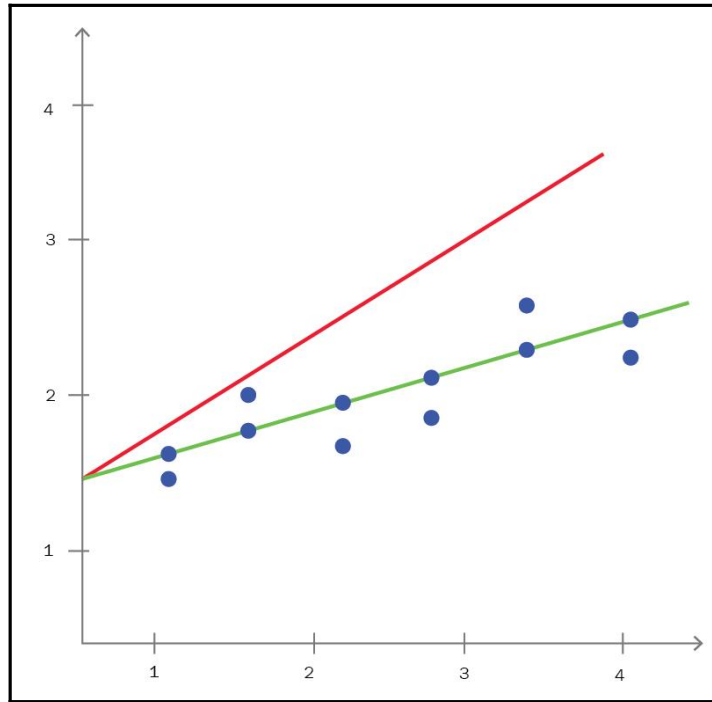


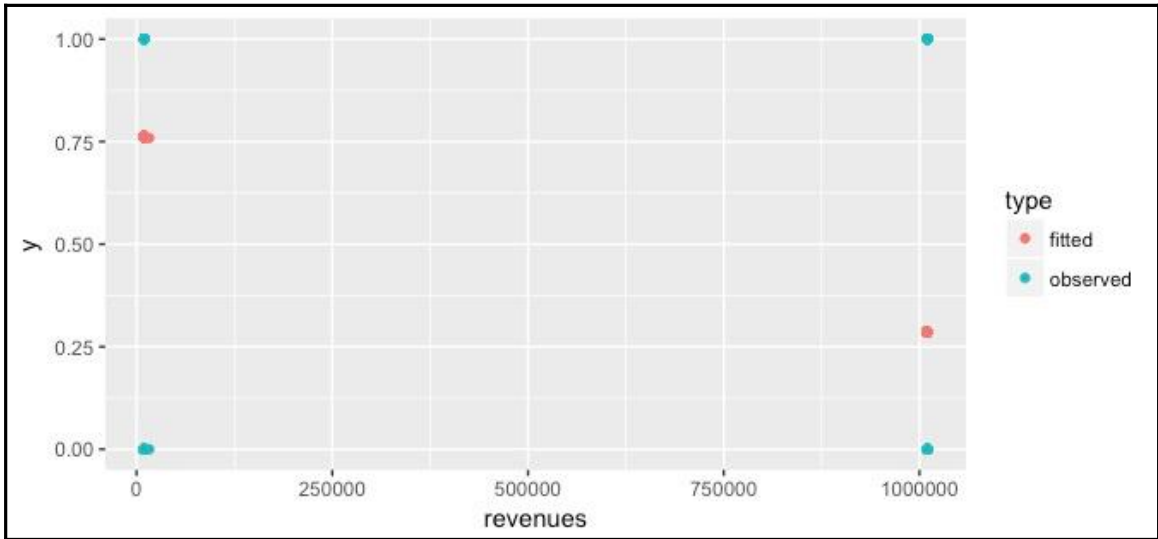


Chapter 7: Our First Guess – a Linear Regression

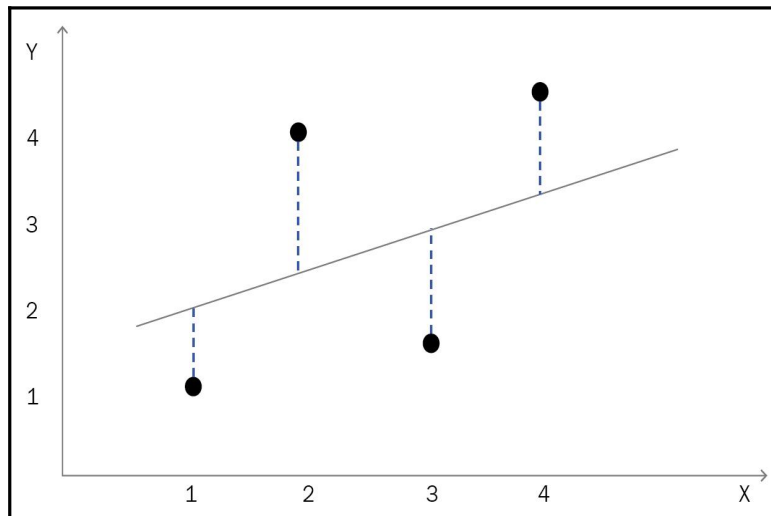
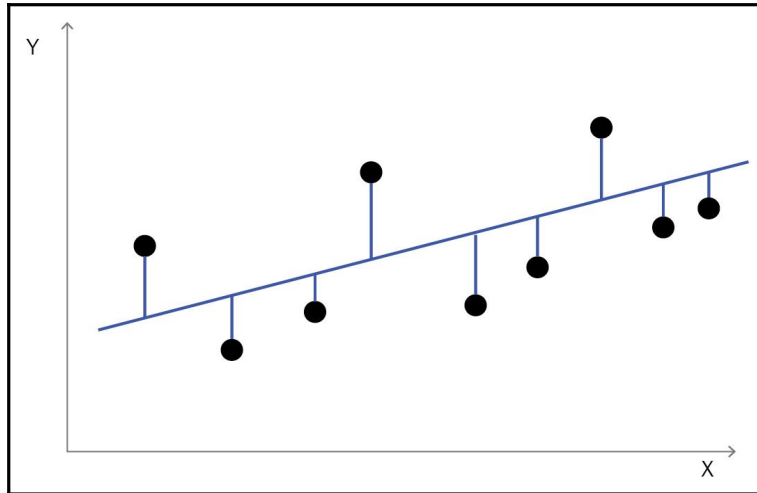


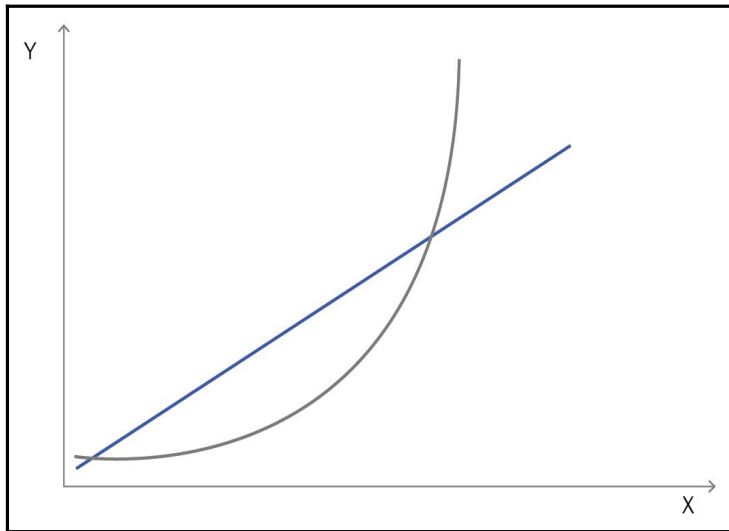
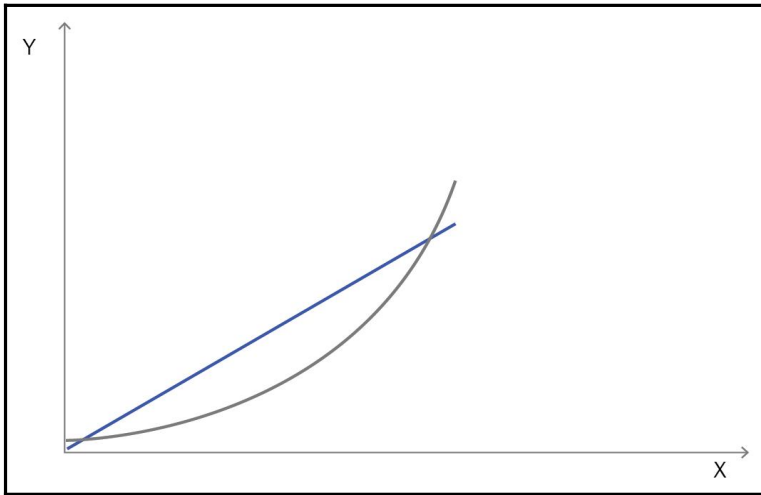




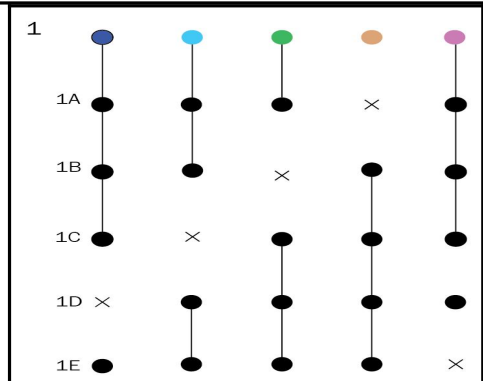
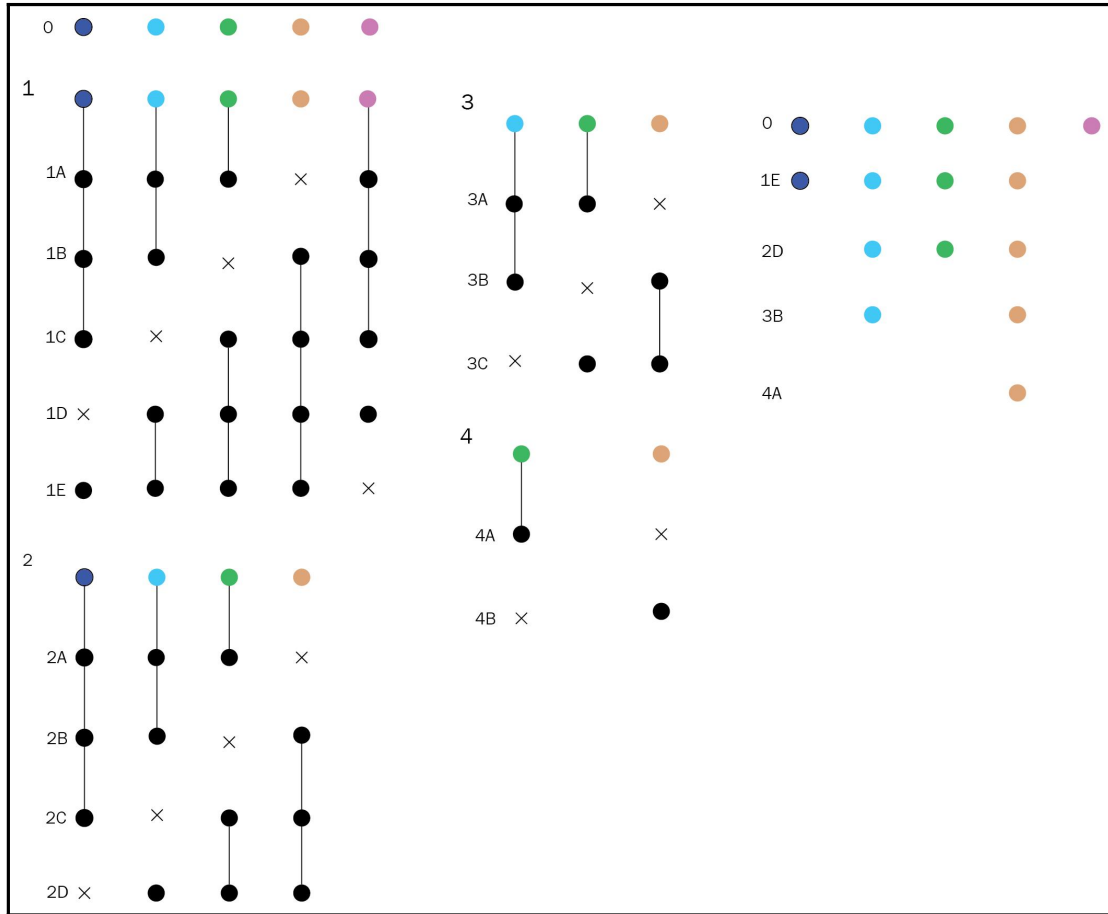


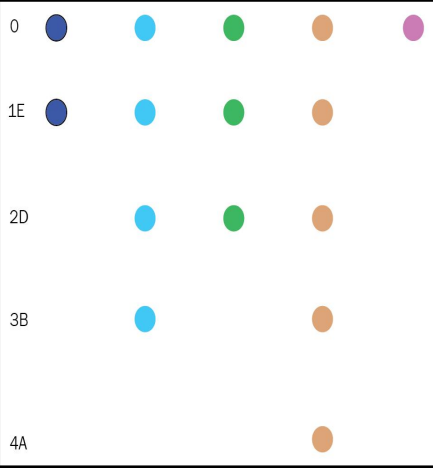
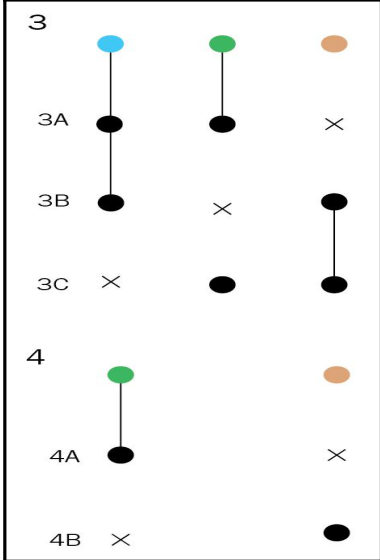
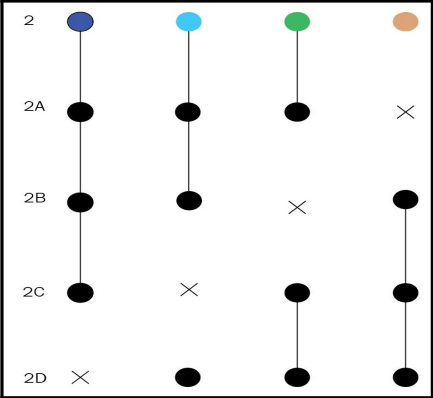
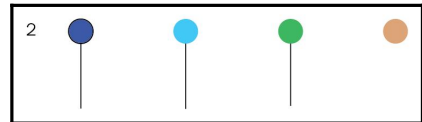
Chapter 8: A Gentle Introduction to Model Performance Evaluation

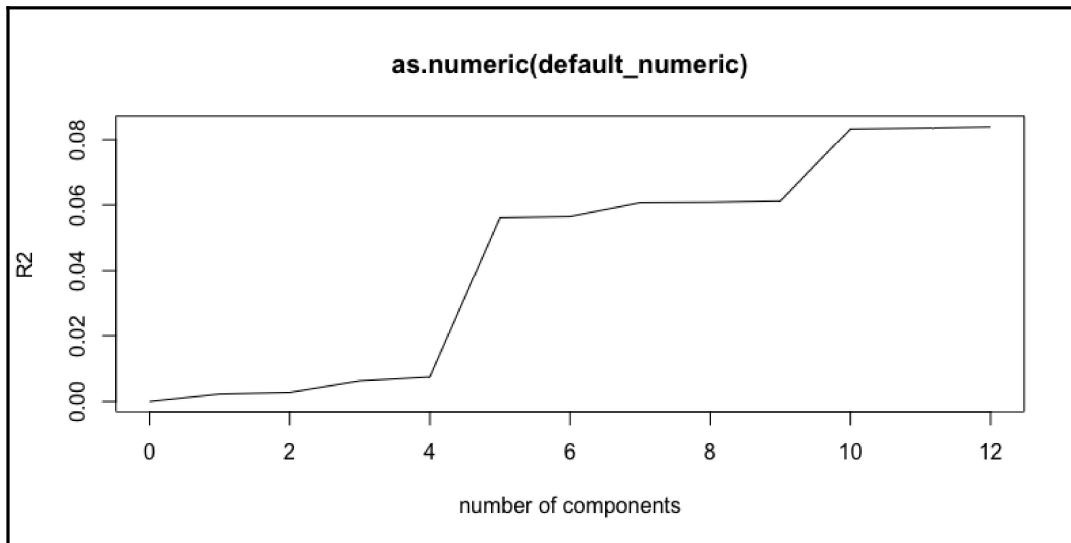
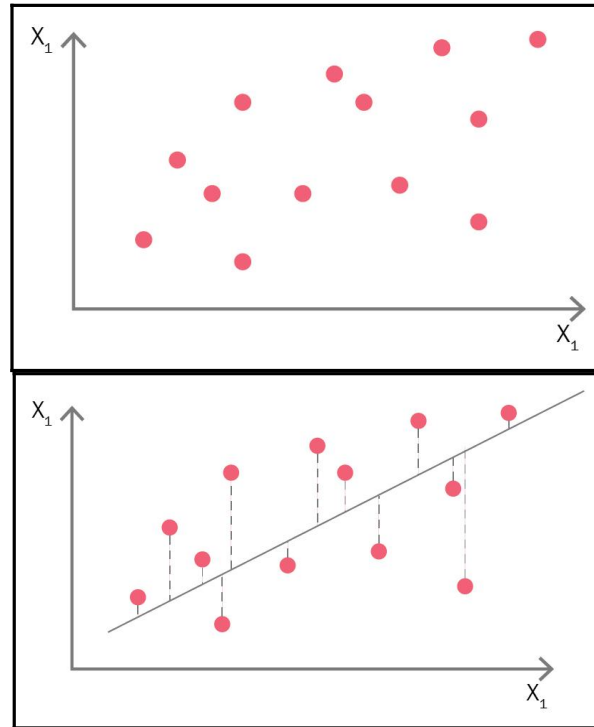




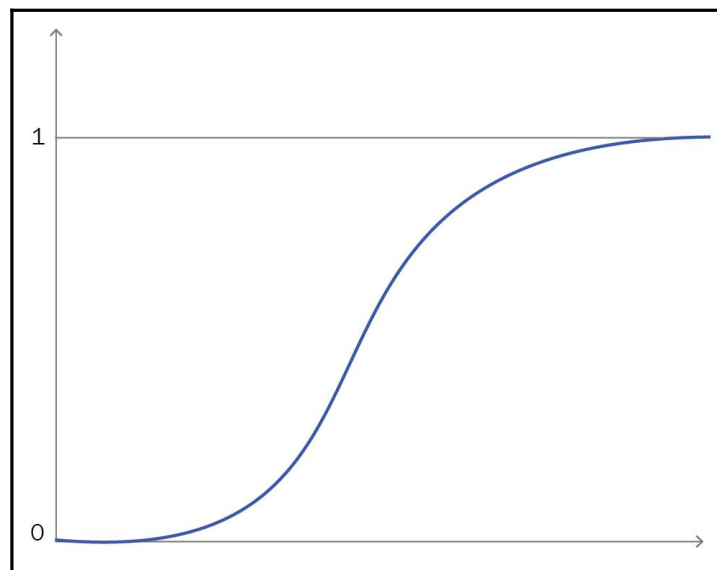
Chapter 9: Don't Give up – Power up Your Regression Including Multiple Variables

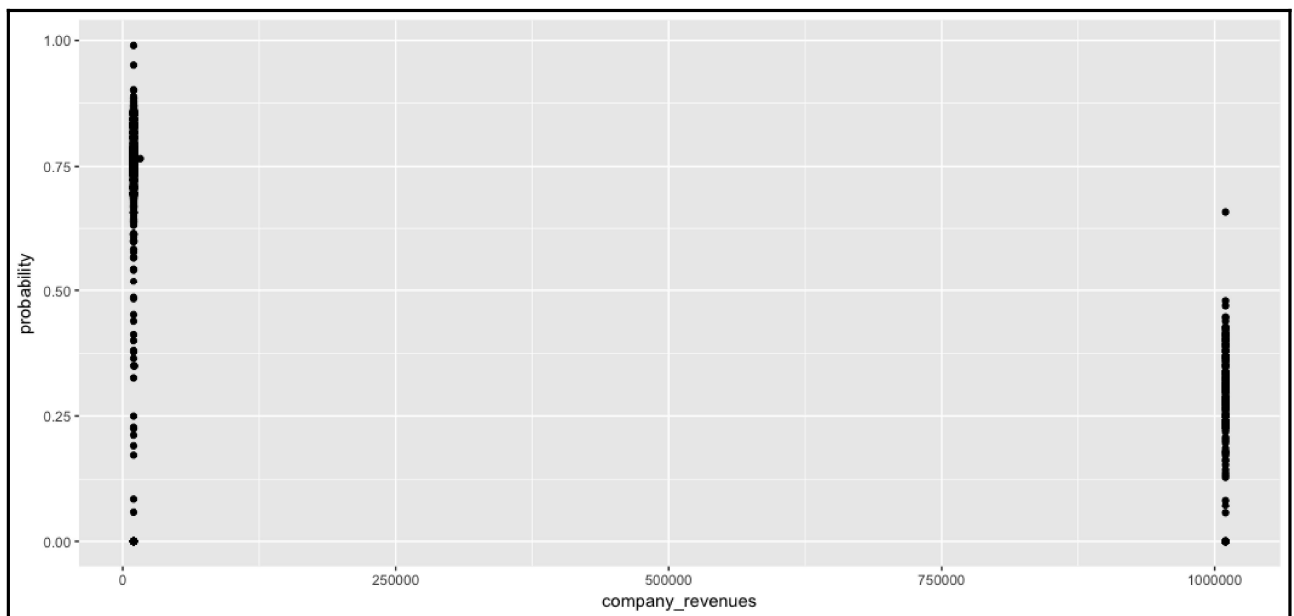
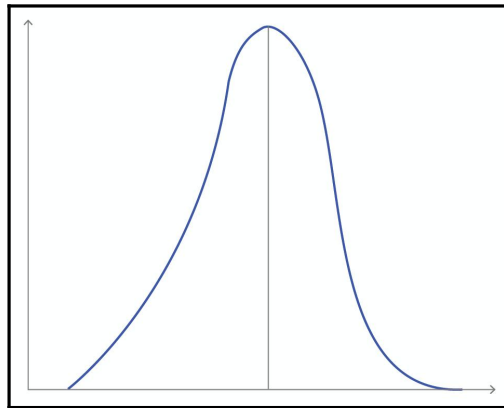
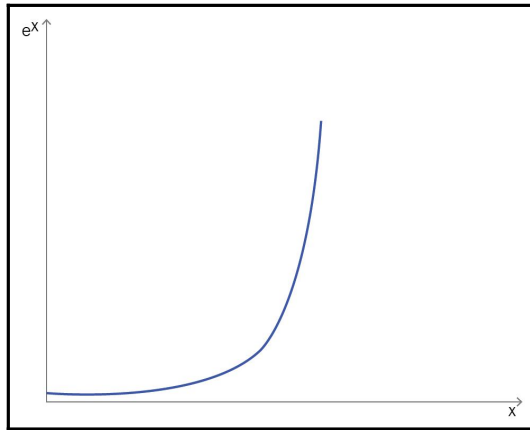


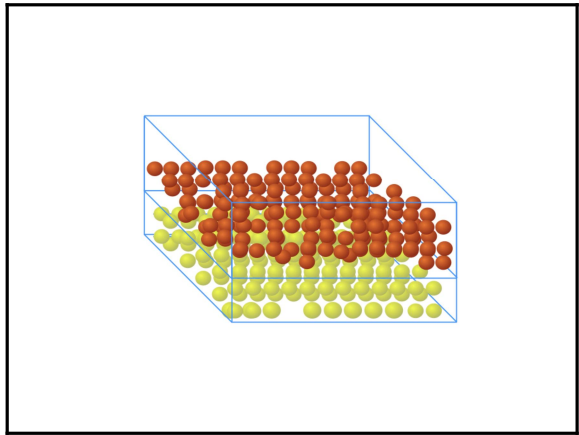
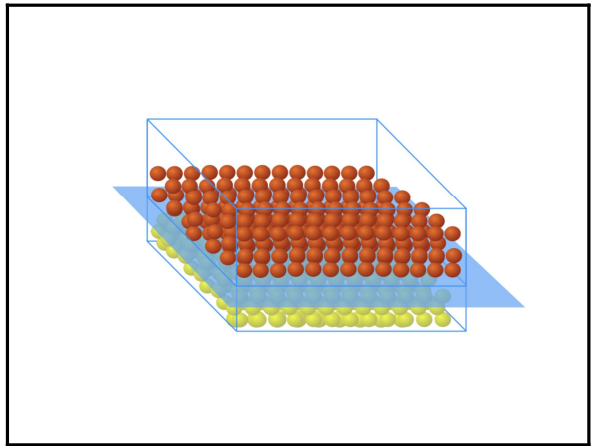
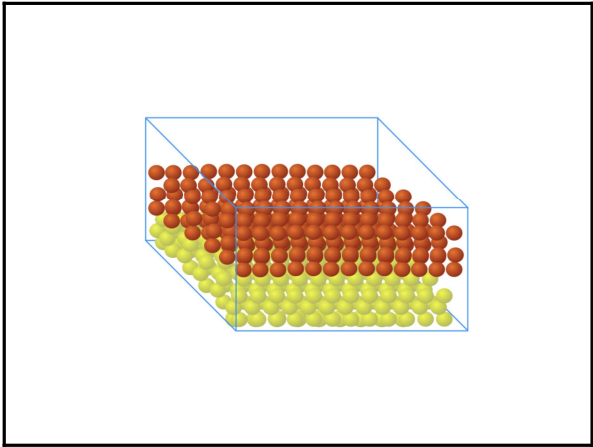


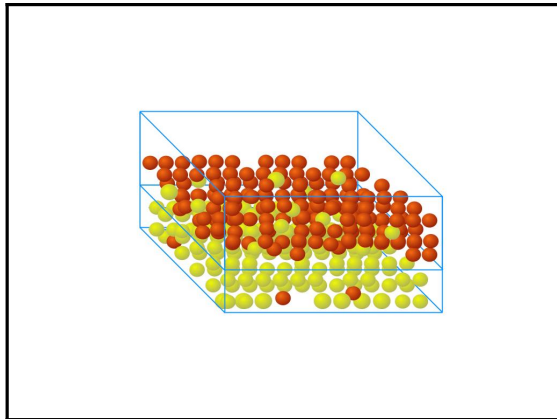
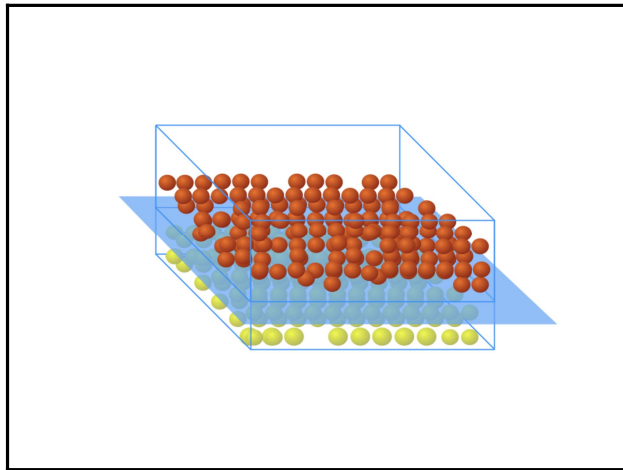


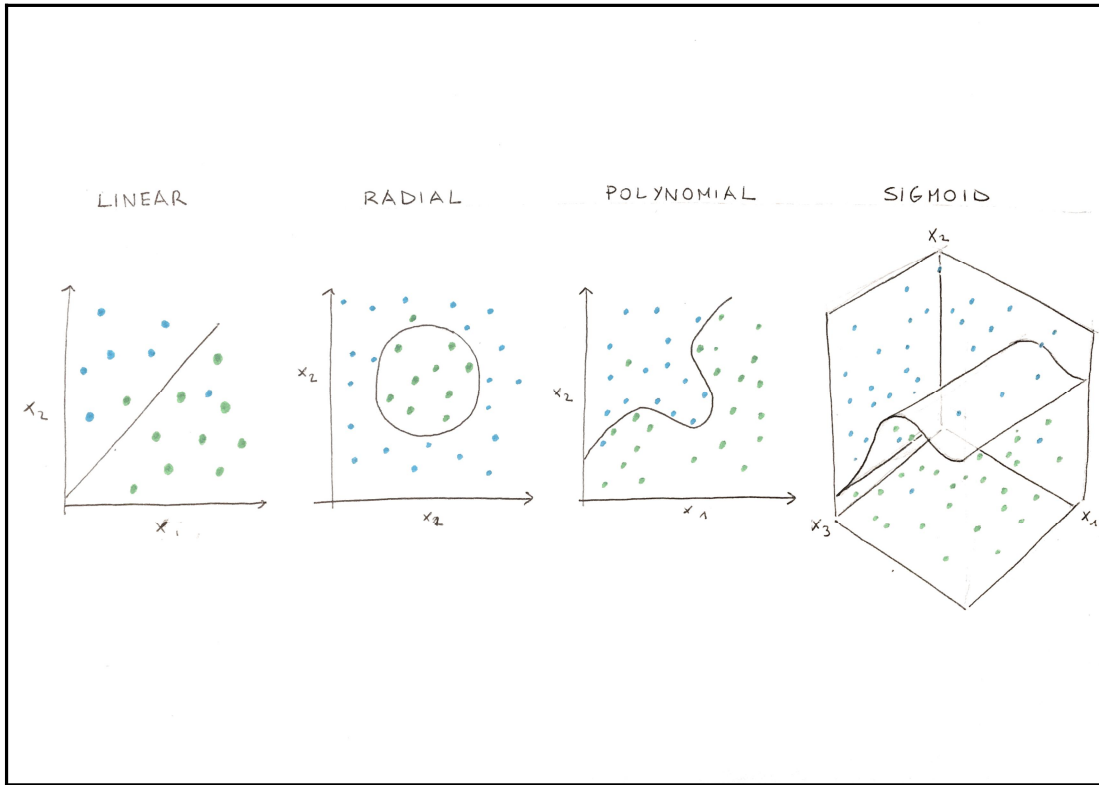
Chapter 10: A Different Outlook to Problems with Classification Models



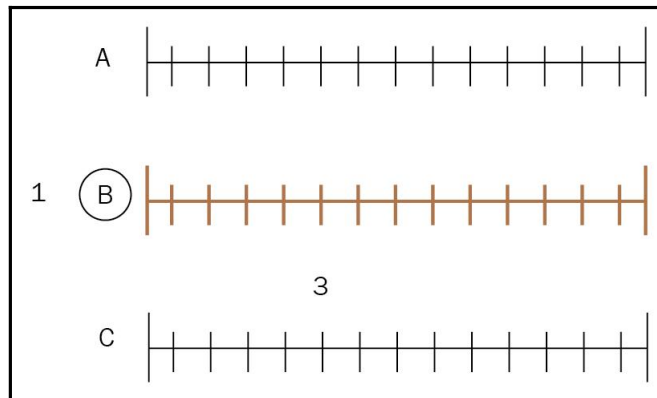


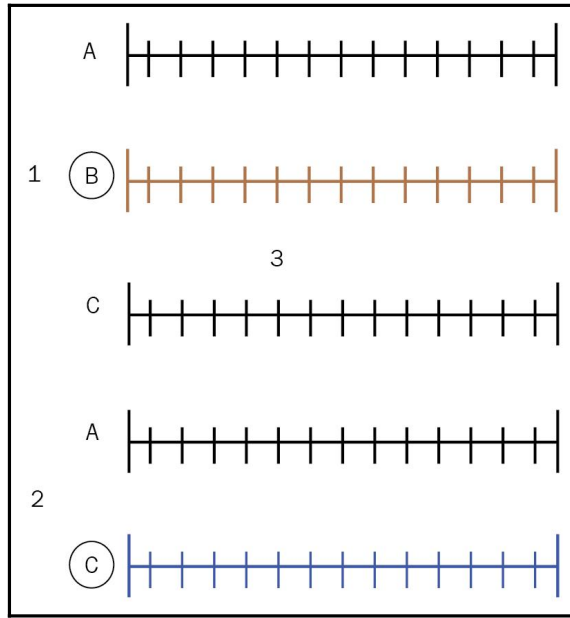


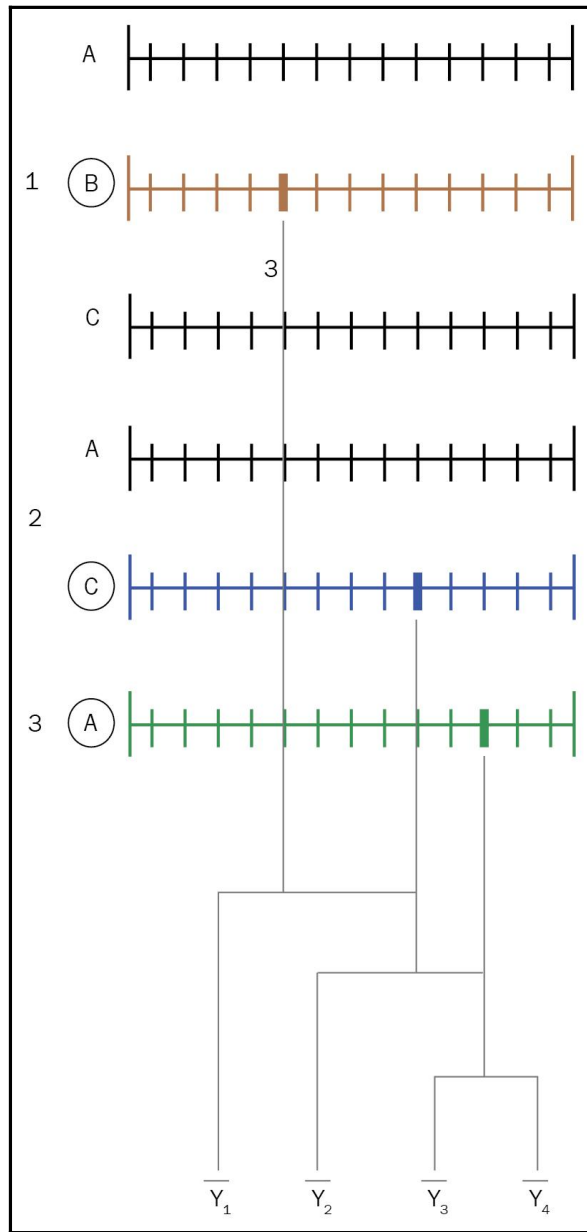


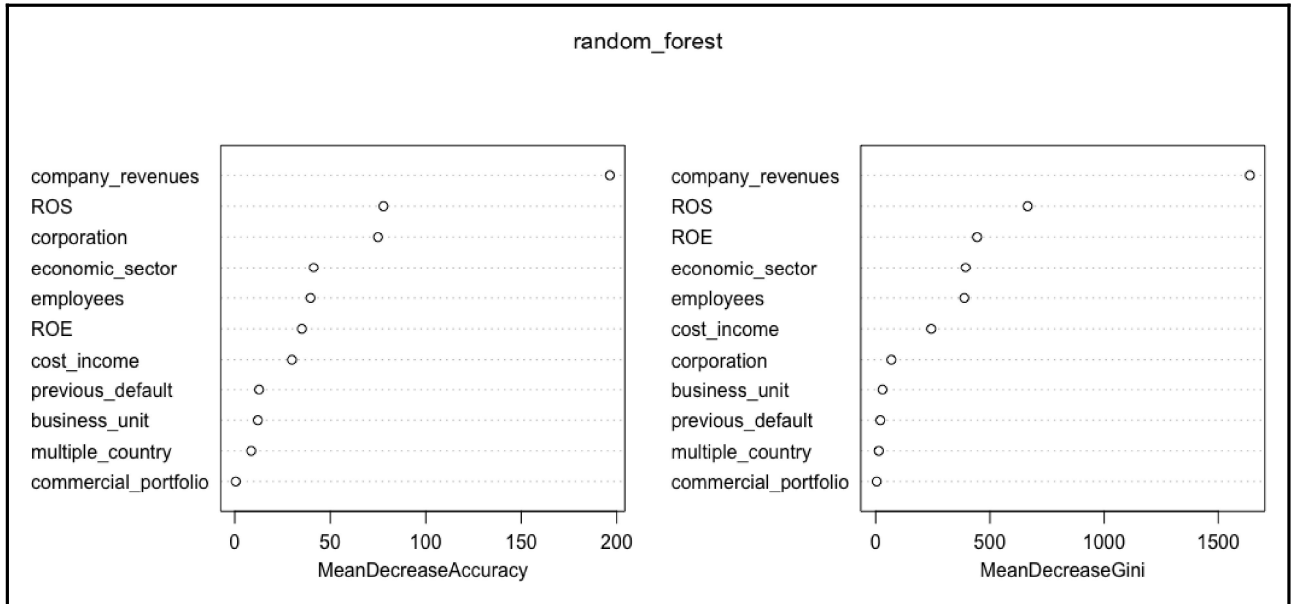
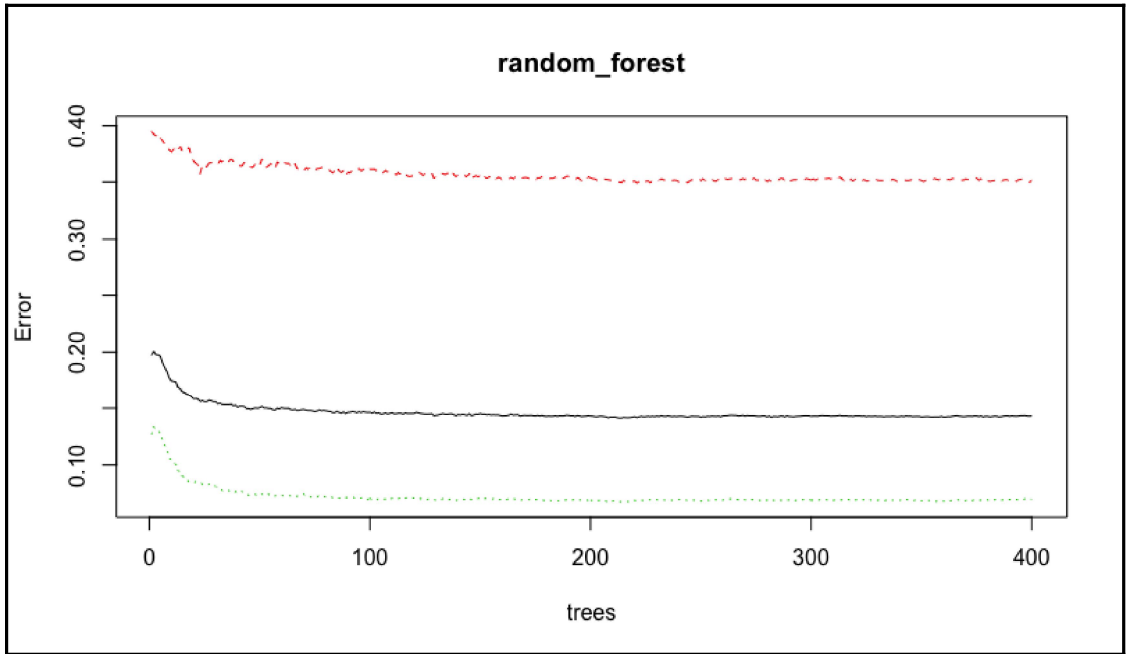


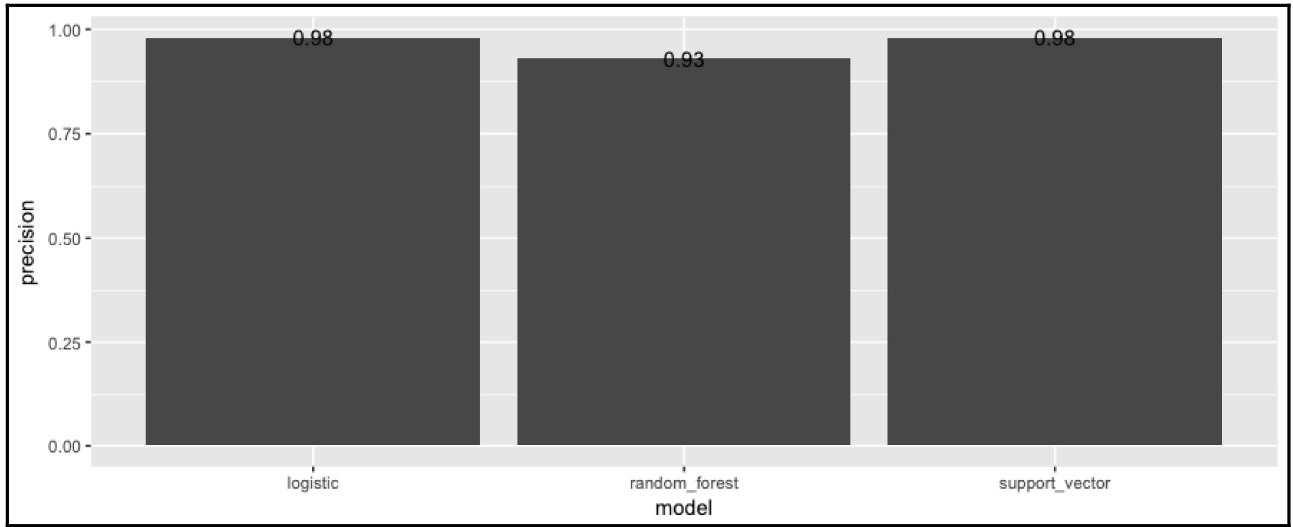
Chapter 11: The Final Clash – Random Forests and Ensemble Learning











Chapter 12: Looking for the Culprit – Text Data Mining with R

