Practical 2: Plotting

Complete this sheet as you work through it. If you run into problems, then ask for help - don't skip sections! Open Rstudio and store any files you download or create in a directory called Practical2.

1 The data

Downloading the data

In today's practical, we are going to examine the big bang data that was collected in Monday's lecture. Use the library command to load the mas1343 package:

```
> library(mas1343)
```

> bb1415=read.table("bb1415.txt",header=TRUE)

Now to load the data from Monday's quiz, type the following:

> bb1415 = read.table("http://www.mas.ncl.ac.uk/~nlf8/ + teaching/mas1343/prac/bb1415.txt",header=TRUE)

Do not enter the + sign. I have done this because the code would not fit on a single line in this handout – you should enter the code in a single line!

This downloads the data from the MAS1343 webpage and stores it in the dataframe bb1415. The argument header=TRUE tells **R** to use the column headings given in the file.

Provided this has all worked smoothly, you should now be able to view the data by typing:

> bb1415

However, the file might be too big for you to view conveniently on the screen. You can get an idea of what the data looks like by using the head function, as was demonstrated in the very first lecture:

```
> head(bb1415)
```

Each row in the dataframe corresponds to a team from Monday's lecture, with their guesses of the ages of the actors/actresses in each of the columns labelled with a character's name. "Watched" takes the value TRUE if most of the team watch The Big Bang Theory; "Gender" takes the values "Male", "Female" or "Both", depending on the make-up of the team. Before you move on, please make sure you familiarise yourself with the dataframe in **R**!

Initial investigations

Since we are going to alter the data set, it's a good idea to create a copy, so assign the dataset to another object – say d:

> d = bb1415

We can inspect the column names using

```
> colnames(d)
```

and we can also look at individual columns, using

> d\$Name

> d\$Shel

The dimensions of d are found using

> dim(d)

Complete: Number of columns: _____ Number of rows: _____

2 Scatter plots

2.1 The basics

Let's start with a scatter plot for Group ID against the Sheldon guesses:

```
> plot(d$ID, d$Shel)
```

The default x- and y-axis labels aren't very good, so we use xlab and ylab arguments, i.e.

> plot(d\$ID, d\$Shel, xlab="Group ID", ylab="Guess")

We can highlight the actual age using abline, i.e.

> abline(h=41)

Other commands for abline are

> abline(h=37, col=1)
> abline(h=38, col=2)
> abline(h=39, col=3,lty=2)

Complete: What does the col argument do?

Complete: What does the lty argument do? _____

Complete: What does the h argument do?

Complete: What happens if you change the h argument to v?

2.2 Separating by Gender

We will now use scatter plots to investigate differences between male and female guesses. First we need to separate the guesses into different categories. The following piece of code extracts the female teams...

```
> x_female = d$ID[d$Gender=="Female"]
> y_female = d$Shel[d$Gender=="Female"]
```

... which we can then plot

> plot(x_female, y_female, pch="F")

Of course we can add in the xlab and ylab arguments.

Complete: Create new variables x_male and y_male. Add these points to your scatter plot, i.e.

> points(x_male, y_male, pch="M")

Make sure you include all teams in your scatter plot. To extend the x-axis, you may have to explicitly state the x-limits:

> ## Change XX & YY to something sensible
> plot(x_female, y_female, pch="F", xlim=c(XX, YY))

Complete: What happens if you omit the pch argument?

Complete: What happens if you change pch="F" to pch=1 or pch=2?

Complete: Do you think that males guess the age of Sheldon better than females?_____

3 Summary Statistics

Use the commands mean, median, quantile and sd to calculate summary statistics for the Sheldon column. *Hint: look at Chapter 3 in your lecture notes!*

Complete: Mean: _____, Median: _____

Complete: Standard deviation: _____, Q1: ____, Q3: _____,

4 Histograms

We will now investigate the distribution of guesses using histograms. To plot a histogram of guesses for Sheldon we use the following code – use the xlab and ylab commands to add labels to the axes, and also the main command to add a main title to your plots.

> hist(d\$Shel)

We can also add in vertical lines to show the mean and median of guesses:

```
> abline(v=mean(d$Shel))
> #What does the lty argument do?
> abline(v=median(d$Shel), col=2, lty=2)
```

We can generate two histograms on the same plot using the par command, e.g.

```
> par(mfrow=c(1, 2))
> hist(d$Shel, main="Sheldon Guesses")
> abline(v=mean(d$Shel), col=2)
> hist(d$Penny, main="Penny Guesses")
> abline(v=mean(d$Penny), col=2)
```

Complete: How are the graphs different? ______ *Hint: Think about summary statistics.*

Usually **R** makes a good initial guess when constructing histograms, but we often have to vary the number of bins to achieve the best result. We can do this using the breaks argument, i.e.

```
> hist(d$Shel, breaks=5)
> #You can also use the "scott" or the "fd" options as in lectures
> #Typically, these are better than the default.
> hist(d$Shel, breaks="scott")
> hist(d$Shel, breaks="fd")
```

Which rule would you use for the Shel data set – i.e. Default, Scott or FD?

5 Boxplots

Let's now investigate the error in age guessing:

|actual age - estimated age|

Sheldon's true age is 41, so the absolute error is:

> shel_err = abs(41-d\$Shel)

To do a boxplot of the error, we use the command:

> boxplot(shel_err)

We can also investigate if people who have watched BB have an advantage when guessing:

```
> #Remember you can use xlab & ylab for labels!
> boxplot(shel_err ~ d$Watched)
> boxplot(shel_err ~ d$Watched, col="bisque") #In colour
```

Do you think watching BB increased guessing precision?

Assignment

Before you begin

Points to note:

- Remember to label your axis correctly.
- For this practical, copy your graphs into Word and type up your comments.
- If you are asked to submit histograms, try different bin number rules.
- Don't go overboard with the number of decimal places you quote.
- If you have more than five pages you will lose marks.
- If you have fewer than five pages you will lose marks.

What I want

The marks for each question are indicated below.

- Page 1: On separate lines I would like:
 - Module Name MAS1343
 - My name Dr Lee Fawcett
 - The practical number *Practical 2*
 - Your name: surname and initial(s). For example, Smith, J.H.
 - Your student id. For example, *b9123456*
- Page 2 [25%]: Choose a Big Bang Theory actor, and plot a scatter plot of Group ID vs Guess.
 - Add a line to indicate the mean guess.
 - Add two other lines to represent the mean ± 2 standard deviations. *Hint: use* sd *to get the standard deviations, then use* abline.
 - Use the lty argument to distinguish between the lines.
- **Page 3** [25%]: Produce a boxplot of the error against Gender for the Big Bang actor you previously chose (not Sheldon!). Do you think that females are better than males at guessing? What about mixed sex groups?
- **Page 4** [25%]: Produce histograms of the error for four different characters. Put the histograms on the same plot. Use the **abline** command to indicate your group's guess on each of the plots. Your scales on the four histograms should be the same.
 - If you can't remember your group, choose the first group on the list.
- Page 5 [25%]: Guesses from two years ago can be found by typing:
 - > data(bb1213)
 - Using some graphs and summary statistics compare this year's guesses to those from two years ago. Do you think the class in 2012/13 are better at guessing than the this year's class? I don't went every possible plot and summary statistic. Just a brief overview.