# R Programming: Worksheet 2

After today you should know:

- the differences between a data frame and a matrix;

- how to access parts of a data frame using `subset()`;

- how to alter the elements of a data frame, including its row names and variable names;

- what a factor is, and how to create one from a character vector;

- how to get boxplots using factors.

1. **Data Frames**

   Load the `MASS` library and take a look at the dataset called `survey`.

   ```
   > library(MASS)
   > head(survey)
   ```

   You can look at the documentation:

   ```
   > ?survey
   ```

   and get a brief summary of each variable:

   ```
   > summary(survey)
   ```

   (a) Find the mean pulse rate of the students. What goes wrong here?

   The vector of pulses stored in `survey` contains some entries which are labelled `NA`. This is used in `R` to represent **missing data**.

   (b) Try looking at the documentation for `mean()` to see how to get around this.

   (c) The ages are recorded as fractions representing a number of months. Change that columns of the data frame so that it contains whole years (the `floor()` command may be useful).

   (d) Find the mean pulse rate for students under 20.

   **Subsetting.** Suppose I want to obtain the records of students who are over 190 cm tall. Since data frames allow me to subset just like a matrix, I might just think of typing:

   ```
   > survey[survey$Height > 190, ]
   ```

   What goes wrong here? Try instead the following:

```
> subset(survey, Height > 190)
```

The subset function ignores missing values, which is usually the behaviour we would prefer. We can also select only some of the fields of the data frame if we prefer:

```
> subset(survey, Height > 190, select = c("Pulse", "Clap"))
```

Recall that the & operator does a point-wise logical 'and' comparison.

```
> subset(survey, (Pulse > 70) & (Smoke == "Heavy"))
```

Similarly, | is for 'or', and ! for 'not'.

```
> with(survey, (Pulse > 70) | (Pulse < 45))
> !(survey$Age > 30)
```

(e) Find the mean age of students who write with their right hand.

(f) What proportion of left handers do not clap with their left hand on top?

(g) Using the **plot()** command, plot the pulse of the subjects against their age. Try subtracting 10 from the age and taking the logarithm (using the function **log()**), to obtain a slightly clearer picture.

2. **More Manipulation**

Load the **hills** dataset from lectures. What happens when you apply the **plot()** command to this data?

```
> plot(hills)
```

You can access the row names of the **hills** data with the function **rownames()**. You can also **change** the rownames this way:

```
> rownames(hills)[1] <- "Redmantle"
> rownames(hills)[1]
```

If at any point you want to revert to the original data frame, type

```
> rm(hills)
```

You can't really edit the original data set in the **MASS** package, only your own copy. Calling **rm()** deletes the copy.

(a) The race climbs are measured in feet. Change the data set so that they are measured in metres (there are 2.54 cm in 1 inch, and 12 inches in 1 foot). Use the **round()** command to set them to the nearest 10 metres.

(b) Obtain a logical vector indicating which races were longer than an hour. How many are this long?

The **which()** command is a useful way to turn a logical vector into a numerical one. For example, try

```
> which(hills$dist < 5)
```

(c) Use `which()` to select the first three races under 5 miles.

(d) Identify (with an `R` command, not just by looking) which of the races involves `"Meall Ant-Suidhe"`. I've no idea how to pronounce Meall Ant-Suidhe, so remove it from the data frame (this might not normally be considered good statistical practice).

(e) Plot the race times against distance. Can you spot an outlier? (think carefully)

(f) In fact there is a race whose time was recorded as one hour larger than it ought to have been. Correct this.

You can turn the data frame into a matrix like this:

```
> hillsMat = as.matrix(hills)
```

(g) Print the two objects `hills` and `hillsMat` by typing their names. Can you see any difference?

(h) Use `is()` to see that they *are* different. What happens if you use `plot()` on the matrix? Compare this behaviour with that of the command `pairs()`.

(i) What happens if you try to turn the `survey` data into a matrix?

3. **Factors**

   (Extension of Exercise 3.5 from Lectures)

   Take a look at the `birthwt` data from the `MASS` package.

   (a) How is race stored in these data? Is this sensible?

   (b) Turn this into a factor with level names as indicated in the documentation.

   (c) Use the `table()` command to count the number of babies of each race.

   (d) Try the following command:

   ```
   > tab = with(birthwt, table(smoke, low))
   ```

   What do the results in `tab` suggest?

4. **Reading In Data**

   I have emailed you a dataset called `sprays.dat`. The data represent insect counts in agricultural experiments treated with different insecticides. Save the file onto your local drive.

   (a) Read the data into `R` using the command

   ```
   > dat = read.table("sprays.dat", header = TRUE)
   ```

   Check that the first row is as you expect.

```
> head(dat)
```

What happens if you omit the `header=TRUE` part?

(b) Look at the two variables in this new data frame; what types do they have?difference between them do you notice? Is their format appropriate?

(c) Find the mean number of insects for each different experimental unit.

(d) Look at the documentation for the command `tapply()`. Now repeat the previous question with a single command.

(e) Try using the command `quantile()` on the counts. Use `tapply()` to find the upper and lower quartiles of the data broken down by spray type.

(f) * Can you write a command which gives you just a vector of the six upper quartiles? [Hint: how can you pass the argument `probs=0.75` to `quantile()` when you are using `tapply()`?]

(g) Use the `plot()` command to produce separate boxplots for each level (as we did in class for the height data).