Part IV

Graphical Presentation of Data

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

• Graphical displays of data can be very useful in showing the main features of a data set.

- Graphical displays of data can be very useful in showing the main features of a data set.
- The appropriate form of graph depends on the nature of the variables being displayed and what aspects are to be shown.

▲ロト ▲母 ▶ ▲ 国 ▶ ▲ 国 ● の Q @

- Graphical displays of data can be very useful in showing the main features of a data set.
- The appropriate form of graph depends on the nature of the variables being displayed and what aspects are to be shown.
- However it should always be borne in mind that the object is to provide a clear and truthful representation of the data, not to distort and not to impress with unnecessary "fancy" features.

▲ロト ▲母 ▶ ▲ 国 ▶ ▲ 国 ● の Q @

- Graphical displays of data can be very useful in showing the main features of a data set.
- The appropriate form of graph depends on the nature of the variables being displayed and what aspects are to be shown.
- However it should always be borne in mind that the object is to provide a clear and truthful representation of the data, not to distort and not to impress with unnecessary "fancy" features.

▲ロト ▲母 ▶ ▲ 国 ▶ ▲ 国 ● の Q @





http://www.fusioncharts.com/explore/pie-doughnut-charts





æ

Export von Bananen in Tonnen von 1994-2005



4.1 Label your %#?* axes!



◆□▶ ◆□▶ ◆目▶ ◆目▶ 三目 - のへで

4.1 Label your %#?* axes!



Rant over!

▲ロト ▲団ト ▲ヨト ▲ヨト 三回 - の々で

4.2 Qualitative data: bar charts

- The most useful way to display qualitative data is usually with a bar chart.
- The length of each bar is proportional to the frequency of the corresponding value of the variable in the sample of data.
- Note that the widths of the bars should be equal to avoid giving a false impression.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

MPAA breakdown



Figure: Barchart of the mpaa ratings for 4847 films.

・ロト ・四ト ・ヨト ・ヨト

臣

MPAA breakdown

>	tabl	<mark>e(</mark> movi	es\$mpa	a)
N	2-17	PG	PG-13	R
	16	526	989	3316

▲口▶ ▲御≯ ★注≯ ★注≯ 二注

<pre>> table(movies\$mpaa)</pre>								
NC-17	PG	PG-13	R					
16	526	989	3316					

Inside the barplot function:

>	<pre>barplot(table(movies\$mpaa), xlab="MPAA Rating",</pre>
+	<pre>ylab="Frequency", border = "black",</pre>
+	<pre>col="mistyrose")</pre>



<pre>> table(movies\$mpaa)</pre>								
NC-17	PG	PG-13	R					
16	526	989	3316					

Inside the barplot function:

> barplot(table(movies\$mpaa), xlab="MPAA Rating", + ylab="Frequency", border = "black", + col="mistyrose")

《曰》 《聞》 《臣》 《臣》 三臣 -

Remember to load the data first!

- > library(mas1343)
- > data(movies)

- To represent the distribution of a sample of values of a continuous variable we can use a histogram.
- The range of values of the variable is divided into intervals, known as *classes*, and the frequencies in classes are represented by columns.
- As the variable is continuous, there are no gaps between neighbouring columns unlike a bar chart.

▲ロト ▲母 ▶ ▲ 国 ▶ ▲ 国 ● 今 Q @

- To represent the distribution of a sample of values of a continuous variable we can use a histogram.
- The range of values of the variable is divided into intervals, known as *classes*, and the frequencies in classes are represented by columns.
- As the variable is continuous, there are no gaps between neighbouring columns unlike a bar chart.
- Note also that, strictly speaking, it is the *area* of the column which is proportional to the frequency, not the height.
- The reason for this is that columns need not be of the same width.
- Computer software tends to use columns of the same width.
- However this default can be overridden in R if you really want to.

▲□▶ ▲□▶ ▲□▶ ▲□▶ = ● のへで



Mean film budget

Figure: Histogram of film ratings and budgets.

・ロト ・四ト ・ヨト ・ヨト

æ

- Figure 4.2 shows histograms of the film budgets.
- When dealing with densities (relative frequency), we can easily work out the height using this formula:

$$\mathsf{Height} = \frac{\mathsf{frequency}}{n \times \mathsf{Bin-width}}$$

When the y-axis is labelled with density or relative frequencies, the area under the histogram is one.

- Figure 4.2 shows histograms of the film budgets.
- When dealing with densities (relative frequency), we can easily work out the height using this formula:

$$\mathsf{Height} = \frac{\mathsf{frequency}}{n \times \mathsf{Bin-width}}$$

When the y-axis is labelled with density or relative frequencies, the area under the histogram is one.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ つく⊙

 Bin widths should be chosen so that you get a good idea of the distribution of the data, without being swamped by random variation.

To generate Figure 4.2 in R we use the following commands:

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○□ のへで

```
> hist(movies$Budget, col="grey",
+ main="Mean film budget", freq=FALSE,
+ xlab="Budget ($)")
```

4.3.1 How many bins should we have?

First we will define the notation we will use:

- n: the sample size;
- k: the number of bins in the histogram;
- h: the bin-width.

Then the number of bins we will use to construct a histogram is:

$$k = \left\lceil \frac{\max(x) - \min(x)}{h} \right\rceil$$
(4.1)

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

where $\lceil \cdot \rceil$ is the ceiling function.

Table 4.1: Sturges' rule

- $k_{ST} = \lceil \log_2 n + 1 \rceil$
- Default
- Tends not to be very good for n > 30.



・ロト ・四ト ・ヨト ・ヨト

æ

Table 4.1: Scotts' rule

•
$$h_{SC} = 3.49 \times s \times n^{-1/3}$$

breaks = "Scott"



<ロ> (四) (四) (三) (三) (三)

æ

Table 4.1: Freedman-Diaconis

- $h_{FR} = 2 \times IQR(x) \times n^{-1/3}$
- breaks = "FD"
- When the distribution is symmetric, this is very similar to Scott's rule.



(ロ)、<()、<()、<()、<()、<()、<()</p>

The FD rule

Suppose we want to calculate the number of bins for the budget movie variable.

Suppose we want to calculate the number of bins for the budget movie variable. For Sturges' rule we have:

Suppose we want to calculate the number of bins for the budget movie variable. For Sturges' rule we have:

 $k_{ST} =$



Suppose we want to calculate the number of bins for the budget movie variable. For Sturges' rule we have:

 $k_{ST} = \lceil \log_2(n) \rceil + 1 =$



Suppose we want to calculate the number of bins for the budget movie variable. For Sturges' rule we have:

$$k_{ST} = \lceil \log_2(n) \rceil + 1 = \lceil \log_2(4847) \rceil + 1 =$$

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 三臣 - の�?

Suppose we want to calculate the number of bins for the budget movie variable. For Sturges' rule we have:

$$k_{ST} = \lceil \log_2(n) \rceil + 1 = \lceil \log_2(4847) \rceil + 1 = 14$$
.



Mean film budget

For Scotts' rule, we first calculate the bin width

For Scotts' rule, we first calculate the bin width

$$h_{SC} =$$

For Scotts' rule, we first calculate the bin width

$$h_{SC} = \frac{3.49 \times s}{n^{1/3}} =$$

For Scotts' rule, we first calculate the bin width

$$h_{SC} = rac{3.49 imes s}{n^{1/3}} = rac{3.49 imes 23039711}{4847^{1/3}} \simeq$$

For Scotts' rule, we first calculate the bin width

$$h_{SC} = \frac{3.49 \times s}{n^{1/3}} = \frac{3.49 \times 23039711}{4847^{1/3}} \simeq 4,751,289$$

For Scotts' rule, we first calculate the bin width

$$h_{SC} = \frac{3.49 \times s}{n^{1/3}} = \frac{3.49 \times 23039711}{4847^{1/3}} \simeq 4,751,289$$

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

then the number of bins is:
For Scotts' rule, we first calculate the bin width

$$h_{SC} = \frac{3.49 \times s}{n^{1/3}} = \frac{3.49 \times 23039711}{4847^{1/3}} \simeq 4,751,289$$

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

$$k_{SC} =$$

For Scotts' rule, we first calculate the bin width

$$h_{SC} = \frac{3.49 \times s}{n^{1/3}} = \frac{3.49 \times 23039711}{4847^{1/3}} \simeq 4,751,289$$

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

$$k_{SC} = \left\lceil \frac{\max(x) - \min(x)}{h_{sc}} \right\rceil =$$

For Scotts' rule, we first calculate the bin width

$$h_{SC} = \frac{3.49 \times s}{n^{1/3}} = \frac{3.49 \times 23039711}{4847^{1/3}} \simeq 4,751,289$$

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

$$k_{SC} = \left\lceil \frac{\max(x) - \min(x)}{h_{sc}} \right\rceil = \left\lceil \frac{2 \times 10^8 - (-1)}{4751289} \right\rceil =$$

For Scotts' rule, we first calculate the bin width

$$h_{SC} = \frac{3.49 \times s}{n^{1/3}} = \frac{3.49 \times 23039711}{4847^{1/3}} \simeq 4,751,289$$

$$k_{SC} = \left\lceil \frac{\max(x) - \min(x)}{h_{SC}} \right\rceil = \left\lceil \frac{2 \times 10^8 - (-1)}{4751289} \right\rceil = 43$$



▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

For the FD rule, the bin width is:

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

For the FD rule, the bin width is:

 $h_{FD} =$

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

For the FD rule, the bin width is:

$$h_{FD} = \frac{2 \times IQR(x)}{n^{1/3}} =$$

For the FD rule, the bin width is:

$$h_{FD} = rac{2 imes IQR(x)}{n^{1/3}} = rac{2 imes (8 imes 10^6 - (-1))}{4847^{1/3}} \simeq$$

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

For the FD rule, the bin width is:

$$h_{FD} = rac{2 imes IQR(x)}{n^{1/3}} = rac{2 imes (8 imes 10^6 - (-1))}{4847^{1/3}} \simeq 945,429$$
 .

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

For the FD rule, the bin width is:

$$h_{FD} = rac{2 imes IQR(x)}{n^{1/3}} = rac{2 imes (8 imes 10^6 - (-1))}{4847^{1/3}} \simeq 945,429$$
 .

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

For the FD rule, the bin width is:

$$h_{FD} = rac{2 imes IQR(x)}{n^{1/3}} = rac{2 imes (8 imes 10^6 - (-1))}{4847^{1/3}} \simeq 945,429$$
 .

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

Hence,

 $k_{FD} =$

For the FD rule, the bin width is:

$$h_{FD} = rac{2 imes IQR(x)}{n^{1/3}} = rac{2 imes (8 imes 10^6 - (-1))}{4847^{1/3}} \simeq 945,429$$
 .

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

$$k_{FD} = \left\lceil \frac{\max(x) - \min(x)}{h_{sc}} \right\rceil$$

For the FD rule, the bin width is:

$$h_{FD} = rac{2 imes IQR(x)}{n^{1/3}} = rac{2 imes (8 imes 10^6 - (-1))}{4847^{1/3}} \simeq 945,429$$
 .

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

$$k_{FD} = \left\lceil \frac{\max(x) - \min(x)}{h_{sc}} \right\rceil = \left\lceil \frac{2 \times 10^8 - (-1)}{945429} \right\rceil =$$

For the FD rule, the bin width is:

$$h_{FD} = rac{2 imes IQR(x)}{n^{1/3}} = rac{2 imes (8 imes 10^6 - (-1))}{4847^{1/3}} \simeq 945,429$$
 .

$$k_{FD} = \left\lceil \frac{\max(x) - \min(x)}{h_{sc}} \right\rceil = \left\lceil \frac{2 \times 10^8 - (-1)}{945429} \right\rceil = 212$$





- A box and whisker plot, sometimes simply called a 'boxplot', is another way to represent continuous data.
- This kind of plot is particularly useful for comparing two or more groups, by placing the boxplots side-by-side.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで



Figure: Box and whisker plots of (a) film length (b) film length split according to the mpaa rating.

▲□▶ ▲圖▶ ▲厘▶ ▲厘

æ

To do this in R we use the following commands:

```
> par(mfrow=c(1, 2))
> boxplot(movies$Length, ylab="Film length",
+ col="bisque")
> boxplot(movies$Length ~ movies$mpaa,
+ ylab="Film length",
+ col="bisque")
```



> boxplot(movies\$Length ~ movies\$mpaa + movies\$Romance)

4.4.1 Boxplot example 1

For the data set

0.1	0.1	0.2	0.4	0.4	0.8	0.8	0.8
0.9	0.9	1.0	1.4	1.6	2.0	2.4	3.5

Table: An example data set.

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 三臣 - の�?

construct a boxplot.

4.4.1 Solution

First we calculate the median and quartiles

Median	1 st quartile	3 rd quartile	IQR
0.85	0.4	1.55	1.15

Table: Summary statistics for the first example data in Table 4.3.

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 三臣 - の�?

Solution

To calculate the outliers and whiskers, we first calculate:

$$W_L =$$
 lower quartile $-1.5IQR = -1.325$

and

$$W_U$$
 = upper quartile + 1.5/QR = 3.275

Since $3.5 > W_U = 3.275$, this means that 3.5 is an outlying point. Since we have no points less than W_L , the lower whisker is the smallest data point, i.e. 0.1. The upper whisker is $\max(x^*)$ where x^* does not include any outlying points, i.e. 2.4.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

4.4.1 Boxplot example 1



▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

For the following data set, construct a boxplot

9.0	32.8	33.0	34.9	35.4	39.7	41.6	42.0
42.3	43.2	46.9	49.2	51.6	51.7	55.0	81.0

Table: An example data set.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Median	1 st quartile	3 rd quartile	IQR	W _L	W _U
42.15	35.025	51.00	15.975		

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 三臣 - の�?

Median	1 st quartile	3 rd quartile	IQR	W _L	W _U
42.15	35.025	51.00	15.975		

$$W_L =$$

Median	1 st quartile	3 rd quartile	IQR	W_L	W _U
42.15	35.025	51.00	15.975		

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 三臣 - の�?

 W_L = lower quartile - 1.5/QR =

Median	1 st quartile	3 rd quartile	IQR	WL	W_U
42.15	35.025	51.00	15.975	11.0625	

$$W_L$$
 = lower quartile - 1.5/QR = 11.0625

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 三臣 - の�?

Median	1 st quartile	3 rd quartile	IQR	WL	W_U
42.15	35.025	51.00	15.975	11.0625	

$$W_L$$
 = lower quartile - 1.5*IQR* = 11.0625
 W_U =

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 三臣 - の�?

Median	1 st quartile	3 rd quartile	IQR	WL	W _U
42.15	35.025	51.00	15.975	11.0625	

 W_L = lower quartile - 1.5/QR = 11.0625

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○□ ○○○

$$W_U$$
 = Upper quartile + 1.5/QR =

Median	1 st quartile	3 rd quartile	IQR	WL	W _U
42.15	35.025	51.00	15.975	11.0625	74.9625

 W_L = lower quartile - 1.5/QR = 11.0625

$$W_U$$
 = Upper quartile + 1.5/QR = 74.9625

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○□ ○○○

9.0	32.8	33.0	34.9	35.4	39.7	41.6	42.0
42.3	43.2	46.9	49.2	51.6	51.7	55.0	81.0

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

• From table 4.5, we see that $9.0 < W_L = 11.0625$.

9.0	32.8	33.0	34.9	35.4	39.7	41.6	42.0
42.3	43.2	46.9	49.2	51.6	51.7	55.0	81.0

• From table 4.5, we see that $9.0 < W_L = 11.0625$. So the lower whisker extends to the point 32.8 and 9.0 is an outlier.

9.0	32.8	33.0	34.9	35.4	39.7	41.6	42.0
42.3	43.2	46.9	49.2	51.6	51.7	55.0	81.0

• From table 4.5, we see that $9.0 < W_L = 11.0625$. So the lower whisker extends to the point 32.8 and 9.0 is an outlier.

(日) (四) (문) (문) (문)

• From table 4.5, we see that $81 > W_U = 74.9625$.

9.0	32.8	33.0	34.9	35.4	39.7	41.6	42.0
42.3	43.2	46.9	49.2	51.6	51.7	55.0	81.0

- From table 4.5, we see that $9.0 < W_L = 11.0625$. So the lower whisker extends to the point 32.8 and 9.0 is an outlier.
- From table 4.5, we see that $81 > W_U = 74.9625$. So the upper whisker extends to the point 55.0 and 81 is an outlier.

《曰》 《聞》 《臣》 《臣》 三臣 …

4.4.2 Boxplot example 2



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Part V

Control statements and functions

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで
This function takes in a single argument x and returns x^2 :

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 三臣 - の�?

```
> Fun1 = function(x) {
+ return (x*x)
+ }
```

The key elements in the function call are:

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

• The word function;

The key elements in the function call are:

- The word function;
- The brackets () which enclose the **argument** list. This list may be empty.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

The key elements in the function call are:

- The word function;
- The brackets () which enclose the **argument** list. This list may be empty.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

• A sequence of statements in curly braces $\{ \}$.

The key elements in the function call are:

- The word function;
- The brackets () which enclose the **argument** list. This list may be empty.

▲ロト ▲母 ▶ ▲ ヨ ▶ ▲ ヨ ● つんで

- A sequence of statements in curly braces $\{ \}$.
- A return statement.

We 'call' Fun1 in the following manner:

> F	un1(5)		
[1]	25		

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

We 'call' Fun1 in the following manner:

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

> Fun1(5)
[1] 25

> y = Fun1(10)

We 'call' Fun1 in the following manner:

> Fun1(5) [1] 25

> y = Fun1(10)

> y [1] 100

▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のへで

We 'call' Fun1 in the following manner:

> Fun1(5)
[1] 25

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○□ のへで

> y = Fun1(10)

> y [1] 100

> z = c(1, 2, 3, 4)
> Fun1(z)
[1] 1 4 9 16

Of course, the old saying 'Garbage in, Garbage out' is true:

```
#Incorrect function calls.
> Fun1()
Error in Fun1() : argument "x" is missing, with no default
> Fun1("5")
Error in x * x : non-numeric argument to binary operator
```

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○□ のへで

The error messages do give you an idea of what went wrong.

What is the value of y?

```
> fun1 = function(x) {
+     return(2*x)
+ }
> y = fun1(2)
```

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

- 5
- 10
- 2
- 4

What is the value of y?

```
> x = 5
> fun1 = function(x) {
+         return(2*x)
+ }
> (y= fun1(2))
[1] 4
```

• 5

• 10

• 2



What is the value of y?

> fun1 = function(x) {
+ z = 3
+ return(x*z)
+ }
> y = fun1(2)

▲ロト ▲圖ト ▲画ト ▲画ト 三国 - のへで

• 2

• 3

• 4

What is the value of y?

```
> fun1 = function(x) {
+  z = 3
+  return(x*z)
+ }
> (y = fun1(2))
[1] 6
```

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

• 2

• 3

• 4

What is the value of y?

> fun1 = function(x) {
+ z = 3
+ return(x*z)
+ }
> k = 5
> y = fun1(k)

▲ロト ▲圖ト ▲画ト ▲画ト 三国 - のへで

What is the value of y?

```
> fun1 = function(x) {
+   z = 3
+   return(x*z)
+ }
> k = 5
> (y = fun1(k))
[1] 15
```

▲ロト ▲圖ト ▲画ト ▲画ト 三国 - のへで

Other variations to this simple function are:

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

```
> Fun2 = function(x=1) {
+ return (x*x)
+ }
```

Other variations to this simple function are:

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

```
> Fun2 = function(x=1) {
+ return (x*x)
+ }
```

> Fun2() [1] 1

Other variations to this simple function are:

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

```
> Fun2 = function(x=1) {
+ return (x*x)
+ }
```

> Fun2() [1] 1

> Fun2(4)
[1] 16

```
> Fun3 = function(x, y) {
+ return (x*y)
+ }
```

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

```
> Fun3 = function(x, y) {
+ return (x*y)
+ }
```

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

> Fun3(3, 4)
[1] 12

What is the value of y?

```
> fun1 = function(x=3) {
+     return(2*x)
+ }
> y = fun1(2)
```

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

- 3
- 4
- 6

What is the value of y?

```
> fun1 = function(x=3) {
+         return(2*x)
+ }
> (y = fun1(2))
[1] 4
```

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

```
• 3
```

• 4

What is the value of y?

```
> fun1 = function(x=3) {
+ return(2*x)
+ }
> y = fun1()
```

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

• 3

• 4

What is the value of y?

```
> fun1 = function(x=3) {
+   return(2*x)
+ }
> (y = fun1())
[1] 6
```

▲ロト ▲圖ト ▲画ト ▲画ト 三国 - のへで

- 3
- 4
- 6

What is the value of y?

```
> fun1 = function(x=3) {
+     x = 2
+     return(2*x)
+ }
> y = fun1()
```

▲ロト ▲圖ト ▲画ト ▲画ト 三国 - のへで

- 3
- 4
- 6

What is the value of y?

```
> fun1 = function(x=3) {
+     x = 2
+     return(2*x)
+ }
> (y = fun1())
[1] 4
```

▲ロト ▲圖ト ▲画ト ▲画ト 三国 - のへで

• 3

• 4

Here the function below takes in a vector, plots a histogram and returns a vector containing the mean and standard deviation:

```
> Investigate = function(values) {
+    hist(values)
+    m_std = c(mean(values), sd(values))
+    return(m_std)
+ }
```

5.1.2 A more useful function

Once we have created our function, we can put it to good use:

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

> Investigate(movies\$Rating)

[1] 5.522715 1.451864

5.1.2 A more useful function

Once we have created our function, we can put it to good use:

> Investigate(movies\$Rating)

[1] 5.522715 1.451864

> Investigate(movies\$Length)
[1] 100.87807 17.34152

5.1.2 A more useful function

Once we have created our function, we can put it to good use:

> Investigate(movies\$Rating)

[1] 5.522715 1.451864

> Investigate(movies\$Length)

[1] 100.87807 17.34152

> Investigate(movies\$Budget)
[1] 10286893 23039711

- + ロ ト + 個 ト + 画 ト + 画 - りへの

• When we call a function, R first looks for *local* variables, then *global* variables.

・ロト ・御 ト ・ ヨト ・ ヨト … ヨ

• For example, Fun4 uses a global variable:

```
> blob = 5
> Fun4 = function() {
+ return(blob)
+ }
```

- When we call a function, R first looks for *local* variables, then *global* variables.
- For example, Fun4 uses a global variable:

```
> blob = 5
> Fun4 = function() {
+ return(blob)
+ }
```

```
> Fun4()
[1] 5
```

・ロト ・四ト ・ヨト ・ヨト - ヨー

But Fun5 uses a local variable:

```
> blob = 5
> Fun5 = function() {
+     blob = 6
+     return(blob)
+ }
```

▲ロト ▲母 ▶ ▲ ヨ ▶ ▲ ヨ ● つんで

But Fun5 uses a local variable:

```
> blob = 5
> Fun5 = function() {
+     blob = 6
+     return(blob)
+ }
```

▲ロト ▲母 ▶ ▲ ヨ ▶ ▲ ヨ ● つんで

```
> Fun5()
[1] 6
```
5.1.3 Variable scope

But Fun5 uses a local variable:

```
> blob = 5
> Fun5 = function() {
+     blob = 6
+     return(blob)
+ }
```

> Fun5()		
[1] 6		
> blob		
[1] 5		

▲ロト ▲母 ▶ ▲ ヨ ▶ ▲ ヨ ● つんで

What is the value of y?

```
> x = 1
> fun1 = function(x) {
+     x = 2
+     return(x)
+ }
> y= fun1(3)
```

▲ロト ▲圖ト ▲画ト ▲画ト 三国 - のへで

• 1

• 2

What is the value of y?

```
> x = 1
> fun1 = function(x) {
+     x = 2
+     return(x)
+ }
> (y= fun1(3))
[1] 2
```

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

• 1

• 2

What is the value of y?

```
> x = 1
> fun1 = function(x=3) {
+    x = 2
+    return(x)
+ }
> y= fun1()
```

▲ロト ▲圖ト ▲画ト ▲画ト 三国 - のへで

• 1

• 2

What is the value of y?

```
> x = 1
> fun1 = function(x=3) {
+    x = 2
+    return(x)
+ }
> (y= fun1())
[1] 2
```

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

• 1

• 2

What is the value of y?

```
> x = 1
> fun1 = function(x=3) {
+ return(2*x)
+ }
> x = fun1(x)
> y = fun1(x)
```

▲ロト ▲圖ト ▲画ト ▲画ト 三国 - のへで

What is the value of y?

```
> x = 1
> fun1 = function(x=3) {
+ return(2*x)
+ }
> x = fun1(x)
> (y = fun1(x))
[1] 4
```

▲ロト ▲圖ト ▲画ト ▲画ト 三国 - のへで

5.2 The cat command

• A useful function to help debugging is the cat function.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

• This function is used to print messages to the screen.

5.2 The cat command

- A useful function to help debugging is the cat function.
- This function is used to print messages to the screen.
- For example,

```
> x = 5
> cat(x, "\n")
5
> (y = cat(x, "\n"))
5
NULL
```

• We will use the cat function in the next section.

5.3 Conditionals

• Conditional statements are features of a programming language which perform different computations or actions depending on whether a condition evaluates to TRUE or FALSE.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

• They are used in almost all computer programs.

The basic structure of an if statement is:



where expr is evaluated to be either TRUE or FALSE.



The following example illustrates if statements in R:

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

> x = 5
> y = 5
> if(x<5) {
+ y = 0
+ }</pre>

The following example illustrates if statements in R:

> x = 5
> y = 5
> if(x<5) {
+ y = 0
+ }</pre>

> v		
· ,		
[1] 5	5	

In this code chunk, x < 5 evaluates to be FALSE so the following brackets are not evaluated.

We test for greater than in a similar manner:

> x = 5 > y = 5 > if(x > 0) { + y = 0 + } > y [1] 0

Here x > 0 evaluates to be TRUE so, y is set equal to 0. If we wanted to test for equality with zero, then we would use =>.

We can link together a number of if statements

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

```
> x = 0
> if(x > 0) {
+ cat("x is greater than zero")
+ } else if(x < 0) {
+ cat("x is less than zero")
+ } else {
+ cat("x must be zero!")
+ cat("\n")
+ }
x must be zero!</pre>
```

The final else is optional.

```
> IsNegative = function(value) {
+ is_pos = FALSE
+ if(value < 0) {
+ is_pos = TRUE
+ }
+ return(is_pos)
+ }</pre>
```

```
> IsNegative = function(value) {
+ is_pos = FALSE
+ if(value < 0) {
+ is_pos = TRUE
+ }
+ return(is_pos)
+ }</pre>
```

```
> IsNegative(1)
[1] FALSE
```

```
> IsNegative = function(value) {
+ is_pos = FALSE
+ if(value < 0) {
+ is_pos = TRUE
+ }
+ return(is_pos)
+ }</pre>
```

```
> IsNegative(1)
[1] FALSE
```

```
> IsNegative(-5.6)
[1] TRUE
```

A more sophisticated function could be

```
> IsGreaterThan = function(value1, value2) {
+ is_greater_than = FALSE
+ if(value1 > value2) {
+ is_greater_than = TRUE
+ }
+ return(is_greater_than)
+ }
```

《曰》 《聞》 《臣》 《臣》 三臣 …

```
> IsGreaterThan(-5, -6)
[1] TRUE
```



```
> IsGreaterThan(-5, -6)
[1] TRUE
```

```
> IsGreaterThan(10, 10)
[1] FALSE
```



What is the value of y?

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

```
> x = 5
> if(x==6) {
+ y = FALSE
+ } else {
+ y = TRUE
+ }
```

TRUE

FALSE

What is the value of y?

```
> x = 5
> if(x==6) {
+  y = FALSE
+ } else {
+  y = TRUE
+ }
> y
[1] TRUE
```

• TRUE

FALSE

What is the value of y?

```
> x = 5
> if(x > 6){
+  y = 0
+ } else if(x >= 5) {
+  y = 1
+ } else {
+  y = 2
+ }
```

012

What is the value of y?

```
> x = 5
> if(x > 6){
+ y = 0
+ }else if(x >= 5) {
+ y = 1
+ } else {
+ y = 2
+ }
> y
[1] 1
```

・ロト ・日 ・ モト ・モト ・ モー うへで

• 0

• 1

What is the value of z?

```
> x = 5
> y = "male"
> if(x > 6 & y == "Female"){
+ z = 0
+ }else if(x < 5 & y == "Male") {
+ z = 1
+ } else {
+ z = 2
+ }</pre>
```

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

• 0

• 1

What is the value of z?

```
> x = 5
> y = "male"
> if(x > 6 & y == "Female"){
+ z = 0
+ }else if(x < 5 & y == "Male") {
 z = 1
+
+ } else {
+ z = 2
+ }
> z
[1] 2
 • 0
 • 1
 • 2
```

- At times we would like to perform some operation on a vector or a data frame.
- Often R has built-in functions that will do this for you, e.g. mean, sd, other times we have to write our own functions.

- For example, suppose we want calculate $\sum_{i=1}^{10} i^2$.
- In R we can use a for loop:

```
> x = 0
> for(i in 1:10) {
+ x = x + i^2
+ }
```

- At times we would like to perform some operation on a vector or a data frame.
- Often R has built-in functions that will do this for you, e.g. mean, sd, other times we have to write our own functions.
- For example, suppose we want calculate $\sum_{i=1}^{10} i^2$.
- In R we can use a for loop:

```
> x = 0
> for(i in 1:10) {
+ x = x + i^2
+ }
```

> x		
[1] 385		

▲ロト ▲母 ▶ ▲ ヨ ▶ ▲ ヨ ● つんで

$$\sum_{j=-5}^{-1} e^{j}/j^2$$
 ,

```
> total = 0
> for(j in -5:-1) {
+ total = total + exp(j)/j^2
+ }
```

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

$$\sum_{j=-5}^{-1} e^j / j^2$$
 ,

```
> total = 0
> for(j in -5:-1) {
+ total = total + exp(j)/j^2
+ }
```

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

> total [1] 0.4086594

A more tricky example is, calculate $\sum e^k / k^2$, for $k = 3, 6, 9, \dots, 21$:

```
> total = 0
> for(i in 1:7) {
+     k = i*3
+     total = total + exp(k)/k<sup>2</sup>
+ }
> total
[1] 3208939
```

What it the value of x?

```
> total = 0
> for(blob in 1:4) {
+   total = total + 1
+ }
> x = total
```

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

What it the value of x?

```
> total = 0
> for(blob in 1:4) {
+ total = total + 1
+ }
> (x = total)
[1] 4
```

▲ロト ▲圖ト ▲画ト ▲画ト 三国 - のへで

What it the value of x?

```
> total = 0
> for(blob in 1:10) {
+ if(blob > 9) {
+ total = total + 1
+ }
+ }
> x = total
```

▲ロト ▲圖ト ▲画ト ▲画ト 三国 - のへで

What it the value of x?

```
> total = 0
> for(blob in 1:10) {
+ if(blob > 9) {
+ total = total + 1
+ }
+ }
> (x = total)
[1] 1
```
Rather than have to constantly write R code to solve the summations in §5.4 we can create a function to solve the general form:

$$\sum_{i=i_{\mathrm{s}}}^{i_{\mathrm{g}}} \frac{e^{i}}{i^{2}} \quad \text{for } i=i_{\mathrm{s}}, \ i_{\mathrm{s}}+j, \ i_{\mathrm{s}}+2j, \ \ldots, i_{\mathrm{g}} \ .$$

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 三臣 - の�?

So in R we have

```
> Summation1 = function(i_s, i_e, j) {
+ total = 0
+ for(i in i_s:(i_e/j)) {
+ k = i*j
+ total = total + exp(k)/k^2
+ }
+ return(total)
+ }
```

So in R we have

```
> Summation1 = function(i_s, i_e, j) {
+ total = 0
+ for(i in i_s:(i_e/j)) {
+     k = i*j
+     total = total + exp(k)/k^2
+     }
+     return(total)
+ }
```

```
> Summation1(-5, -1, 1)
[1] 0.4086594
```

So in R we have

```
> Summation1 = function(i_s, i_e, j) {
+ total = 0
+ for(i in i_s:(i_e/j)) {
+     k = i*j
+     total = total + exp(k)/k^2
+   }
+  return(total)
+ }
```

```
> Summation1(-5, -1, 1)
[1] 0.4086594
```

```
> Summation1(3, 21, 3)
[1] 3208925
```

- We use the apply function when we want to *apply* the same function to every row or column of a data frame.
- For example, suppose we have a data frame with three columns:

```
> (df4 = data.frame(c1 = 1:4, c2 = 4:7, c3 = 2:5))
c1 c2 c3
1 1 4 2
2 2 5 3
3 3 6 4
4 4 7 5
```

- The apply function takes (at least) three arguments.
- The first argument is the data frame, the second the number 1 or 2 indicating row or column and the third a function to apply to each row or column.

<ロト <回ト < 注ト < 注ト = 注

So

> apply(df4, 1, mean)
[1] 2.333333 3.333333 4.333333 5.333333

calculates the mean value of every row,

Standard deviation of every column

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

> apply(df4, 2, sd)
 c1 c2 c3
1.290994 1.290994 1.290994

Suppose one of the columns was non-numeric

```
> (df5 = data.frame(c1 = 1:3, c2 = 4:6,
+ c3 = LETTERS[1:3]))
c1 c2 c3
1 1 4 A
2 2 5 B
3 3 6 C
```

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○□ ○○○

then taking the mean doesn't really make sense:

```
> apply(df5, 1, mean)
[1] NA NA NA
```

Suppose one of the columns was non-numeric

```
> (df5 = data.frame(c1 = 1:3, c2 = 4:6,
+ c3 = LETTERS[1:3]))
c1 c2 c3
1 1 4 A
2 2 5 B
3 3 6 C
```

then taking the mean doesn't really make sense:

```
> apply(df5, 1, mean)
[1] NA NA NA
```

Instead, we remove the column, then calculate the mean:

```
> apply(df5[ ,1:2], 1, mean)
[1] 2.5 3.5 4.5
```

What type of variable is z?

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

> apply(z, 1, mean)

- Data frame
- Double
- String
- Iogical

What type of variable is z?

> apply(z, 1, mean)

- Data frame
- Double
- String
- Iogical

Which of these statements is correct?

> apply(z, 2, median)

- We calculate the median of z
- We calculate the median for each row of z
- We calculate the median for each column of z

《曰》 《聞》 《臣》 《臣》 三臣 …

• We calculate the mean of z

Which of these statements is correct?

> apply(z, 2, median)

- We calculate the median of z
- We calculate the median for each row of z
- We calculate the median for each column of z

《曰》 《聞》 《臣》 《臣》 三臣 …

• We calculate the mean of z

- The function tapply is very useful, but at first glance can be tricky to understand.
- It's best described using an example:

>	tapply(m	ovies\$Length,	<pre>movies\$mpaa, mean)</pre>		
	NC - 17	PG	PG-13	R	
1	10.18750	97.23384 104	.97877	100.18818	

<ロ> (四) (四) (三) (三) (三) (三)

- In the above code, we have calculated the average movie length conditional on its MPAA rating.
- So the average length of a PG movie is 97 minutes and the average NC-17 movie length is 110mins.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

• With tapply we can do we very interesting things.

• For example, in the next piece of code, we plot the average movie rating conditional on its length:

> tapply(movies\$Length, movies\$Rating, mean)[1:6]
 1 1.2 1.3 1.4 1.5 1.6
85.5 93.0 87.5 85.0 67.0 86.0

• For example, in the next piece of code, we plot the average movie rating conditional on its length:

> tapply(movies\$Length, movies\$Rating, mean)[1:6]
 1 1.2 1.3 1.4 1.5 1.6
85.5 93.0 87.5 85.0 67.0 86.0

> rating_by_len = tapply(movies\$Length, movies\$Rating,

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

+ mean)

> plot(names(rating_by_len), rating_by_len)



Figure: Plot of movie length conditional on it's rating.

◆□▶ ◆□▶ ◆□▶ ◆□▶

æ

5.7 Help

- R has a very good help system.
- If you need information about a particular function, say plot.
- Then typing ?plot in a R terminal will bring up the associated help page.
- The internet is another very good source of R help.
- Unfortunately, using Google isn't particularly useful since the letter "R" appears on most web pages!
- However, you can use

```
http://www.rseek.org/
```

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

Using this search engine limits searches to R web-pages.