

# Homework 1

Zhongying Ru 22060321

1. The Iowa data set `iowa.csv` is a toy example that summarises the yield of wheat (bushels per acre) for the state of Iowa between 1930-1962. In addition to yield, year, rainfall and temperature were recorded as the main predictors of yield.
  - a. First, we need to load the data set into R using the command `read_csv()`. Use the help function to learn what arguments this function takes. Once you have the necessary input, load the data set into R and make it a data frame called `iowa.df`.
  - b. How many rows and columns does `iowa.df` have?
  - c. What are the names of the columns of `iowa.df`?
  - d. What is the value of row 5, column 7 of `iowa.df`?
  - e. Display the second row of `iowa.df` in its entirety.

```
# Create a data frame.
iowa.df<-read.csv(file = "data/iowa.csv", sep = ';', header = T) # answer to question1.a
head(iowa.df)
```

```
##   Year Rain0 Temp1 Rain1 Temp2 Rain2 Temp3 Rain3 Temp4 Yield
## 1 1930 17.75 60.2  5.83  69.0  1.49  77.9  2.42  74.4  34.0
## 2 1931 14.76 57.5  3.83  75.0  2.72  77.2  3.30  72.6  32.9
## 3 1932 27.99 62.3  5.17  72.0  3.12  75.8  7.10  72.2  43.0
## 4 1933 16.76 60.5  1.64  77.8  3.45  76.4  3.01  70.5  40.0
## 5 1934 11.36 69.5  3.49  77.2  3.85  79.7  2.84  73.4  23.0
## 6 1935 22.71 55.0  7.00  65.9  3.35  79.4  2.42  73.6  38.4
```

```
class(iowa.df)
```

```
## [1] "data.frame"
```

```
dim(iowa.df)
```

```
## [1] 33 10
```

```
colnames(iowa.df)
```

```
## [1] "Year" "Rain0" "Temp1" "Rain1" "Temp2" "Rain2" "Temp3" "Rain3" "Temp4"
```

```
## [10] "Yield"
```

```
iowa.df[5,7]
```

```
## [1] 79.7
```

```
print(iowa.df[2,])
```

```
##   Year Rain0 Temp1 Rain1 Temp2 Rain2 Temp3 Rain3 Temp4 Yield
## 2 1931 14.76 57.5  3.83   75  2.72  77.2   3.3  72.6  32.9
```

- b. 33 rows, 10 columns.
- c. The column names are “Year”, “Rain0”, “Temp1”, “Rain1”, “Temp2”, “Rain2”, “Temp3”, “Rain3”, “Temp4”, “Yield”.
- d. The value of row 5, column 7 is 79.7

e. 2 1931 14.76 57.5 3.83 75 2.72 77.2 3.3 72.6

---

2. Syntax and class-typing.

- a. For each of the following commands, either explain why they should be errors, or explain the non-erroneous result.

```
vector1 <- c("5", "12", "7", "32")
max(vector1)
sort(vector1)
sum(vector1)
```

The class of the variable *vector1* is *character*. The elements of type *character* are comparable but non-additive. So there is an error in the line `sum(vector1)`: *invalid 'type' (character) of argument*.

- b. For the next series of commands, either explain their results, or why they should produce errors.

```
vector2 <- c("5", 7, 12)
# vector2[2] + vector2[3]
```

```
dataframe3 <- data.frame(z1="5", z2=7, z3=12)
dataframe3[1,2] + dataframe3[1,3]
```

```
list4 <- list(z1="6", z2=42, z3="49", z4=126)
list4[[2]]+list4[[4]]
list4[2]+list4[4]
```

ERROR1. The first element in *vector2* is of the *character* type, and it is non-additive. The first line should be changed to `vector2 <- as.numeric(unlist(c("5", 7, 12)))`. Then the second line can work normally.

ERROR2. The line `list4[2]+list4[4]` contains an error *non-numeric argument to binary operator*. The class of `list4[[2]]` is numeric, but `list4[2]` is of type *list* (by using the function `class(x)`). And *list* is non-additive.

---

3. Working with functions and operators.

- a. The colon operator will create a sequence of integers in order. It is a special case of the function `seq()` which you saw earlier in this assignment. Using the help command `?seq` to learn about the function, design an expression that will give you the sequence of numbers from 1 to 10000 in increments of 372. Design another that will give you a sequence between 1 and 10000 that is exactly 50 numbers in length.
- b. The function `rep()` repeats a vector some number of times. Explain the difference between `rep(1:3, times=3)` and `rep(1:3, each=3)`.

```
# a
s.step372 <- seq(1, 10000, 372)
s.50 <- seq(1, 10000, round((10000-1)/50))
```

```
# b
rep(1:3, times=3) # repeat the whole sequence "1 2 3" for 3 times
```

```
## [1] 1 2 3 1 2 3 1 2 3
```

```
rep(1:3, each=3) # repeat element in the sequence "1 2 3" for 3 times
```

```
## [1] 1 1 1 2 2 2 3 3 3
```

---

MB.Ch1.2. The orings data frame gives data on the damage that had occurred in US space shuttle launches prior to the disastrous Challenger launch of 28 January 1986. The observations in rows 1, 2, 4, 11, 13, and 18 were included in the pre-launch charts used in deciding whether to proceed with the launch, while remaining rows were omitted.

Create a new data frame by extracting these rows from orings, and plot total incidents against temperature for this new data frame. Obtain a similar plot for the full data set.

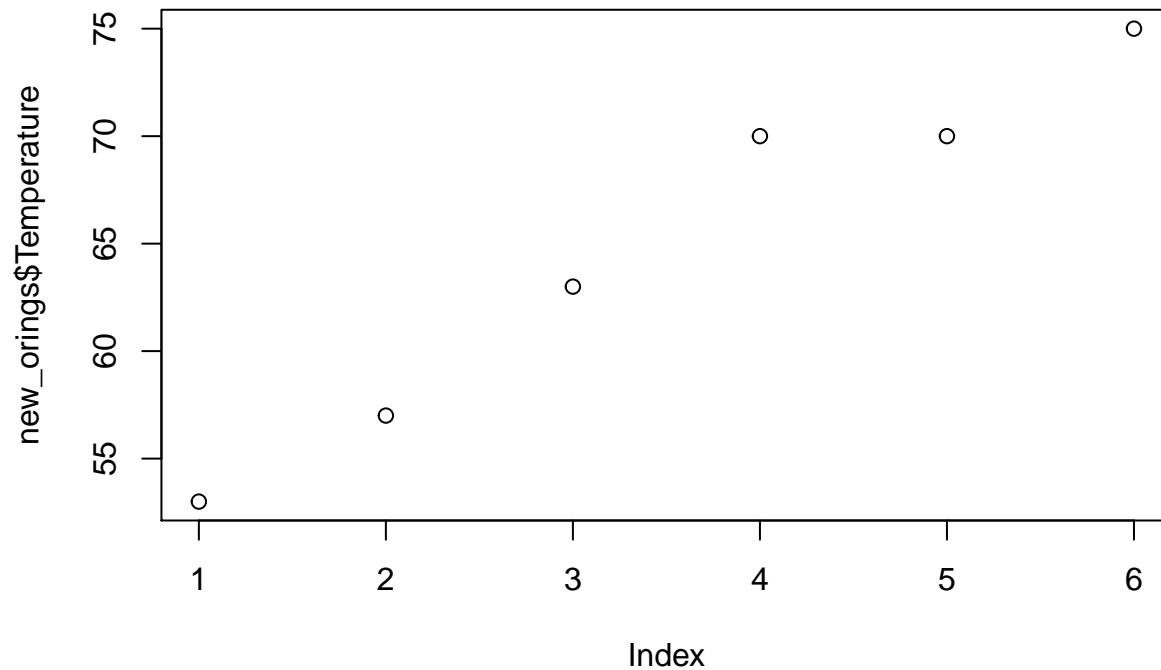
```
orings
```

##	Temperature	Erosion	Blowby	Total
## 1	53	3	2	5
## 2	57	1	0	1
## 3	58	1	0	1
## 4	63	1	0	1
## 5	66	0	0	0
## 6	67	0	0	0
## 7	67	0	0	0
## 8	67	0	0	0
## 9	68	0	0	0
## 10	69	0	0	0
## 11	70	1	0	1
## 12	70	0	0	0
## 13	70	1	0	1
## 14	70	0	0	0
## 15	72	0	0	0
## 16	73	0	0	0
## 17	75	0	0	0
## 18	75	0	2	1
## 19	76	0	0	0
## 20	76	0	0	0
## 21	78	0	0	0
## 22	79	0	0	0
## 23	81	0	0	0

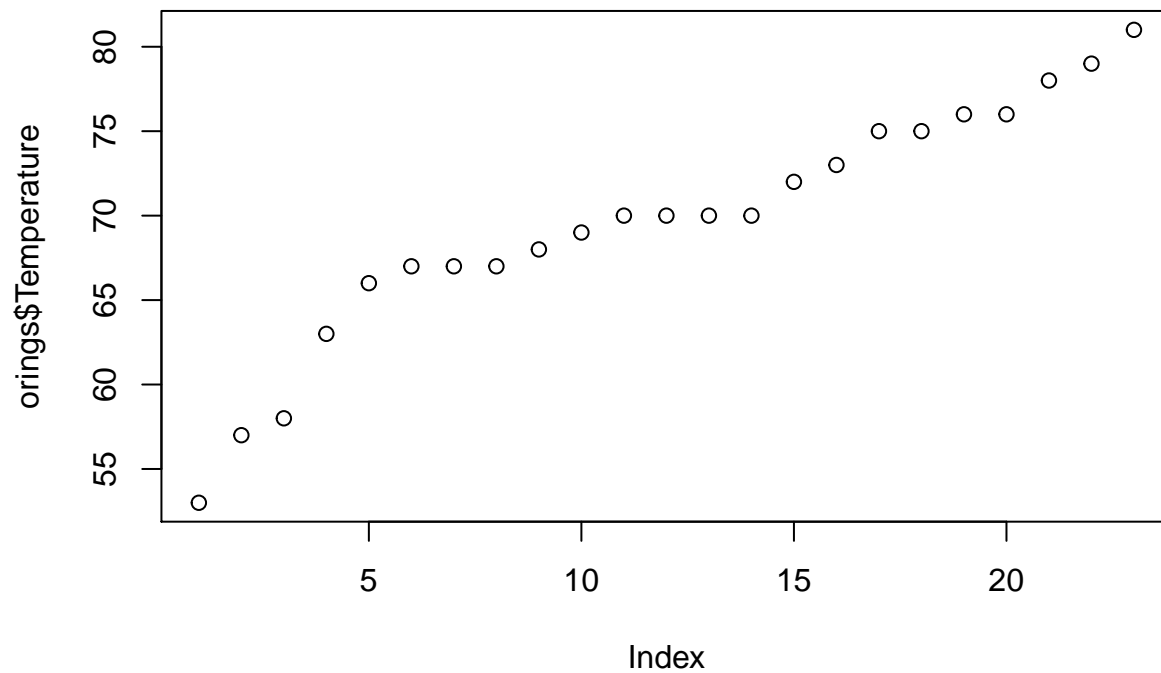
```
new_orings <- orings[c(1,2,4,11,13,18),]  
new_orings
```

##	Temperature	Erosion	Blowby	Total
## 1	53	3	2	5
## 2	57	1	0	1
## 4	63	1	0	1
## 11	70	1	0	1
## 13	70	1	0	1
## 18	75	0	2	1

```
plot(new_orings$Temperature)
```



```
plot(orings$Temperature)
```



\*\*\*

MB.Ch1.4. For the data frame ais (DAAG package)

- (a) Use the function `str()` to get information on each of the columns. Determine whether any of the columns hold missing values.

```
str(ais) # find that no columns hold missing values
```

```
## 'data.frame': 202 obs. of 13 variables:
## $ rcc : num 3.96 4.41 4.14 4.11 4.45 4.1 4.31 4.42 4.3 4.51 ...
## $ wcc : num 7.5 8.3 5 5.3 6.8 4.4 5.3 5.7 8.9 4.4 ...
## $ hc : num 37.5 38.2 36.4 37.3 41.5 37.4 39.6 39.9 41.1 41.6 ...
```

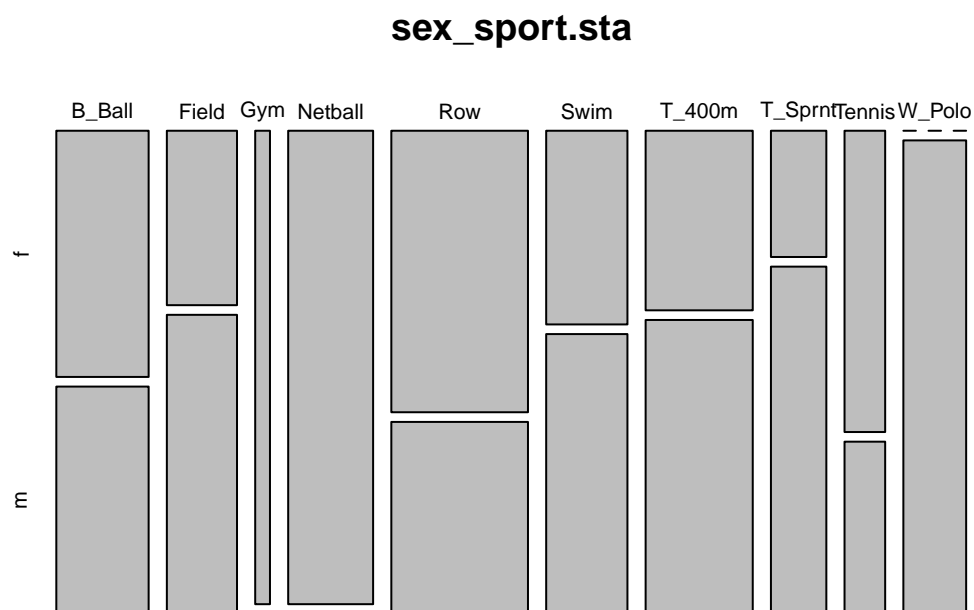
```
## $ hg      : num 12.3 12.7 11.6 12.6 14 12.5 12.8 13.2 13.5 12.7 ...
## $ ferr    : num 60 68 21 69 29 42 73 44 41 44 ...
## $ bmi     : num 20.6 20.7 21.9 21.9 19 ...
## $ ssf     : num 109.1 102.8 104.6 126.4 80.3 ...
## $ pcBfat  : num 19.8 21.3 19.9 23.7 17.6 ...
## $ lbm     : num 63.3 58.5 55.4 57.2 53.2 ...
## $ ht      : num 196 190 178 185 185 ...
## $ wt      : num 78.9 74.4 69.1 74.9 64.6 63.7 75.2 62.3 66.5 62.9 ...
## $ sex     : Factor w/ 2 levels "f","m": 1 1 1 1 1 1 1 1 1 1 ...
## $ sport   : Factor w/ 10 levels "B_Ball","Field",...: 1 1 1 1 1 1 1 1 1 1 ...
```

- (b) Make a table that shows the numbers of males and females for each different sport. In which sports is there a large imbalance (e.g., by a factor of more than 2:1) in the numbers of the two sexes?

```
sex_sport <- ais[, c('sport', 'sex')]
sex_sport.sta <- table(sex_sport$sport, sex_sport$sex)
sex_sport.sta
```

```
##
##           f  m
## B_Ball   13 12
## Field     7 12
## Gym       4  0
## Netball  23  0
## Row      22 15
## Swim      9 13
## T_400m   11 18
## T_Sprnt   4 11
## Tennis    7  4
## W_Polo    0 17
```

```
rate <- sex_sport.sta[,2]/sex_sport.sta[,1]
plot(sex_sport.sta)
```



```
rate[rate > 2 | rate < 0.5]
```

```
##      Gym Netball T_Sprnt W_Polo
```

```
##      0.00      0.00      2.75      Inf
```

As the table and plot show, there is a large sex imbalance (by a factor of more than 2:1 or less than 1:2) in the 4 sports — Gym, Netball, T\_Sprnt, W\_Polo.

MB.Ch1.6. Create a data frame called `Manitoba.lakes` that contains the lake's elevation (in meters above sea level) and area (in square kilometers) as listed below. Assign the names of the lakes using the `row.names()` function.

	elevation	area
Winnipeg	217	24387
Winnipegosis	254	5374
Manitoba	248	4624
SouthernIndian	254	2247
Cedar	253	1353
Island	227	1223
Gods	178	1151
Cross	207	755
Playgreen	217	657

```
temp <- '          elevation  area
Winnipeg      217 24387
Winnipegosis  254  5374
Manitoba       248  4624
SouthernIndian 254  2247
Cedar          253  1353
Island         227  1223
Gods           178  1151
Cross          207   755
Playgreen      217   657'
tb <- read.table(text = temp)
row.names <- c("Winnipeg", "Winnipegosis", "Manitoba", "SouthernIndian", "Cedar", "Island", "Gods", "Cross", "Playgreen")
```

- (a) Use the following code to plot  $\log_2(\text{area})$  versus elevation, adding labeling information (there is an extreme value of area that makes a logarithmic scale pretty much essential):

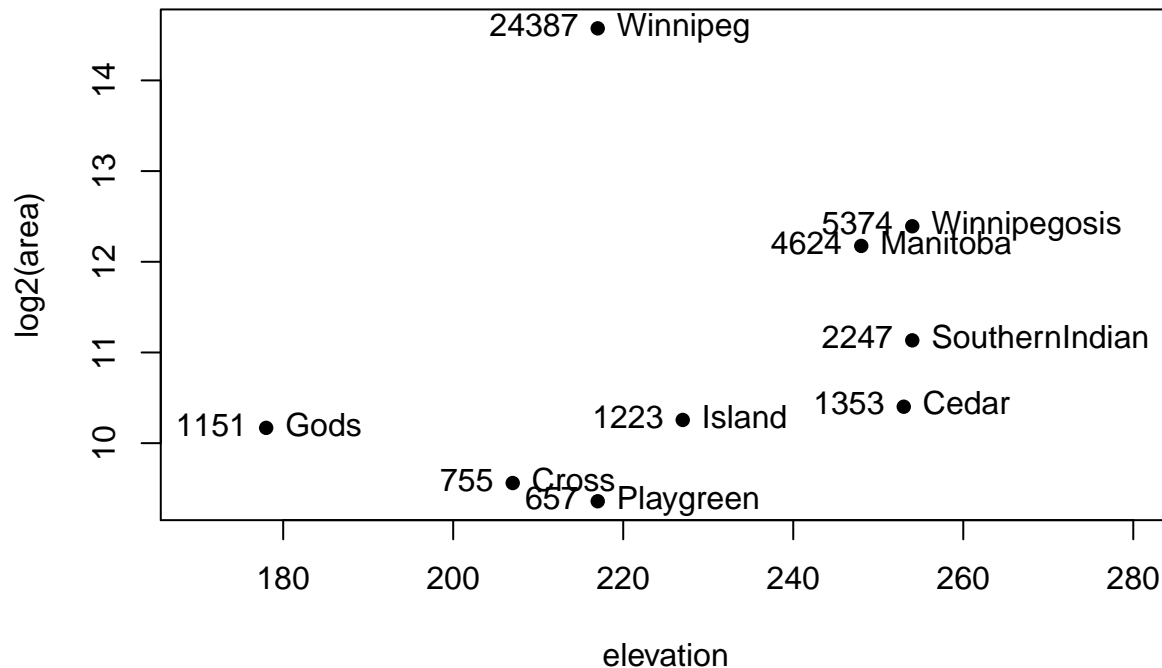
```
attach(Manitoba.lakes) # add to environment
plot(log2(area) ~ elevation, pch=16, xlim=c(170,280))
# NB: Doubling the area increases log2(area) by 1.0
text(log2(area) ~ elevation, labels=row.names(Manitoba.lakes), pos=4)
text(log2(area) ~ elevation, labels=area, pos=2) # pos 1↓ 2← 3↑ 4→
title("Manitoba's Largest Lakes")
```

```
## Warning in title("Manitoba's Largest Lakes"): conversion failure on 'Manitoba's
## Largest Lakes' in 'mbcsToSbcs': dot substituted for <e2>
```

```
## Warning in title("Manitoba's Largest Lakes"): conversion failure on 'Manitoba's
## Largest Lakes' in 'mbcsToSbcs': dot substituted for <80>
```

```
## Warning in title("Manitoba's Largest Lakes"): conversion failure on 'Manitoba's
## Largest Lakes' in 'mbcsToSbcs': dot substituted for <99>
```

## Manitoba...s Largest Lakes



Devise captions that explain the labeling on the points and on the y-axis. It will be necessary to explain how distances on the scale relate to changes in area.

A proper caption is “Elevation - log2(area) of Manitoba’s Largest Lakes”.

- (b) Repeat the plot and associated labeling, now plotting area versus elevation, but specifying `log="y"` in order to obtain a logarithmic y-scale.

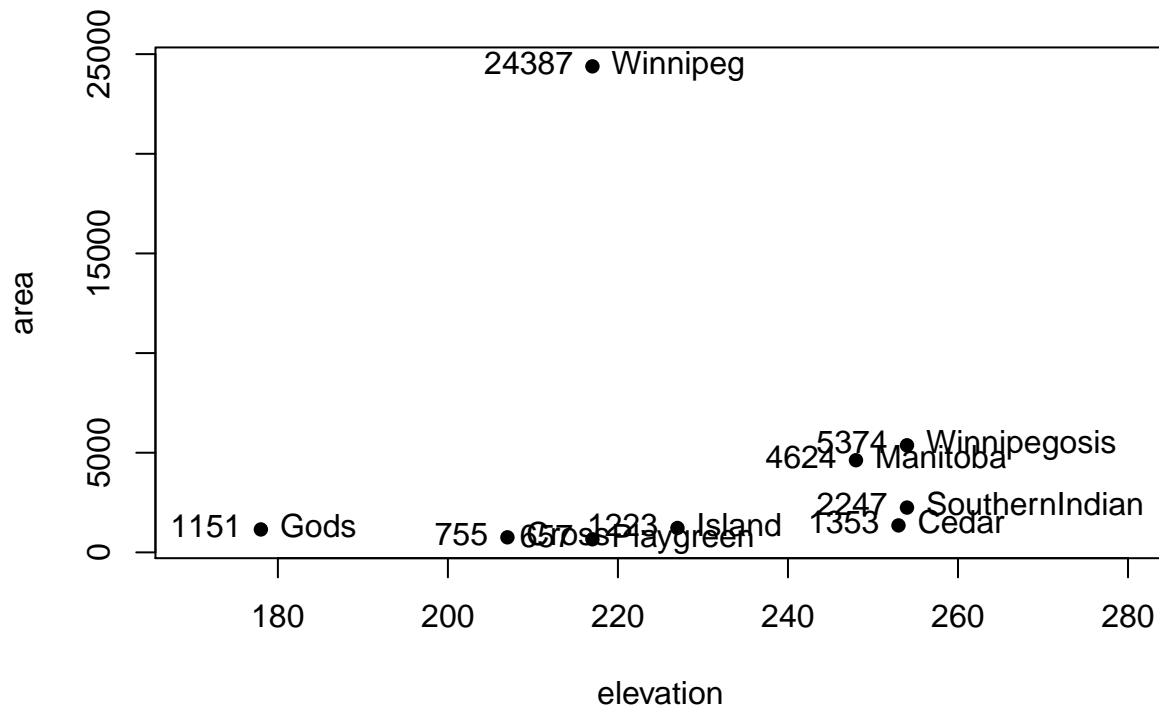
```
plot(area ~ elevation, pch=16, xlim=c(170,280), ylog=T)
text(area ~ elevation, labels=row.names(Manitoba.lakes), pos=4, ylog=T)
text(area ~ elevation, labels=area, pos=2, ylog=T)
title("Manitoba's Largest Lakes")
```

```
## Warning in title("Manitoba's Largest Lakes"): conversion failure on 'Manitoba's
## Largest Lakes' in 'mbcsToSbcs': dot substituted for <e2>
```

```
## Warning in title("Manitoba's Largest Lakes"): conversion failure on 'Manitoba's
## Largest Lakes' in 'mbcsToSbcs': dot substituted for <80>
```

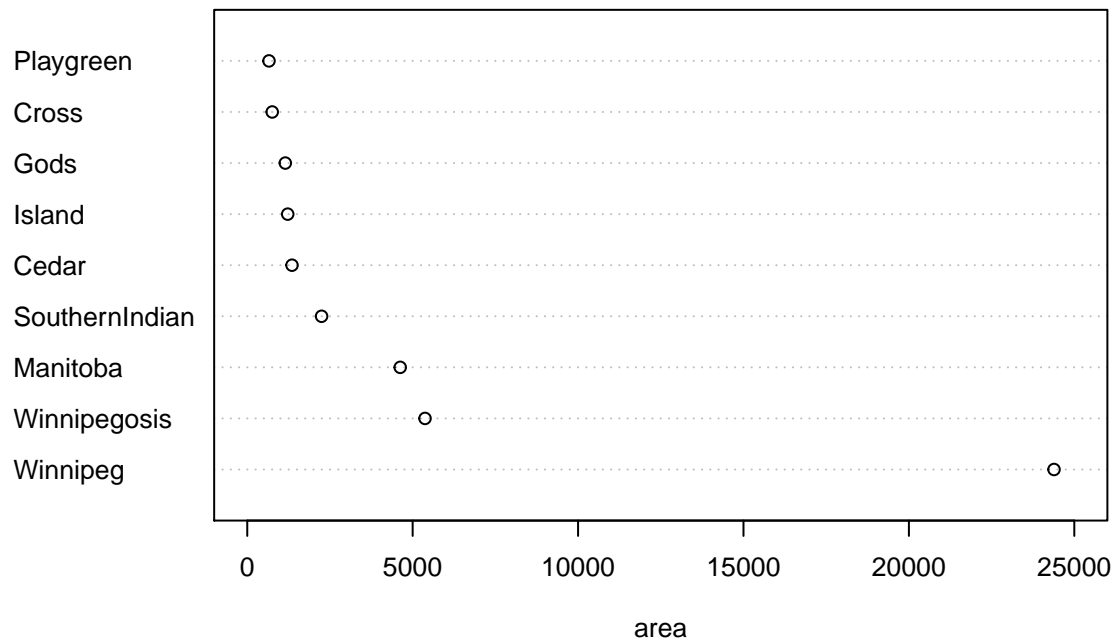
```
## Warning in title("Manitoba's Largest Lakes"): conversion failure on 'Manitoba's
## Largest Lakes' in 'mbcsToSbcs': dot substituted for <99>
```

## Manitoba...s Largest Lakes



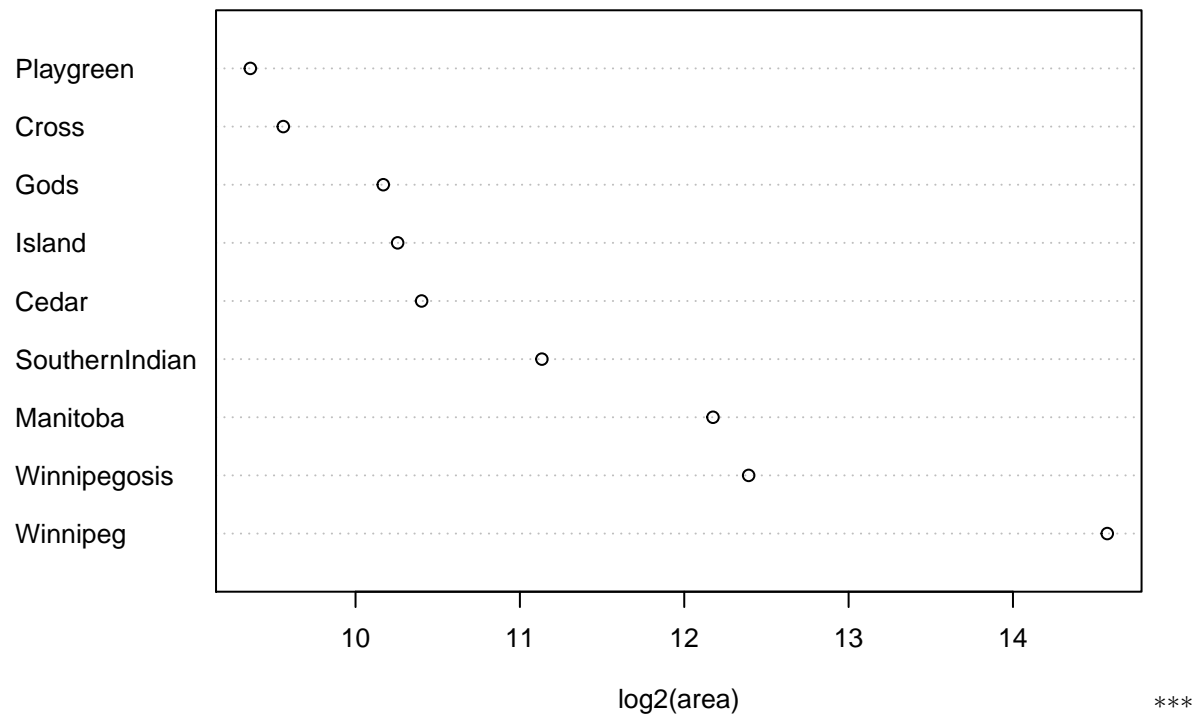
MB.Ch1.7. Look up the help page for the R function `dotchart()`. Use this function to display the areas of the Manitoba lakes (a) on a linear scale, and (b) on a logarithmic scale. Add, in each case, suitable labeling information.

```
dotchart(area, xlim = c(0, 25000), labels=row.names(Manitoba.lakes), xlab = "area", cex = .8)
```



```
dotchart(log2(area), labels=row.names(Manitoba.lakes), xlab = "log2(area)", cex = .8)
```





MB.Ch1.8. Using the `sum()` function, obtain a lower bound for the area of Manitoba covered by water.

```
sum(Manitoba.lakes[,2])
```

```
## [1] 41771
```

The lower bound for the area of Manitoba covered by water is 41771.