

Data Cleaning

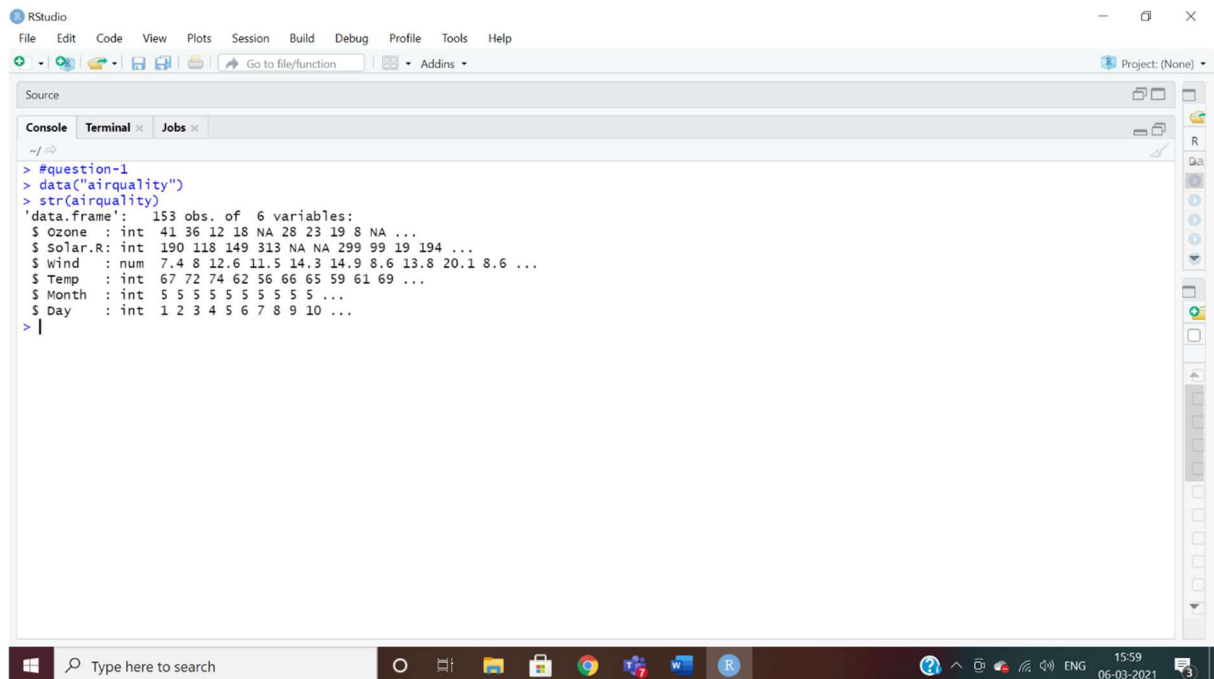
In this study, we will use **airquality** dataset.

- Print the structure of the dataset

Answer: `data("airquality")`

`str(airquality)`

OUTPUT



The screenshot shows the RStudio interface with the console window open. The following R commands were executed in the console:

```
> #question-1
> data("airquality")
> str(airquality)
```

The output of the `str(airquality)` command is displayed below the commands:

```
'data.frame': 153 obs. of 6 variables:
 $ Ozone : int 41 36 12 18 NA 28 23 19 8 NA ...
 $ Solar.R: int 190 118 149 313 NA NA 299 99 19 194 ...
 $ wind : num 7.4 8 12.6 11.5 14.3 14.9 8.6 13.8 20.1 8.6 ...
 $ Temp : int 67 72 74 62 56 66 65 59 61 69 ...
 $ Month : int 5 5 5 5 5 5 5 5 5 ...
 $ Day : int 1 2 3 4 5 6 7 8 9 10 ...
```

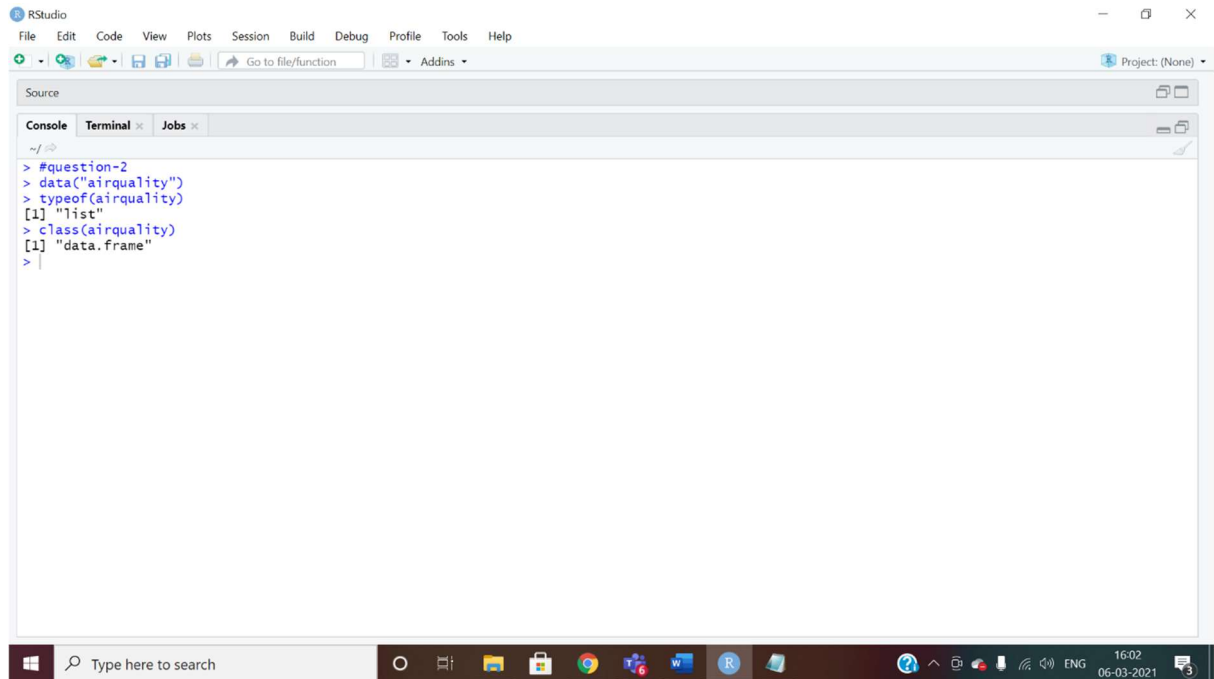
- What is the datatype of the dataset?

Answer: `data("airquality")`

`typeof(airquality)`

`class(airquality)`

OUTPUT



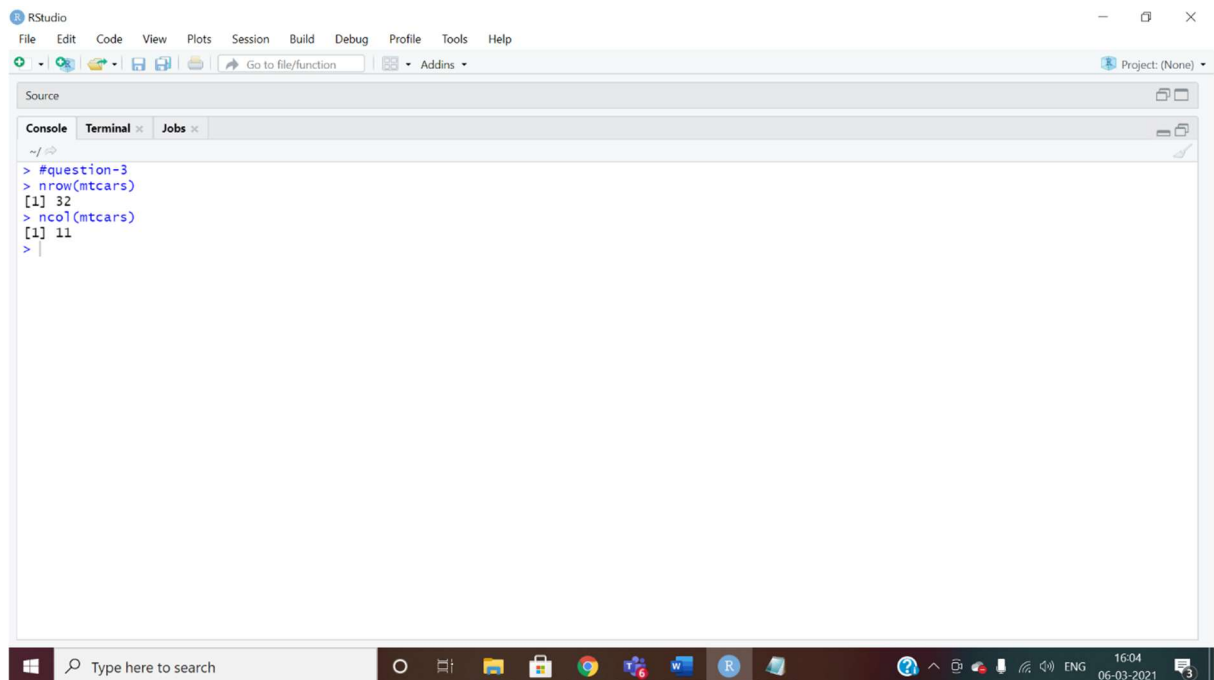
· How many columns and rows are there in the dataset??

Answer:

nrow(airquality)

ncol(airquality)

OUTPUT

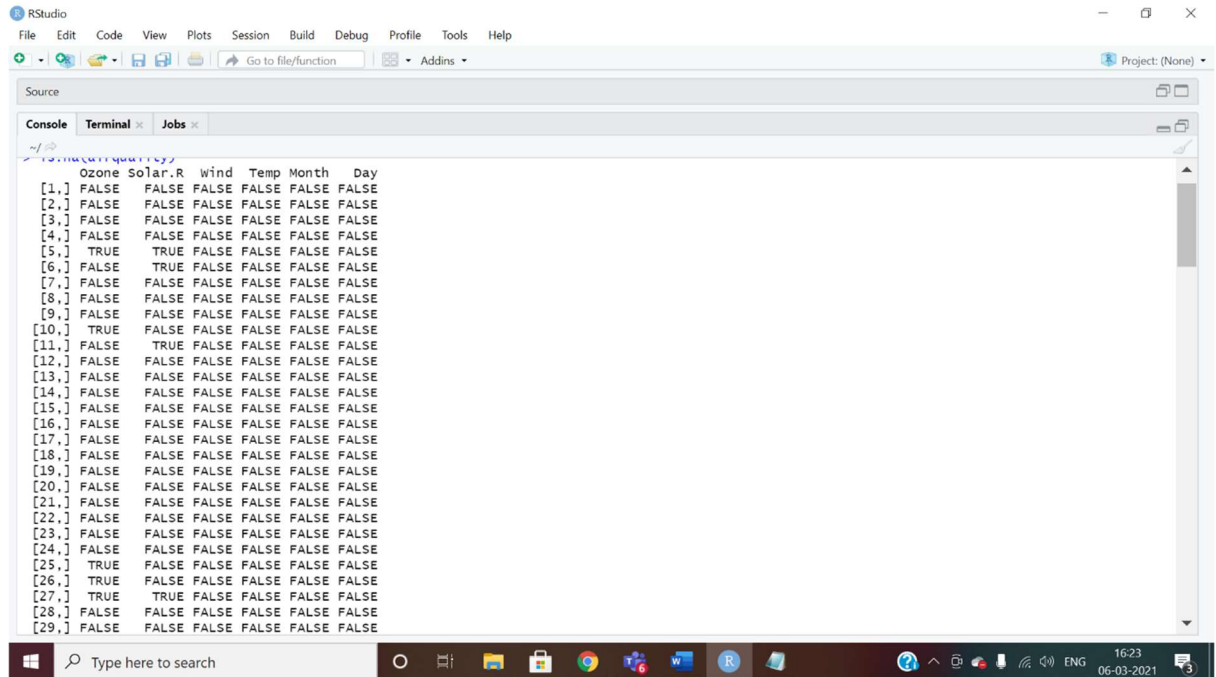


- Use the function `is.na()` to find whether any missing values are in the dataset `airquality`

Answer:

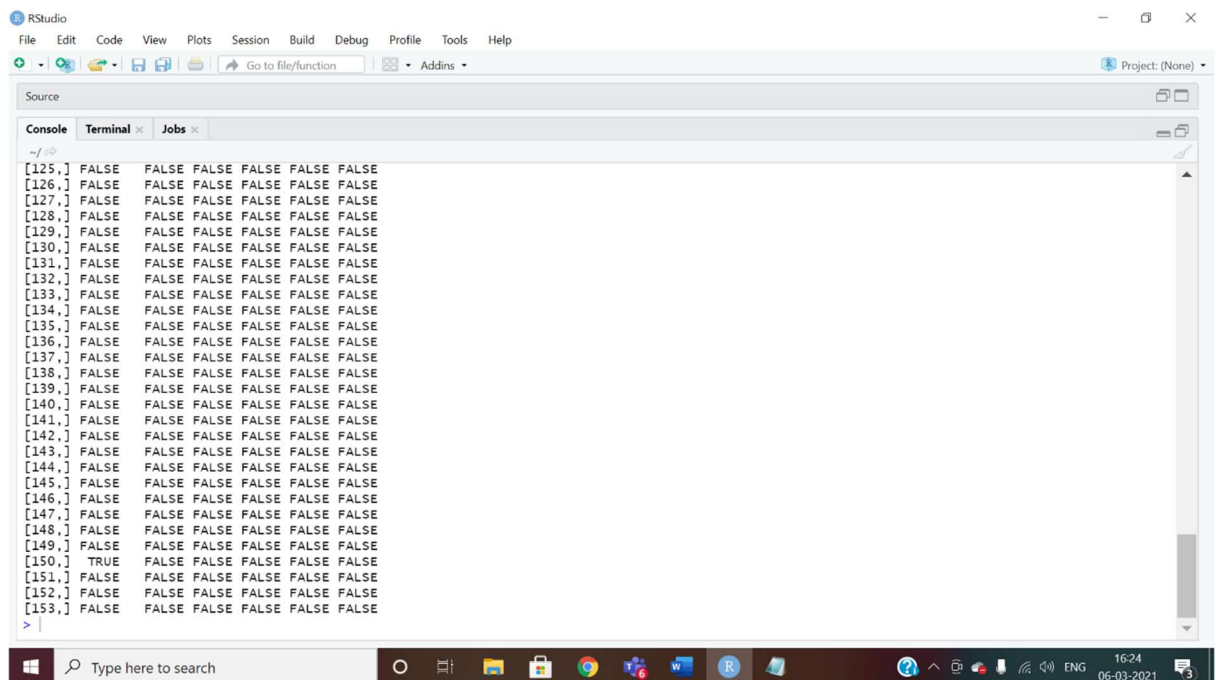
`is.na(airquality)`

OUTPUT



RStudio interface showing the console output of the command `is.na(airquality)`. The output is a matrix with 29 rows and 6 columns. The columns are labeled `Ozone`, `Solar.R`, `wind`, `Temp`, `Month`, and `Day`. The rows are indexed from [1.] to [29.]. The output shows the presence of missing values (TRUE) for `Ozone` at rows 5, 10, 25, 26, 27, and 28, and for `Solar.R` at rows 5, 6, 10, 11, 25, 26, 27, and 28. All other cells are FALSE.

	Ozone	Solar.R	wind	Temp	Month	Day
[1.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[2.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[3.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[4.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[5.]	TRUE	TRUE	FALSE	FALSE	FALSE	FALSE
[6.]	FALSE	TRUE	FALSE	FALSE	FALSE	FALSE
[7.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[8.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[9.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[10.]	TRUE	FALSE	FALSE	FALSE	FALSE	FALSE
[11.]	FALSE	TRUE	FALSE	FALSE	FALSE	FALSE
[12.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[13.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[14.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[15.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[16.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[17.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[18.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[19.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[20.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[21.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[22.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[23.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[24.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[25.]	TRUE	FALSE	FALSE	FALSE	FALSE	FALSE
[26.]	TRUE	FALSE	FALSE	FALSE	FALSE	FALSE
[27.]	TRUE	TRUE	FALSE	FALSE	FALSE	FALSE
[28.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[29.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE



RStudio interface showing the console output of the command `is.na(airquality)` for rows 125 to 153. The output is a matrix with 29 rows and 6 columns. The columns are labeled `Ozone`, `Solar.R`, `wind`, `Temp`, `Month`, and `Day`. The rows are indexed from [125.] to [153.]. The output shows the presence of missing values (TRUE) for `Ozone` at rows 150 and 151, and for `Solar.R` at rows 150 and 151. All other cells are FALSE.

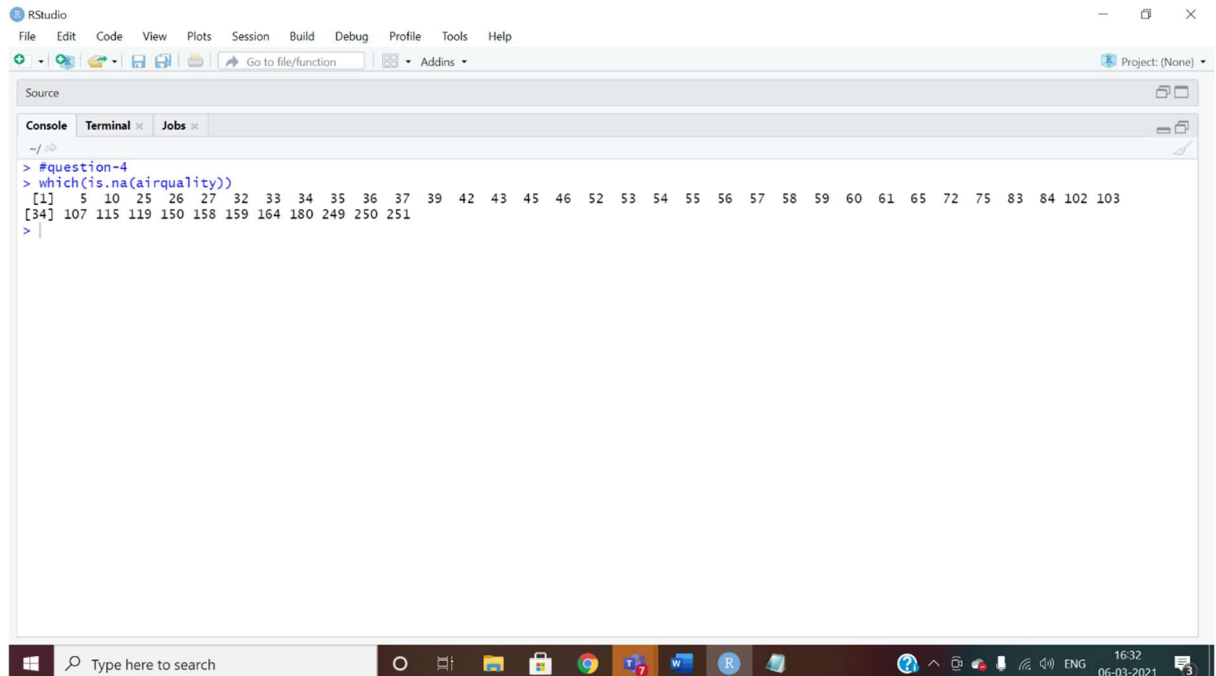
	Ozone	Solar.R	wind	Temp	Month	Day
[125.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[126.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[127.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[128.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[129.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[130.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[131.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[132.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[133.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[134.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[135.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[136.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[137.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[138.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[139.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[140.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[141.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[142.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[143.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[144.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[145.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[146.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[147.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[148.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[149.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[150.]	TRUE	FALSE	FALSE	FALSE	FALSE	FALSE
[151.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[152.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
[153.]	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE

- Print the indices of the missing values in the dataset `airquality` in row major representation

Answer

`which(is.na(airquality))`

OUTPUT



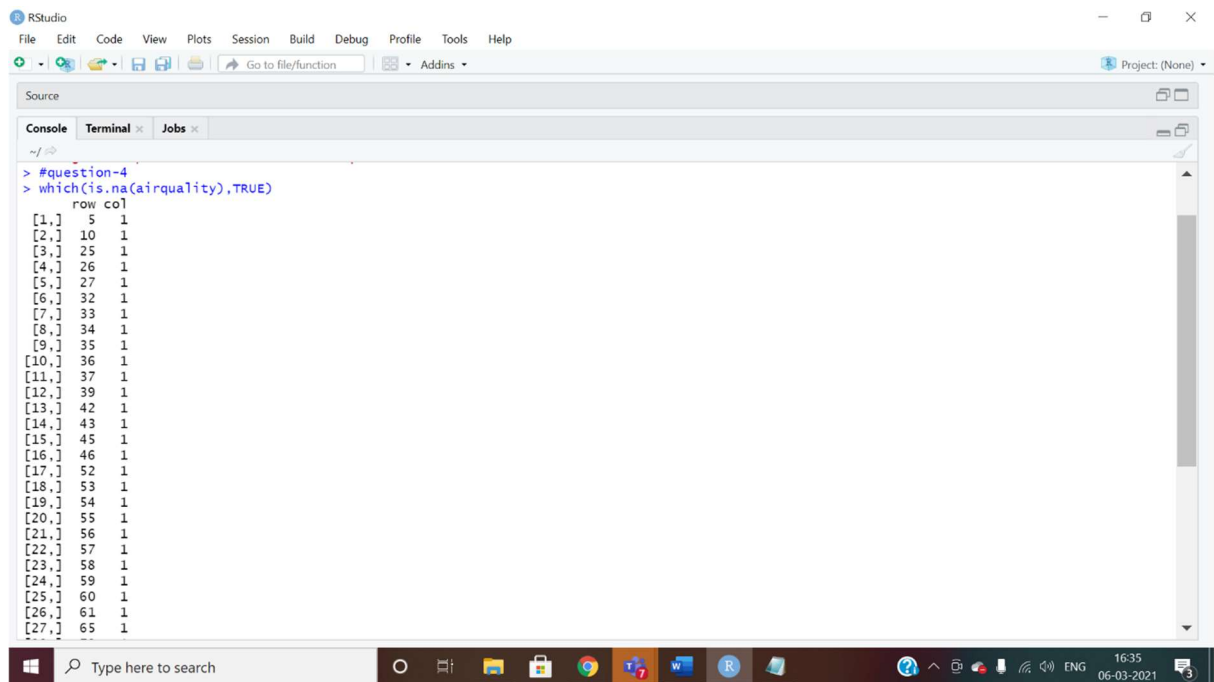
```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Source
Console Terminal Jobs
~/
> #question-4
> which(is.na(airquality))
[1] 5 10 25 26 27 32 33 34 35 36 37 39 42 43 45 46 52 53 54 55 56 57 58 59 60 61 65 72 75 83 84 102 103
[34] 107 115 119 150 158 159 164 180 249 250 251
>
```

- Print indices of the missing values in row and column number wise (Hint: Use function `which()` and argument `arr.ind = TRUE`)

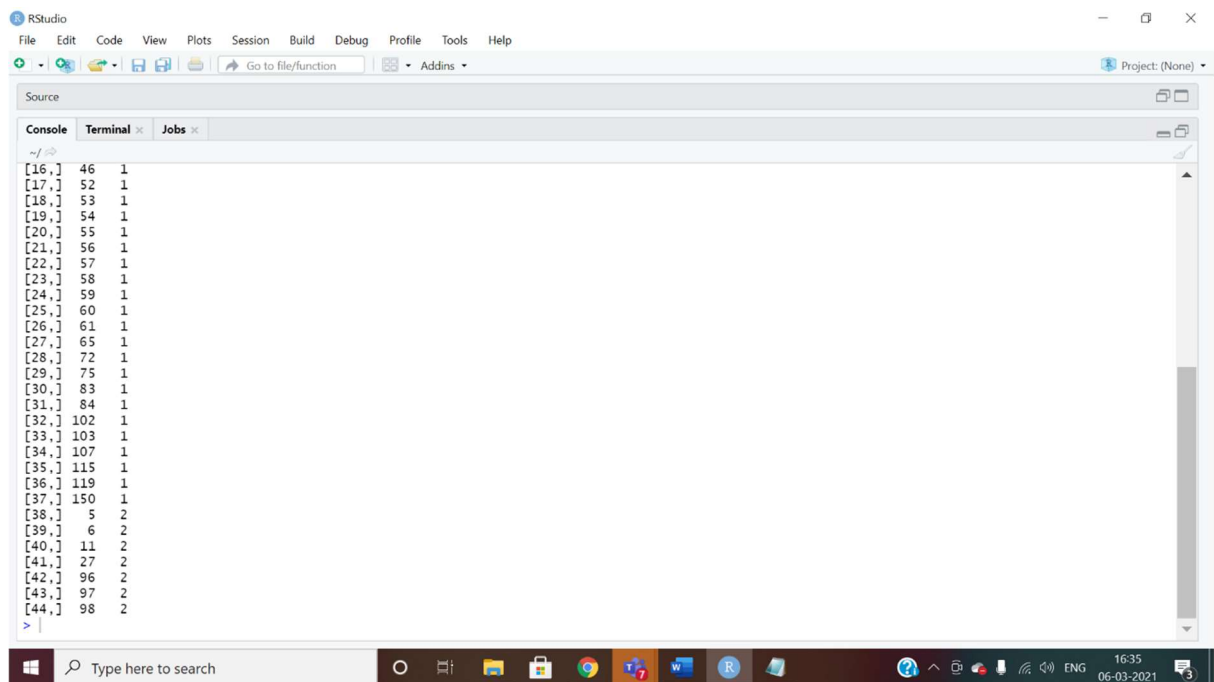
Answer:

`which(is.na(airquality),TRUE)`

OUTPUT



```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins
Project: (None)
Source
Console Terminal Jobs
~#
> #question-4
> which(is.na(airquality), TRUE)
      row col
[1,]    5  1
[2,]   10  1
[3,]   25  1
[4,]   26  1
[5,]   27  1
[6,]   32  1
[7,]   33  1
[8,]   34  1
[9,]   35  1
[10,]  36  1
[11,]  37  1
[12,]  39  1
[13,]  42  1
[14,]  43  1
[15,]  45  1
[16,]  46  1
[17,]  52  1
[18,]  53  1
[19,]  54  1
[20,]  55  1
[21,]  56  1
[22,]  57  1
[23,]  58  1
[24,]  59  1
[25,]  60  1
[26,]  61  1
[27,]  65  1
```



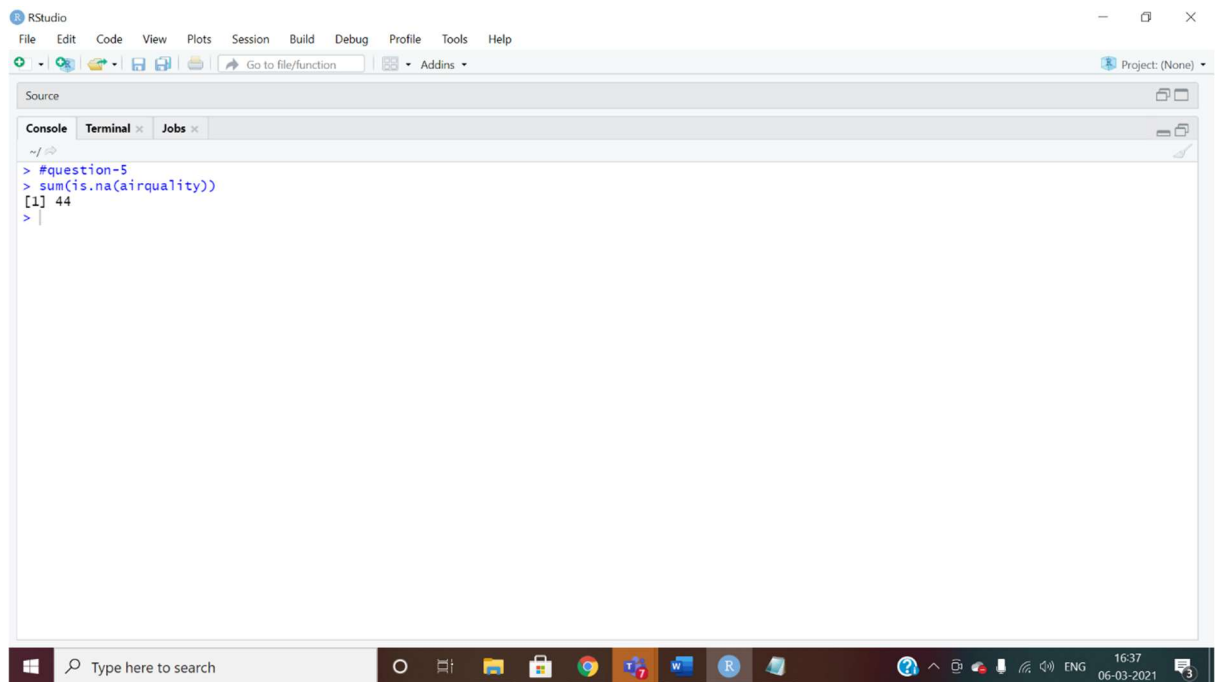
```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins
Project: (None)
Source
Console Terminal Jobs
~#
[16,]  46  1
[17,]  52  1
[18,]  53  1
[19,]  54  1
[20,]  55  1
[21,]  56  1
[22,]  57  1
[23,]  58  1
[24,]  59  1
[25,]  60  1
[26,]  61  1
[27,]  65  1
[28,]  72  1
[29,]  75  1
[30,]  83  1
[31,]  84  1
[32,] 102  1
[33,] 103  1
[34,] 107  1
[35,] 115  1
[36,] 119  1
[37,] 150  1
[38,]    5  2
[39,]    6  2
[40,]   11  2
[41,]   27  2
[42,]   96  2
[43,]   97  2
[44,]   98  2
>
```

• How many missing values are in the dataset airquality?

Answer:

`sum(is.na(airquality))`

OUTPUT

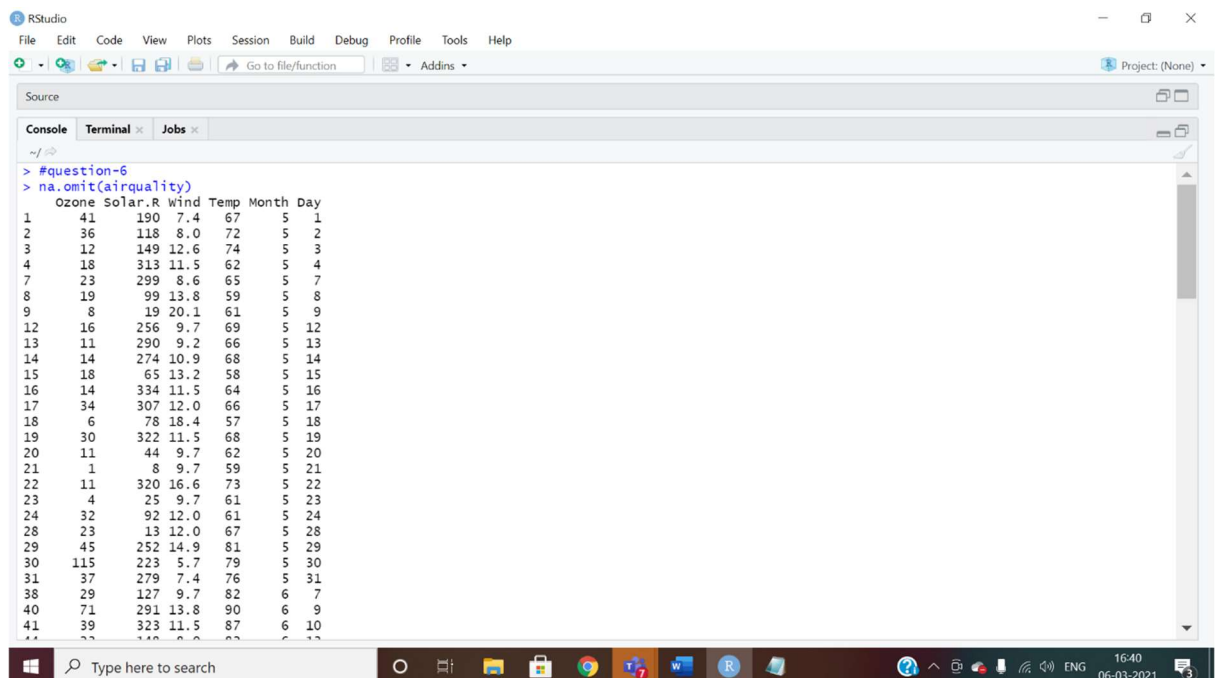


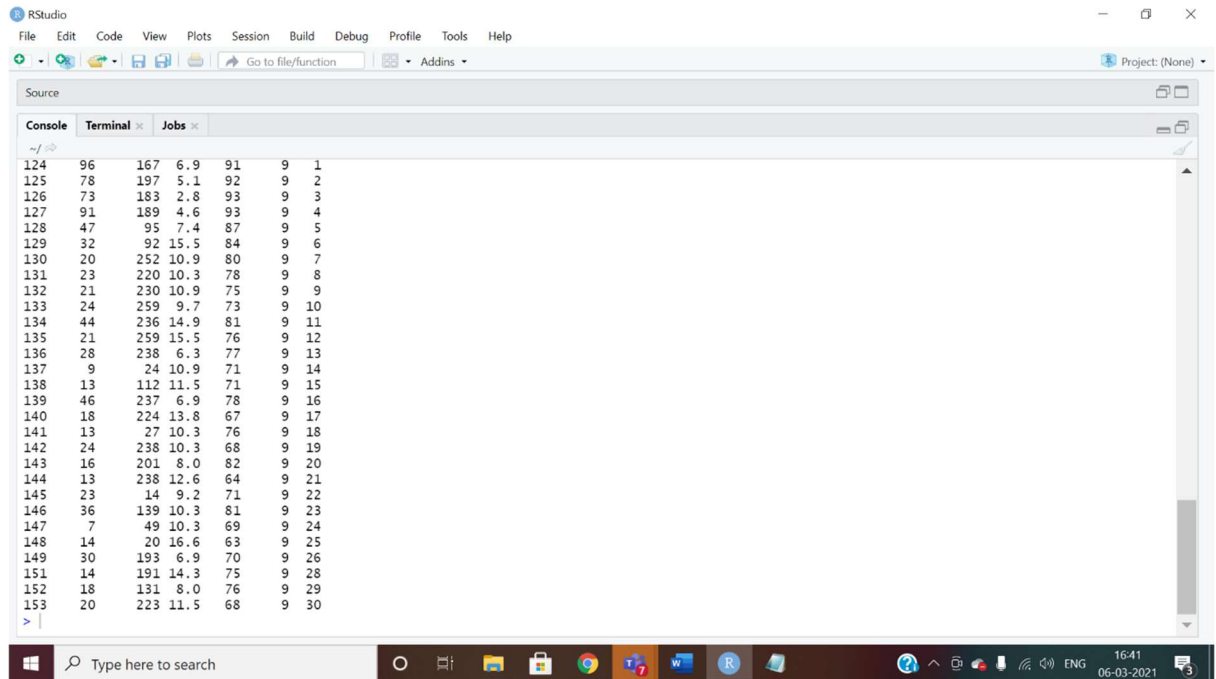
· How would you omit all rows containing missing values?

Answer:

`na.omit(airquality)`

OUTPUT



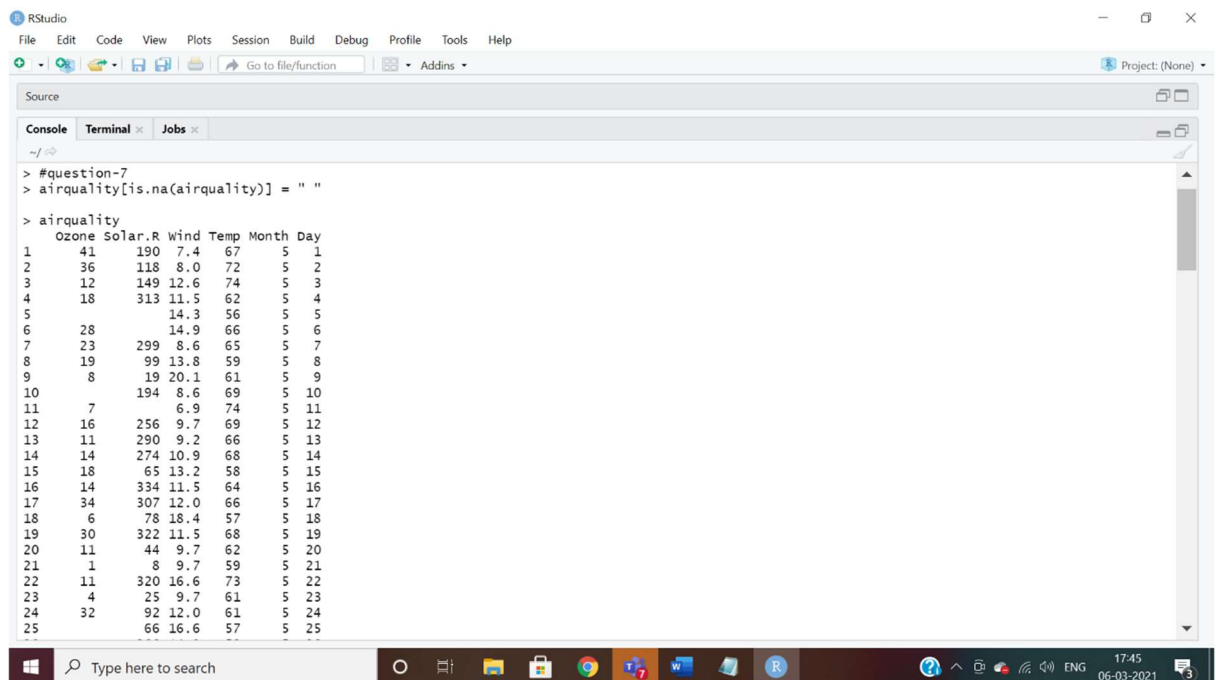


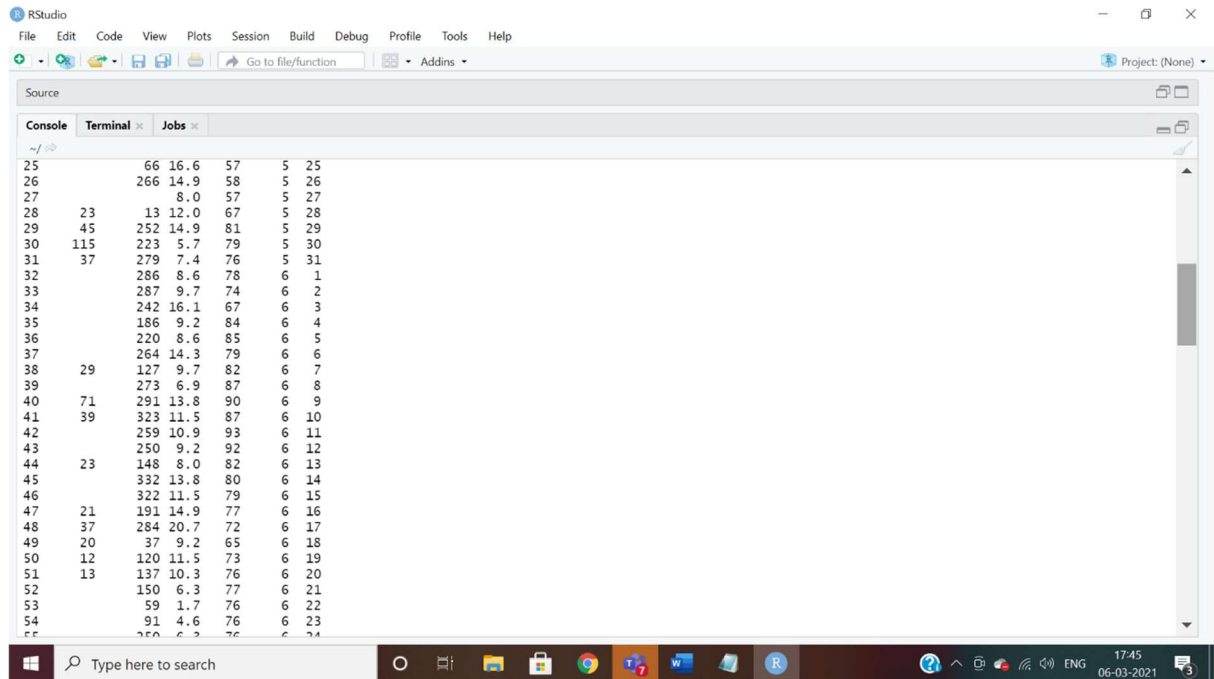
• Print the records without missing values in the dataset `airquality` using the function

Answer:

```
airquality[is.na(airquality)] = " "
```

OUTPUT



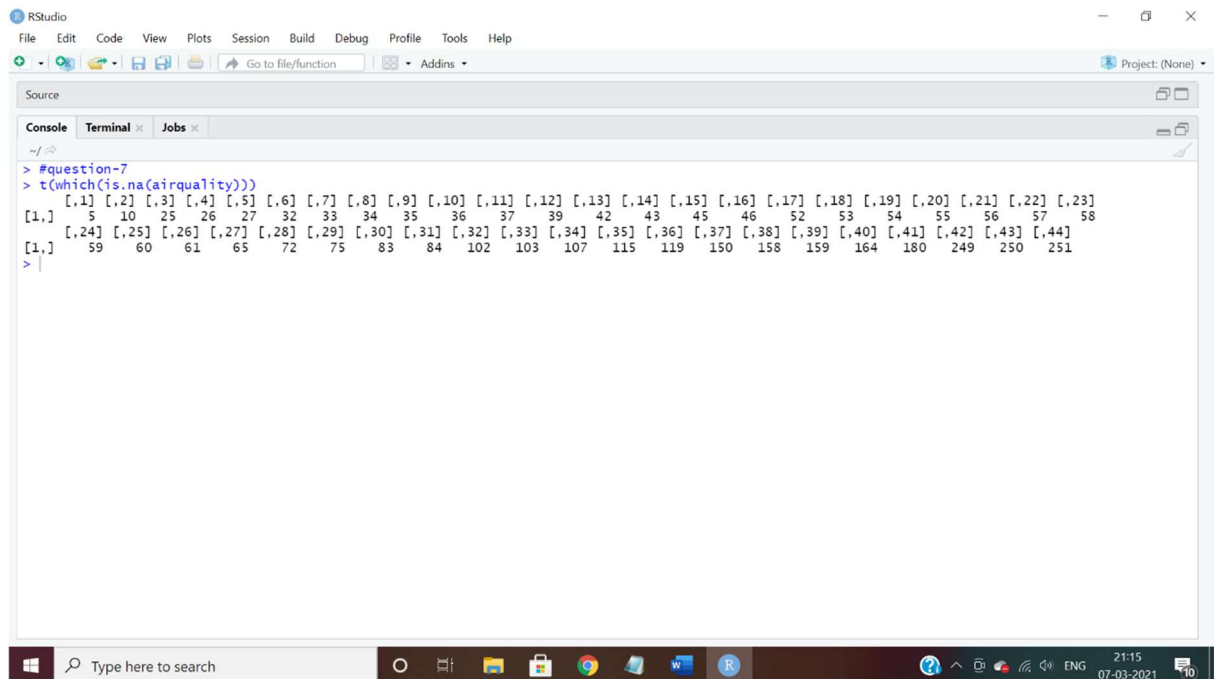


• Print the indices of the missing values in the dataset `airquality` in column major representation

Answer:

`t(which(is.na(airquality)))`

OUTPUT

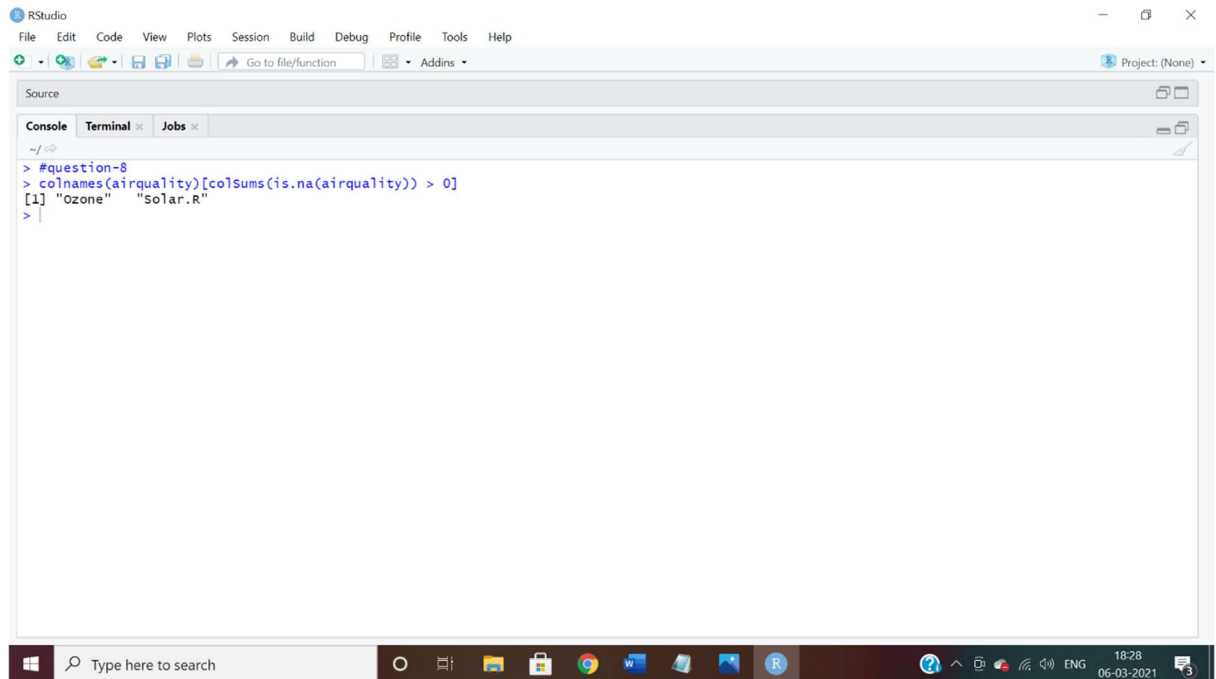


• Print Names of the Columns which contains Missing Values in the dataset airquality

Answer:

```
colnames(airquality)[colSums(is.na(airquality)) > 0]
```

OUTPUT



The screenshot shows the RStudio interface. The console window displays the following R code and its output:

```
> #question-8
> colnames(airquality)[colSums(is.na(airquality)) > 0]
[1] "ozone" "solar.R"
```

The output indicates that the columns 'ozone' and 'solar.R' contain missing values.

Data Recording

Consider a numeric vector `x <- c(3,4,5,6,7,8)`

Write a command to recode the values less than 6 with zero in the vector `x`

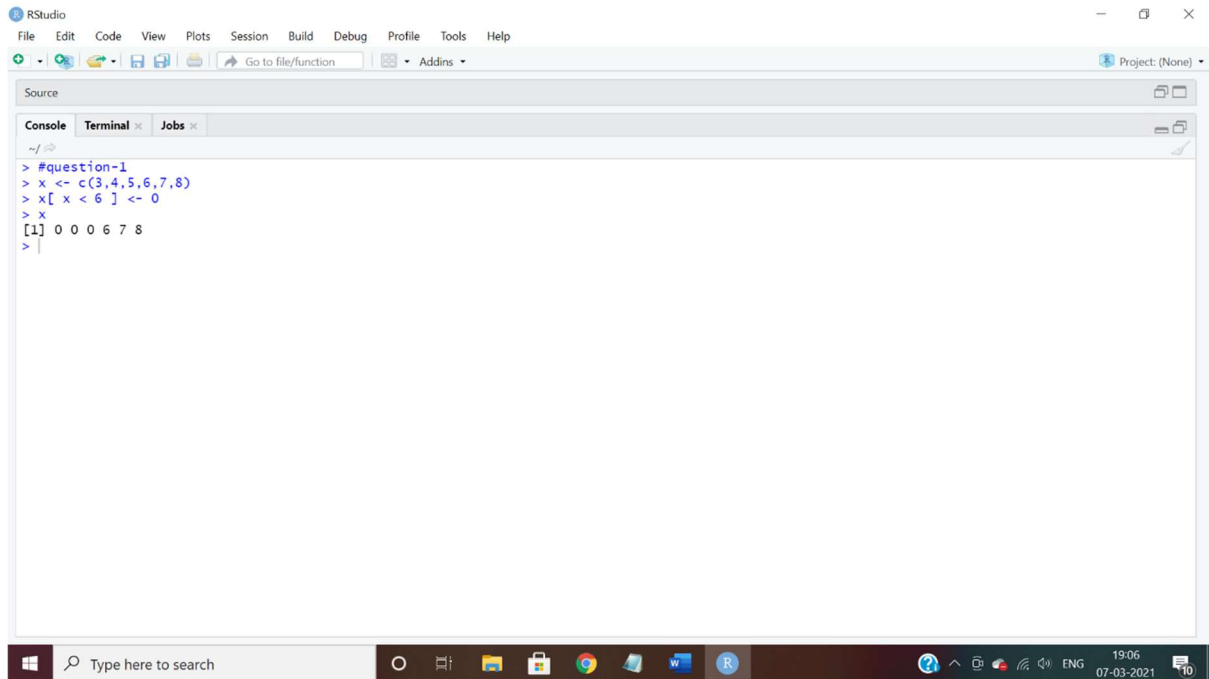
Answer:

```
x <- c(3,4,5,6,7,8)
```

```
x[ x < 6 ] <- 0
```

```
x
```

OUTPUT



RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins

Project: (None)

Source

Console Terminal Jobs

```
~/  
> #question-1  
> x <- c(3,4,5,6,7,8)  
> x[x < 6] <- 0  
> x  
[1] 0 0 0 6 7 8  
>
```

Type here to search

19:06 07-03-2021

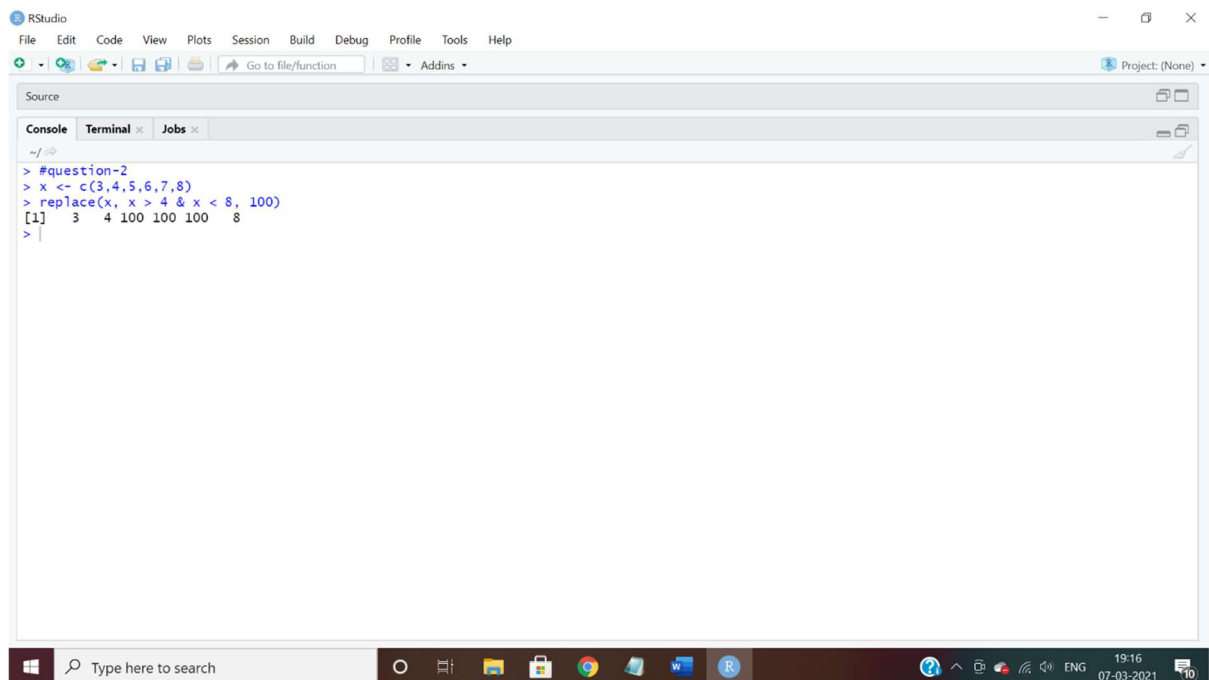
- Write a command to recode the values between 4 and 8 with 100

Answer:

```
x <- c(3,4,5,6,7,8)
```

```
replace(x, x > 4 & x < 8, 100)
```

OUTPUT



RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins

Project: (None)

Source

Console Terminal Jobs

```
~/  
> #question-2  
> x <- c(3,4,5,6,7,8)  
> replace(x, x > 4 & x < 8, 100)  
[1] 3 4 100 100 100 8  
>
```

Type here to search

19:16 07-03-2021

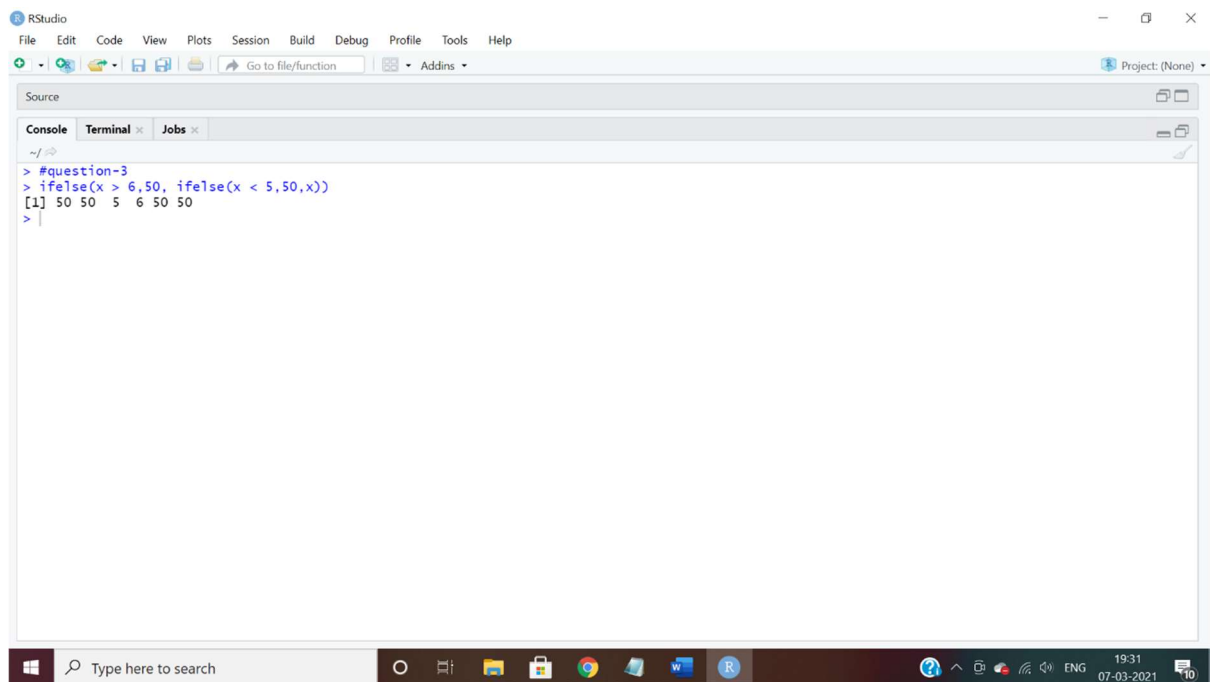
- Write a command to recode the values that are less than 5 or greater than 6 with 50

Answer:

```
x <- c(3,4,5,6,7,8)
```

```
ifelse(x > 6,50, ifelse(x < 5,50,x))
```

OUTPUT



- Write a command to recode the values less than 6 with NA in the vector x

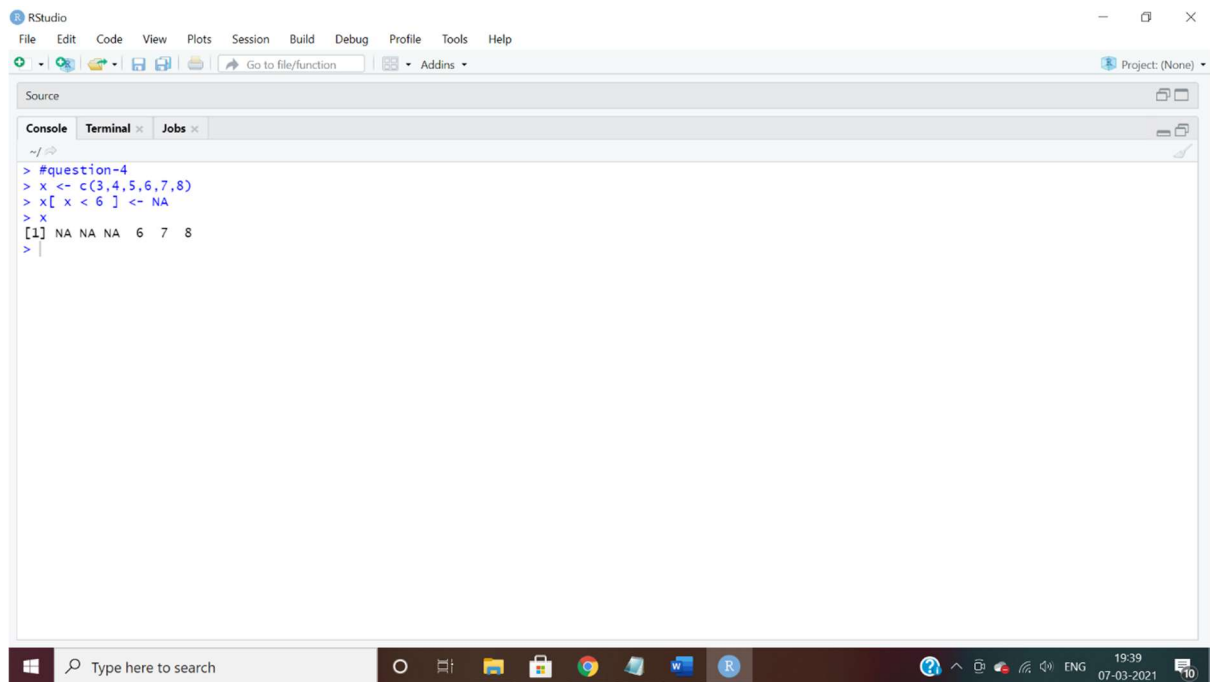
Answer:

```
x <- c(3,4,5,6,7,8)
```

```
x[ x < 6 ] <- NA
```

```
x
```

OUTPUT



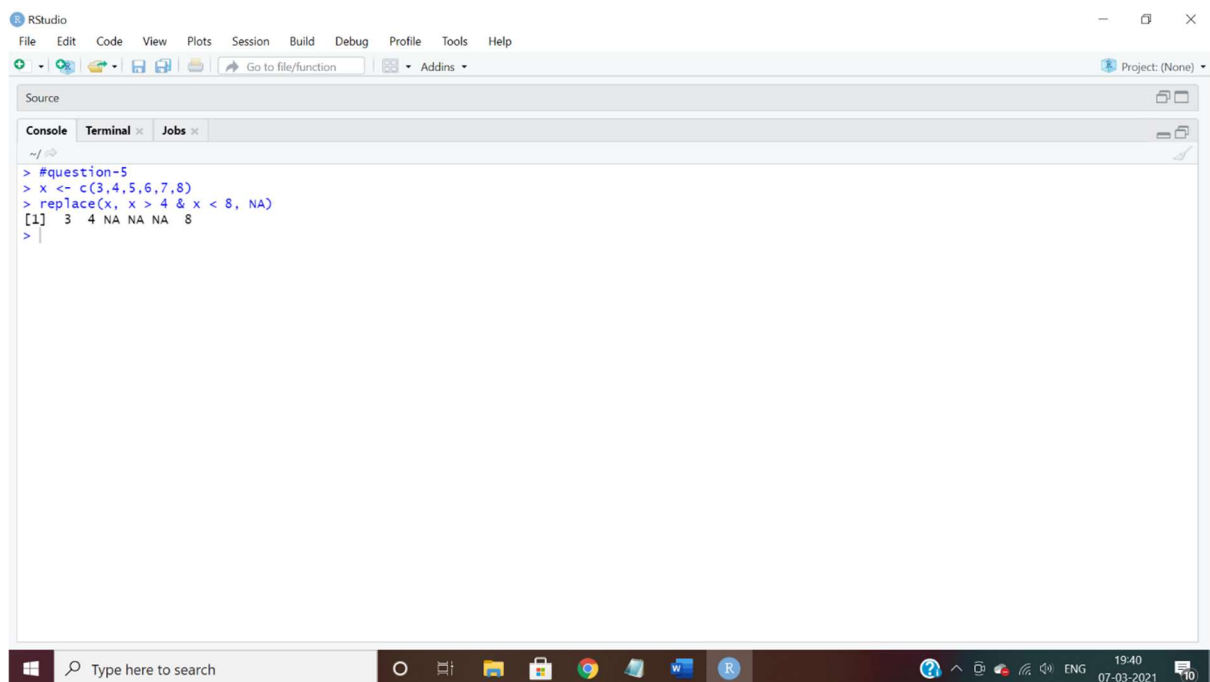
- Write a command to recode the values between 4 and 8 with NA

Answer:

```
x <- c(3,4,5,6,7,8)
```

```
replace(x, x > 4 & x < 8, NA)
```

OUTPUT



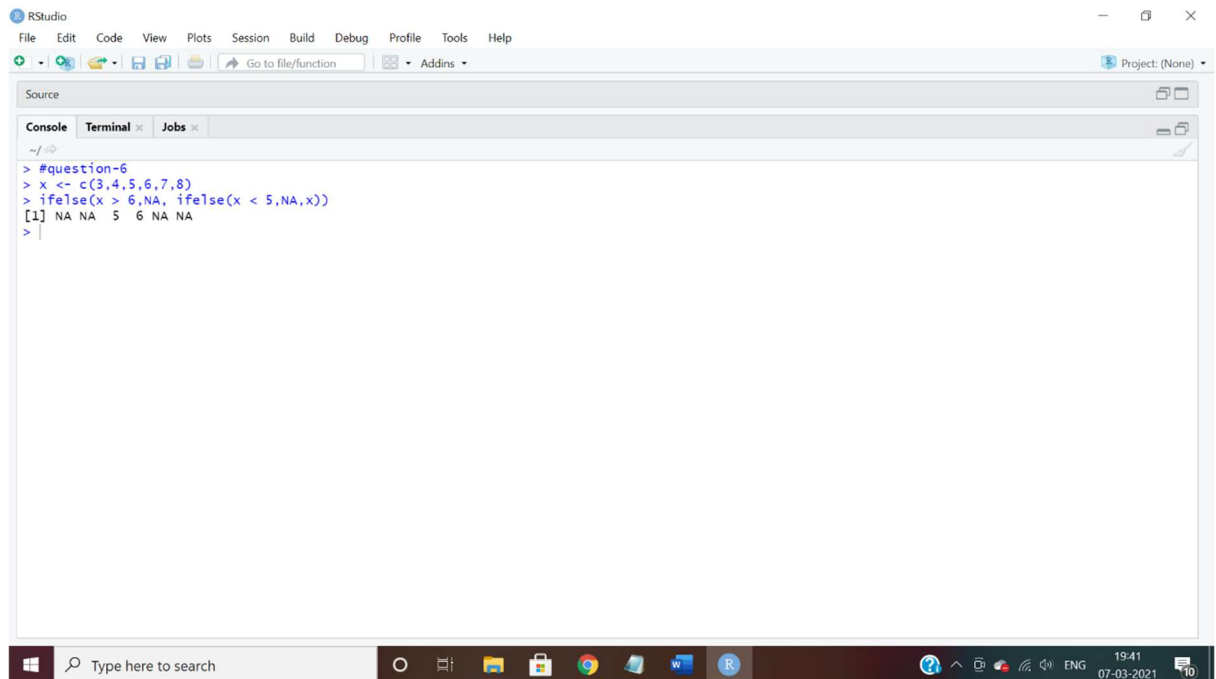
- Write a command to recode the values that are less than 5 or greater than 6 with NA

Answer:

```
x <- c(3,4,5,6,7,8)
```

```
ifelse(x > 6,NA, ifelse(x < 5,NA,x))
```

OUTPUT



The screenshot shows the RStudio interface with the console pane active. The following R code is entered and executed:

```
> #question-6  
> x <- c(3,4,5,6,7,8)  
> ifelse(x > 6,NA, ifelse(x < 5,NA,x))  
[1] NA NA 5 6 NA NA  
>
```

The output of the code is displayed in the console: [1] NA NA 5 6 NA NA. The RStudio window title is 'RStudio' and the project is '(None)'. The Windows taskbar at the bottom shows the date and time as 19:41 on 07-03-2021.