# Introduction to R

**School of Computer Science
University of Windsor**

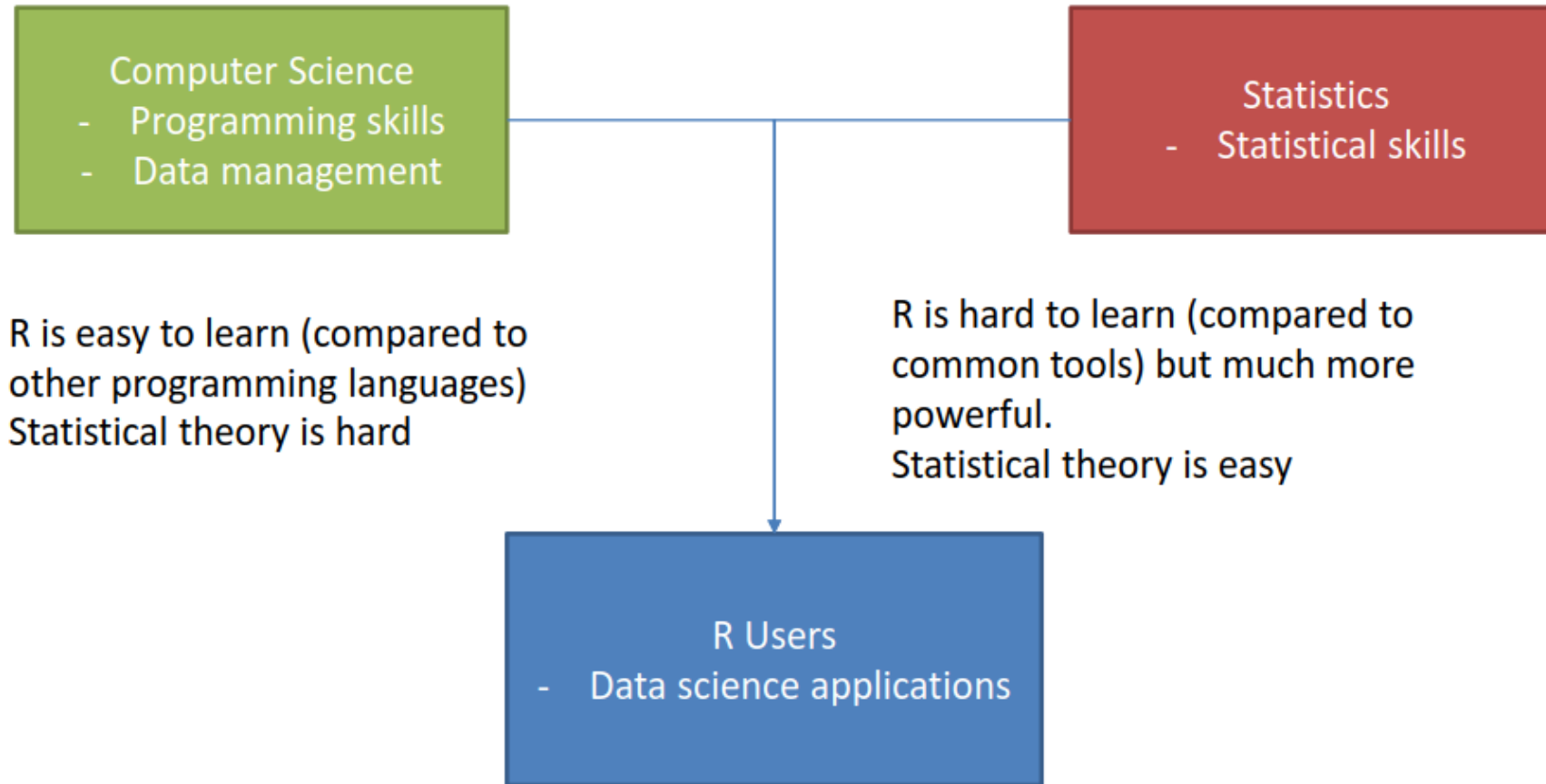**Dr Khan**

# Lecture Content

- Introduction to R
  - GUI,
  - Data Import & Export
  - Attributes & Data Types,
  - Descriptive Statistics

University of Windsor

# Introduction to R

- R is a programming language and software framework for statistical analysis and graphics,

# Installing R & R Studio

- You should install R & R Studio as follow:
  1. Install R version 3.3.1 from: **https://www.r-project.org/** (The R programming Language)
  2. Install Free AGPL Rstudio (GUI to R) **https://www.rstudio.com/products/rstudio/download2/**
  3. Your computer should access the internet during all R sessions

- The Website contains for R documentation is: **https://cran.r-project.org/doc/manuals/r-release/R-intro.html**

- https://www.youtube.com/watch?v=cX532N_XLIs

University of Windsor

# R Studio

# IDE

**Console**

- Where you type commands and receive text output.

**Script Window**

- Script files are text files used to store scripts of R commands. Multiple can be open at once.

- Source runs an entire file.

- Run runs a highlighted selection.

- Write multiline code, including functions, in a script file and then run them from there.

University of Windsor
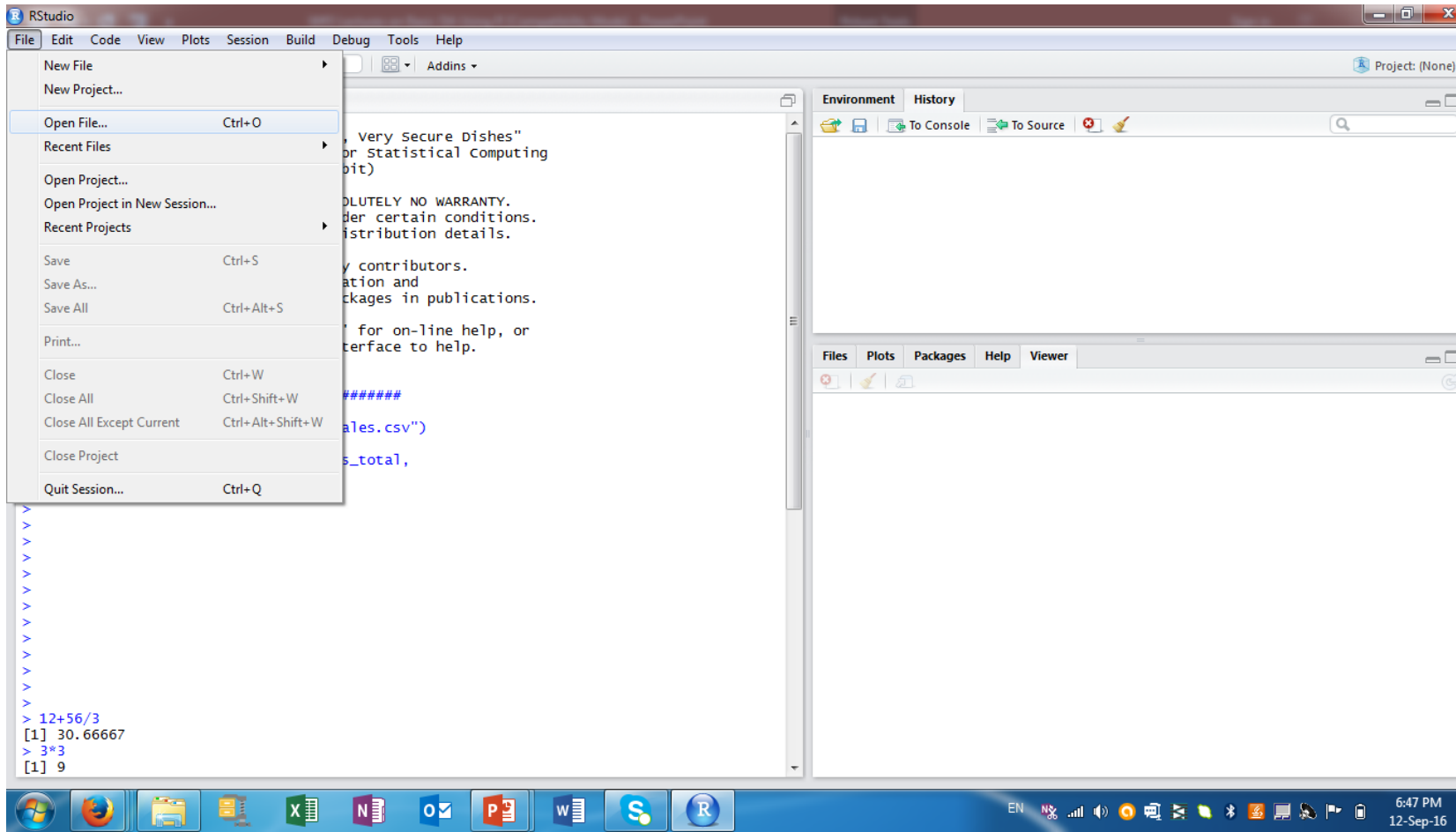
# IDE

## Environment & History

- Environment – Display the objects (including functions) present in the environment.

- Shows you the names of all the data objects (like vectors, matrices, and data frames) that you've defined in your current R session. You can also see information like the number of observations and rows in data objects.

- History – Display commands previously entered into the console.

## Files, Plots, Packages, Help & Viewer Window

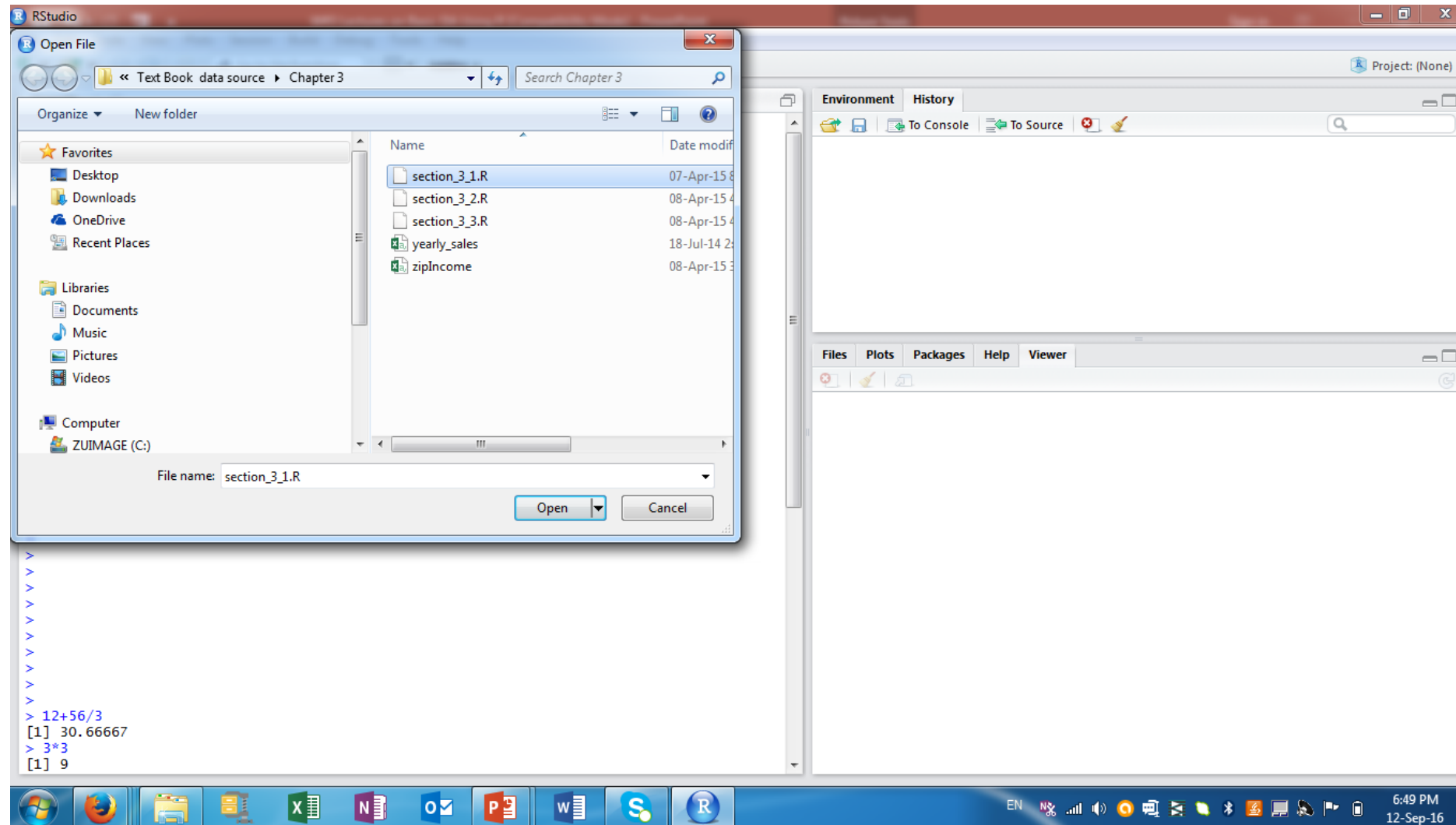- Files – Navigate your computer's file system. Double clicking a file will open it in the script window.

- Plots – Basic graphic output. Export graphics using the export button.

- Packages – Manage packages.

- Help – Displays help information.

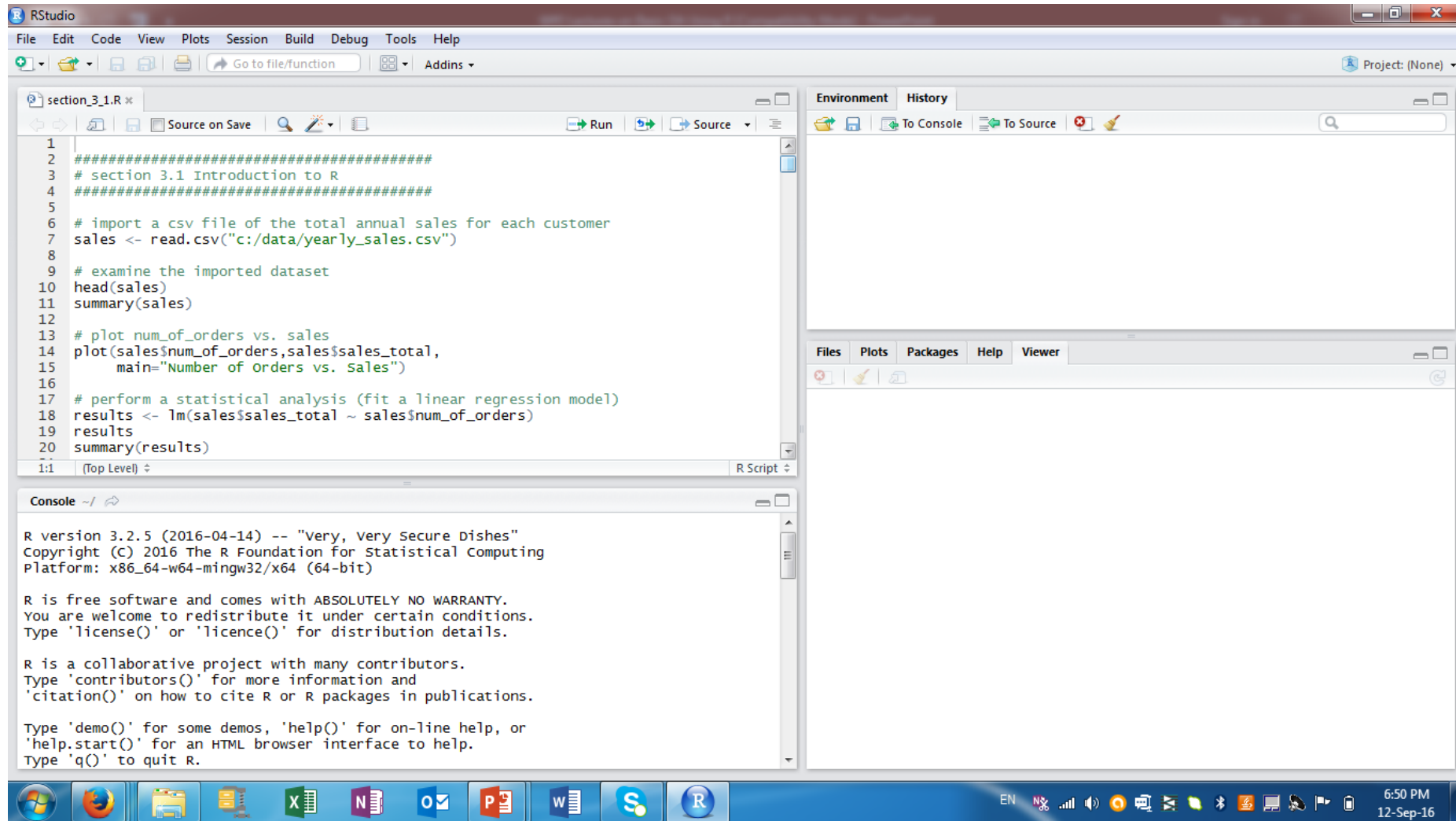- Viewer – Used to view local web content, web graphics and local web applications. We will not use it.

University of Windsor

# Open File/Project

# R code

# Script Window

# Introduction to R

```r
# import a csv file of the total annual sales for each customer
sales <- read.csv("c:/data/yearly_sales.csv")
# examine the imported dataset
head(sales)
summary(sales)
# plot num_of_orders vs. sales
plot(sales$num_of_orders,sales$sales_total, main="Number of Orders vs. Sales")
# Get the working directory
getwd()
# Set the working directory
setwd("D:/Users/Z10596/Desktop/R_files")

# Add a column for the average sales per order
sales$per_order <- sales$sales_total/sales$num_of_orders

# export data as tab delimited without the row names
write.table(sales,"sales_modified.txt", sep="\t", row.names = FALSE)
```
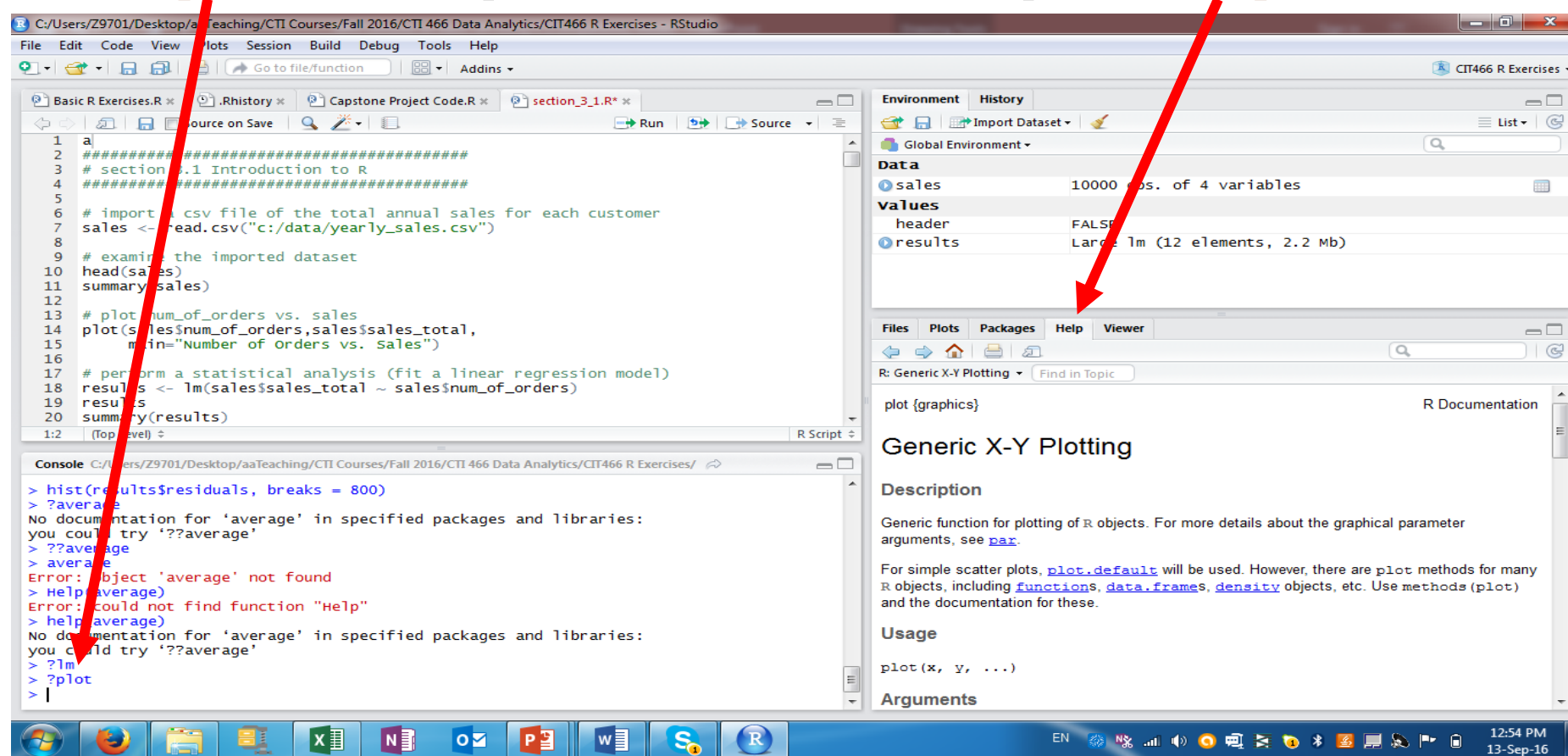
University of Windsor

# Accessing Help in R Studio

You can either use **help(R function)** or use **? R command/function**

Below **?plot** asks R to explain what Plot means and response in **Help Window**

# Import CSV Data Set file

**sales <- read.csv("c:/data/yearly_sales.csv")** means Import *yearly_sales.csv* dataset file and **(<- )** means save it into a file called ***Sales***



*Read-csv* imports the *Yearly_sales.csv* file and save it into the file *Sales*

# Head () Function

**Head (Sales)** *function by default list the six Records of Sales as shown below*

# Summary() Function
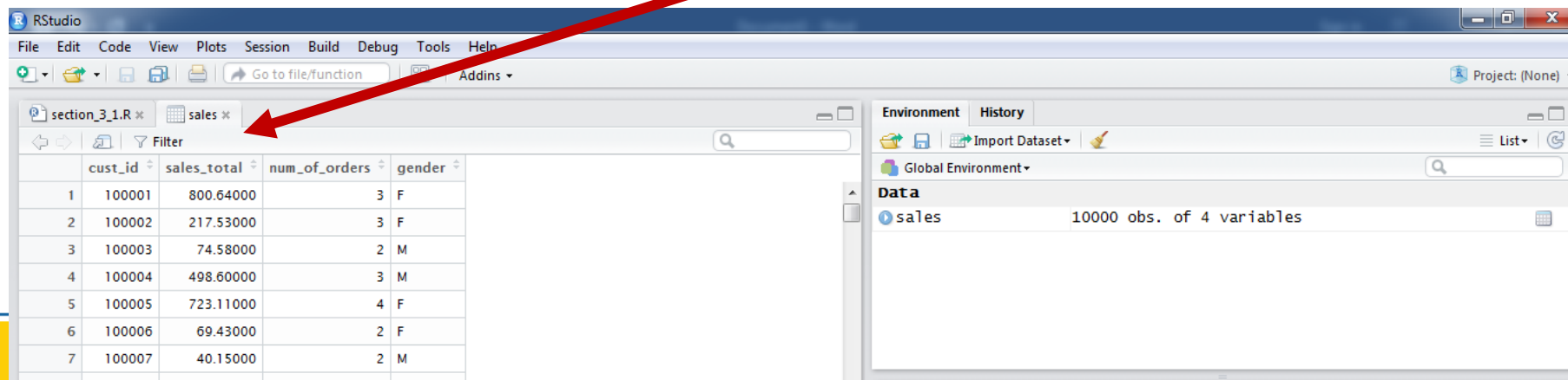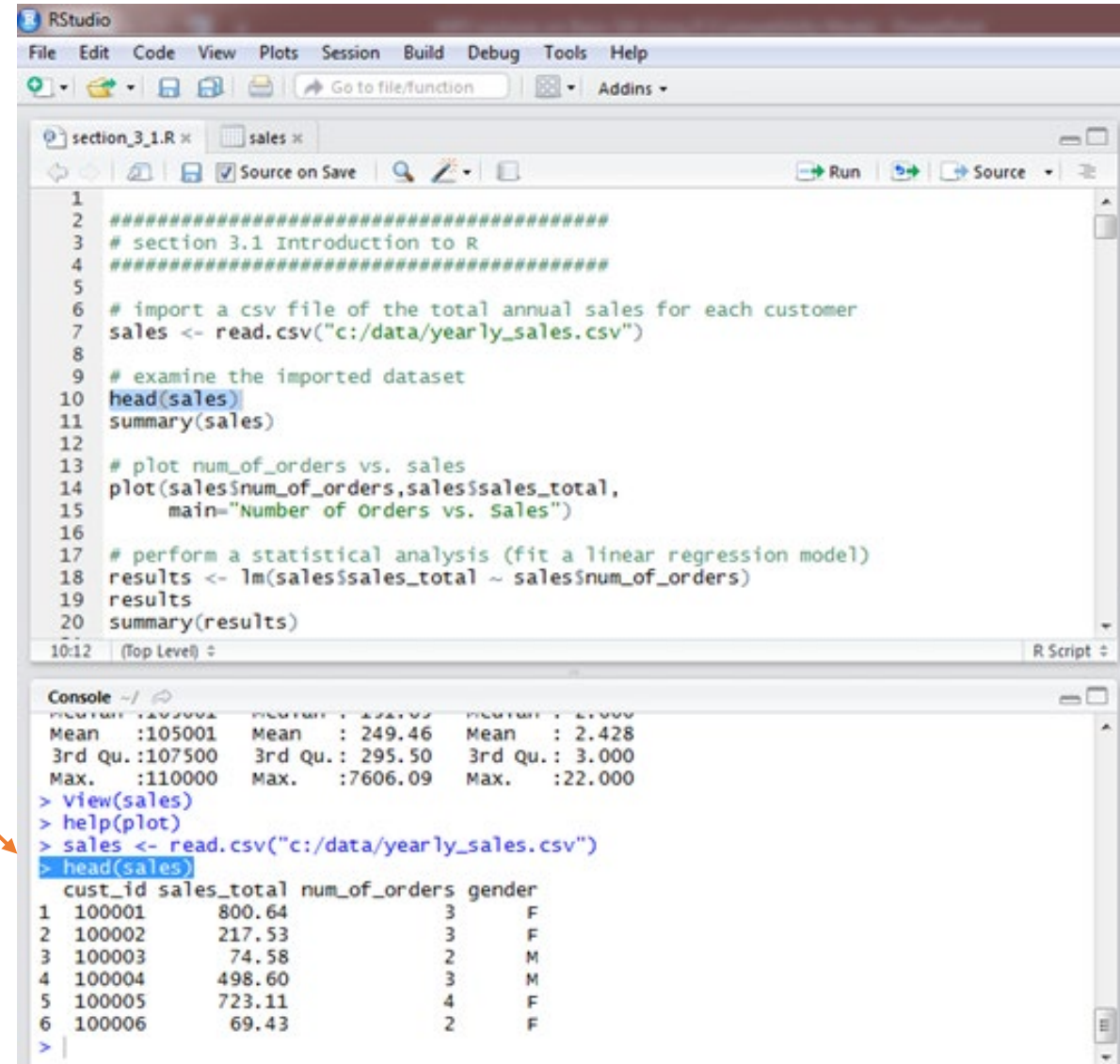
- **Summary()**
  Function provides
  some descriptive
  statistics such as
  Means and
  Median, etc.

```
 1
 2  #########################################
 3  # section 3.1 Introduction to R
 4  #########################################
 5
 6  # import a csv file of the total annual sales for each customer
 7  sales <- read.csv("c:/data/yearly_sales.csv")
 8
 9  # examine the imported dataset
10  head(sales)
11  summary(sales)
12
13  # plot num_of_orders vs. sales
14  plot(sales$num_of_orders,sales$sales_total,
15       main="Number of Orders vs. Sales")
16
17  # perform a statistical analysis (fit a linear regression model)
18  results <- lm(sales$sales_total ~ sales$num_of_orders)
19  results
20  summary(results)
```

11:15    (Top Level) ⬦                                          R Script

```
Console ~/ ⌂

        cust_id    sales_total   num_of_orders   gender
1   100001      800.64              3           F
2   100002      217.53              3           F
3   100003       74.58              2           M
4   100004      498.60              3           M
5   100005      723.11              4           F
6   100006       69.43              2           F
> summary(sales)
    cust_id        sales_total        num_of_orders      gender
 Min.   :100001  Min.   :  30.02   Min.   : 1.000   F:5035
 1st Qu.:102501  1st Qu.:  80.29   1st Qu.: 2.000   M:4965
 Median :105001  Median : 151.65   Median : 2.000
 Mean   :105001  Mean   : 249.46   Mean   : 2.428
 3rd Qu.:107500  3rd Qu.: 295.50   3rd Qu.: 3.000
 Max.   :110000  Max.   :7606.09   Max.   :22.000
>
```
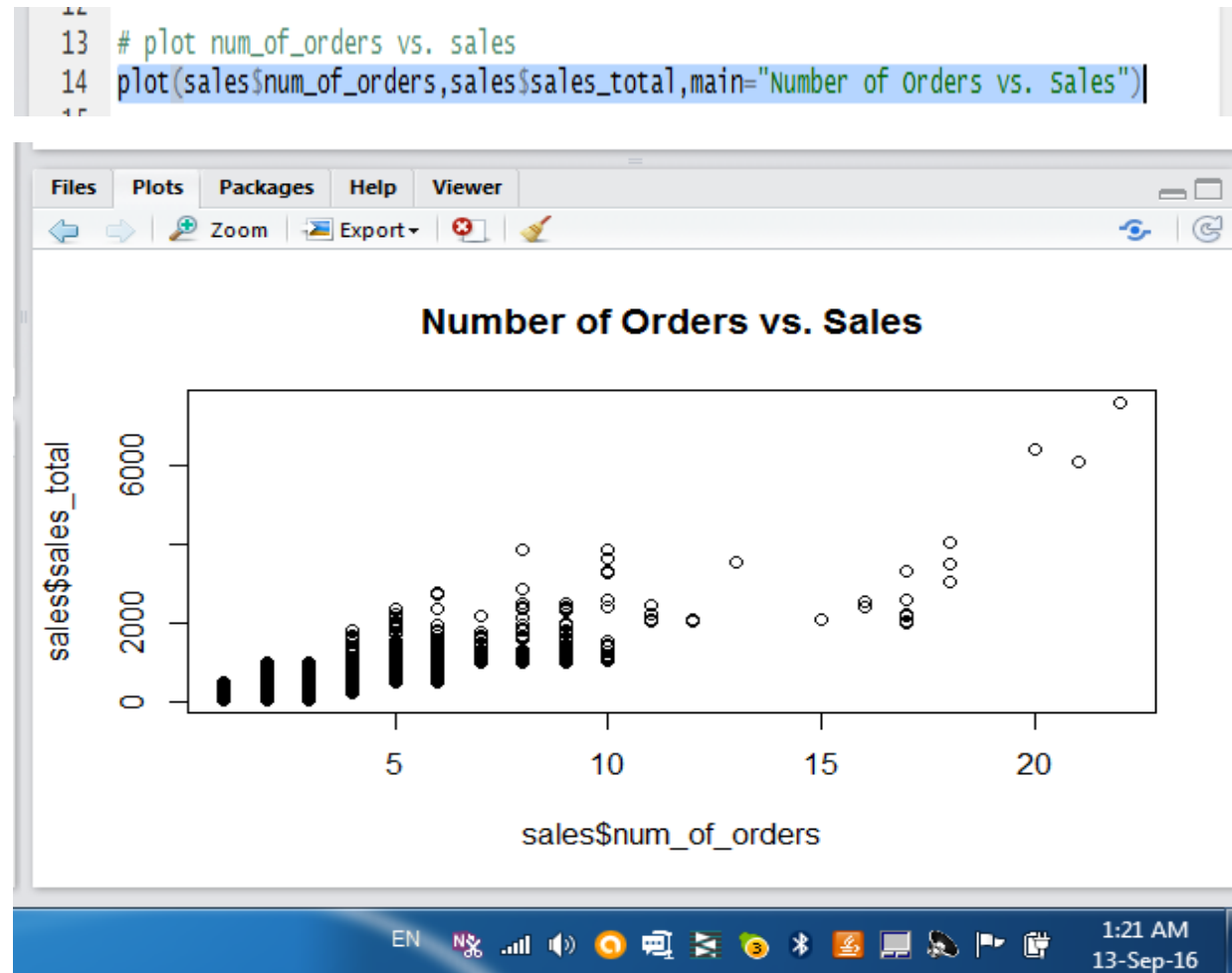
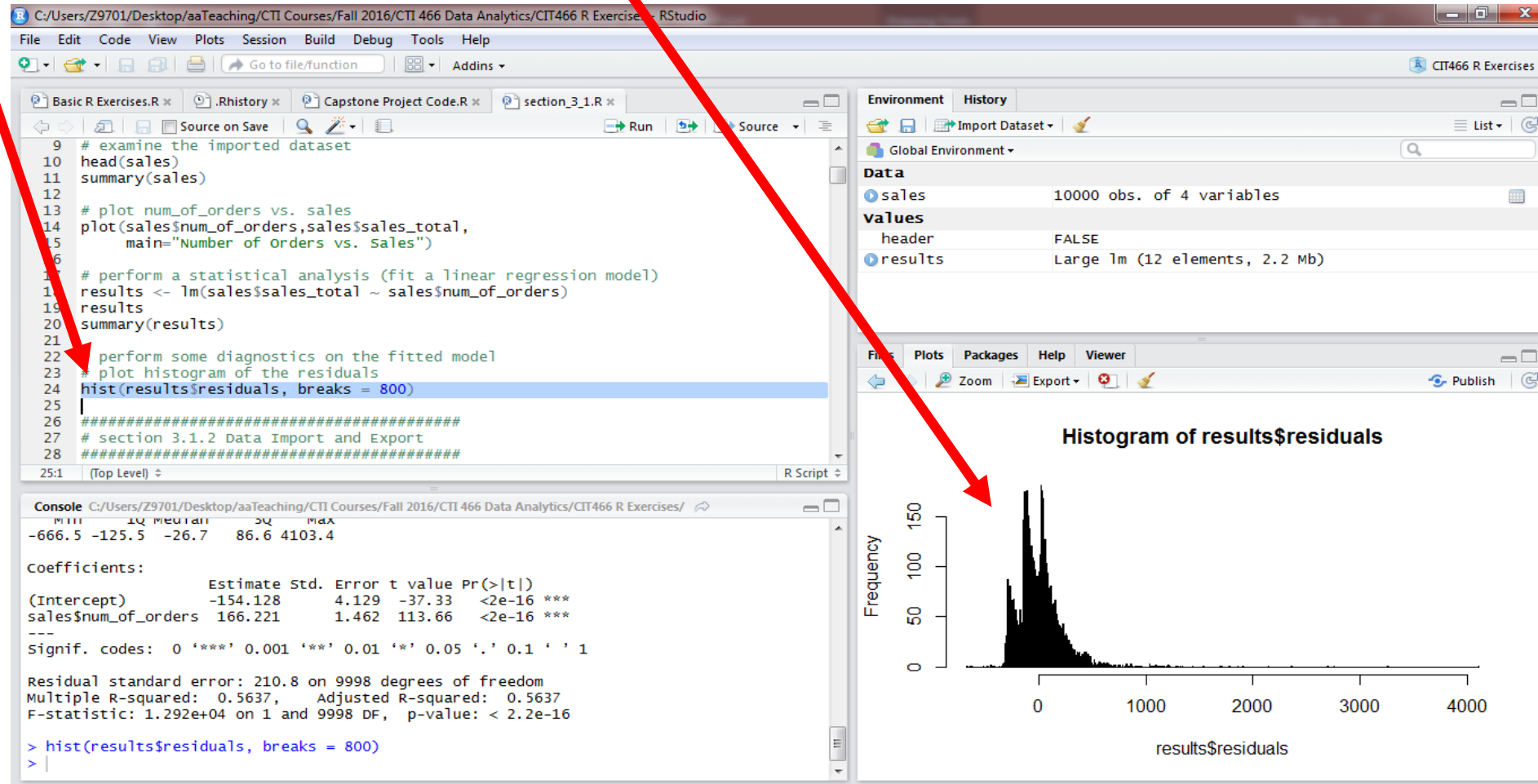University of Windsor

# Plot () function

- Plotting a dataset's content can provide information about the relationships between the various column,

- In this example, Plot() function generate a scatterplot of the number of orders (*Sales$sum_of_orders*) against the annual sales (*Sales$sales_toltal*)

- NB: **$** selects an attribute of a table e.g. *sum_of_orders* attribute of *Sales* Table

# Performing diagnostic on the model

**Hist()** function draws histogram of Residual Results to analyze the model. Here we have large residuals

# Data Import and Export

# Example of CSV files



**Weights3468.csv - Notepad**

```
"Date","Weight""Wed Jun 30 08:00:01 GMT 20
Jun 29 08:00:01 GMT 2010","180.2""Mon Jun
2010","180.2""Sun Jun 27 08:00:01 GMT 2010
Jun 26 08:00:01 GMT 2010","180.2""Fri Jun
2010","180.2""Thu Jun 24 08:00:01 GMT 2010
Jun 23 08:00:01 GMT 2010","180.2""Tue Jun
2010","181.4""Mon Jun 21 08:00:01
Jun 20 08:00:01 GMT 2010","181.4"
2010","181.4""Fri Jun 18 08:00:01
Jun 17 08:00:01 GMT 2010","181.4""
2010","181.4""Tue Jun 15 08:00:01
Jun 14 08:00:01 GMT 2010","181.4"
2010","181.4""Sat Jun 12 08:00:01
Jun 11 08:00:01 GMT 2010","180.0"
2010","180.0""Wed Jun 09 08:00:01
Jun 08 08:00:01 GMT 2010","180.0"
2010","180.0""Sun Jun 06 08:00:01
Jun 05 08:00:01 GMT 2010","178.2"
2010","178.2""Thu Jun 03 08:00:01
Jun 02 08:00:01 GMT 2010","178.2"
2010","178.2""Mon May 31 08:00:01
May 30 08:00:01 GMT 2010","178.2"
2010","178.2""Fri May 28 08:00:01
May 27 08:00:01 GMT 2010","178.2"
2010","178.2""Tue May 25 08:00:01
May 24 08:00:01 GMT 2010","178.2"
2010","178.2""Sat May 22 08:00:01
May 21 08:00:01 GMT 2010","178.2""
```

**widgets.csv - Notepad**

```
Widget1, blue, £10
Widget2, red, £12
Widget3, green, £14
Widget4, black, £16
Widget5, white, £18
```

**Text file - Notepad**

```
Index
One
Print Runs (x1000)
Page numbers
Orders (x1000)
 1    1    2800      22.
 2    1    2670      14.
 3    1    2800      37.
 4    1    2784      15.
 5    1    2800      38.
 6    1    2620     172.
 7    1    2620     249.
 8    1    2470      84.
 9    1    2620     242.
10    1    2475     100.
11    1    2620     114.
```

**yearly_sales - Notepad**

```
cust_id,sales_total,num_of_orders,gender
100001,800.64,3,F
100002,217.53,3,F
100003,74.58,2,M
100004,498.6,3,M
100005,723.11,4,F
100006,69.43,2,F
100007,40.15,2,M
100008,58.61,2,M
100009,364.63,2,F
100010,44.31,2,M
100011,216.41,1,F
100012,157.92,2,F
100013,289.58,1,M
100014,1044.4,7,M
100015,82.3,3,M
```

University of Windsor

# Usage of read.csv function

*read.csv()* converts Comma Separated Values (CSV) file into formatted Column & Row table and upload into R aeropospace as shown below

# Data Import and Export

- In the annual Sales example the dataset was imported using *read.csv* as follow:  **sales <- read.csv("c:/data/yearly_sales.csv")**

- To simplify multiple files with long path names, the ***setwd()*** function can be used to set the working directory for subsequent import and export as follows:

  **setwd("c:/data/")**

  **sales <- read.csv("yearly_sales.csv")**

- Other import function include ***read.table()*** and ***read.delim()*** *function are also used* to import CSV files like *yearly_Sales.csv* or other common files such as TXT.

- There are also two additional R functions: ***read.csv2()*** and **read.delim2()**

# Main Differences between R Import Functions

| Function | Headers | Separators | Decimal Points |
|---|---|---|---|
| read.table() | FALSE | " " | "." |
| read.csv() | TRUE | "," | "." |
| read.csv2() | TRUE | ";" | "," |
| read.delim() | TRUE | "\t" | "." |
| read.delim2 | TRUE | "\t" | "." |

University of Windsor

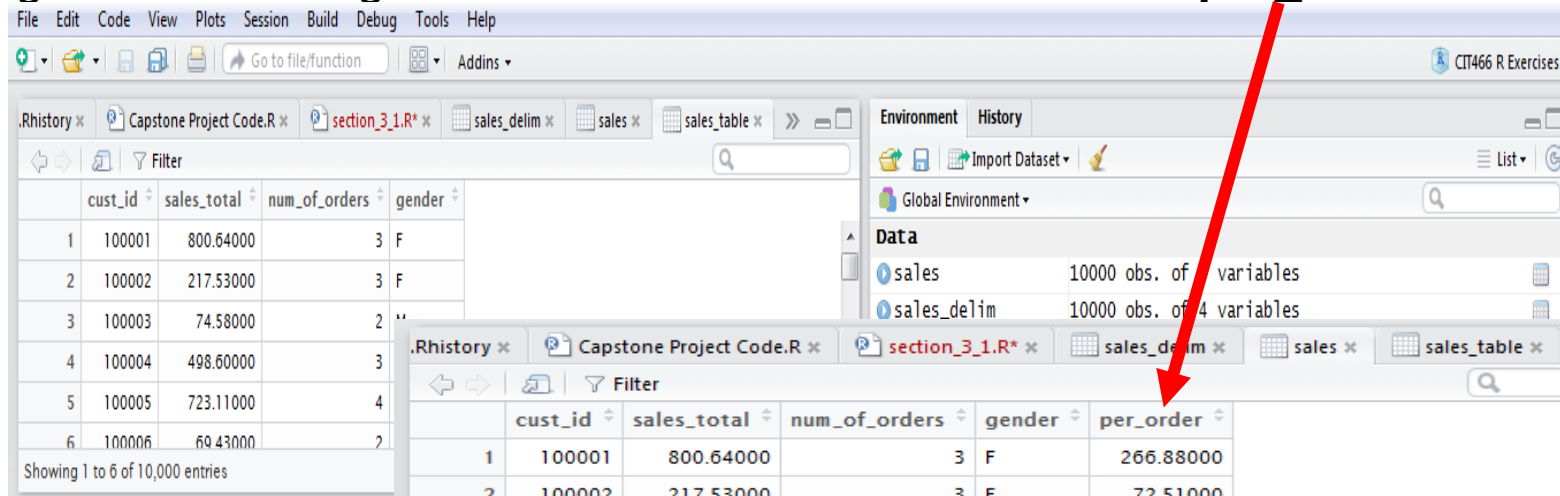# R Export Functions

- The analogue R functions are *write.table(), write.csv()* and *write.csv2() enable exporting of R data sets to an external file*

- *Example below show making change to Sales file and exporting it*

```
37
38   # add a column for the average sales per order
39   sales$per_order <- sales$sales_total/sales$num_of_orders
40   # export data as tab delimited without the row names
41   write.table(sales,"sales_modified.txt", sep="\t", row.names=FALSE)
```

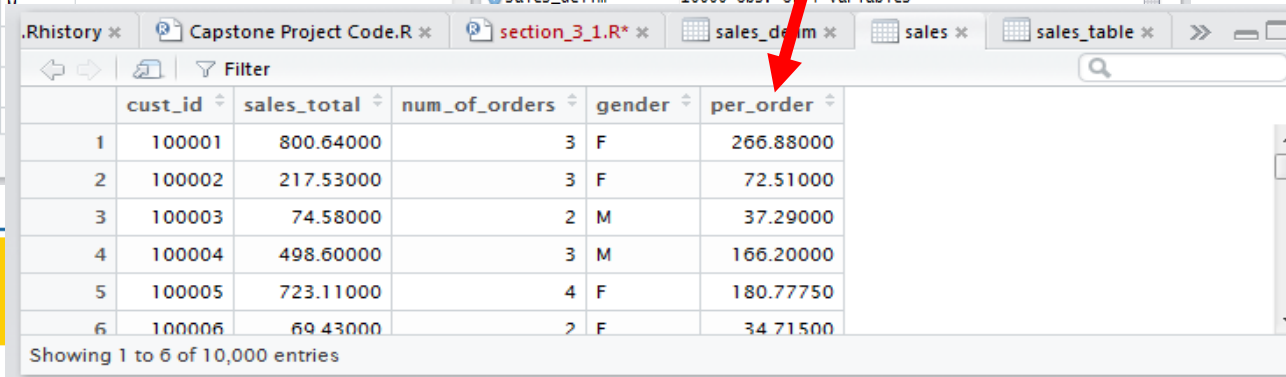This will give the following Sales table with an additional column *per_order:*



Before

After

# Any Questions

University of Windsor