

Marketing Analytics

**Technology: Statistical Analysis Software & R
R Basics**

Stephan Sorger

www.stephansorger.com

Disclaimer:

- All images such as logos, photos, etc. used in this presentation are the property of their respective copyright owners and are used here for educational purposes only
- Some material adapted from: Sorger, “Marketing Analytics: Strategic Models and Metrics”

Statistical Analysis Software: Introduction

Topic	Definition
Definition	Software designed for in-depth analysis Unlike MS Excel (general purpose spreadsheet)
Origins	SAS conceived in 1966 by Anthony J. Barr Placed statistical procedures in formatted file framework
Uses	Advanced statistical techniques Nonlinear functions; Multiple regression; Conjoint
Advantages	Powerful; Accurate; Specific tools
Disadvantages	Command line interface; steep learning curve Very expensive

Statistical Analysis Software: Major Suppliers

Criteria	SAS	SPSS	R
Market	Fortune 500	Universities	Universities
Focus	Power	Ease of use	Price
User	Power user	Student	Price-sensitive
Origins	Industry	Education	Open Source
Learning	Difficult	Moderate	Moderate
Cost	\$86,600/yr+	\$16,000/yr+	Free
UI	Command Line	Point & Click	Command Line
Database	32,768 var.	1 file at a time	
Graphics	SAS/Graph	High quality	Different packages
Analogy	Microsoft	Apple	Linux

UCLA, Statistical Software Packages Comparison, ats.ucla.edu:

http://www.ats.ucla.edu/stat/mult_pkg/compare_packages.htm

MineQuest Business Analytics, "Cost of Licensing WPS 3.0 vs. SAS 9.3." February 2013.

http://www.minequest.com/downloads/Pricing_Comparisons_Between_WPS_and_SAS.pdf

IBM SPSS Statistics website, "Buy IBM SPSS Statistics Now"

<http://www-01.ibm.com/software/analytics/spss/products/statistics/buy-now.html>

R: Introduction

Topic	Description
Description	Free statistical computing and graphics software package Widely used among statisticians and data miners Increased popularity in 2010 - on
History	Started in 1993 as implementation of S programming language (1976) R developed by Ross Ihaka and Robert Gentleman “R” from <u>R</u> oss & <u>R</u> obert, as well as play on “S”
Functions	R includes many functions, which can be expanded through packages
Data	Can handle multiple simultaneous data sets, unlike Excel Data types: scalars, vectors, matrices, data frames, and lists Vectors: numerical, character, logical
Commercial	Revolution Analytics offers enterprise version (\$)

References:

1. Venables, W.N., Smith, D.M., “An Introduction to R.” Version 3.0.1. May 16, 2013.

<http://www.cran.r-project.org/doc/manuals/R-intro.pdf>

R: Basics

Topic	Description
Commands	Based on UNIX; case sensitive Commands separated by “;” or by newline <CR>
Comments	#Hashtags to indicate comments
Prompt	> #system is waiting for you to type something Traditional version not menu-driven, unlike consumer software
Arithmetic	> 5 + 4 [1] 9 #system returns the sum of 5 + 4, which is 9
Assignment (=)	> x <- 3 # assign the number “3” to the object “x”; similar to “=” sign
Help	2 ways to get help; Example: Get help with “read.csv” command ?(read.csv) help(read.csv)

R: Basics

Topic	Description
Functions	R features a rich set of functions c() : Function c Statistics functions: mean(x); median(x); range(x); etc. Arithmetic functions: 4^2; log (10); sqrt (16)
Vector	> x <- v(1, 2, 3) # assign v, a vector of numbers to the object x
Matrix	> y <- matrix(c(1, 2, 3, 4, 5, 6), 2, 3 # create 2 x 3 matrix
Print	Ask R to print out numbers inside an object, such as a vector by printing it > print (x) # ask R to print out x > x # Or, you can just type the variable and hit return
Plot	Ask R to plot out lines based on a dataset by plotting the data > plot(data)
Small subset	R is a large, complex language. We cover only a small % in this class

R: Getting Started

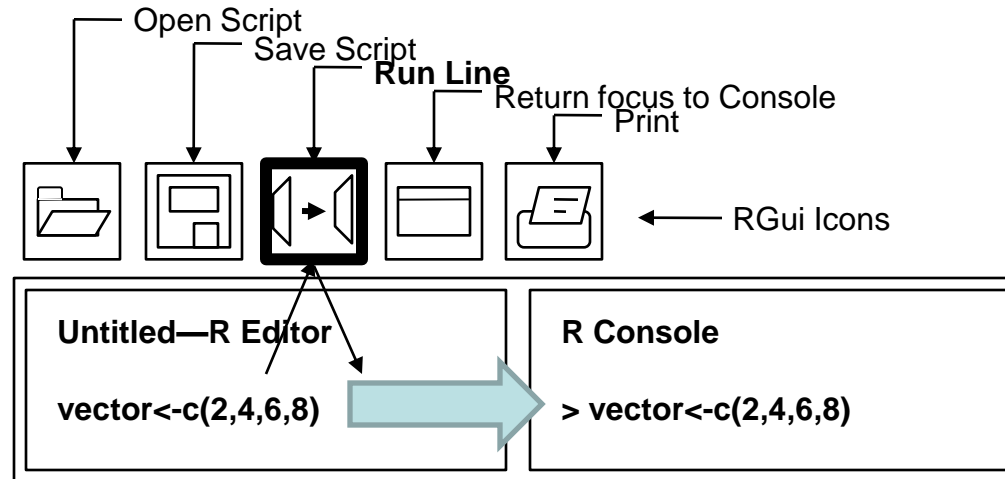
Topic	Description
Download R	Windows: http://cran.r-project.org/bin/windows/base/ Mac: http://cran.r-project.org/bin/macosx/
Launch R	Double-click to launch Will see prompt in “R Console” >

R Console

>

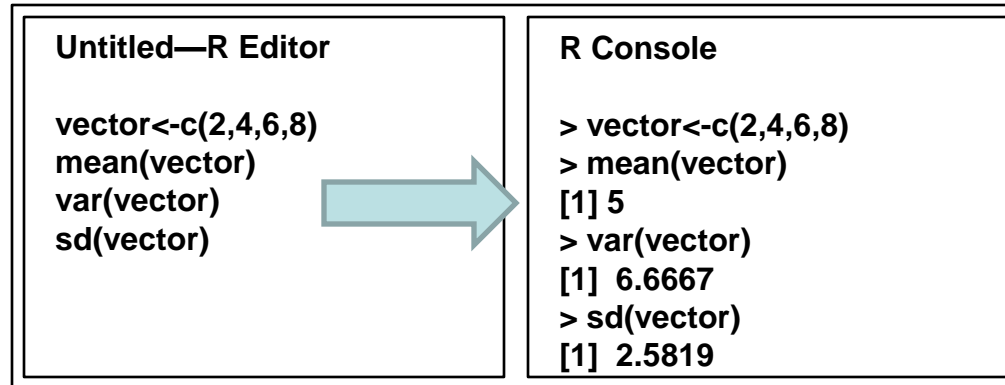
R Editor and R Console: Layout

Topic	Description
R Console	Traditional R interface; Command Line Interface
R Editor	Attempt at easier to use User Interface
New Script	To open Editor, click on Select File > New Script
Run Line	Execute (run) line: Highlight line on R editor; Click on “Run Line”



R Editor: Typical Usage

Topic	Description
Statistics	Find statistics mean(vector) <RUN LINE> (mean) var(vector) <RUN LINE> (variance) sd(vector) <RUN LINE> (standard deviation)
Run Line	Select Run Line icon to move to R Console and execute



Loading Data into R

Topic	Description
1. Preparation	Remove introductory content; First line should be data headers Save Excel file as Comma Separated Values (CSV)
2. Directory	Optional: Set up working directory for dataset; allows shorter filepaths Windows: See “Windows Explorer help” for more info Mac: See “Finder help” for more info
3. Filename	Need complete filename Example: “C:\My Documents\Folder A\Filename.csv” Alternative 1: Right click to see filename Alternative 2: Find filename in Windows Explorer (Windows); Finder (Mac) Alternative 3: Drag csv file and drop into R Console; Will show filename

Loading Data into R

Topic	Description
4. Read CSV data	<code>Datafile <- read.csv("C:\\My Documents\\Desktop\\Filename.csv", header=T)</code>
5. Check data	Print out dataset to ensure it was loaded correctly <code>print(Datafile)</code> : will print out entire datafile; OK for small datasets <code>str(Datafile)</code> : Shows structure of Datafile; "data.frame: 4 obs. of 4 variables" <code>summary(Datafile)</code> : Shows summary: Min; Max; Mean; Median
6. Run regression	<code>lm</code> : Regression analysis in R; stands for Linear Model <code>lm(Dependent~Independent+Independent, Dataset)</code>
7. Interpret Results	Compare results obtained with R with those from Microsoft Excel

Example: Causal Analysis Forecast for Real Estate

Step	Description
------	-------------

1. Preparation	Remove introductory information; First row = header row
----------------	---

“Save As” CSV

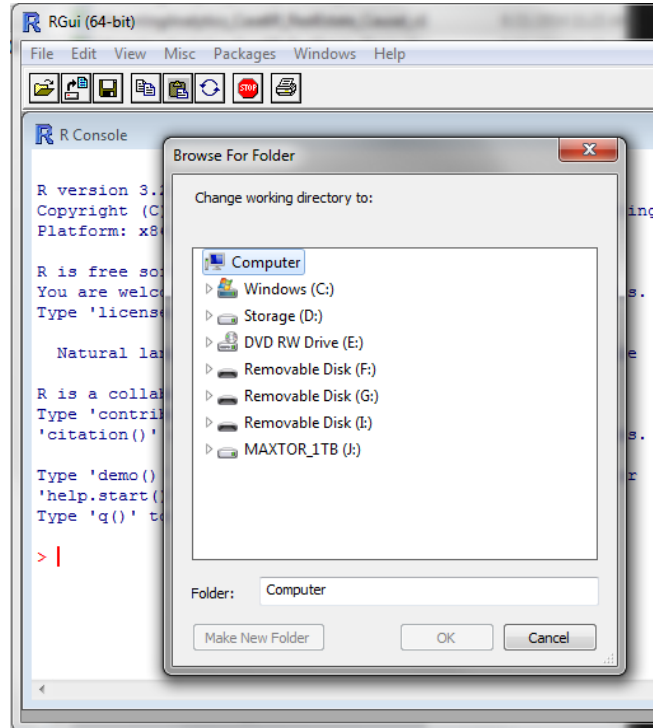
	A	B	C	D	E	F	G	H	I	J	K	L
1	Price	House	Lot									
2		6.0	6.9	42.7								
3		5.8	8.0	36.6								
4		5.6	8.0	44.0								
5		3.5	3.8									
6		3.4	6.1									
7		3.4	4.3									
8		2.7	3.8									
9		2.6	5.0									
10		2.6	3.6									
11		2.3	3.1									
12		2.3	3.9									
13		2.3	3.2									
14		1.9	3.5									
15		1.9	3.4									
16		1.9	3.2									
17		1.6	3.3									
18		1.6	2.3									
19		1.5	2.3									
20												
21												
22												
23												
24												
25												
26												
27												
28												
29												
30												
31												
32												

Example: Causal Analysis Forecast for Real Estate

Step	Description
2. Directory	Optional: Can set up working directory”

In R, select
File → Change dir...

then select where you
want to put R files



Example: Causal Analysis Forecast for Real Estate

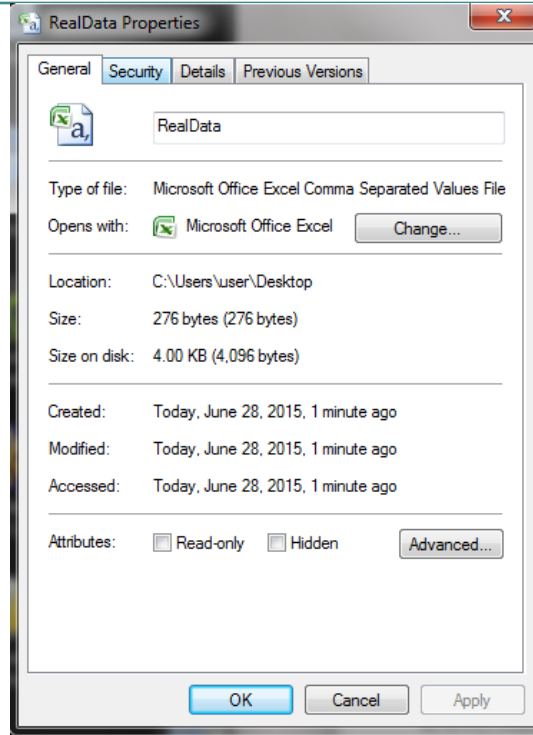
Step	Description
3. Filename	“C:\\Users\\user\\Desktop\\RealData.csv”

Windows:

Right-click on file
to get file properties;
will show full filename
under “Location”

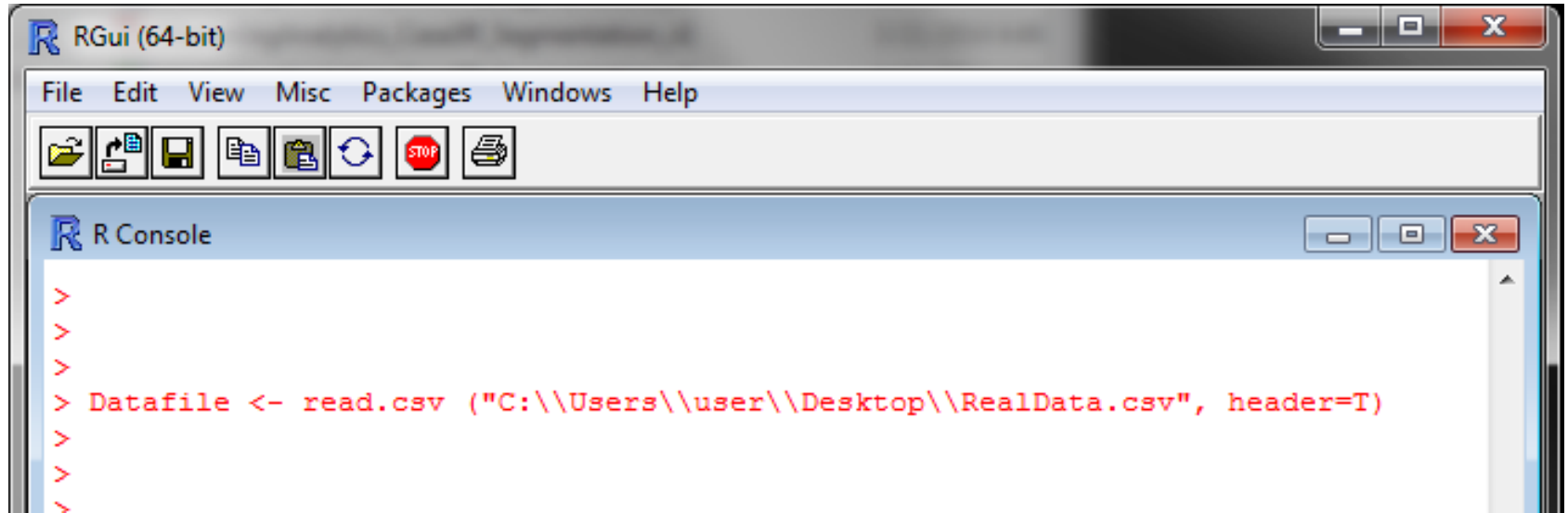
Mac:

Check Finder
to find full filename
OR:
Drag file into R



Example: Causal Analysis Forecast for Real Estate

Step	Description
4. Read Data	<code>Datafile <- read.csv("C:\\Users\\user\\Desktop\\RealData.csv", header=T)</code>



The screenshot shows the RGui (64-bit) window. The menu bar includes File, Edit, View, Misc, Packages, Windows, and Help. Below the menu bar is a toolbar with icons for file operations and execution. The R Console window is open, displaying the following code in red text:

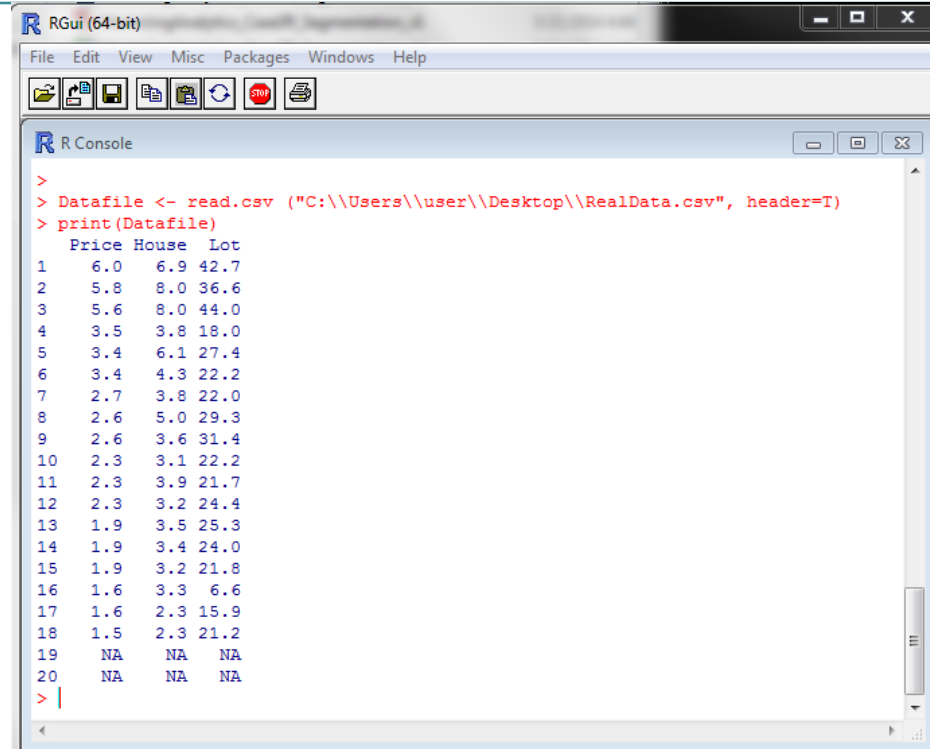
```
>  
>  
>  
> Datafile <- read.csv ("C:\\Users\\user\\Desktop\\RealData.csv", header=T)  
>  
>  
>
```

Alternative: Set up working directory

Example: Causal Analysis Forecast for Real Estate

Step	Description
5. Check Data	print (Datafile) ; check if dataset looks OK

For large datasets,
ask R to provide
summary data
instead of printing
out entire dataset

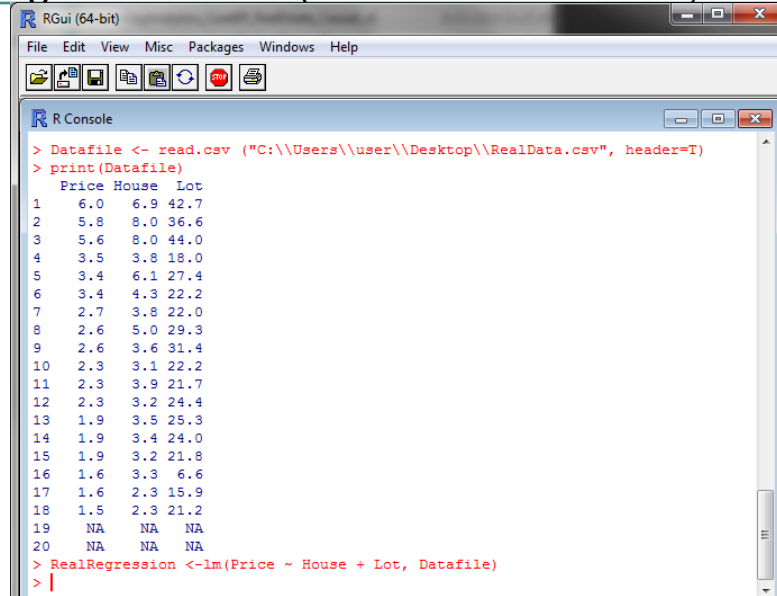


```
RGui (64-bit)
File Edit View Misc Packages Windows Help
[Icons]
R Console
>
> Datafile <- read.csv ("C:\\Users\\user\\Desktop\\RealData.csv", header=T)
> print(Datafile)
  Price House Lot
1    6.0   6.9 42.7
2    5.8   8.0 36.6
3    5.6   8.0 44.0
4    3.5   3.8 18.0
5    3.4   6.1 27.4
6    3.4   4.3 22.2
7    2.7   3.8 22.0
8    2.6   5.0 29.3
9    2.6   3.6 31.4
10   2.3   3.1 22.2
11   2.3   3.9 21.7
12   2.3   3.2 24.4
13   1.9   3.5 25.3
14   1.9   3.4 24.0
15   1.9   3.2 21.8
16   1.6   3.3  6.6
17   1.6   2.3 15.9
18   1.5   2.3 21.2
19   NA    NA  NA
20   NA    NA  NA
> |
```

Example: Causal Analysis Forecast for Real Estate

Step	Description
6. Run Regression	$\text{lm}(\text{Dependent} \sim \text{Independent} + \text{Independent}, \text{Dataset})$ Dependent variable: Price; Independent variable: House; Lot Equation: $\text{Price} = c_1 + c_2 * (\text{House Size}) + c_3 * (\text{Lot Size})$ <code>RealRegression <- lm(Price ~ House + Lot, Datafile)</code>

Find tilde symbol “ ~ ”
at upper left of keyboard,
to left of number “1”

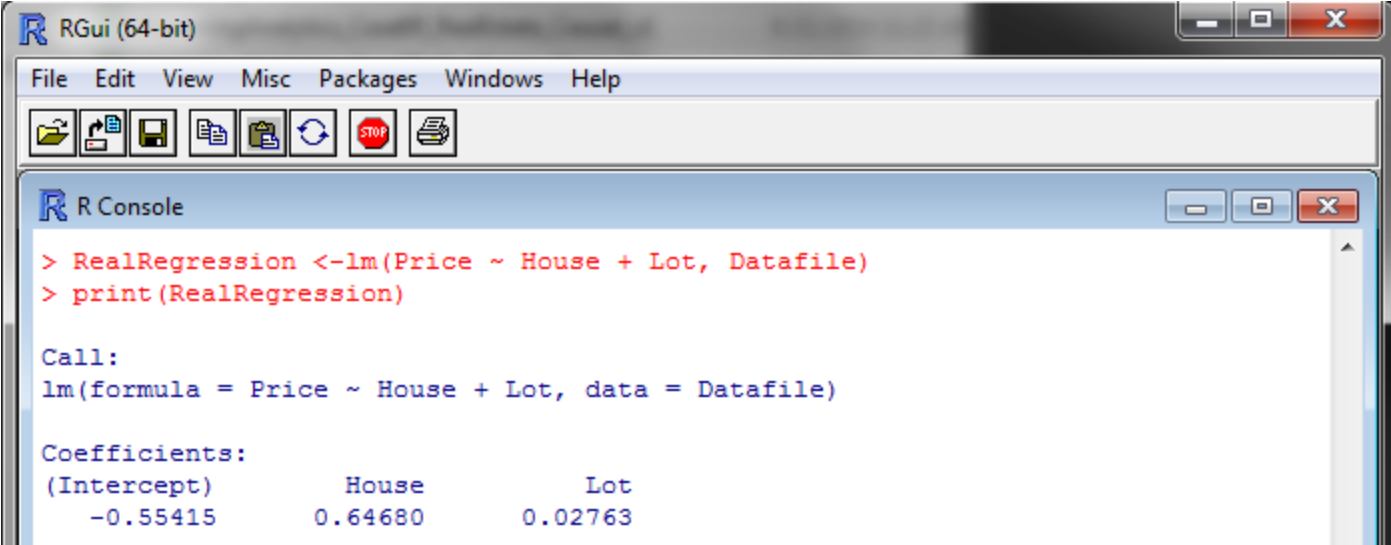


```
RGui (64-bit)
File Edit View Misc Packages Windows Help
R Console
> Datafile <- read.csv("C:\\Users\\user\\Desktop\\RealData.csv", header=T)
> print(Datafile)
  Price House Lot
1   6.0   6.9 42.7
2   5.8   8.0 36.6
3   5.6   8.0 44.0
4   3.5   3.8 18.0
5   3.4   6.1 27.4
6   3.4   4.3 22.2
7   2.7   3.8 22.0
8   2.6   5.0 29.3
9   2.6   3.6 31.4
10  2.3   3.1 22.2
11  2.3   3.9 21.7
12  2.3   3.2 24.4
13  1.9   3.5 25.3
14  1.9   3.4 24.0
15  1.9   3.2 21.8
16  1.6   3.3  6.6
17  1.6   2.3 15.9
18  1.5   2.3 21.2
19  NA    NA  NA
20  NA    NA  NA
> RealRegression <- lm(Price ~ House + Lot, Datafile)
> |
```

Example: Causal Analysis Forecast for Real Estate

Topic	Description												
7. Interpret Results	Compare results from R with those from Excel												
	<table border="1"><thead><tr><th>Method</th><th>Coefficient</th><th>House Size</th><th>Lot Size</th></tr></thead><tbody><tr><td>Excel</td><td>-0.554</td><td>+0.646</td><td>+0.027</td></tr><tr><td>R</td><td>-0.55415</td><td>+0.64680</td><td>+0.02763</td></tr></tbody></table>	Method	Coefficient	House Size	Lot Size	Excel	-0.554	+0.646	+0.027	R	-0.55415	+0.64680	+0.02763
Method	Coefficient	House Size	Lot Size										
Excel	-0.554	+0.646	+0.027										
R	-0.55415	+0.64680	+0.02763										

R results agree well with those of Excel



```
RGui (64-bit)
File Edit View Misc Packages Windows Help
R Console
> RealRegression <-lm(Price ~ House + Lot, Datafile)
> print(RealRegression)

Call:
lm(formula = Price ~ House + Lot, data = Datafile)

Coefficients:
(Intercept)      House         Lot
   -0.55415      0.64680      0.02763
```