[View All](#) [View Keyframes](#)

## Overview of R

- Over the past decade, companies have been using a variety of statistical packages, including SPSS, SAS, Minitab, Stata, and others.
- Recently, companies have begun to migrate to a powerful statistical and data mining package called **R**.
- R is an open source software.
  - Free
  - Cost of other packages has driven new users to R.
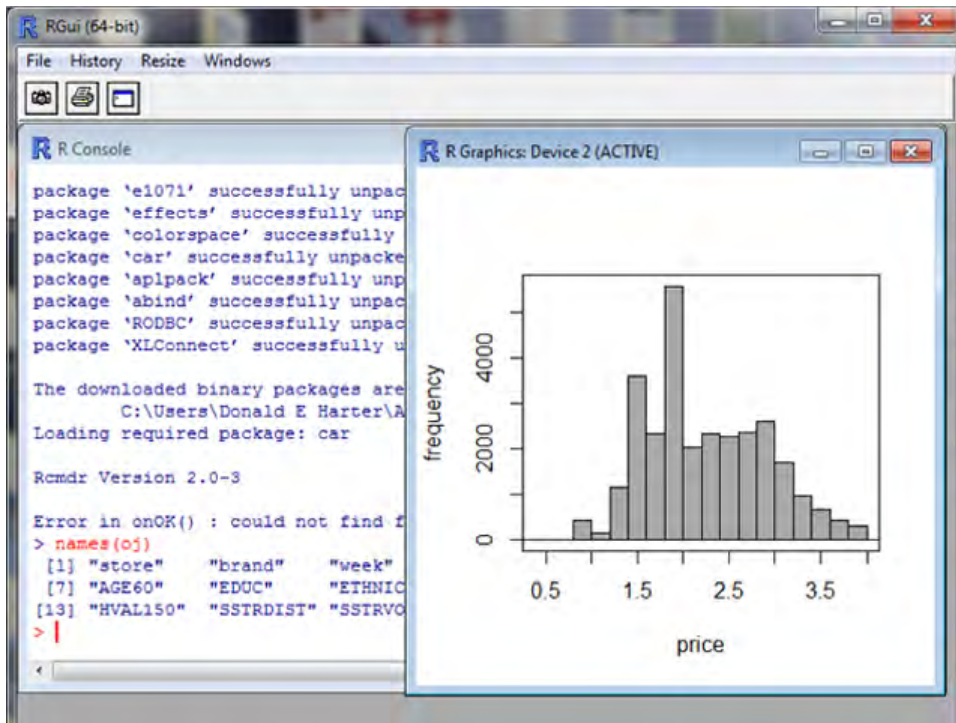- R is a language, but we will mostly operate it using a graphical user interface.

7 of 7
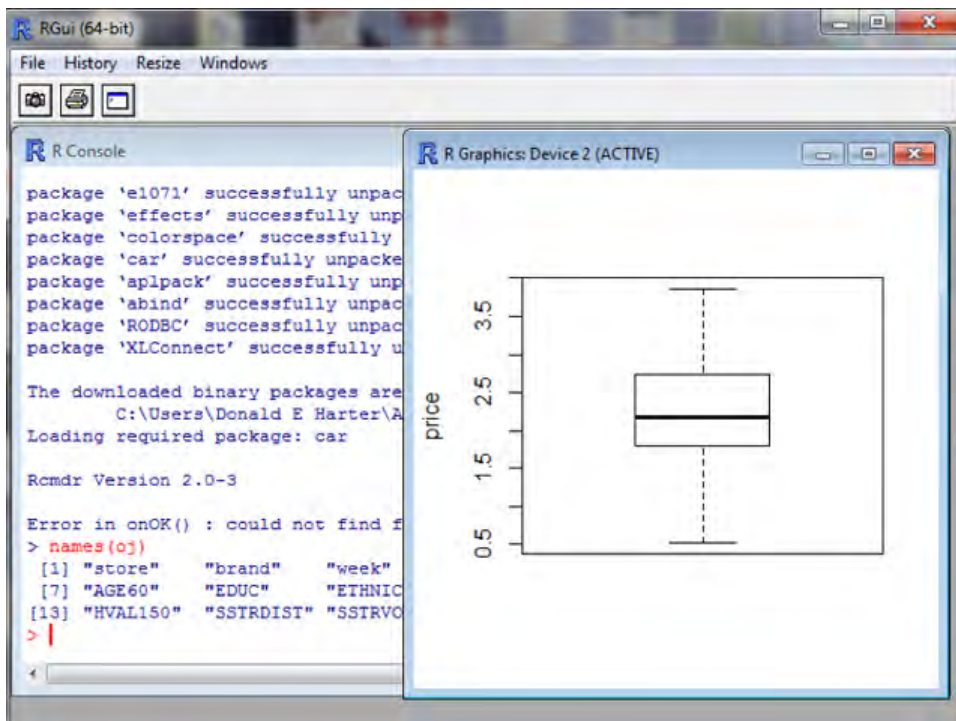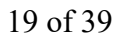
# Visualization

- R has the ability to generate histograms, boxplots, scatterplots, and XY plots similar to Excel.
- Graphs from R can differentiate by brand, product category, or other variable.
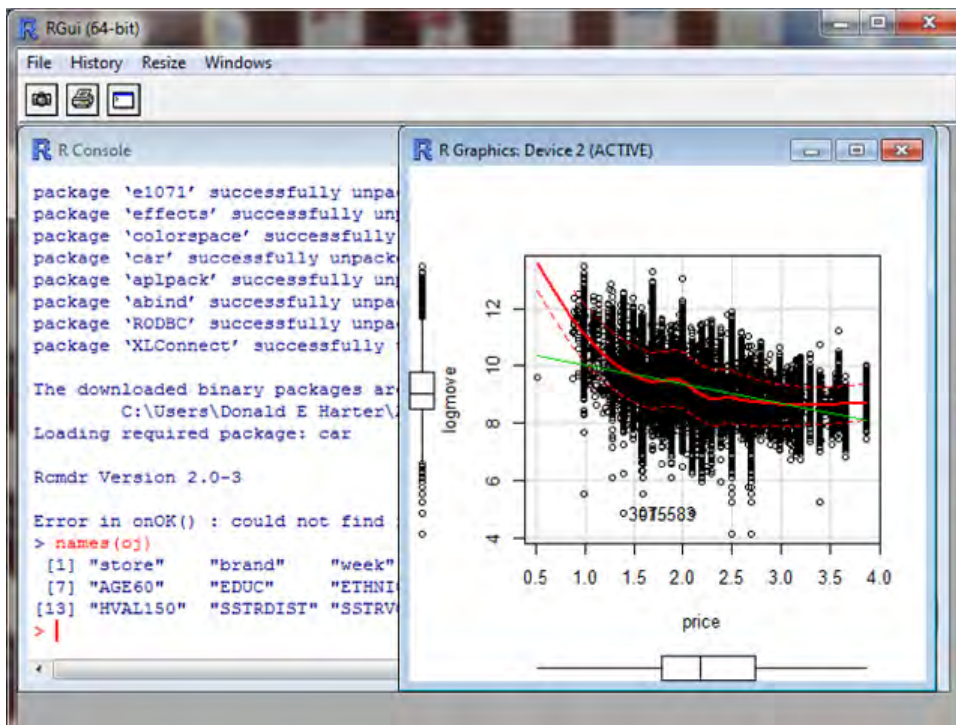  - Gives greater insight into patterns

4 of 39



5 of 39

## Histograms and Boxplots

- **Histograms** show frequency of data within intervals or bins
  - Easy to develop in Excel or R.
- **Boxplots** identify:
  - Maximum and minimum—whiskers
  - 25%-ile and 75%-ile—box
  - Median—center bar
- Boxplots in R can also differentiate by brands.
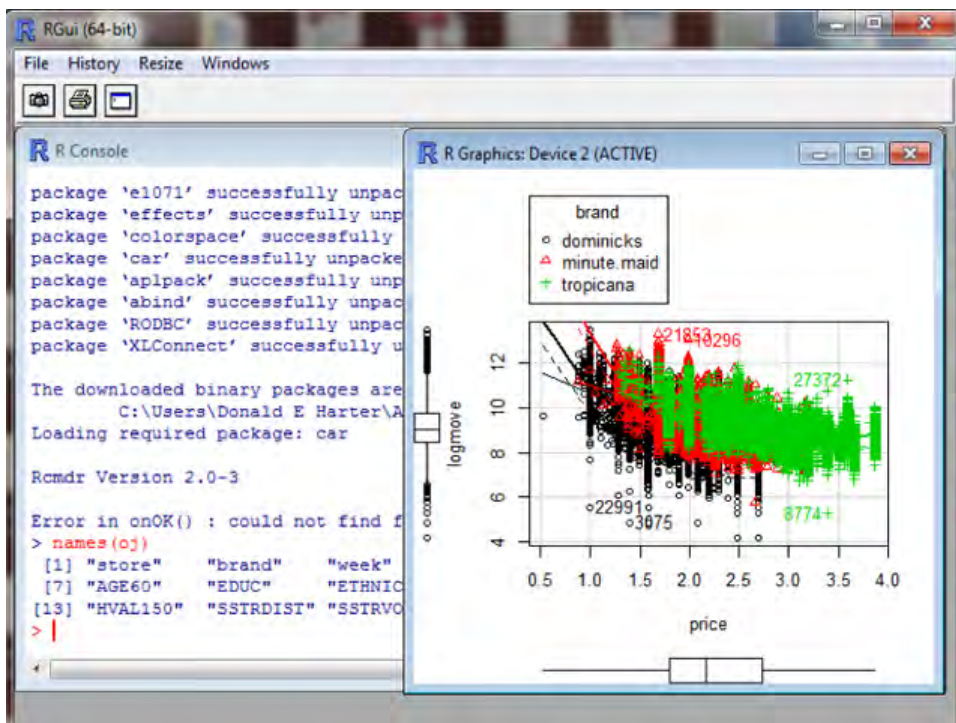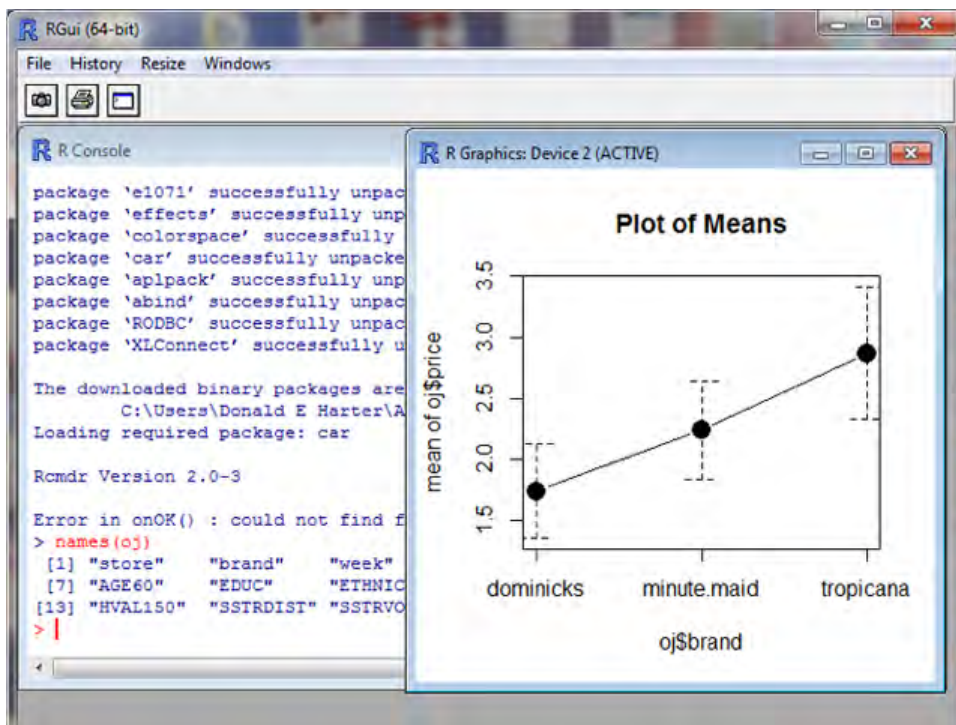
24 of 39

## Scatterplots

- **Scatterplots** offer not only trend line, but price by price.
- From example:
    - Solid red line—average sales by price
    - Green line—trend line
    - Dotted red line—standard deviation
- Scatterplots of log sales by price can also be categorized by brand.
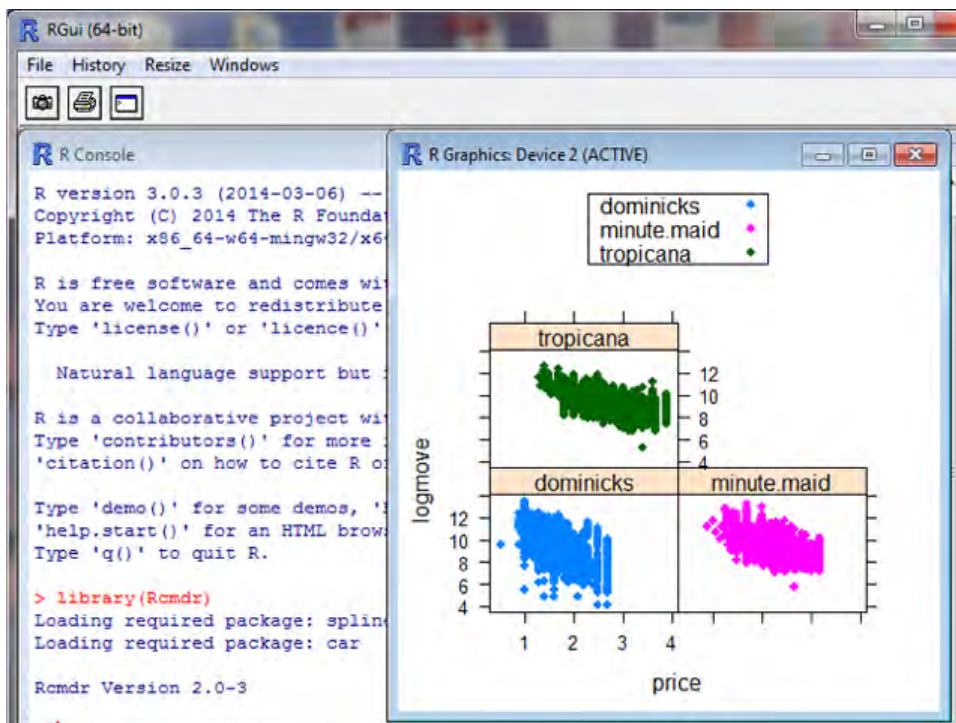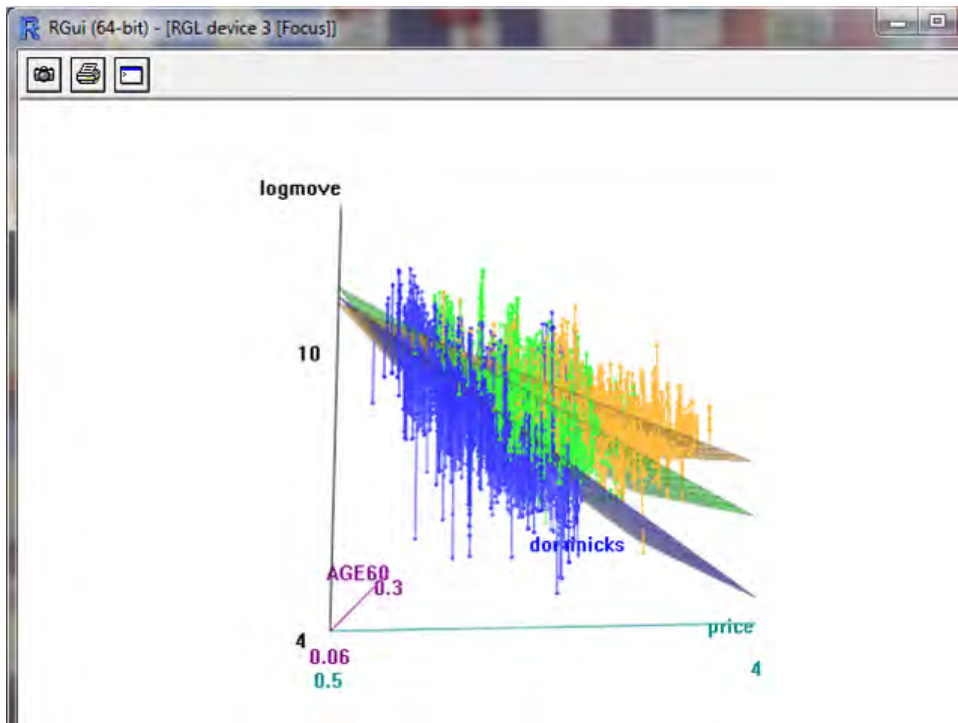
30 of 39

## Plot of Means and XY Plots

- **Plot of means** is where the average price by brand is plotted.
- **XY plots** display the effect of price on sales per brand by separating and placing plots next to each other.

5 of 13

## 3-D Visualization

- R can perform more sophisticated 3-D graphing than Excel.
    - Plot the planes—helps find patterns
    - Rotate the graphs—offers different senses of parameters.
- 3-D allows you to see regression lines and planes to identify patterns.

13 of 13

<u>View All</u> <u>View Keyframes</u>



4 of 13

## Statistical Summaries

- Similar to Excel, R can calculate descriptive statistics.
- R has the additional ability to calculate descriptive statistics by brand.
- Stratify and characterize data while performing analysis.
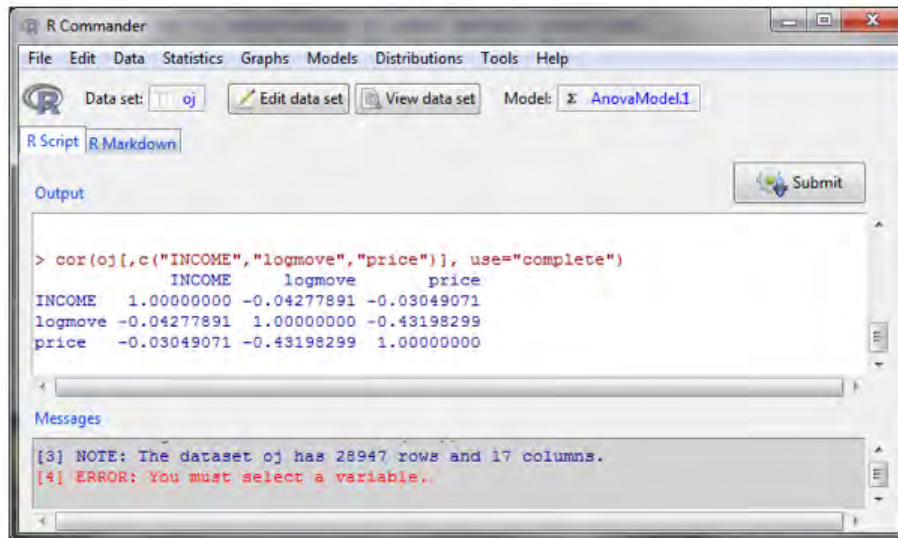
13 of 13

## Correlations

- Measure how two variables are related
- Positively correlated: one variable increases, the other one also increases
- Negatively correlated: one variable increases, the other one decreases

4 of 11



5 of 11

8 of 11

## Analysis of Variance

- Store managers want to compare product placement positions, sales, and prices of different products or brands.
- Analysis of variance (ANOVA) is used to compare the averages for different categories or brands.
- ANOVA calculates the average of each item, combines with standard deviation/variance, and determines if they are statistically different from one another.
- ANOVA is a valuable technique for product placement.

11 of 11

```
R Commander                                                    [_][□][x]
File  Edit  Data  Statistics  Graphs  Models  Distributions  Tools  Help
(R)  Data set: [TT oj]  [/ Edit data set] [Q View data set]  Model: Σ RegModel.5
[R Script] [R Markdown]
plot(confint(.Pairs))
par(old.oma)
remove(.Pairs)
.cluster <-  KMeans(model.matrix(~-1 + AGE60 + price, oj), centers = 2,
    iter.max = 10, num.seeds = 10)

Output                                                      [≈ Submit]

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept) 11.987186   0.208142  57.591  < 2e-16 ***
AGE60        1.709162   0.087691  19.491  < 2e-16 ***
INCOME      -0.145482   0.019211  -7.573 3.76e-14 ***
price       -0.688144   0.008279 -83.123  < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9117 on 28943 degrees of freedom
Multiple R-squared:  0.2002, Adjusted R-squared:  0.2002
F-statistic:  2416 on 3 and 28943 DF,  p-value: < 2.2e-16

Messages
[3] NOTE: The dataset oj has 28947 rows and 17 columns.
[4] ERROR: You must select a variable.
```
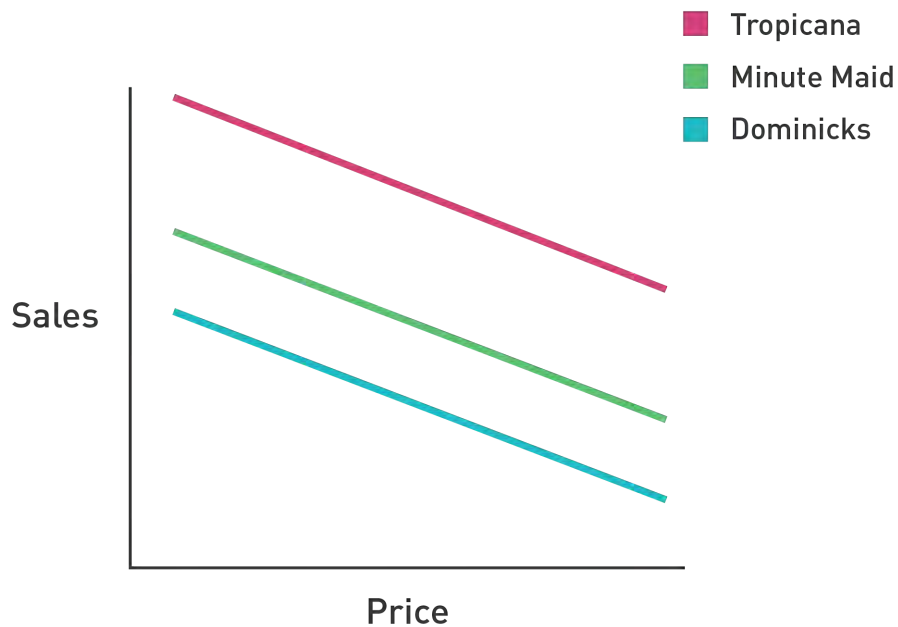
4 of 9

## Regression

- R can perform linear regression.
  - Output is similar to Excel's multivariate regression output
- Analysis produces coefficients, statistical significance, and R-squared values.
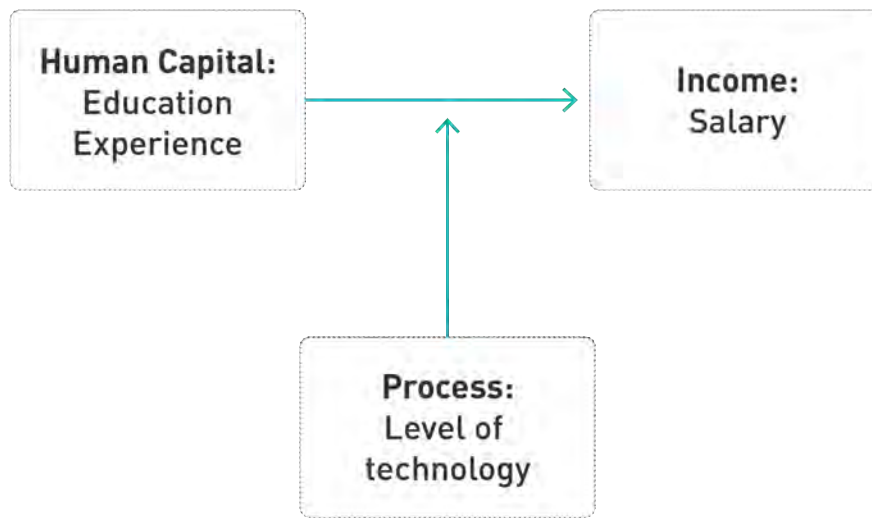  - Sets foundation to view more complex regressions

9 of 9

5 of 10

## Regression with Dummy Variables

- Dummy variables change the intercept of some items to be different from others.
  - E.g., one brand will always be higher in sales than another
  - Reflected in different intercepts or heights of the regression line
- This allows you to look at price differentiation in different product markets.
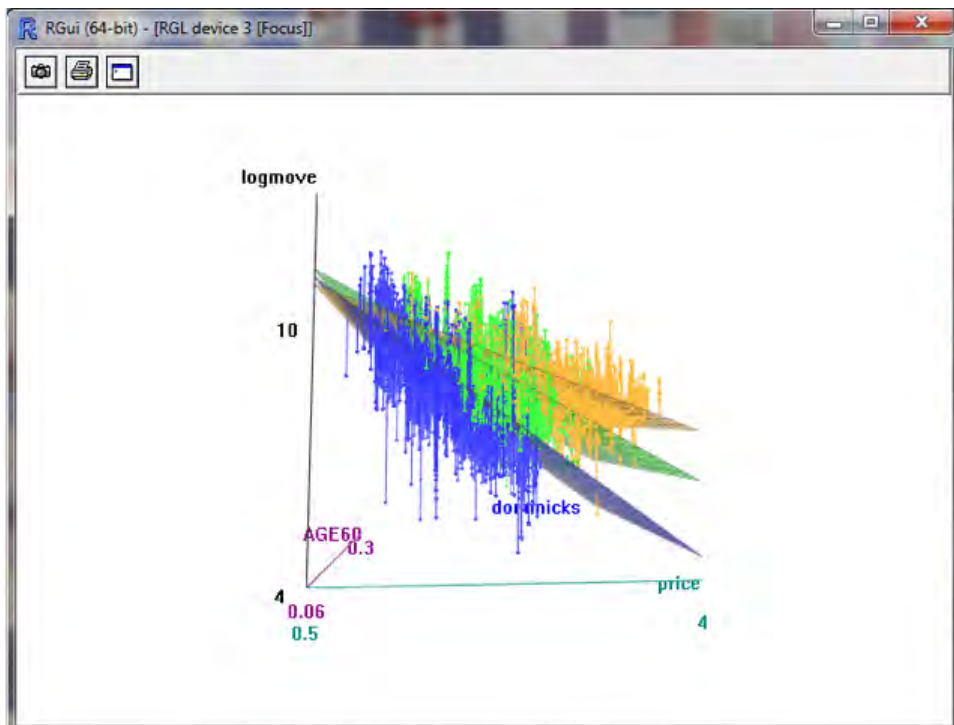
10 of 10

10 of 19

## Regression with Moderating Effects

- Dummy variables can only change the intercepts.
- Moderating effects allow the slope of the line to change.
  - Recognizes that two variables might interact, magnifying their effects
- Education and experience both have a positive effect on income.
  - When an experienced, educated professional learns technology, the technology is the moderating effect.
- **Moderating effect:** An interaction of variables that leverage each other
- The moderating effect acts as a catalyst to accelerate the effect of certain variables.

18 of 19

19 of 19