



# Business Analytics and Data Science

## INTRODUCTION

# Know your Instructor



- Author "[R for Business Analytics](#)"
- Author "[R for Cloud Computing](#)"
- Founder "[Decisionstats.com](#)"
- University of Tennessee, Knoxville MS (courses in statistics and computer science)
- MBA (IIM Lucknow, India - 2003)
- B.Engineering (DCE 2001)

<http://linkedin.com/in/ajayohri>

# Classroom Rules

- From Instructor
- From Audience
  - mobile phones should be kindly switched off
    - Yes, this includes Whatsapp
  - Ask Questions at end of session
  - Take Notes
  - Please Take Notes

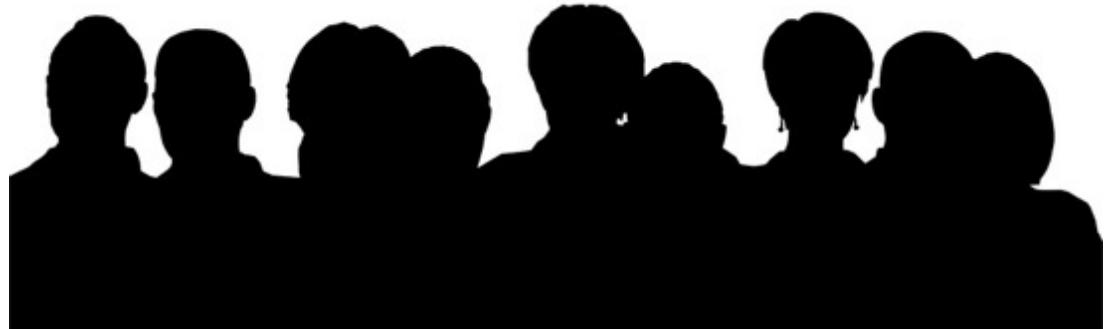


# Introduce Yourself

Name

Education Degree from Institute

Work Ex in Years in Domain



# Introduce Yourself

Name

Education Degree from Institute

Work Ex in Years in Domain

What expectations from this training



# Expectations

How Data Science can help your career ?



# Support Team

Madhuresh

# Introduction to Data Science

## Course Outline

### Class 1 : Introduction

- Understanding Business Analytics
- Introduction to RFM Analysis and LTV Analysis
- Refresher in Statistics
- Basics of Data Driven Decision Making
- Installation of R, Rtools, R Studio, R packages and GUIs
- Using RStudio and Using GUIs

### Quiz 1

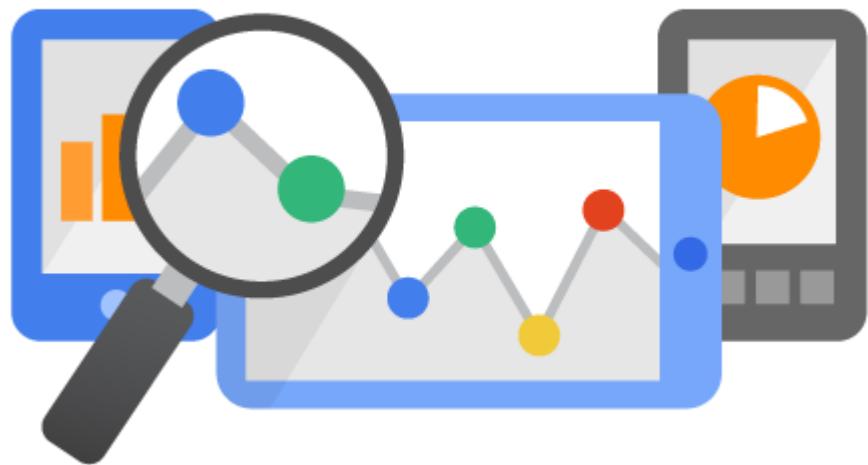
### Class 2 : Data Input

- Data Input using standard data sources
- Data Input for large data (RDBMS)
- Using SQL from within R
- Web Scraping
- Specialized Packages for other data sources

### Quiz 2

### Dataset - Iris

# Introduction to Data Science



# Information Ladder

The **information ladder** was created by education professor Norman Longworth to describe the stages in human learning.

According to the ladder, a learner moves through the following progression to construct “wisdom” from “data”

Data →

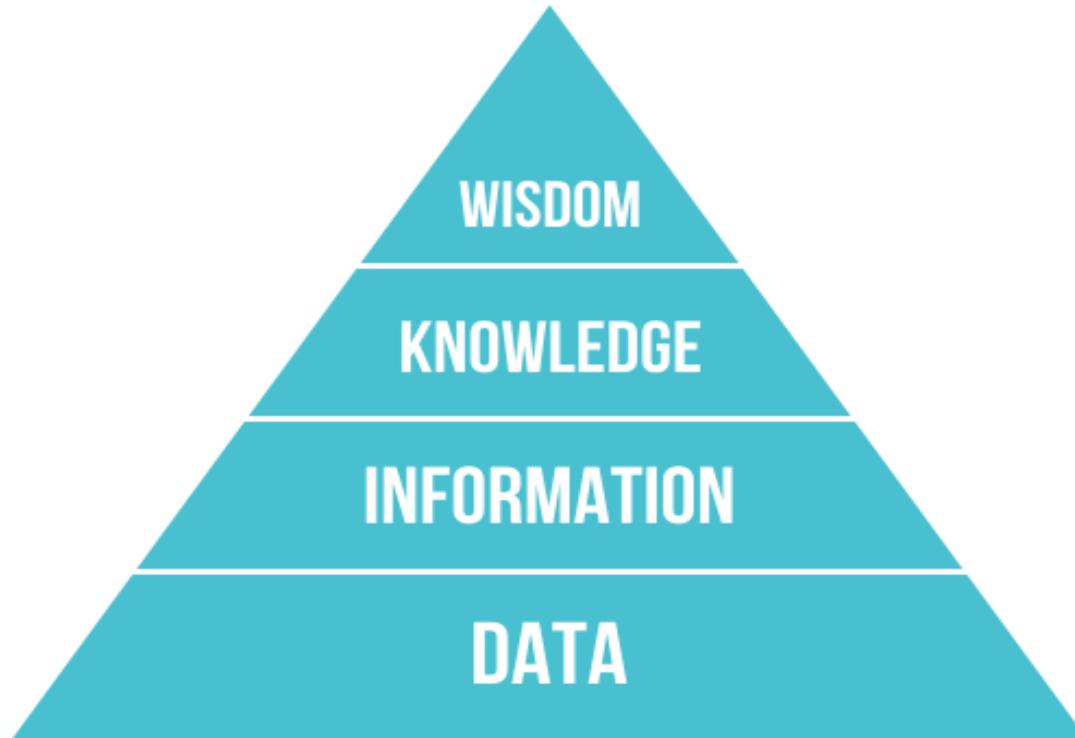
Information →

Knowledge →

Understanding →

Insight →

Wisdom



# Basics of Data Science

**Data Science** is the extraction of knowledge from data,<sup>[1][2]</sup> which is a continuation of the field data mining and predictive analytics, also known as knowledge discovery and data mining (KDD). It employs techniques and theories drawn from many fields within the broad areas of mathematics, statistics, information theory and information technology, including signal processing, probability models, machine learning, statistical learning, data mining, database, data engineering, pattern recognition and learning, visualization, predictive analytics, uncertainty modeling, data warehousing, data compression, computer programming, and high performance computing. Methods that scale to Big Data are of particular interest in data science, although the discipline is not generally considered to be restricted to such data. The development of machine learning, a branch of artificial intelligence used to uncover patterns in data from which predictive models can be developed, has enhanced the growth and importance of data science.

CONFUSING?

[http://en.wikipedia.org/wiki/Data\\_science](http://en.wikipedia.org/wiki/Data_science)

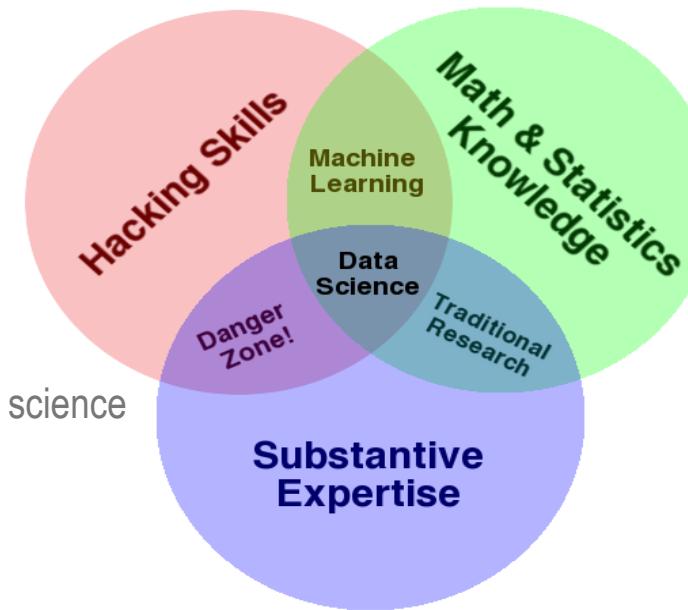
# Basics of Data Science

the culture of academia, which does not reward researchers for understanding technology.

DANGER ZONE- this overlap of skills gives people the ability to create what appears to be a legitimate analysis without any understanding of how they got there or what they have created

Being able to manipulate text files at the command-line, understanding vectorized operations, thinking algorithmically; these are the hacking skills that make for a successful data hacker.

data plus math and statistics only gets you machine learning, which is great if that is what you are interested in, but not if you are doing data science



<http://drewconway.com/zia/2013/3/26/the-data-science-venn-diagram>

# Business Intelligence

**Business intelligence** (BI) is an umbrella term that includes the applications, infrastructure and tools, and best practices that enable access to and analysis of information to improve and optimize decisions and performance.

The key general categories of business intelligence tools are:

- [Spreadsheets](#)
- [Reporting and querying software](#): tools that extract, sort, summarize, and present selected data
- [OLAP](#): Online analytical processing
- [Digital dashboards](#)
- [Data mining](#)
- [Data warehousing](#)
- [Local information systems](#)



Definition – study of business data using statistical techniques and programming for creating decision support and insights for achieving business goals

Predictive- To predict the future.

Descriptive- To describe the past.

# So what is a Data Scientist ?

a data scientist is simply a data analyst living in **california**

# What is a Data Scientist

a data scientist is simply a person who can

**write code**

**understand statistics**

**derive insights from data**

# Oh really, is this a Data Scientist ?

A data scientist is simply a person who can

**write code = in** R,Python,Java, SQL, Hadoop (Pig,HQL,MR) etc

**= for** data storage, querying, summarization, visualization

**= how** efficiently, and in time (fast results?)

**= where** on databases, on cloud, servers

**and** understand **enough** statistics

**to** derive **insights** from data

**so** **business** can make **decisions**

# Guide for Data Scientists

<http://www.kdnuggets.com/2014/05/guide-to-data-science-cheat-sheets.html>

By Ajay Ohri, May 2014.

Over the past few years, as the buzz and apparently the demand for data scientists has continued to grow, people are eager to learn how to join, learn, advance and thrive in this seemingly lucrative profession. As someone who writes on analytics and occasionally teaches it, I am often asked - How do I become a data scientist?

Adding to the complexity of my answer is data science seems to be a multi-disciplinary field, while the university departments of statistics, computer science and management deal with data quite differently.

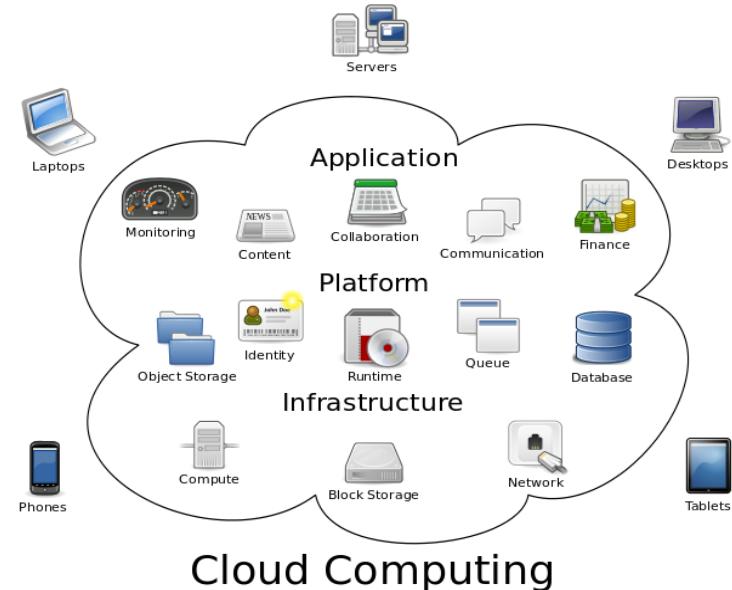
But to cut the marketing created jargon aside, a data scientist is simply a person who can write code in a few languages (primarily R, Python and SQL) for data querying, manipulation , aggregation, and visualization using enough statistical knowledge to give back actionable insights to the business for making decisions.

<http://www.slideshare.net/ajayohri/cheat-sheets-for-data-scientists>

- Business Analytics
  - Understanding what solution business needs
- Data Science
  - Primarily R programming skills
  - Some Applied Statistical Methods
  - Exposure to new domains and techniques

# Cloud Computing

1. the practice of using a network of remote servers hosted on the Internet to store, manage, and process data, rather than a local server or a personal computer.



<http://csrc.nist.gov/publications/nistpubs/800-145/SP800-145.pdf>

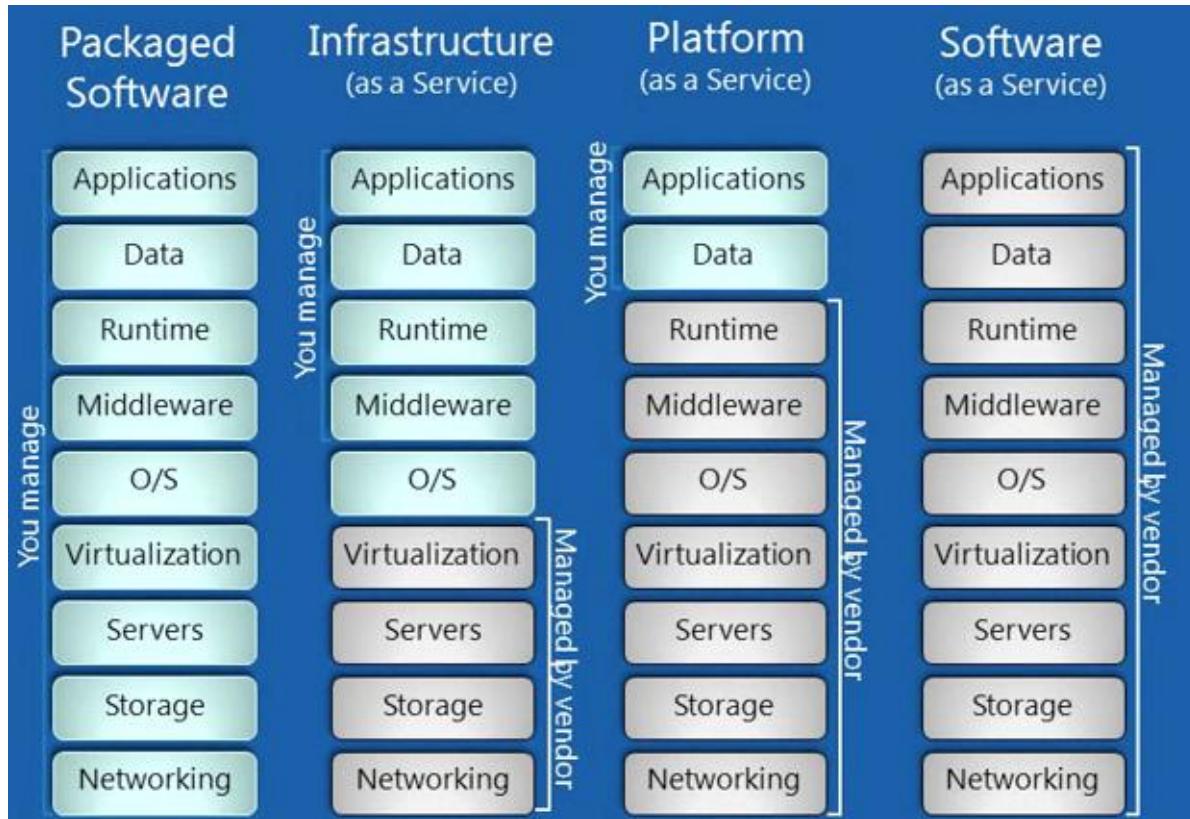
# Cloud Computing

1. the practice of using a network of remote servers hosted on the Internet to store, manage, and process data, rather than a local server or a personal computer.



<http://csrc.nist.gov/publications/nistpubs/800-145/SP800-145.pdf>

# Cloud Computing



# LTV Analytics

Life Time Value (LTV) will help us answer 3 fundamental questions:

1. Did you pay enough to acquire customers from each marketing channel?
2. Did you acquire the best kind of customers?
3. How much could you spend on keeping them sweet with email and social media?

## Breaking Down LTV Further

### LTV WILL BE DIFFERENT FOR DIFFERENT KINDS OF CUSTOMERS

Step 2 in this graphic is intended to help you determine LTV as a total average (an average of all your customers). To do this, companies will typically average the data from randomly chosen customers (as shown in Step 1 above). Sometimes it's helpful to break down the average further and perform separate LTV calculations for different kinds of customers. Try and segment your customer base by total purchases over a long time period, and it will help you determine the LTV of a "good" customer versus an "average" one. This type of analysis will help you determine how much more you should pay in order to acquire a "good" customer. See chart below.

#### INVESTING IN "GOOD" CUSTOMERS

Companies should be worried about the lasting impact of "buying cheap customers." How likely are these customers to buy another product, or hang around for a few years? Sometimes it pays to invest in "good" customers. "Good" customers might cost more to acquire, but they'll likely be more profitable as well.

Let's say that the LTV of an "average" customer is \$8,000, and the LTV of a "good" customer is \$10,000. By subtracting the two LTVs, you can see that you might expect to pay \$2,000 more to acquire "good" customers.



## Customer Satisfaction Boosts LTV

One of the most effective ways to boost LTV is to increase customer satisfaction. Research has found that a 5% increase in customer retention can increase profits by 25% to 95%. The same study found that it costs six to seven times more to gain a new customer than to keep an existing one.

# LTV Analytics

**Questions.** Fill in the yellow boxes and the spreadsheet will take care of the rest.

	Best Customers	Average Customers
<b>Acquisition Cost</b> . How much did you pay to acquire these customers?	£40.00	£12.00
<b>Average order value</b> . How much do they spend per order?	£92.00	£70.00
<b>Orders per year?</b> Quite simply, How many orders do they place per year?	5	2
<b>Retention?</b> How many years will they be customers for?	3	2
<b>Net profit?</b> What is the net profit percentage of goods sold?	10%	10%

**Answers.** These cells will be magically calculated based on the values you put in the table on the left

	Best Customers	Average Customers
Lifetime Gross Revenue	£1,380.00	£280.00
Life Time Net Profit	£98.00	£16.00

<http://www.kaushik.net/avinash/analytics-tip-calculate-ltv-customer-lifetime-value/>

# LTV Analytics

Questions. Fill in the yellow boxes and the spreadsheet will do the rest.		Year 1	Year 2	Year 3	Year 4	Year 5
Screen	Full Screen					
<b>Customer Segment.</b> How many of a specific group of customers will you start with?						
<b>Acquisition Cost?</b> How much did you pay for each new customer? We won't use this figure - see note to explain						
<b>Retention Rate.</b> What % of customers will you keep from one year to the next?					75%	80%
<b>Total Orders.</b> How many orders/sales per customer per year? They may place more in future years		3	3	4	4	5
<b>Average order value.</b> How much is each sale or order worth, and will this rise over time?		£60.00	£65.00	£70.00	£75.00	£80.00
<b>Net Profit.</b> What % of each order is left after all costs have been accounted for?		10%	12%	12%	15%	15%
<b>Discount Rate.</b> This recognises our money <b>could</b> be better spent on something else - see note to explain		Some companies include this in their LTV calculations, especially where the investment is high over a long time period. Just set all the fields to 1 if you'd prefer to ignore it!			0.729	0.656
Answers. These cells will be magically calculated based on the values you put in the table above.		Year 1	Year 2	Year 3	Year 4	Year 5
<b>Total Customers.</b> The number of customers at the start of each year from the original segment		3,000	1,800	1,170	819	614
<b>Total Revenue per Customer.</b> This is the total revenue per year for individual customers		£180	£195	£280	£300	£400
<b>Total Revenue.</b> Annual revenue generated by all the customers in that year		£5,40,000	£3,51,000	£3,27,600	£2,45,700	£2,45,700
<b>Cumulative Revenue.</b> The revenue generated from the (remaining) original customers every year		£5,40,000	£8,91,000	£12,18,600	£14,64,300	£17,10,000
<b>Annual Net profit per customer.</b> Simply, the profit each customer generates in that year.		£18.00	£23.40	£33.60	£45.00	£60.00
<b>Total Net Profit.</b> Profit generated by all the original customers in that year.		£54,000	£42,120	£39,312	£36,855	£36,855
<b>Profit at Net Present Value.</b> The profit made each year, even if we offset a better way of spending it!		£54,000	£37,908	£31,843	£26,867	£24,177
<b>Cumulative Net Profit at NPV.</b> The profit generated in successive years from the original customers.		£54,000	£96,120	£1,35,432	£1,72,287	£2,09,142
<b>Individual LTV at NPV.</b> The cumulative amount of net profit each original customer is worth each year.		£18.00	£32.04	£45.14	£57.43	£69.71

# LTV Analytics

Download the zip file from [http://www.kaushik.net/avinash/avinash\\_Ltv.zip](http://www.kaushik.net/avinash/avinash_Ltv.zip)

Do the class exercise based on numbers given by instructor

Give a brief supporting statement on analysis

# LTV Analytics :Another Approach

## Step 1: Average Your Variables

CUSTOMER EXPENDITURES PER VISIT



NUMBER OF VISITS PER WEEK (THE "PURCHASE CYCLE")



AVG. CUSTOMER VALUE PER WEEK (EXPENDITURES × VISITS, IN USD)



<https://blog.kissmetrics.com/how-to-calculate-lifetime-value/>

# LTV Analytics

## Step 2: Calculate Lifetime Value (LTV)

### CONSTANTS

**t** The Average Customer Lifespan (how long someone remains a customer). In the case of Starbucks, the average customer lifespan is 20 years.

**r** Customer Retention Rate. The percentage of customers, who, over a given period of time, repurchase, when compared to an equal and preceding period of time. Starbucks: 75%.

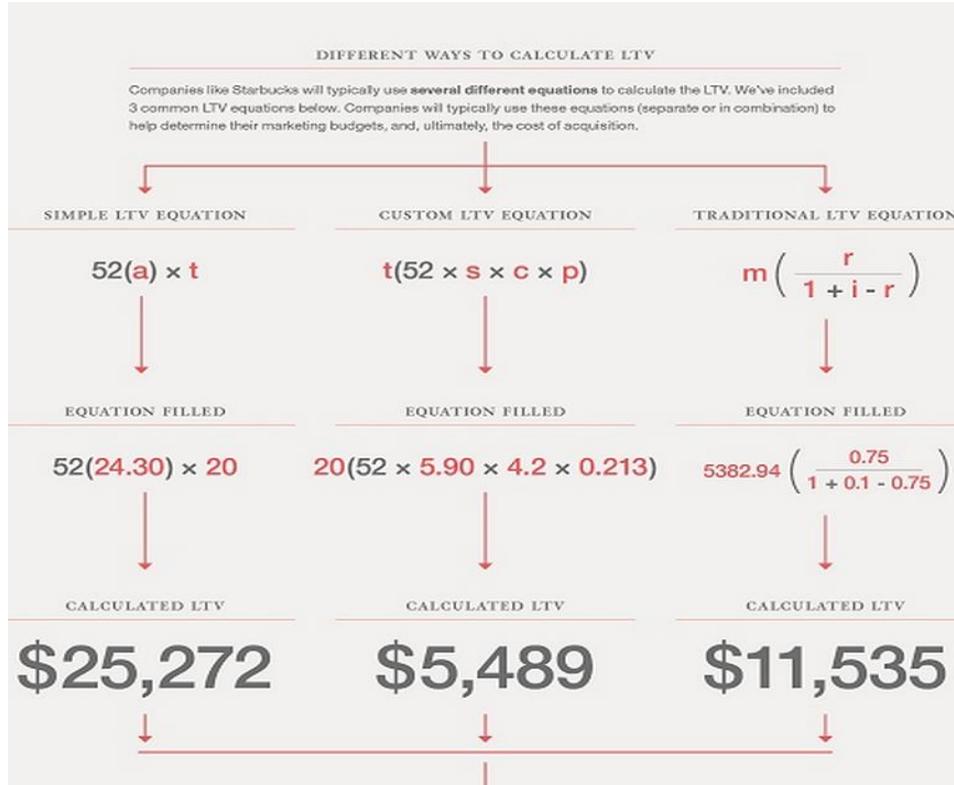
**p** Profit Margin per Customer. Starbucks: 21.3%.

**i** The Rate of Discount. The "rate of discount" is the interest rate used in discounted cash flow analysis to determine the present value of future cash flows. Usually this number falls between 8% and 15%. Starbucks: 10%.

**m** Avg. Gross Margin per Customer Lifespan. Starbucks has a profit margin of 21.3% (see constant "p"). If the average customer spends \$25,272 (see the "Simple LTV Equation" results below) during their time as a customer ("t"), Starbucks has gross margin per customer lifespan of \$5302.94.

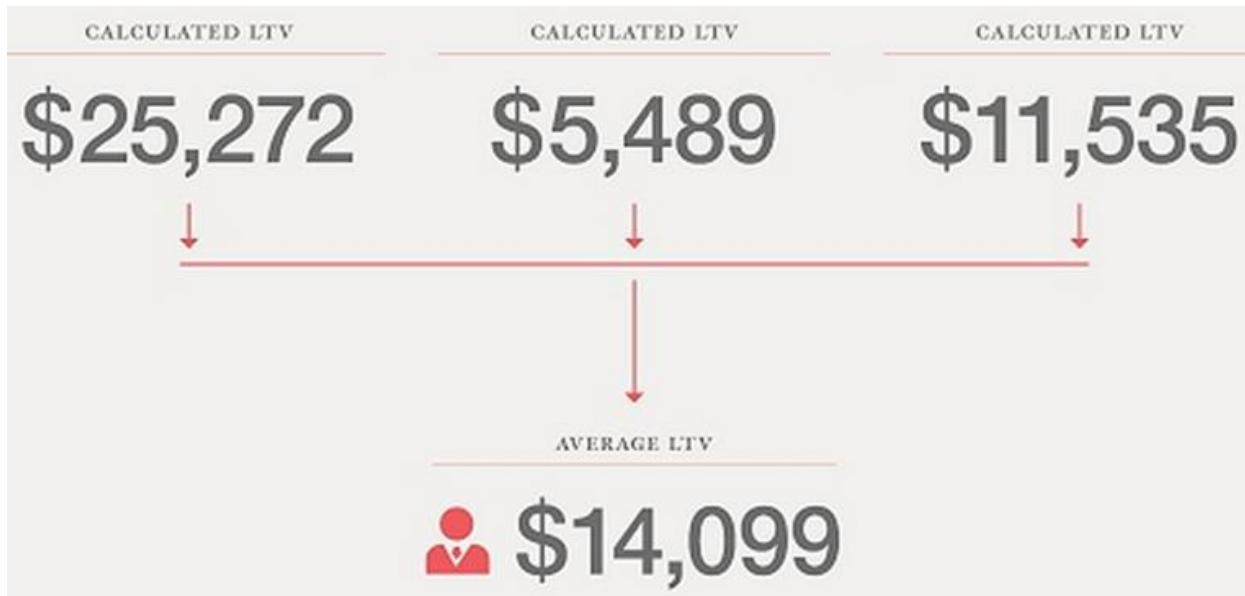
<https://blog.kissmetrics.com/how-to-calculate-lifetime-value/>

# LTV Analytics



<https://blog.kissmetrics.com/how-to-calculate-lifetime-value/>

# LTV Analytics



<https://blog.kissmetrics.com/how-to-calculate-lifetime-value/>

The **Pareto principle** (also known as the **80–20 rule**, the **law of the vital few**, and the **principle of factor sparsity**) states that, for many events, roughly 80% of the effects come from 20% of the causes

- 80% of a company's profits come from 20% of its customers
- 80% of a company's complaints come from 20% of its customers
- 80% of a company's profits come from 20% of the time its staff spend
- 80% of a company's sales come from 20% of its products
- 80% of a company's sales are made by 20% of its sales staff

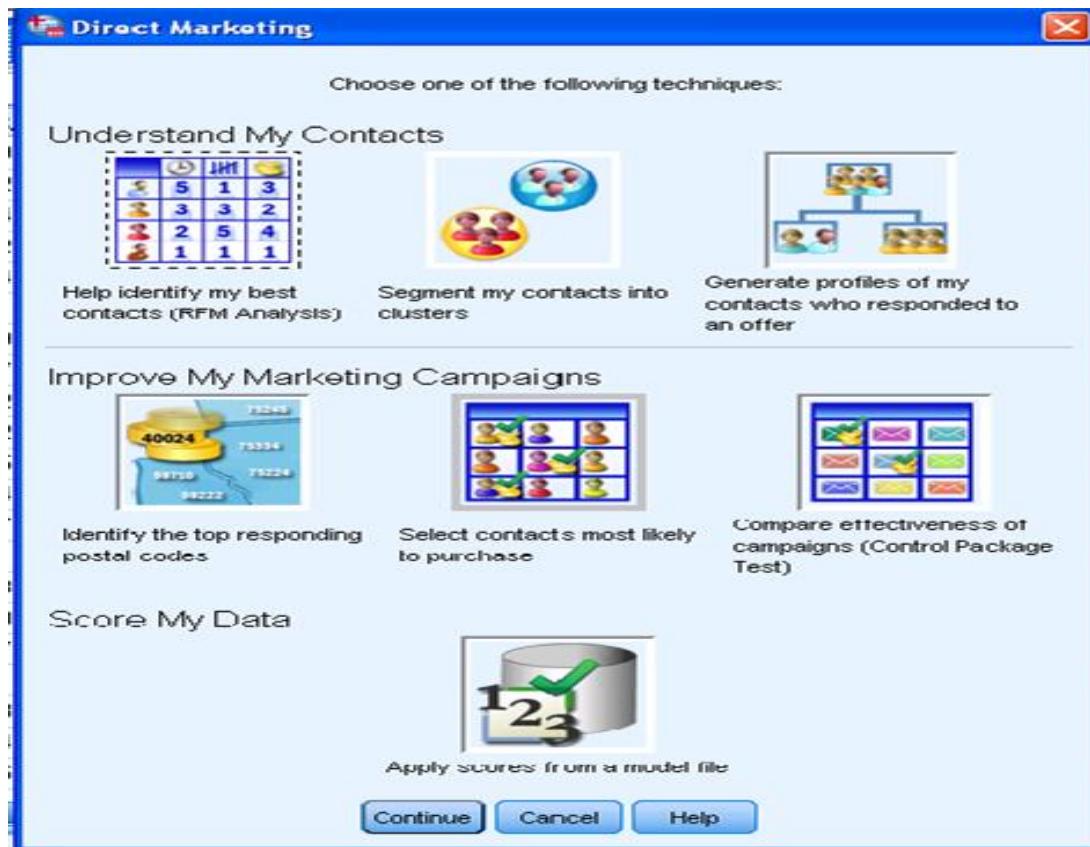
Several criminology studies have found 80% of crimes are committed by 20% of criminals.

# RFM Analysis

RFM is a method used for analyzing customer value.

- Recency - *How recently did the customer purchase?*
- Frequency - *How often do they purchase?*
- Monetary Value - *How much do they spend?*

# RFM Analysis



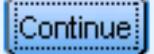
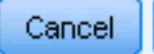
# RFM Analysis

 **RFM Analysis: Data Format** 

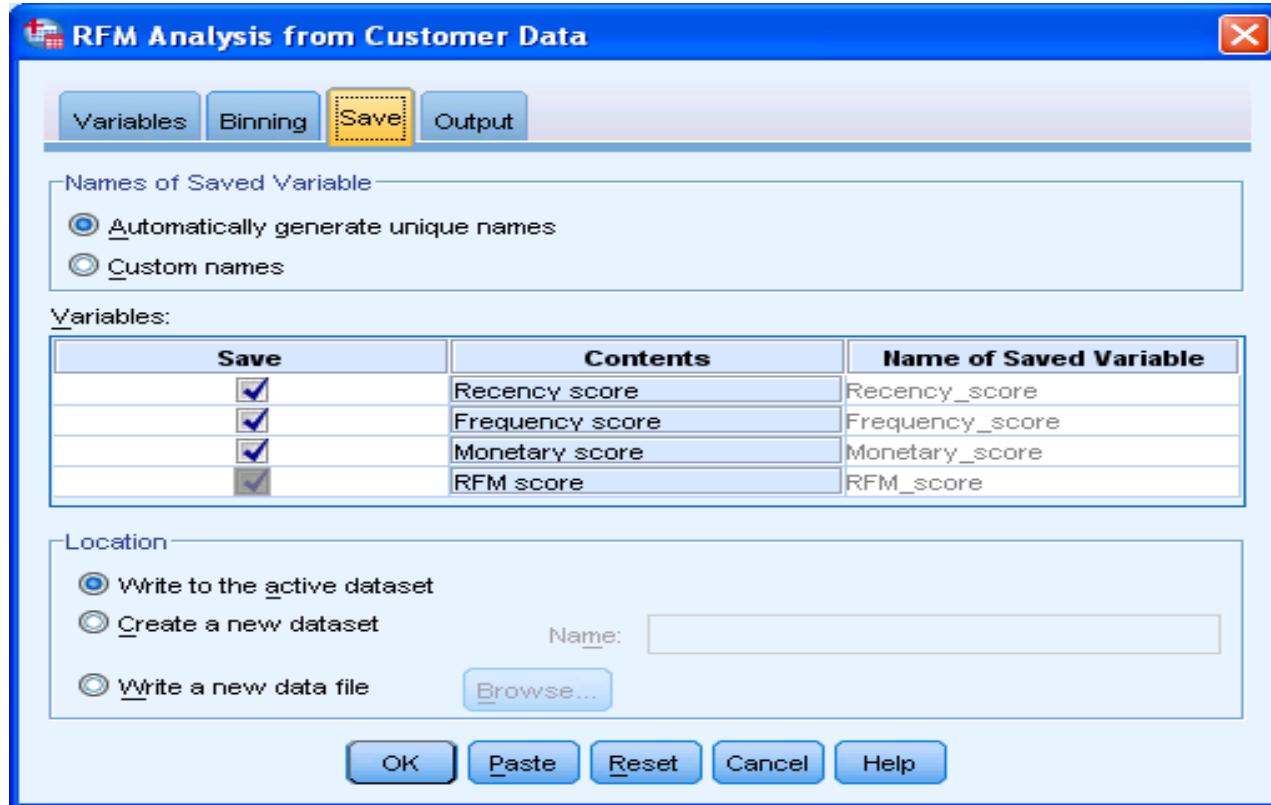
My data are:

**Transaction data**  
Each row contains data for one transaction. For the analysis, transactions will be combined by customer identifiers.

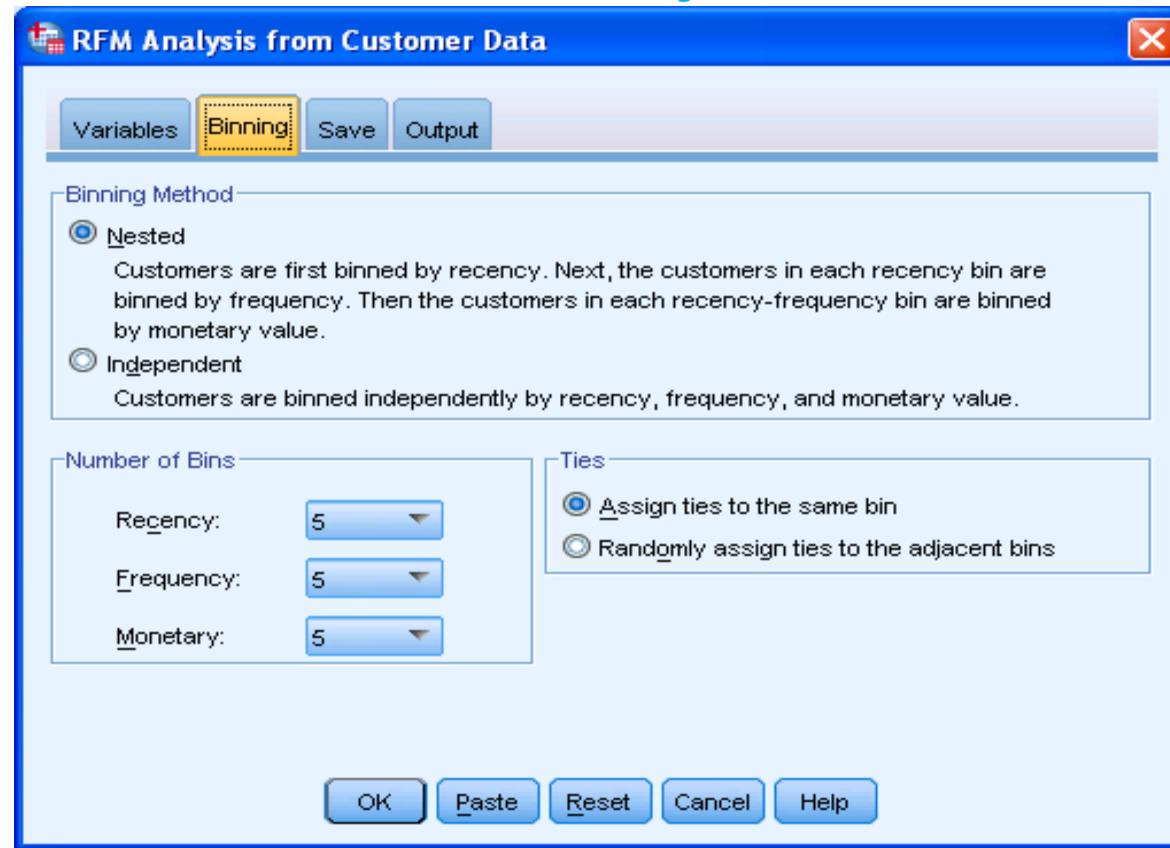
**Customer data**  
Each row contains data for one customer. The data have already been combined by customer over transactions.

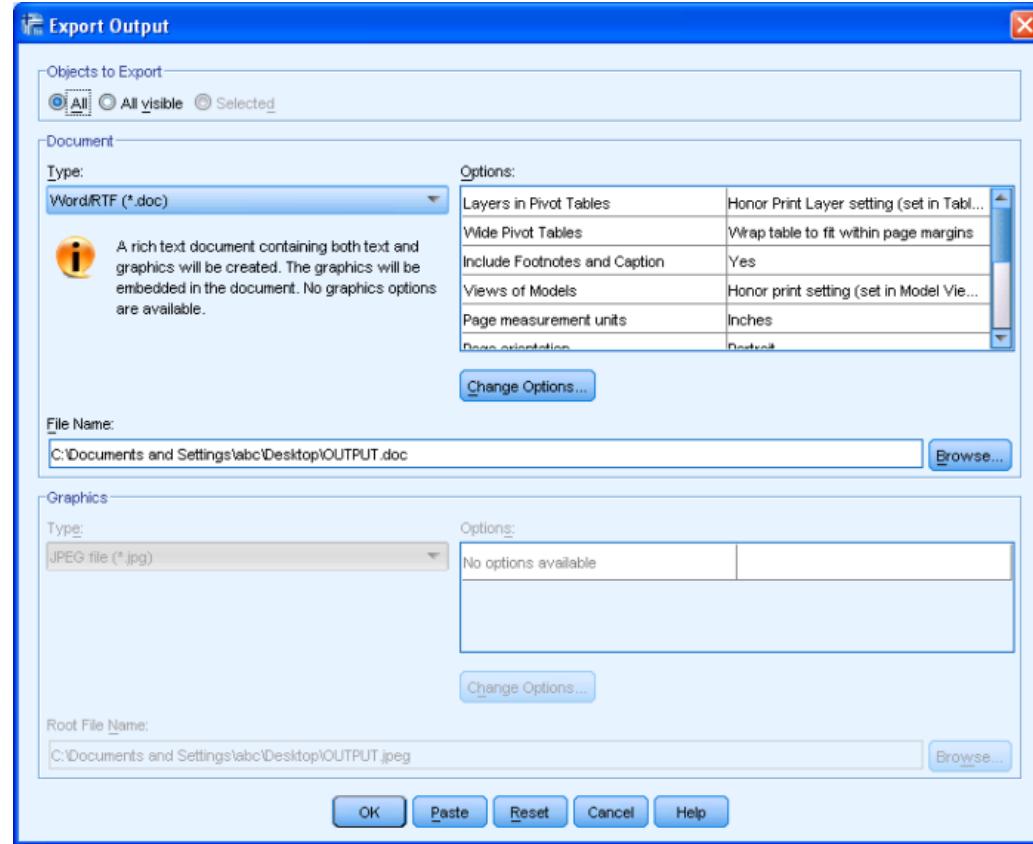
# RFM Analysis



# RFM Analysis



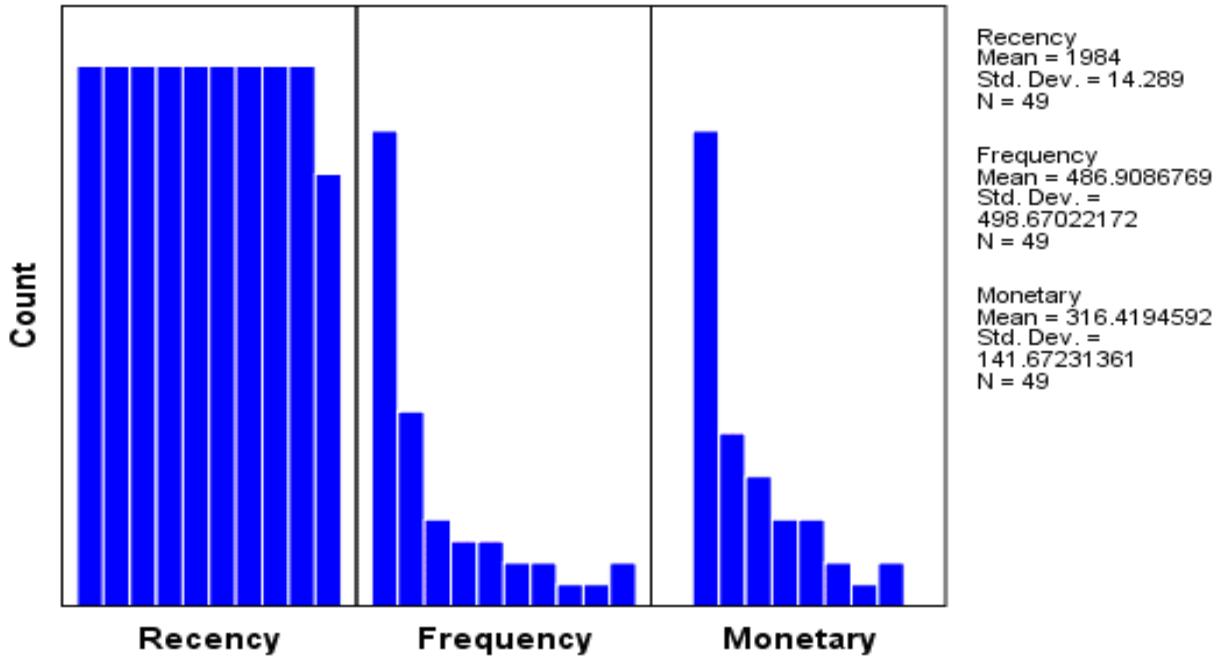
# RFM Analysis



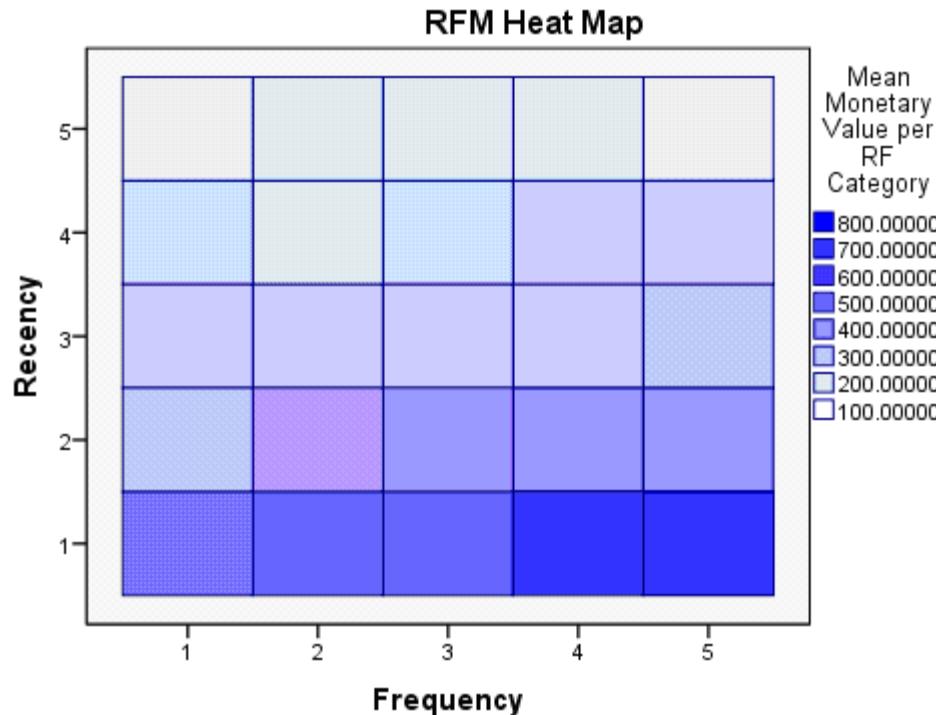
## Using SPSS 19 - example

# RFM Analysis

RFM  
Histograms

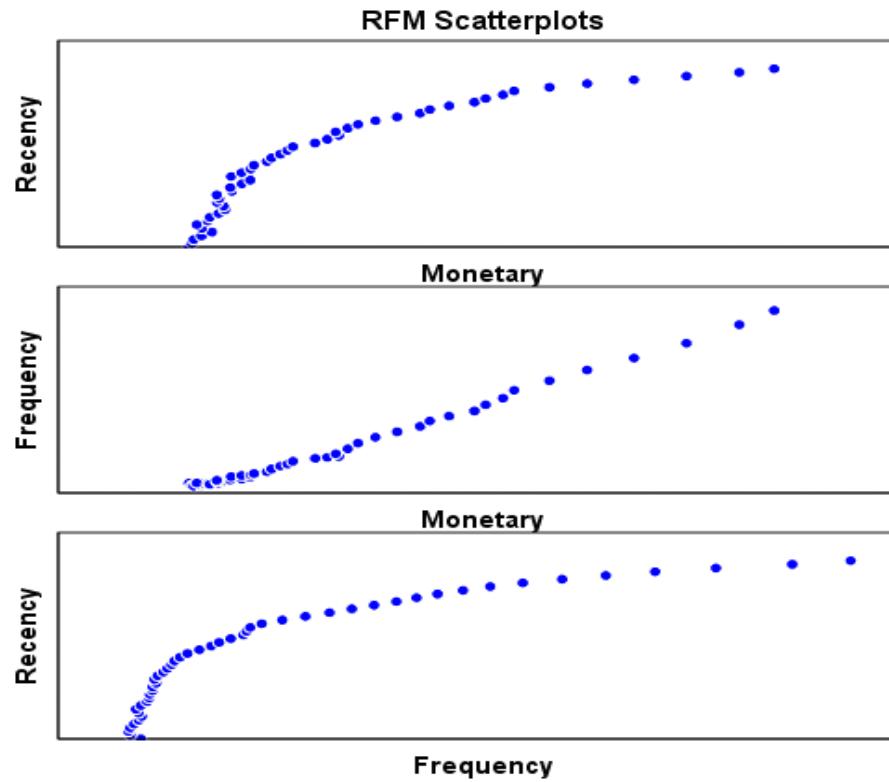


# RFM Analysis

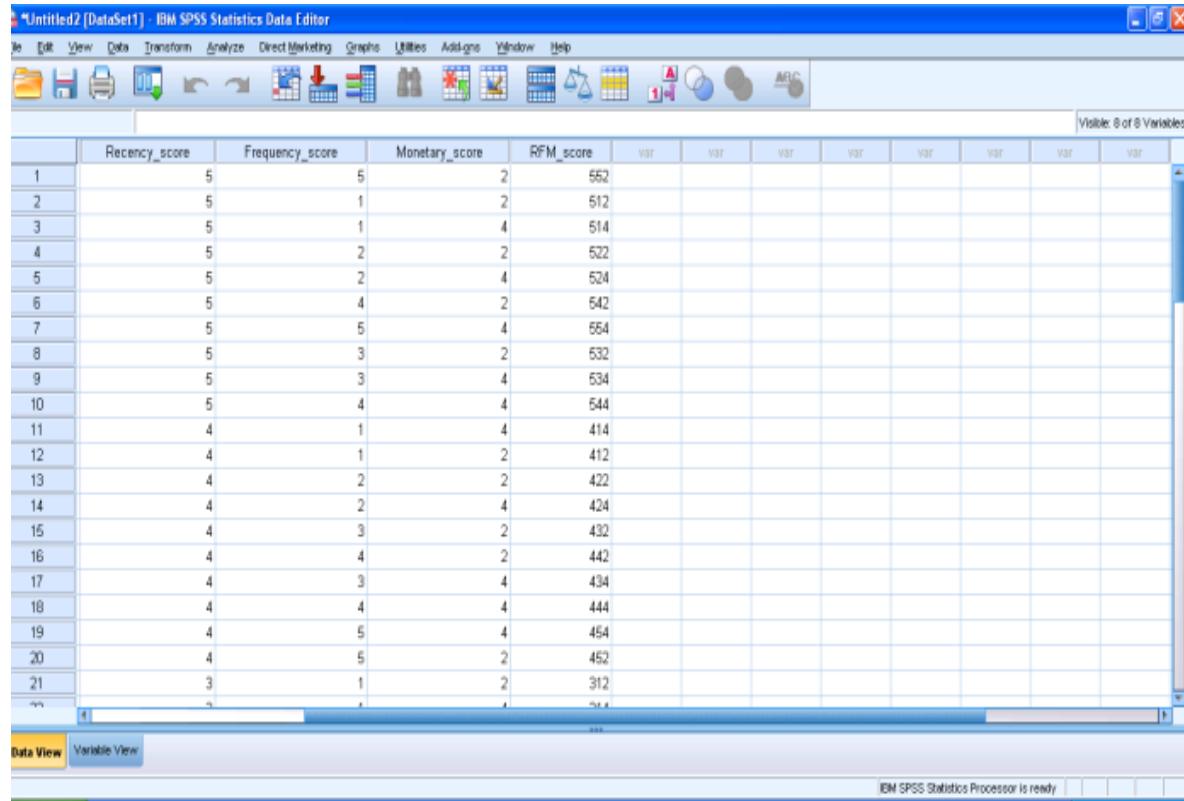


# RFM Analysis

Using SPSS 19 - example



# RFM Analysis



The screenshot shows the IBM SPSS Statistics Data Editor window titled "Untitled2 [DataSet1] - IBM SPSS Statistics Data Editor". The menu bar includes File, Edit, View, Data, Transform, Analyze, Direct Marketing, Graphs, Utilities, Addins, Window, and Help. The toolbar contains various icons for file operations, data manipulation, and analysis. The data view displays a table with 22 rows and 11 columns. The columns are labeled: Recency\_score, Frequency\_score, Monetary\_score, RFM\_score, var, var, var, var, var, var, and var. The "RFM\_score" column contains the calculated RFM scores for each row. The "var" columns are empty. The status bar at the bottom indicates "IBM SPSS Statistics Processor is ready".

	Recency_score	Frequency_score	Monetary_score	RFM_score	var						
1	5	5	2	552							
2	5	1	2	512							
3	5	1	4	514							
4	5	2	2	522							
5	5	2	4	524							
6	5	4	2	542							
7	5	5	4	554							
8	5	3	2	532							
9	5	3	4	534							
10	5	4	4	544							
11	4	1	4	414							
12	4	1	2	412							
13	4	2	2	422							
14	4	2	4	424							
15	4	3	2	432							
16	4	4	2	442							
17	4	3	4	434							
18	4	4	4	444							
19	4	5	4	454							
20	4	5	2	452							
21	3	1	2	312							
22	5	1	4	514							

Using SPSS 19 - example

# RFM Analysis

RFM is a method used for analyzing customer value.

- Recency - *How recently did the customer purchase?*
- Frequency - *How often do they purchase?*
- Monetary Value - *How much do they spend?*

A method

- Recency = 10 - the number of months that have passed since the customer last purchased
- Frequency = number of purchases in the last 12 months (maximum of 10)
- Monetary = value of the highest order from a given customer (benchmarked against \$10k)

Alternatively, one can create categories for each attribute. For instance, the Recency attribute might be broken into three categories: customers with purchases within the last 90 days; between 91 and 365 days; and longer than 365 days. Such categories may be arrived at by applying business rules, or using a data mining technique, to find meaningful **breaks**.

A commonly used shortcut is to use deciles. One is advised to look at distribution of data before choosing breaks.

# Refresher in Statistics

## Mean

Arithmetic Mean- the sum of the values divided by the number of values.

The [geometric mean](#) is an average that is useful for sets of positive numbers that are interpreted according to their product and not their sum (as is the case with the arithmetic mean) e.g. rates of growth.

## Median

the **median** is the number separating the higher half of a data sample, a population, or a probability distribution, from the lower hal

## Mode-

The "mode" is the value that occurs most often.

# Refresher in Statistics

## Range

the **range** of a set of data is the difference between the largest and smallest values.

## Variance

mean of squares of differences of values from mean

## Standard Deviation

square root of its variance

## Frequency

a **frequency distribution** is a table that displays the **frequency** of various outcomes in a sample.

# Distributions

## Bernoulli

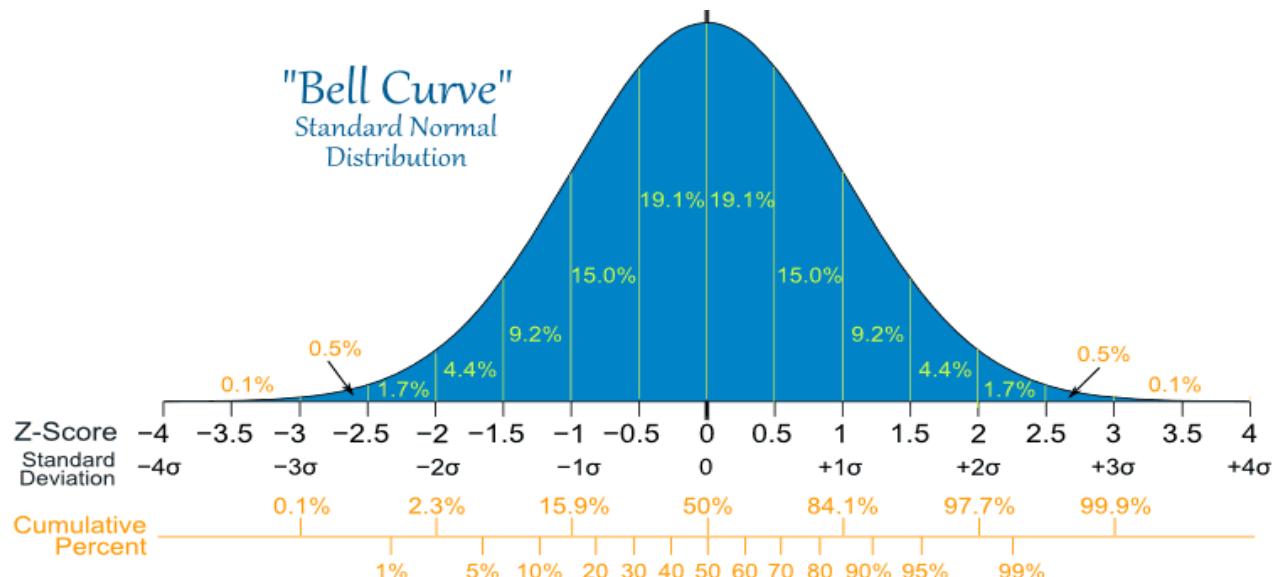
$$p \qquad \qquad q = 1 - p$$

Distribution of a random variable which takes value 1 with success probability and value 0 with failure probability. It can be used, for example, to represent the toss of a coin

# Distributions

## Normal

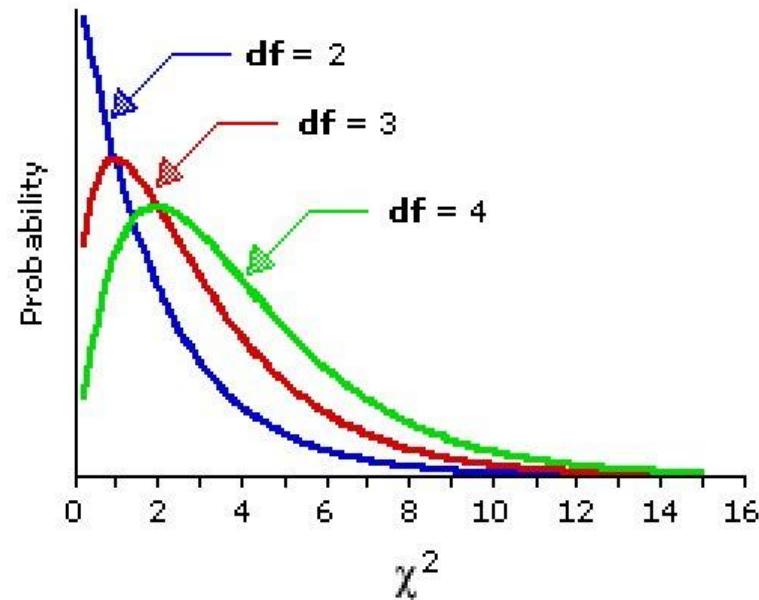
The simplest case of a normal distribution is known as the *standard normal distribution*. This is a special case where  $\mu=0$  and  $\sigma=1$ ,



# Distributions

## Chi Square

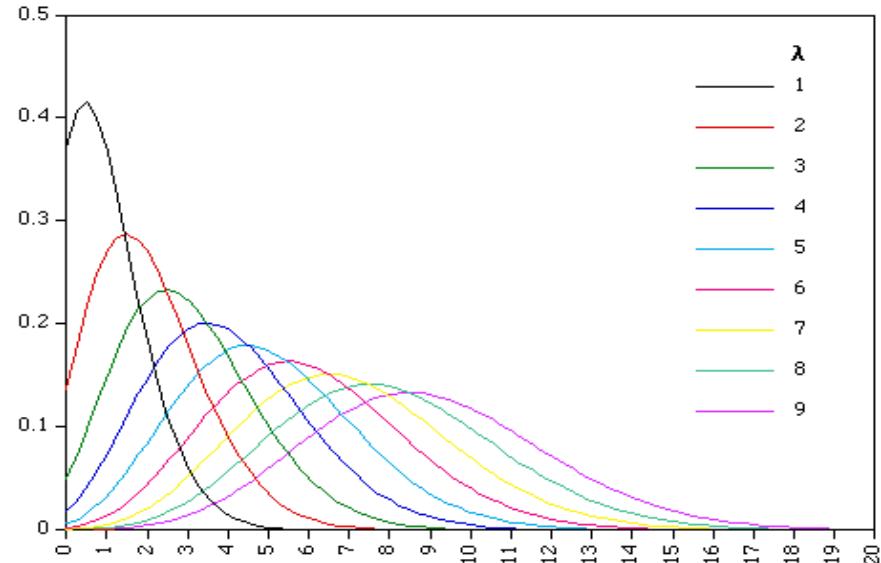
the distribution of a sum of the squares of  $k$  independent standard normal random variables.



# Distributions

## Poisson

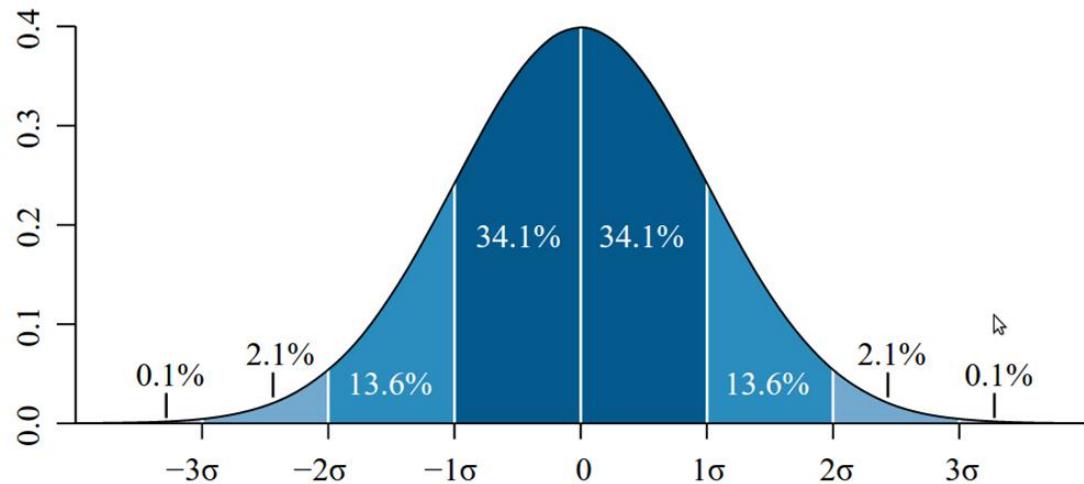
a discrete probability distribution that expresses the probability of a given number of events occurring in a fixed interval of time and/or space if these events occur with a known average rate and independently of the time since the last event



# Refresher in Statistics

## Probability Distribution

The probability density function (pdf) of the normal distribution, also called Gaussian or "bell curve", the most important continuous random distribution. As notated on the figure, the probabilities of intervals of values correspond to the area under the curve.



# Refresher in Statistics

## Central Limit Theorem –

In probability theory, the **central limit theorem (CLT)** states that, given certain conditions, the arithmetic mean of a sufficiently large number of iterates of independent random variables, each with a well-defined expected value and well-defined variance, will be approximately normally distributed, regardless of the underlying distribution.

# Introduction to Modeling

# Hypothesis testing

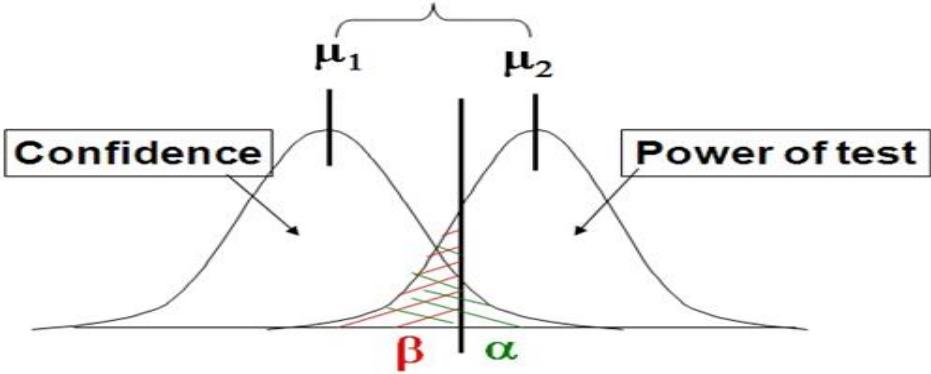
Hypothesis testing is the use of statistics to determine the probability that a given hypothesis is true. The usual process of hypothesis testing consists of four steps.

1. Formulate the null hypothesis (commonly, that the observations are the result of pure chance) and the alternative hypothesis (commonly, that the observations show a real effect combined with a component of chance variation).
2. Identify a test statistic that can be used to assess the truth of the null hypothesis.
3. Compute the P-value, which is the probability that a test statistic at least as significant as the one observed would be obtained assuming that the null hypothesis were true. The smaller the -value, the stronger the evidence against the null hypothesis.
4. Compare the -value to an acceptable significance value (sometimes called an alpha value). If , that the observed effect is statistically significant, the null hypothesis is ruled out, and the alternative hypothesis is valid.

<http://mathworld.wolfram.com/HypothesisTesting.html>

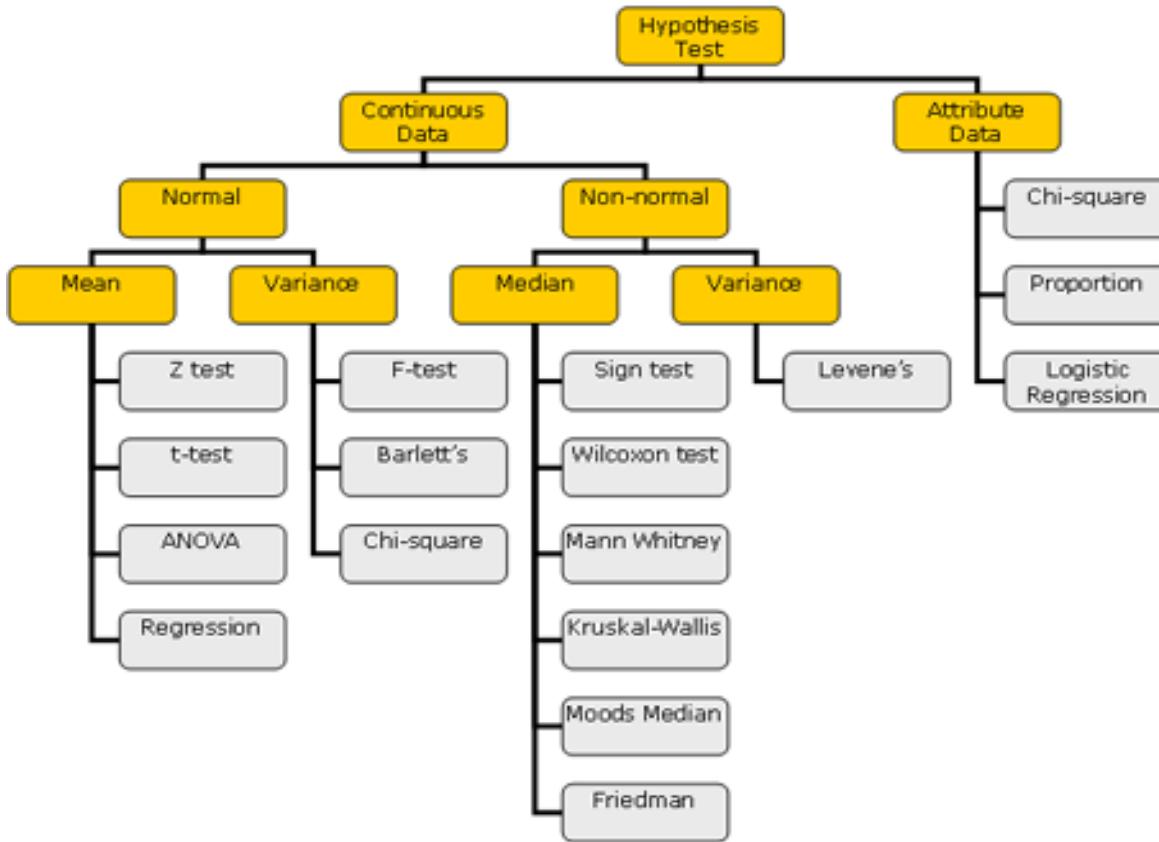
# Hypothesis testing

		Reality	
		$H_0$	$H_a$
$H_0$		Correct decision	Type II error $\beta$ -risk Producer's risk
$H_a$		Type I error $\alpha$ -risk Consumer's risk	Correct decision
Decision			



The Truth (unknown to the researcher)		
The Researcher's Decision	The Null Hypothesis is True	The Null Hypothesis is False
Reject the Null Hypothesis	Type I Error	Correct Decision
Fail to Reject the Null Hypothesis	Correct Decision	Type II Error

# Hypothesis testing



# Hypothesis testing

Comparison of MEANS	Degrees of Freedom	Application	Assumptions	Test Statistic
One Sample Z-Test	Not Applicable	Testing the difference of a sample mean, $\bar{x}$ , with a known population mean, $\mu$ (fixed mean, historical mean, or targeted mean)	Normal distribution Known population $\sigma$ .	$Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$
One Sample t-test	$n-1$	Testing the difference of one sample mean, $\bar{x}$ , with a known population mean, $\mu$ (fixed mean, historical mean, or targeted mean)	Normal distribution Population standard deviation, $\sigma$ , is unknown.	$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$
Two Sample t-test	$n_1 + n_2 - 2$	Testing difference of two sample means when population variances unknown but <u>considered equal</u>	Normal Distribution Requires standard pooled deviation calculation, $s_p$	$t = \frac{\bar{x}_1 - \bar{x}_2}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$
Paired t-test	$n - 1$	Testing two sample means when their respective population standard deviations are unknown but considered equal. Data recorded in pairs and each pair has a difference, $d$ .	Normal Distribution Two dependent samples Always two-tailed test $s_d$ = standard deviation of the differences of all samples	$t = \frac{\bar{d} \sqrt{n}}{s_d}$
One-Way ANOVA	$n_1 - 1$ & $n_2 - 1$	Testing the difference of three or more population means	Normal Distribution $s_1^2$ and $s_2^2$ represent sample variances	$F = \frac{(s_1)^2}{(s_2)^2}$

# Hypothesis testing

R Data Miner - [Rattle]

Execute New Open Save Report Export Stop Quit

Data Explore Test Transform Cluster Associate Model Evaluate Log

Two-Sample Tests:  Kolmogorov-Smirnov  Wilcoxon Rank-Sum  T-test  F-test

Paired Two-Sample Tests:  Correlation  Wilcoxon Signed Rank

Sample 1:  Sample 2:  Group By Target: No Target

**Statistical Tests**

These tests apply to two samples. The paired two sample tests assume that we have two samples or observations, and that we are testing for a change, usually from one time period to another.

**Distribution of the Data**

\* Kolmogorov-Smirnov      Non-parametric      Are the distributions the same?  
\* Wilcoxon Signed Rank      Non-parametric      Do paired samples have the same distribution?

**Location of the Average**

\* T-test      Parametric      Are the means the same?  
\* Wilcoxon Rank-Sum      Non-parametric      Are the medians the same?

**Variation in the Data**

\* F-test      Parametric      Are the variances the same?

**Correlation**

\* Correlation      Pearsons      Are the values from the paired samples correlated?

# Data Mining

**Data Mining** is an analytic process designed to explore **data** (usually large amounts of **data** - typically business or market related - also known as "big **data**") in search of consistent patterns and/or systematic relationships between variables, and then to validate the findings by applying the detected patterns

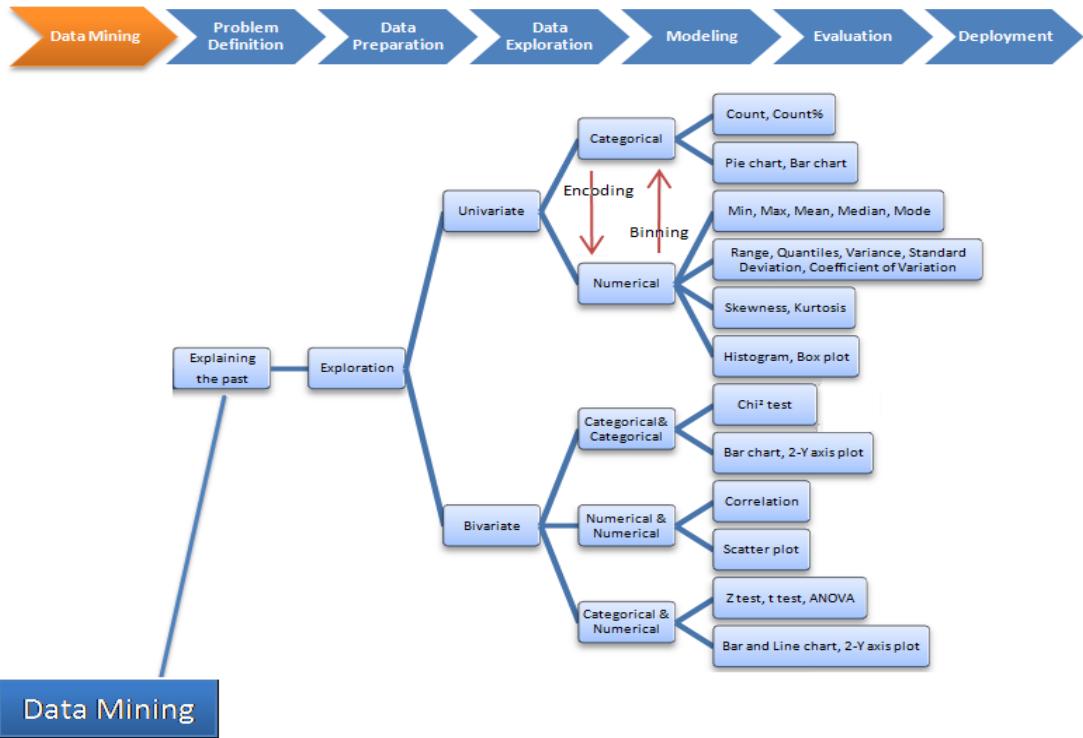
# Data Mining Map

Copyright © 2010-2015, [Dr. Saed Sayad](#)

## An Introduction to Data Mining

Source-

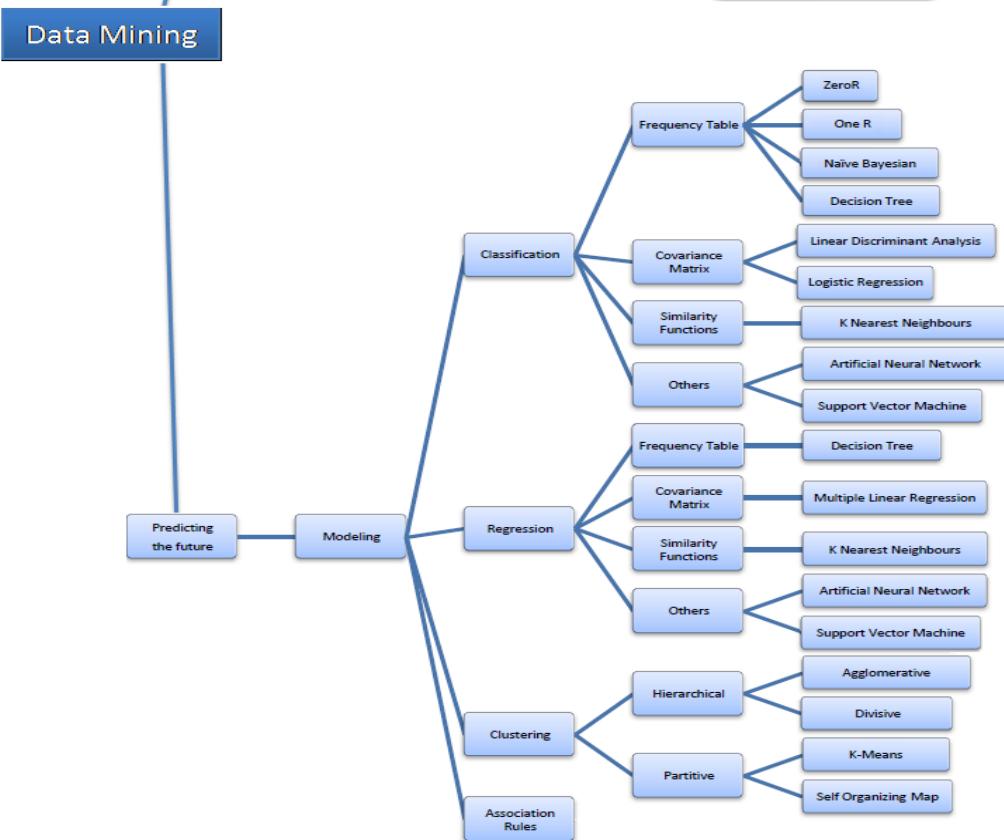
<http://www.saedsayad.com/>



# Data Mining Map

Source-

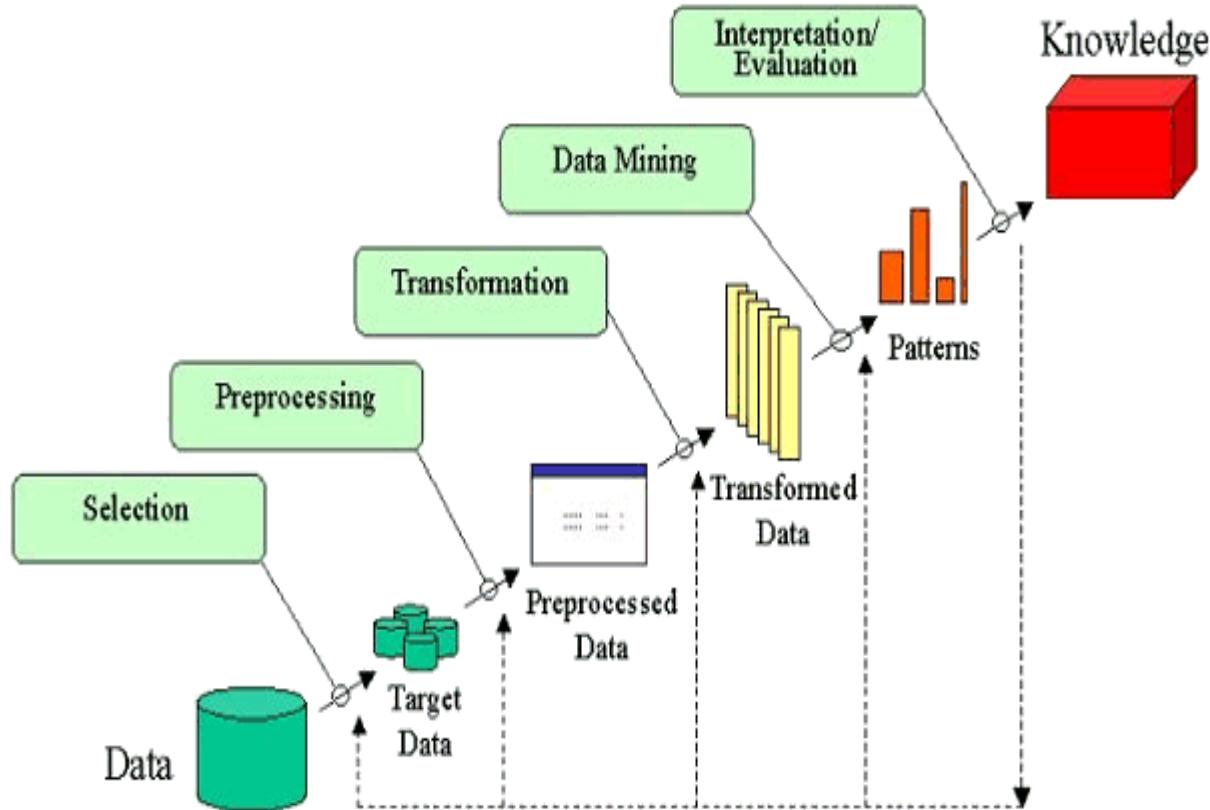
<http://www.saedsayad.com/>



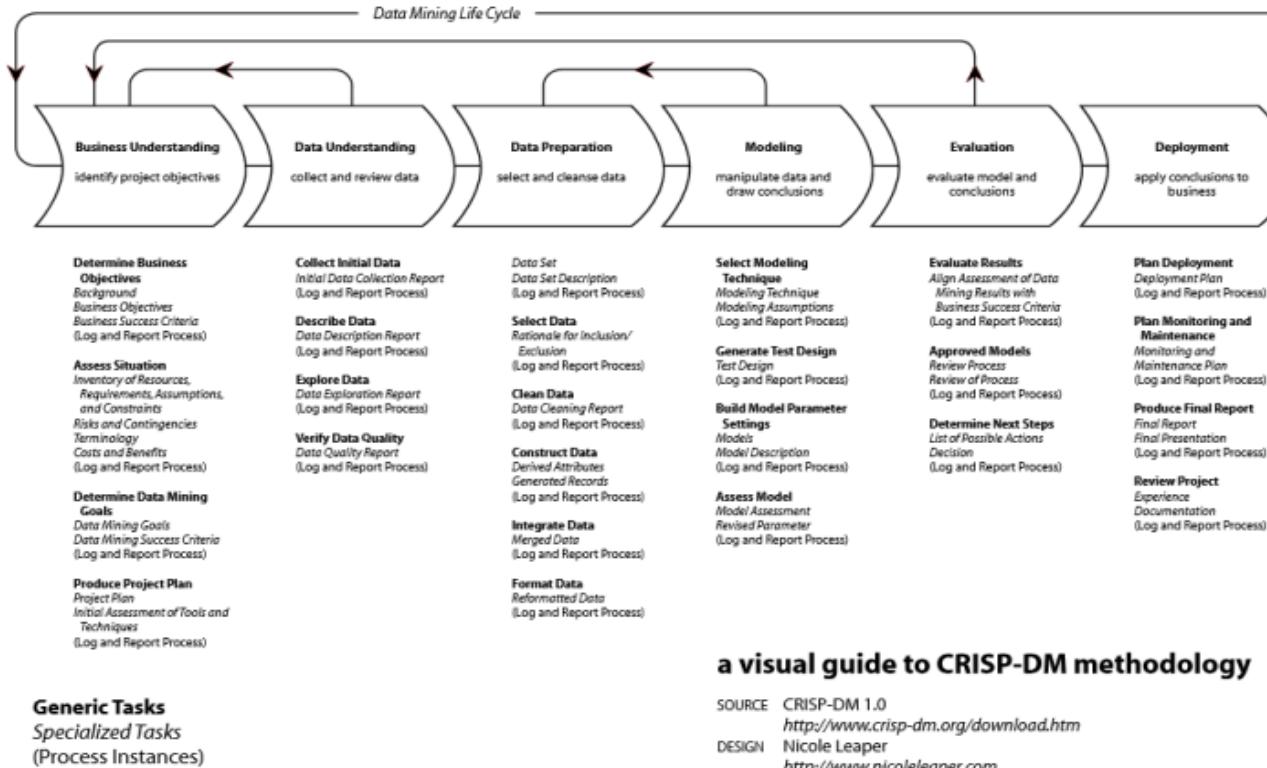
# Examples of Data Mining

- which items sell well together in retail ( market basket)
- which products sell well together on a website ( association analysis)
- which customers are likely to buy a new credit card (regression)

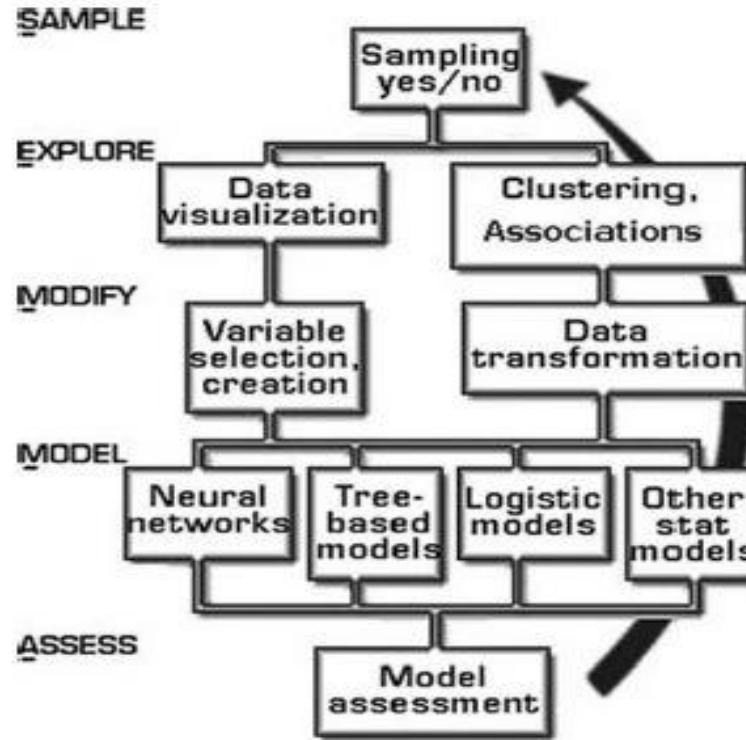
# KDD



# CRISP DM



# SEMMA



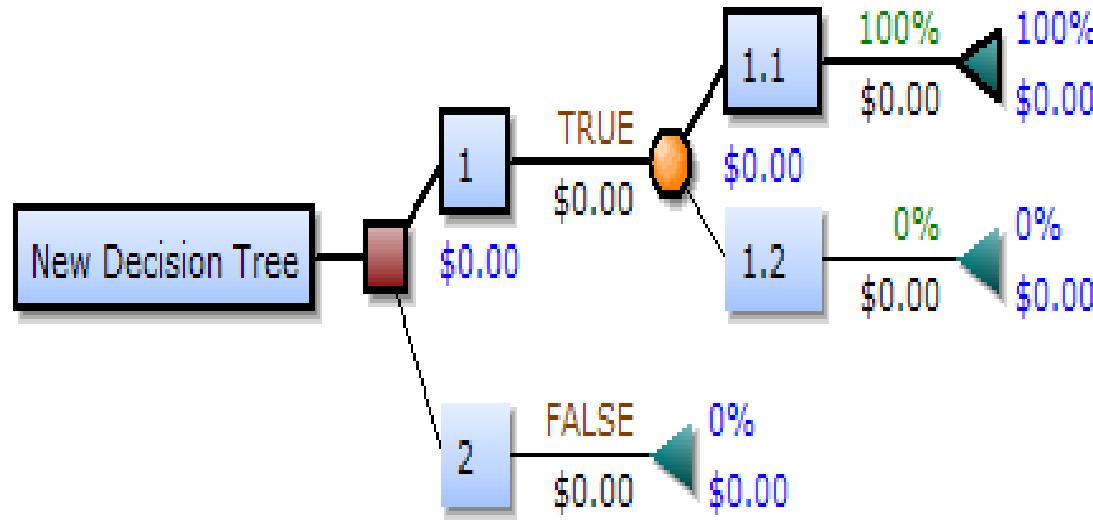
# Machine Learning

**Machine learning** is a subfield of computer science that evolved from the study of pattern recognition and computational learning theory in artificial intelligence.

**Machine learning** explores the construction and study of algorithms that can learn from and make predictions on data.

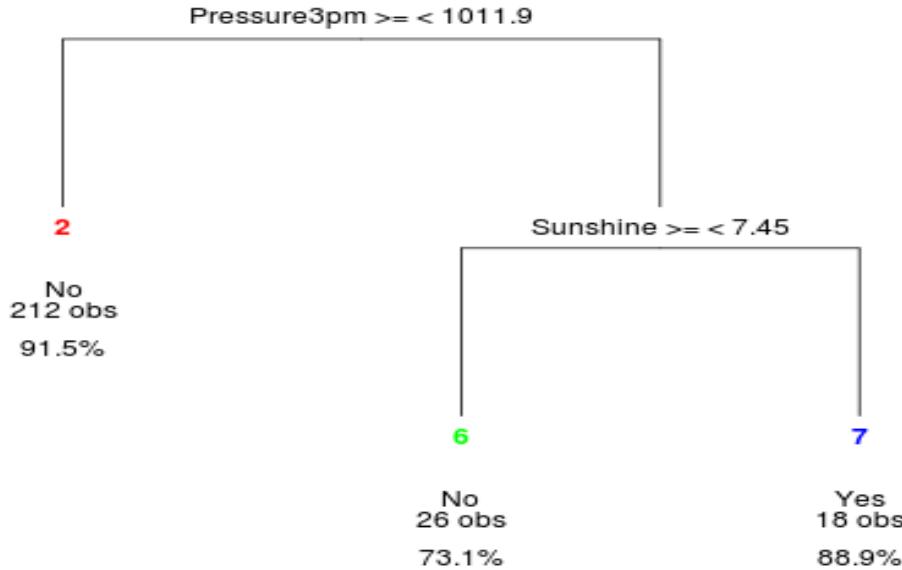
- Supervised learning. The computer is presented with example inputs and their desired outputs, given by a "teacher", and the goal is to learn a general rule that maps inputs to outputs.
- Unsupervised learning, no labels are given to the learning algorithm, leaving it on its own to find structure in its input. Unsupervised learning can be a goal in itself (discovering hidden patterns in data) or a means towards an end.
- In classification, a supervised way, inputs are divided into two or more classes, and the learner must produce a model that assigns unseen inputs to one (or multi-label classification) Spam filtering, where the inputs are email (or other) messages and the classes are "spam" and "not spam".
- In regression, also a supervised problem, the outputs are continuous rather than discrete.
- In clustering, a set of inputs is to be divided into groups. Unlike in classification, the groups are not known beforehand, making this typically an unsupervised task.

# Decision Trees



# Decision Trees

Decision Tree rpart() weather \$ RainTomorrow



Rattle 2012-Mar-12 09:53:45 gjw

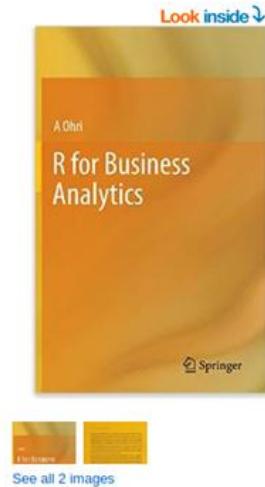
# Association Analysis

As an unsupervised learning technique it has delivered considerable benefit in areas ranging from the traditional shopping basket analysis to the analysis of who bought what other books or who watched what other videos, and in areas including health care, telecommunications, and so on

from

<http://handsondatascience.com/ARulesO.pdf>

# An Example of Data Mining



See all 2 images

**R for Business Analytics** Hardcover – Import, 11 Sep 2012

by A Ohrni (Author)

1 customer review

See all 2 formats and editions

Kindle Edition

₹ 1,673.49

Hardcover

₹ 3,865.16

Read with our free app

10 New from ₹ 3,285.86

EMI Available. Options ▾

Delivery to pincode 110001 - Delhi : within 1 - 2 weeks. Details

*R for Business Analytics* looks at some of the most common tasks performed by business analysts and helps the user navigate the wealth of information in R and its 4000 packages. With this information the reader can select the packages that can help process the analytical tasks with minimum effort and maximum usefulness. The use of Graphical User Interfaces (GUI) is emphasized in this book to further cut down and bend the famous learning curve in learning R. This book is aimed to help you kick-start with analytics including chapters on data visualization, code examples on web analytics and social media analytics, clustering, regression models, text mining, data mining models and forecasting. The book tries to expose the reader to a breadth of business analytics topics without burying the user in needless

Read more

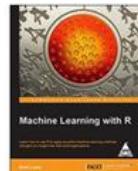
## Customers Who Bought This Item Also Bought



R for Everyone: Advanced Analytics and Graphics



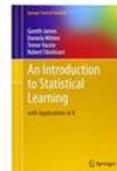
Applied Predictive Modeling



Machine Learning With R: Learn How to Use R to...



Data Smart: Using Data Science to Transform...



An Introduction to Statistical Learning:...

Page 1 of 4

Share



100% Purchase Protection  
Genuine Products | Secure Payments  
Easy Returns

₹ 3,865.16 + FREE Delivery

Inclusive of all taxes

Sold and fulfilled by B2A UK (4.4 out of 5 | 1,934 ratings).



Add to Cart



Buy Now

Add to Wish List

## Other Sellers on Amazon

₹ 4,250.00

+ ₹ 300.00 Delivery charge  
Sold by: A1websites

₹ 4,875.56

+ FREE Delivery  
Sold by: B2A US

₹ 4,937.00

+ ₹ 300.00 Delivery charge  
Sold by: uRead-shop

10 New from ₹ 3,285.86

# An Example of Data Mining

## Examples

[https://en.wikipedia.org/wiki/Apriori\\_algorithm](https://en.wikipedia.org/wiki/Apriori_algorithm)

transaction ID	milk	bread	butter	beer	diapers
1	1	1	0	0	0
2	0	0	1	0	0
3	0	0	0	1	1
4	1	1	1	0	0
5	0	1	0	0	0

# Clustering

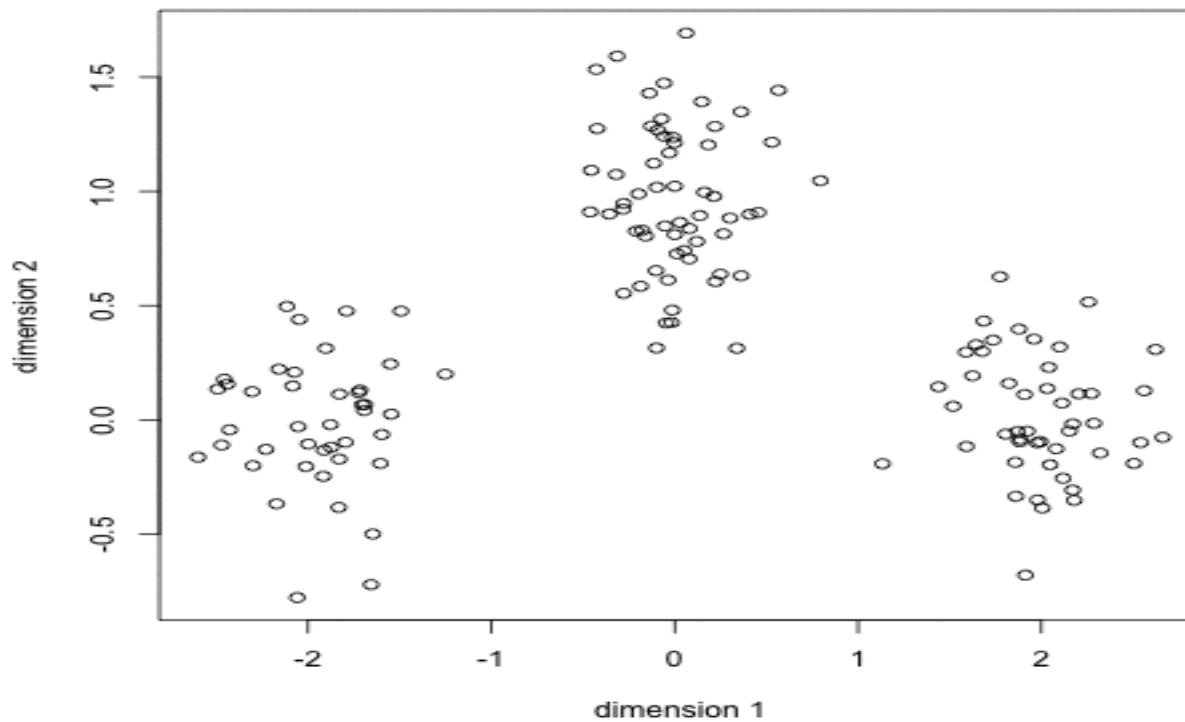
**Cluster analysis** or **clustering** is the task of grouping a set of objects in such a way that objects in the same group (called a **cluster**) are more similar (in some sense or another) to each other than to those in other groups (**clusters**).

**k-means clustering** aims to partition n observations into **k clusters** in which each observation belongs to the **cluster** with the nearest **mean**, serving as a prototype of the**cluster**. This results in a partitioning of the data space into Voronoi cells

<http://shabal.in/visuals/kmeans/1.html>

# Clustering

**step 0**

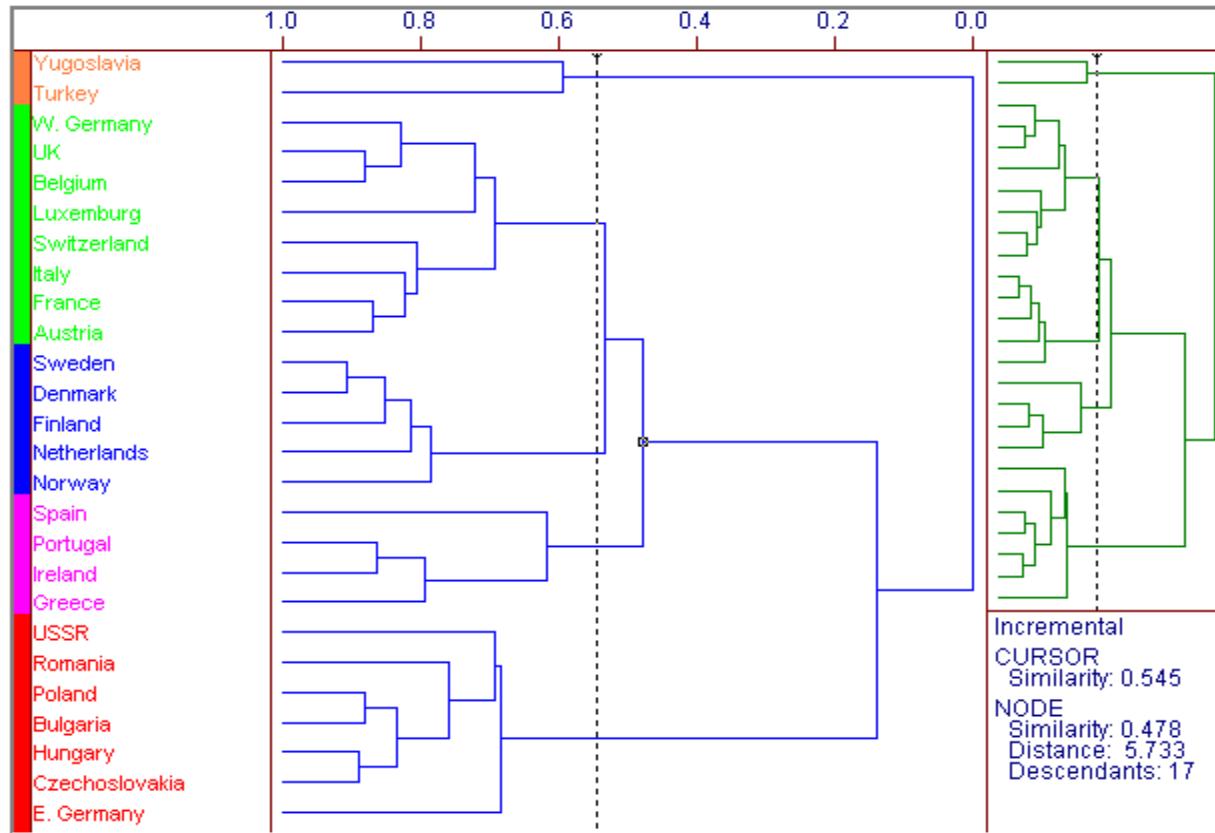


# Clustering

**hierarchical clustering** (also called **hierarchical cluster analysis** or **HCA**) is a method of cluster analysis which seeks to build ahierarchy of clusters. Strategies for hierarchical clustering generally fall into two types: [\[1\]](#)

- **Agglomerative:** This is a "bottom up" approach: each observation starts in its own cluster, and pairs of clusters are merged as one moves up the hierarchy.
- **Divisive:** This is a "top down" approach: all observations start in one cluster, and splits are performed recursively as one moves down the hierarchy.

# Clustering



# Regression

**regression analysis** is a statistical process for estimating the relationships among variables. It includes many techniques for modeling and analysing several variables, when the focus is on the relationship between a dependent variable and one or more independent variables.

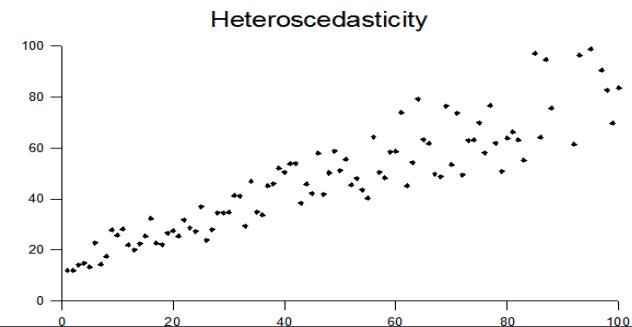
$$y = a + bx$$

$$y = a + bx + cy$$

$$\ln(p / 1-p) = a + bx$$

# Regression

1. In statistics, **multicollinearity** (also collinearity) is a phenomenon in which two or more predictor variables in a multiple regression model are highly correlated, meaning that one can be linearly predicted from the others with a non-trivial degree of accuracy. A multiple regression model with correlated predictors can indicate how well the entire bundle of predictors predicts the outcome variable, but it may not give valid results about any individual predictor
2. **heteroscedasticity**(also spelled heteroskedasticity) refers to the circumstance in which the variability of a variable is unequal across the range of values of a second variable that predicts it.
3. The p-value for each term tests the null hypothesis that the coefficient is equal to zero (no effect). A low p-value (< 0.05) indicates that you can reject the null hypothesis.



# Text Mining

1. Retrieving Text
2. Transforming Text to corpus
3. Cleaning Text (lowercase, punctuation, numbers, commonly used words (stop words))
4. Stemming Words
5. Building a Document-Term Matrix
6. Frequent Terms and Associations
7. Word Cloud

<http://www.rdatamining.com/examples/text-mining>

# Sentiment Analysis

**Sentiment analysis** (also known as **opinion mining**) refers to the use of natural language processing, text analysis and computational linguistics to identify and extract subjective information in source materials.

Generally speaking, sentiment analysis aims to determine the attitude of a speaker or a writer with respect to some topic or the overall contextual polarity of a document

example- <http://www.slideshare.net/ajayohri/twitter-analysis-by-kaify-rais>

# Sentiment Analysis

A sentiment analysis model is used to analyze a text string and classify it with one of the labels that you provide; for example, you could analyze a tweet to determine whether it is positive or negative, or analyze an email to determine whether it is happy, frustrated, or sad.

## R package "sentiment"

Another interesting option that we can use to do our sentiment analysis is by utilizing the R package [sentiment](#) by Timothy Jurka. This package contains two handy functions serving our purposes:

### **classify\_emotion**

This function helps us to analyze some text and classify it in different types of emotion: anger, disgust, fear, joy, sadness, and surprise.

### **classify\_polarity**

In contrast to the classification of emotions, the `classify_polarity` function allows us to classify some text as positive or negative.

example- <https://sites.google.com/site/miningtwitter/questions/sentiment/sentiment>

# Social Network Analysis

**Social network analysis** (SNA) is a strategy for investigating **social** structures through the use of**network** and graph theories. It characterizes networked structures in terms of nodes (individual actors, people, or things within the **network**) and the ties or edges (relationships or interactions) that connect them.

The NSA has been performing social network analysis on [Call Detail Records](#) (CDRs), also known as [metadata](#), since shortly after the [September 11 Attacks](#)

Social Network Analysis to Optimize Tax Enforcement Effort -The South African Revenue Service

<http://aiselaisnet.org/cgi/viewcontent.cgi?article=1579&context=amcis2012>

*Irish Tax & Customs Authority*

[http://www.sas.com/en\\_ie/customers/irish-tax-and-customers.html](http://www.sas.com/en_ie/customers/irish-tax-and-customers.html)

# Social Network Analysis

Bridge: An individual whose weak ties fill a structural hole, providing the only link between two individuals or clusters. It also includes the shortest route when a longer one is unfeasible due to a high risk of message distortion or delivery failure.[\[18\]](#)

Centrality: Centrality refers to a group of metrics that aim to quantify the "importance" or "influence" (in a variety of senses) of a particular node (or group) within a network.

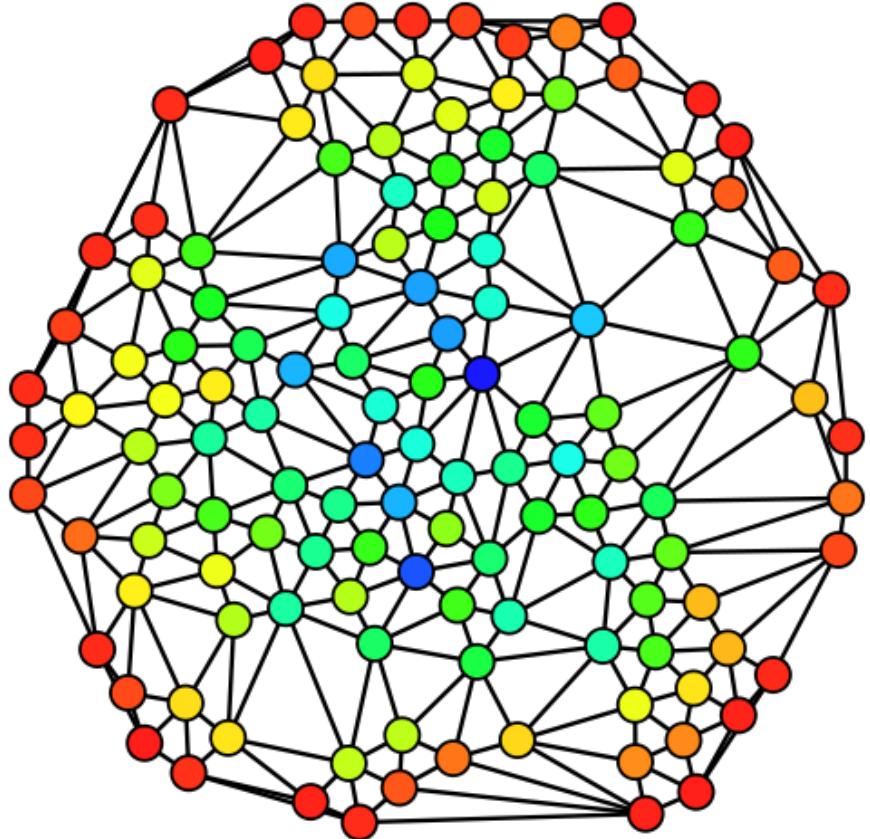
Density: The proportion of direct ties in a network relative to the total number possible.[\[25\]](#)[\[26\]](#)

Distance: The minimum number of ties required to connect two particular actors, as popularized by [Stanley Milgram's small world experiment](#) and the idea of 'six degrees of separation'.

Mutuality/Reciprocity: The extent to which two actors reciprocate each other's friendship or other interaction.[\[16\]](#)

Network Closure: A measure of the completeness of relational triads.

# Social Network Analysis



Hue (from red=0 to blue=max) indicates each node's [betweenness centrality](#).

# Social Network Analysis

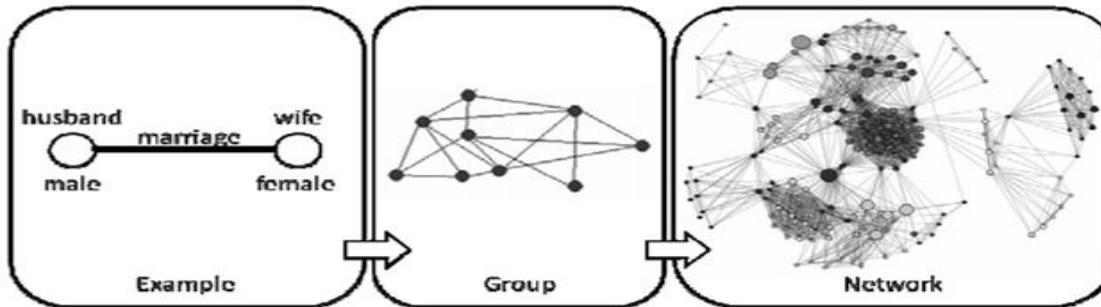


Figure 1. Social Network Analysis

## Managing Tax Compliance through Decision Support Systems

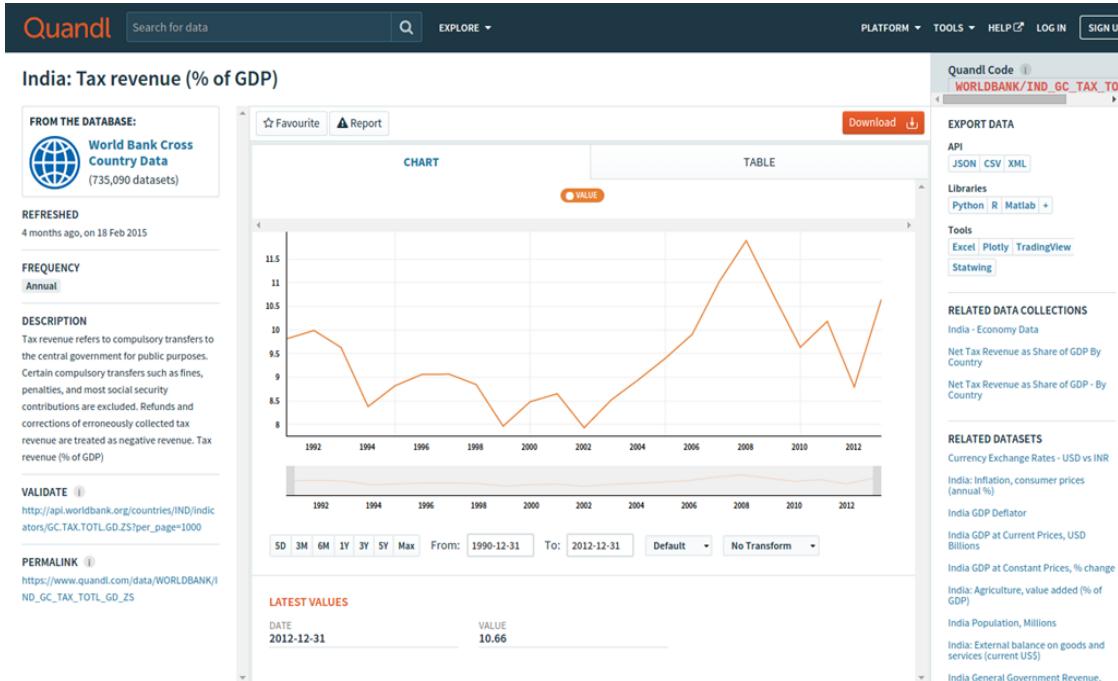
Much like any other organization, tax compliance can be managed at a strategic, operational and tactical level, as presented in Figure 2 (OECD, 2008:10). Strategic compliance management considers the tax system in its entirety whereas operational compliance management focuses on whole taxpayer segments. Tactical compliance management considers targeted individuals, or groups of which social structures such as marriage, employee and employer relationships and tax consultant and taxpayer relationships are but a few examples. The different DSS defined by Power (2002:13-16) can be associated with the types of compliance-management levels. Knowledge driven DSS are associated with strategic management, data driven DSS with operational management, and model driven DSS with tactical management. Model driven DSS is often used to conduct SNA, and DSS tools such as Analyst Notebook and SAS are widely recognized as industry leaders in this domain.

Types of DSS	Management Level	Compliance Monitoring Focus
Knowledge Driven	Strategic	Whole of tax system
Data Driven	Operational	Whole of tax product
Model Driven	Tactical	Whole of taxpayer segment Targeted compliance risk issues Targeted individuals/ groups

Figure 2. Compliance Management and Decision Support Systems (Derived from OECD, 2008:10; Power, 2002:13-16)

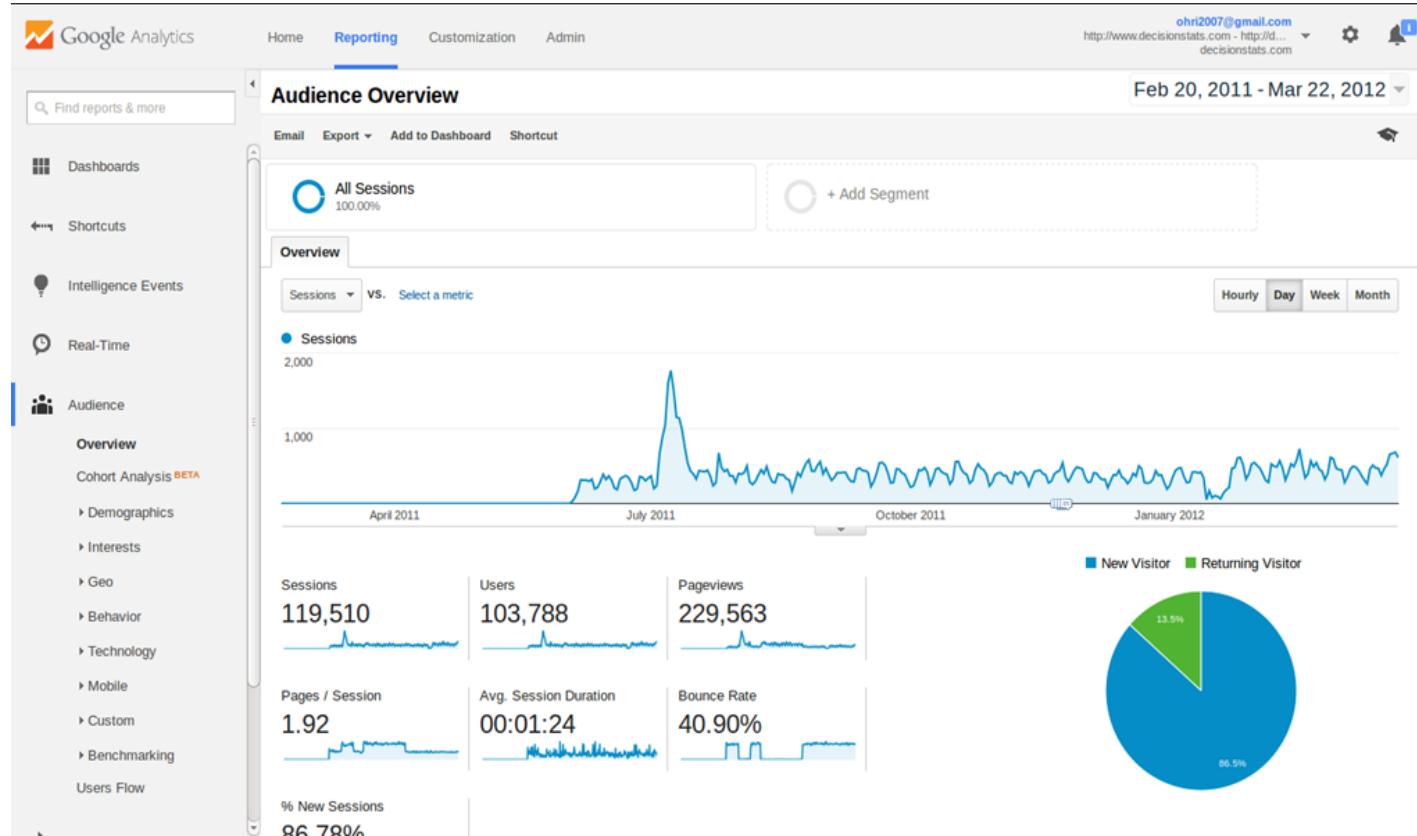
## RESEARCH OBJECTIVE

# Web Data for Time Series



[https://www.quandl.com/data/WORLDBANK/IND\\_GC\\_TAX\\_TOTL\\_GD\\_ZS-India-Tax-revenue-of-GDP](https://www.quandl.com/data/WORLDBANK/IND_GC_TAX_TOTL_GD_ZS-India-Tax-revenue-of-GDP)

# Introduction to Web Analytics



March 23, 2015, 7:50 am

DECISION STATS (WP.com)

Have you tried the upgraded stats page?

Show Me

Days Weeks Months

Views Visitors

Summaries →



122  
Visitors

195  
Views

Best ever  
2,993  
views

565,032  
views

817  
comments

#### VIEWS BY COUNTRY

Today Yesterday

Summaries →

Country	Views
India	58
United States	34
Oman	14
Australia	10
United Kingdom...	9
Finland	6



#### TOP POSTS & PAGES

Today Yesterday

Summaries →

Title	Views
Home page / Archives	41
Cricinfo StatsGuru Database for Statistical and Graphic...	11
Installing Scala on CentOS	11
Windows 7 Error : Verify that the file exists and that you ...	9
Top 10 Graphical User Interfaces in Statistical Software	8
Running R on Amazon EC2	6

# New: Cohort Analysis

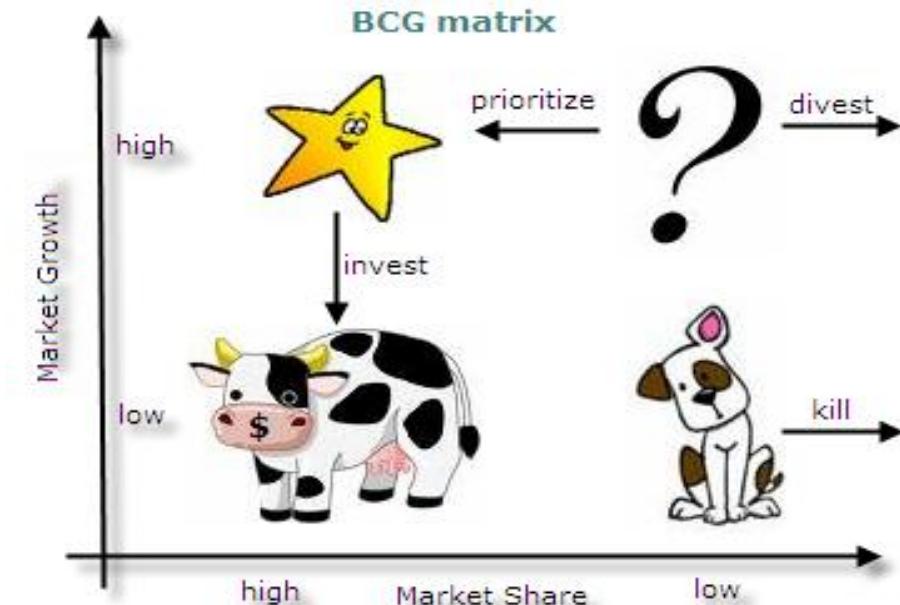
**Cohort analysis** is a subset of [behavioral analytics](#) that takes the data from a given eCommerce platform, web application, or online game and rather than looking at all users as one unit, it breaks them into related groups for analysis. These related groups, or [cohorts](#), usually share common characteristics or experiences within a defined timespan.

User signed up		User signed in by Months												
Time	People	1	2	3	4	5	6	7	8	9	10	11	12	> 12
August 2010	1,021	25.6%	6.0%	5.4%	5.8%	3.3%	2.9%	3.8%	2.9%	2.9%	1.1%	1.6%	1.9%	0.6%
September 2010	1,016	28.0%	8.1%	5.0%	5.7%	4.5%	3.7%	2.4%	3.3%	2.9%	2.2%	1.6%	0.8%	-
October 2010	973	26.6%	6.7%	4.5%	5.4%	4.6%	3.3%	3.1%	2.4%	2.6%	2.2%	0.4%	-	-
November 2010	1,386	28.2%	5.0%	5.3%	4.7%	4.4%	3.0%	3.0%	2.5%	1.7%	0.8%	-	-	-
December 2010	1,652	23.4%	6.6%	3.9%	3.5%	3.0%	2.1%	2.0%	2.0%	0.7%	-	-	-	-
January 2011	1,523	26.3%	6.6%	4.3%	3.9%	3.4%	2.2%	2.4%	0.3%	-	-	-	-	-
February 2011	1,405	28.5%	7.9%	6.5%	5.9%	3.6%	2.9%	0.9%	-	-	-	-	-	-
March 2011	1,312	30.0%	8.7%	7.2%	5.7%	4.7%	1.5%	-	-	-	-	-	-	-
April 2011	1,137	30.2%	8.6%	6.3%	5.1%	1.5%	-	-	-	-	-	-	-	-
May 2011	1,260	28.7%	7.7%	5.6%	2.5%	-	-	-	-	-	-	-	-	-
June 2011	1,155	29.2%	6.6%	2.2%	-	-	-	-	-	-	-	-	-	-
July 2011	1,003	26.5%	2.2%	-	-	-	-	-	-	-	-	-	-	-

	People	Weeks later											
		1	2	3	4	5	6	7	8	9	10	11	12
Oct 7, 2013	44	27.27%	20.45%	22.73%	18.18%	15.91%	11.36%	6.82%	13.64%	13.64%	9.09%	6.82%	2.27%
Oct 14, 2013	50	24.00%	14.00%	24.00%	14.00%	6.00%	14.00%	14.00%	12.00%	6.00%	2.00%	0.00%	-
Oct 21, 2013	49	26.53%	20.41%	16.33%	8.16%	6.12%	12.24%	12.24%	8.16%	6.12%	0.00%	-	-
Oct 28, 2013	43	16.28%	11.63%	11.63%	11.63%	11.63%	11.63%	11.63%	2.33%	0.00%	-	-	-
Nov 4, 2013	69	21.74%	11.59%	7.25%	11.59%	13.04%	5.80%	2.90%	0.00%	-	-	-	-
Nov 11, 2013	62	20.97%	14.52%	16.13%	11.29%	4.84%	3.23%	0.00%	-	-	-	-	-
Nov 18, 2013	83	13.25%	13.25%	13.25%	8.43%	1.20%	1.20%	-	-	-	-	-	-
Nov 25, 2013	74	17.57%	13.51%	8.11%	2.70%	0.00%	-	-	-	-	-	-	-
Dec 2, 2013	97	17.53%	12.37%	1.03%	1.03%	-	-	-	-	-	-	-	-
Dec 9, 2013	62	24.19%	6.45%	1.61%	-	-	-	-	-	-	-	-	-
Dec 16, 2013	40	10.00%	5.00%	-	-	-	-	-	-	-	-	-	-
Dec 23, 2013	16	6.25%	-	-	-	-	-	-	-	-	-	-	-

# Strategy Refresher- Business Strategy Models

## BCG Matrix for Product Lines



# Strategy Refresher- Business Strategy Models

## Porter's Model for Industries



# Strategy Refresher- Business Strategy Models

## Grenier's Theory



# More Strategy Anyone?

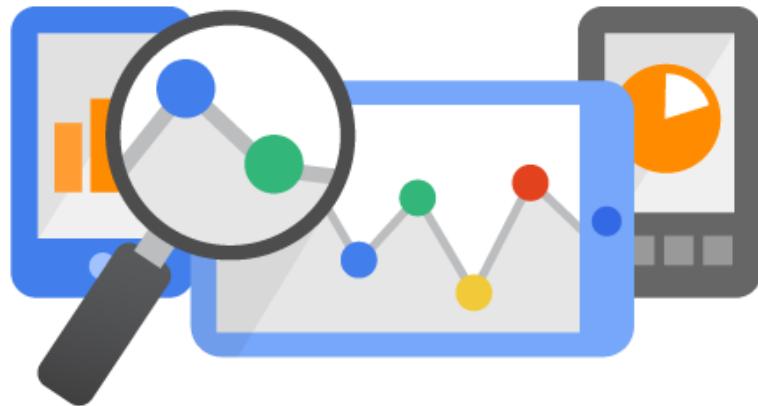
1. Porters 5 forces Model- To analyze industries
2. Business Canvas
3. BCG Matrix- To analyze Product Portfolios
4. Porters Diamond Model- To analyze locations
5. McKinsey 7 S Model- To analyze teams
6. Gernier Theory- To analyze growth of organization
7. Herzberg Hygiene Theory- To analyze soft aspects of individuals
8. Marketing Mix Model- To analyze marketing mix.

<http://decisionstats.com/2013/12/19/business-strategy-models/>

# Introduction

The R language is now the leading language for analytics and statistics on this planet. This R training starts with R language basics and covers basics of analytics

This course provides hands-on experience to execute analytics using the R language. There will be many challenging tasks and focused practicals for the learners.



# Requirements

## Installations

R

<http://cran.r-project.org/>

<http://www.rstudio.com/products/rstudio/download/>

RStudio

R Packages rattle Rcmdr Deducer

# Episode 1

# Learning Objectives

- learn about R
- install R and it's packages

# What will you learn from this lesson

- Installation of R, Rtools, R Studio, R packages and GUIs
- Using RStudio and Using GUIs

# Data Driven Decision Making

- using data and trending historical data
- validating assumptions if any
- using champion challenger to test scenarios
- using experiments
- use baselines
- continuous improvement
  - customer experiences
  - costs
  - revenues

If you can't measure it, you can't manage it -Peter Drucker

# What is R



## The R environment

R is an integrated suite of software facilities for data manipulation, calculation and graphical display. It includes

- an effective data handling and storage facility,
- a suite of operators for calculations on arrays, in particular matrices,
- a large, coherent, integrated collection of intermediate tools for data analysis,
- graphical facilities for data analysis and display either on-screen or on hardcopy, and
- a well-developed, simple and effective programming language which includes conditionals, loops, user-defined recursive functions and input and output facilities.

<http://www.r-project.org/about.html>

# Statistical Software Landscape

SAS

Python (Pandas)

IBM SPSS

R

Julia

Clojure

Octave

Matlab

JMP

E views



# Using R with other software

<https://rforanalytics.wordpress.com/useful-links-for-r/using-r-from-other-software/>

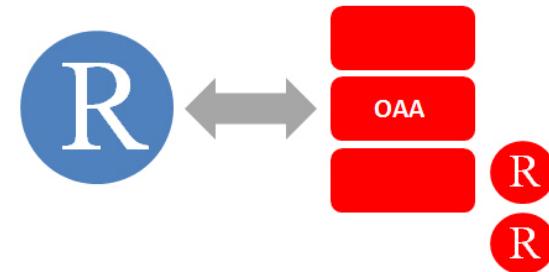
Tableau <http://www.tableausoftware.com/new-features/r-integration>

Qlik <http://qliksolutions.ru/qlikview/add-ons/r-connector-eng/>

Oracle R <http://www.oracle.com/technetwork/database/database-technologies/r/r-enterprise/overview/index.html>

Rapid Miner <https://rapid-i.com/content/view/202/206/lang,en/#>

JMP <http://blogs.sas.com/jmp/index.php?/archives/298-JMP-Into-R!.html>



# Using R with other software

<https://rforanalytics.wordpress.com/useful-links-for-r/using-r-from-other-software/>

SAS/IML <http://www.sas.com/technologies/analytics/statistics/iml/index.html>

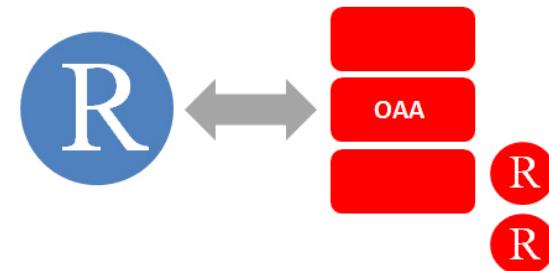
Teradata <http://developer.teradata.com/applications/articles/in-database-analytics-with-teradata-r>

Pentaho <http://bigdatatechworld.blogspot.in/2013/10/integration-of-rweka-with-pentaho-data.html>

IBM SPSS [https://www14.software.ibm.com/webapp/iwm/web/signup.do?source=ibm-](https://www14.software.ibm.com/webapp/iwm/web/signup.do?source=ibm-analytics&S_PKG=ov18855&S_TACT=M161003W&dynform=127&lang=en_US)

[analytics&S\\_PKG=ov18855&S\\_TACT=M161003W&dynform=127&lang=en\\_US](#)

TIBCO TERR <http://spotfire.tibco.com/discover-spotfire/what-does-spotfire-do/predictive-analytics/tibco-enterprise-runtime-for-r-terr>



# Some Advantages of R

open source

free

large number of algorithms and packages esp for statistics

flexible

very good for data visualization

superb community

rapidly growing

can be used with other software



# Some Disadvantages of R

- in memory (RAM) usage
- steep learning curve
- some IT departments frown on open source
- verbose documentation
- tech support
- evolving ecosystem for corporates



# Solutions for Disadvantages of R

- in memory (RAM) usage → specialized packages, in database computing
- steep learning curve → TRAINING !!!
- some IT departments frown on open source → TRAINING and education!
- verbose documentation → CRAN View , R Documentation
- tech support → expanding pool of resources
- evolving ecosystem for corporates → getting better with MS et al

# R used by Government

- In the early days of the [Deepwater Horizon disaster](#), NIST used uncertainty analysis in R to harmonize spill estimates from various sources, and to provide ranges of estimates to other agencies and the media.
- Before new drugs are allowed on the market, the FDA works with pharmaceutical companies to verify safety and efficacy through clinical trials. Despite a [false perception](#) that only commercial software may be used, many pharmaceutical companies are now using open-source R to [analyze data from clinical trials](#).
- The National Weather Service uses R for research and development of [models to predict river flooding](#).
- The newly-formed [Consumer Financial Protection Bureau](#) -- freed from the restrictions of a legacy IT infrastructure -- is championing the use of open-source technologies in government.
- Local governments are also building data-based applications. The SF Estuary Institute [uses R and Google Maps](#) to provide a [tool to track pollution](#) in the San Francisco Bay area.

[http://gsnmagazine.com/node/26483?c=cyber\\_security](http://gsnmagazine.com/node/26483?c=cyber_security)

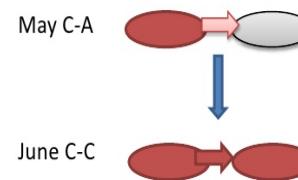
# R used by Telecom

- Churn using

## Social Network Analysis

<http://www.slideshare.net/dataspora/social-network-analysis-for-telecoms>

**Results: A Customer With a Canceller in Their Network Churns at Twice the Rate**



Types of Connections (Edges)

reality	expected by chance	delta
X	Y	2.0

In essence, we are asking whether being connected to another canceller has any effect on one's rate of cancellation. It turns out that it does.

And if we only look at voluntary port-outs, we see that customers churn at 3x the rate.

# R used by Insurance

a few more insurance related packages:

- [ChainLadder](#) – Reserving methods in R. The package provides Mack-, Munich-, Bootstrap, and Multivariate-chain-ladder methods, as well as the LDF Curve Fitting methods of Dave Clark and GLM-based reserving models.
- [cplm](#) – Monte Carlo EM algorithms and Bayesian methods for fitting Tweedie compound Poisson linear models
- [lossDev](#) – A Bayesian time series loss development model. Features include skewed-t distribution with time-varying scale parameter, Reversible Jump MCMC for determining the functional form of the consumption path, and a structural break in this path; by Christopher W. Laws and Frank A. Schmid
- [actuar](#): Loss distributions modelling, risk theory (including ruin theory), simulation of compound hierarchical models and credibility theory check out the [actuar](#) package by C. Dutang, V. Goulet and M. Pigeon.
- [favir](#): Formatted Actuarial Vignettes in R. FAViR lowers the learning curve of the R environment. It is a series of peer-reviewed Sweave papers that use a consistent style.
- [mondate](#): R packackge to keep track of dates in terms of months
- [lifecontingencies](#) – Package to perform actuarial evaluation of life contingencies

and

[Introduction to R for Actuaries](#) by Nigel de Silva

and <http://www.rininsurance.com/>

# R in Finance

R/Finance [home](#) [agenda](#) [register](#) [travel](#) [committee](#)

Friday, May 29th, 2015

08:00 - 09:00 Optional Pre-Conference Tutorials

Ross Bennett: PortfolioAnalytics: Advanced Moment Estimation & Optimization ([pdf](#))

Kris Boudt: High-frequency Price Data Analysis in R ([pdf](#))

Dirk Eddelbuettel: Hands-on Introduction to Rcpp ([pdf](#))

Guy Yollin: Getting Started with Quantstrat

Maria Belianina: An Introduction to OneTick

09:00 - 09:30 Registration (2nd floor Inner Circle) & Continental Breakfast (3rd floor by Sponsor Tables)

Transition between seminars

09:30 - 09:35 Kickoff

09:35 - 09:40 Sponsor Introduction

09:40 - 10:30 **Emanuel Derman:** Understanding the World

10:30 - 10:54 **John Burkett:** Portfolio Optimization: Price Predictability, Utility Functions, Computational Methods, and Applications ([pdf](#))

**Kyle Balkissoon:** A Framework for Integrating Portfolio-level Backtesting with Price and Quantity Information ([html](#))

**Anthony Tsou:** Implementation of Quality Minus Junk

**Ilya Kipnis:** Flexible Asset Allocation With Stepwise Correlation Rank ([pptm](#))

10:54 - 11:20 Break

11:20 - 11:40 **Sanjiv Das:** Efficient Rebalancing of Taxable Portfolios ([pdf](#))

11:40 - 12:00 **Marjan Wauters:** Characteristic-based equity portfolios: economic value and dynamic style allocation ([pdf](#))

12:00 - 12:20 **Bernhard Pfaff:** The sequel of cccp: Solving cone constrained convex programs

12:20 - 13:40 Lunch

13:40 - 14:00 **Markus Gesmann:** Communicating risk - a perspective from an insurer ([pdf](#))

14:00 - 14:20 **Doug Martin:** Nonparametric vs Parametric Shortfall: What are the Differences?

<http://www.rinfinance.com/>

# R in Finance

<http://cran.r-project.org/web/views/Finance.html>

This CRAN Task View contains a list of packages useful for empirical work in Finance, grouped by topic.

- The Rmetrics suite of packages comprises [fArma](#), [fAsianOptions](#), [fAssets](#), [fBasics](#), [fBonds](#), [timeDate](#) (formerly: [fCalendar](#)), [fCopulae](#), [fExoticOptions](#), [fExtremes](#), [fGarch](#), [fImport](#), [fNonlinear](#), [fOptions](#), [fPortfolio](#), [fRegression](#), [timeSeries](#) (formerly: [fSeries](#)), [fTrading](#), [fUnitRoots](#) and contains a very large number of relevant functions for different aspect of empirical and computational finance.
- The [RQuantLib](#) package provides several option-pricing functions as well as some fixed-income functionality from the QuantLib project to R.
- The [quantmod](#) package offers a number of functions for quantitative modelling in finance as well as data acquisition, plotting and other utilities.
- The [portfolio](#) package contains classes for equity portfolio management; the [portfolioSim](#) builds a related simulation framework. The [backtest](#) offers tools to explore portfolio-based hypotheses about financial instruments. The [stockPortfolio](#) package provides functions for single index, constant correlation and multigroup models. The [pa](#) package offers performance attribution functionality for equity portfolios.
- The [PerformanceAnalytics](#) package contains a large number of functions for portfolio performance calculations and risk management.

# R in Finance

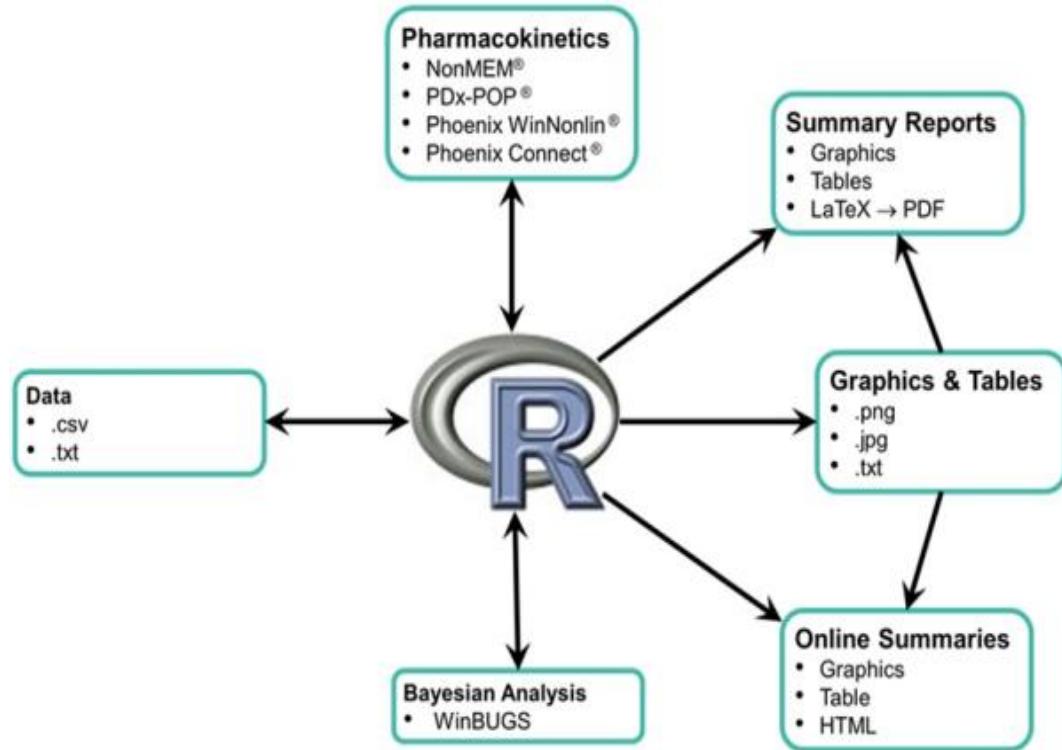
<http://blog.revolutionanalytics.com/2013/08/r-drug-development-and-the-fda.html>

*Opening the Doors to Open Source Programming in Drug Development.*

*R: Regulatory Compliance and Validation Issues A Guidance Document for the Use of R in Regulated Clinical Trial Environments* in which he concluded that useR 2012 FDA statistician Jea Brodsky presented a [poster](#) described how FDA scientists “use R on a daily basis” and have themselves written R packages for use at various stages in the drug submission process.

*Open Source Software in the Biopharma Industry: Challenges and Opportunities,*

# R in Pharma



<http://web.quanticate.com/bid/102741/Using-the-Statistical-Programming-Language-R-in-the-Pharma-Industry>

# R in Pharma

<http://cran.r-project.org/web/views/ClinicalTrials.html>

This task view gathers information on specific R packages for design, monitoring and analysis of data from clinical trials. It focuses on including packages for clinical trial design and monitoring in general plus data analysis packages for a specific type of design.

## Design and Monitoring

- **TrialSize** This package has more than 80 functions from the book *Sample Size Calculations in Clinical Research* (Chow & Wang & Shao, 2007, 2nd ed., Chapman & Hall/CRC).
- **ad** This Package contains calculations for adaptive seamless designs using early outcomes for treatment selection.
- **bcopt** This package implements a wide variety of one- and two-parameter Bayesian CRM designs. The program can run interactively, allowing the user to enter outcomes after each cohort has been recruited, or via simulation to assess operating characteristics.
- **blockrand** creates randomizations for block random clinical trials. It can also produce a PDF file of randomization cards.
- **conf.design** This small package contains a series of simple tools for constructing and manipulating confounded and fractional factorial designs.
- **CRTSize** This package contains basic tools for the purpose of sample size estimation in cluster (group) randomized trials. The package contains traditional power-based methods, empirical smoothing (Rotondi and Donner, 2009), and updated meta-analysis techniques (Rotondi and Donner, 2011).
- **dfrm** This package provides functions to run the CRM and TITE-CRM in phase I trials and calibration tools for trial planning purposes.
- **experiment** contains tools for clinical experiments, e.g., a randomization tool, and it provides a few special analysis options for clinical trials.
- **Efr2** This package creates regular and non-regular Fractional Factorial designs. Furthermore, analysis tools for Fractional Factorial designs with 2-level factors are offered (main effects and interaction plots for all factors simultaneously, cube plot for looking at the simultaneous effects of three factors, full or half normal plot, alias structure in a more readable format than with the built-in function alias). The package is currently subject to intensive development. While much of the intended functionality is already available, some changes and improvements are still to be expected.
- **GroupSeq** performs computations related to group sequential designs via the alpha spending approach, i.e., interim analyses need not be equally spaced, and their number need not be specified in advance.
- **gsDesign** derives group sequential designs and describes their properties.
- **ld98and** from **Hmisc** computes and plots group sequential stopping boundaries from the Lan-DeMets method with a variety of  $\alpha$ -spending functions using the ld98 program from the Department of Biostatistics, University of Wisconsin written by DM Reboussin, DL DeMets, KM Kim, and KKG Lan.
- **ldbounds** uses Lan-DeMets Method for group sequential trial; its functions calculate bounds and probabilities of a group sequential trial.
- **longpower** The longpower package contains functions for computing power and sample size for linear models of longitudinal data based on the formula due to Liu and Liang (1997) and Diggle et al (2002). Either formula is expressed in terms of marginal model or Generalized Estimating Equations (GEE) parameters. This package contains functions which translate pilot mixed effect model parameters (e.g. random intercept and/or slope) into marginal model parameters so that the formulas of Diggle et al or Liu and Liang formula can be applied to produce sample size calculations for two sample longitudinal designs assuming known variance.
- **PPS** generates predicted interval plots, simulates and plots confidence intervals of an effect estimate given observed data and a hypothesis about the distribution of future data.
- **PowerTOST** contains functions to calculate power and sample size for various study designs used for bioequivalence studies. See function `known.designs()` for study designs covered. Moreover the package contains functions for power and sample size based on 'expected' power in case of uncertain (estimated) variability. Added are functions for the power and sample size for the ratio of two means with normally distributed data on the original scale (based on Fieller's confidence ('fiducial') interval).
- **pwr** has power analysis functions along the lines of Cohen (1988).
- **PvcrGSD** is a set of tools to compute power in a group sequential design.
- **qtlDesign** provides tools for the design of QTL experiments.
- **seqmon** is computes the probability of crossing sequential efficacy and futility boundaries in a clinical trial. It implements the Armitage-McPherson and Rowe Algorithm using the method described in Schoenfeld (2001).

## Design and Analysis

- Package **AGSDest** This package provides tools and functions for parameter estimation in adaptive group sequential trials.
- Package **clinfun** has functions for both design and analysis of clinical trials. For phase II trials, it has functions to calculate sample size, effect size, and power based on Fisher's exact test, the operating characteristics of a two-stage boundary, Optimal and Minimax 2-stage Phase II designs given by Richard Simon, the exact 1-stage Phase II design and can compute a stopping rule and its operating characteristics for toxicity monitoring based repeated significance testing. For phase III trials, it can calculate sample size for group sequential designs.

# Companies using R

from <http://www.revolutionanalytics.com/companies-using-r>

ANZ, the fourth largest bank in Australia, using R for credit risk analysis

Bank of America uses R for reporting.

The Consumer Financial Protection Bureau uses R for data analysis.

## Facebook

Facebook and R:

- Analysis of Facebook Status Updates
- Facebook's Social Network Graph
- How Google and Facebook are using R
- Predicting Colleague Interactions with R

# Pre Requisites

- Installation of R

<http://cran.rstudio.com/bin/windows/base/>



[CRAN  
Mirrors](#)  
[What's new?](#)  
[Task Views](#)  
[Search](#)

[About R](#)  
[R Homepage](#)  
[The R Journal](#)

[Software](#)  
[R Sources](#)  
[R Binaries](#)  
[Packages](#)  
[Other](#)

[Documentation](#)  
[Manuals](#)  
[FAQs](#)  
[Contributed](#)

R-3.1.1 for Windows (32/64 bit)

[Download R 3.1.1 for Windows](#) (54 megabytes, 32/64 bit)  
[Installation and other instructions](#)  
[New features in this version](#)

If you want to double-check that the package you have downloaded exactly matches the package distributed by R, you can compare the [md5sum](#) of the .exe to the [true fingerprint](#). You will need a version of md5sum for windows: both [graphical](#) and [command line versions](#) are available.

## Frequently asked questions

- [How do I install R when using Windows Vista?](#)
- [How do I update packages in my previous version of R?](#)
- [Should I run 32-bit or 64-bit R?](#)

Please see the [R FAQ](#) for general information about R and the [R Windows FAQ](#) for Windows-specific information.

## Other builds

- Patches to this release are incorporated in the [r-patched.snapshot.build](#)
- A build of the development version (which will eventually become the next major release of R) is available in the [r-devel.snapshot.build](#)
- [Previous releases](#)

Note to webmasters: A stable link which will redirect to the current Windows binary release is  
[CRAN MIRROR->bin/windows/base/release.htm](#).

Last change: 2014-07-10, by Duncan Murdoch

- R Studio

- R Packages

# Pre Requisites

- Installation of R
  - Rtools
- R Studio
- R Packages

<http://cran.rstudio.com/bin/windows/Rtools/>



[CRAN](#)  
[Mirror](#)  
[What's new?](#)  
[Task Views](#)  
[Search](#)

[About R](#)  
[R Home Page](#)  
[The R Journal](#)

[Software](#)  
[R Sources](#)  
[R-Banners](#)  
[Packages](#)  
[Other](#)

[Documentation](#)  
[Manuals](#)  
[FAQs](#)  
[Contributed](#)

## Building R for Windows

This document is a collection of resources for building packages for R under Microsoft Windows, or for building R itself (version 1.9.0 or later). The original collection was put together by Prof. Brian Ripley; it is currently being maintained by Duncan Murdoch.

The authoritative source of information for tools to work with the current release of R is the "R Administration and Installation" manual. In particular, please read the ["Windows Tools" appendix](#).

### Tools Downloads

With the change to gcc 4.2.1, some of the tools for 32 bit compiles became incompatible with obsolete versions of R. Since then we have been maintaining one actively updated version of the tools, and other "frozen" snapshots of them. We recommend that users use the latest release of Rtools with the latest release of R.

The current version of this file is recorded here: [VERSION.txt](#)

Download	R compatibility	Frozen?
<a href="#">Rtools31.exe</a>	R 3.0.x to 3.1.x	No
<a href="#">Rtools30.exe</a>	R >= 2.15.1 to R 3.0.x	Yes
<a href="#">Rtools215.exe</a>	R >= 2.14.1 to R 2.15.1	Yes
<a href="#">Rtools214.exe</a>	R 2.13.x or R 2.14.x	Yes
<a href="#">Rtools213.exe</a>	R 2.13.x	Yes
<a href="#">Rtools212.exe</a>	R 2.12.x	Yes
<a href="#">Rtools211.exe</a>	R 2.10.x or R 2.11.x	Yes
<a href="#">Rtools210.exe</a>	R 2.9.x or 2.10.x	Yes
<a href="#">Rtools29.exe</a>	R 2.8.x or R 2.9.x	Yes
<a href="#">Rtools28.exe</a>	R 2.7.x or R 2.8.x	Yes
<a href="#">Rtools27.exe</a>	R 2.6.x or R 2.7.x	Yes
<a href="#">Rtools26.exe</a>	R 2.6.x, R 2.5.x or (untested) earlier	Yes

The change history to the Rtools is [below](#).

### Tools for 64 bit Windows builds

Rtools 2.12 and later include both 32 bit and 64 bit tools.

# Pre Requisites

- Installation of R
  - RTools



Products Resources Pricing About Us Blog



Download RStudio

Home / Overview / RStudio / Download RStudio

RStudio is a set of integrated tools designed to help you be more productive with R. It includes a console, syntax-highlighting editor that supports direct code execution, as well as tools for plotting, history, debugging and workspace management.

If you run R on a Linux server and want to enable users to remotely access RStudio using a web browser please download RStudio Server.

**Do you need support or a commercial license?**  
Check out our commercial offerings

**Download RStudio Desktop v0.98.1074 — Release Notes**

RStudio requires R 2.11.1 (or higher). If you don't already have R, you can download it here.

Let's stay in touch. Give us your email and we'll keep you in the loop.

\* Email  Submit

## Installers for ALL Platforms

Installers	Size	Date	MD5
RStudio 0.98.1074 - Windows XP/Vista/7/8	45 MB	2011-10-14	74d7bc76ec04287fac79cdada5dfaa8dd
RStudio 0.98.1074 - Mac OS X 10.6+ (64-bit)	38.4 MB	2014-10-14	f01c43f29af679400c0faeae7ee33fbfc
RStudio 0.98.1074 - Debian 6+ (Ubuntu 10.04+) (32-bit)	54.3 MB	2014-10-14	759d865a59b22b28202a5a0025e77278
RStudio 0.98.1074 - Debian 6+ (Ubuntu 10.04+) (64-bit)	86.1 MB	2014-10-14	077a31c714a7df2afffr3r0573d044ah8

- R Packages

# Pre Requisites

- Installation of R

- Rtools

- R Studio

<http://www.rstudio.com/products/rstudio/download/>

- R Packages

about eight packages supplied with the R distribution and many more are available through the CRAN family of Internet sites covering a very wide range of modern statistics.

The screenshot shows the R Documentation homepage. At the top, there's a search bar and navigation links for Discussion, About, and Rd documentation package. Below that is a section for Top Ranked CRAN Packages, showing a list of packages with their names, package IDs, and download counts. The main area features logos for CRAN, Bioconductor, and GitHub, followed by a search form with fields for All Fields, Package Name, Function Name, Title, Description, and Author(s), and a green Start search button. A note at the bottom says "Rdocumentation is a tool that helps you easily find and browse the documentation of all current and some past packages on CRAN. Click on the search bar at the top left for instant search or fill out the forms below for advanced search!"

#	Package	#
1	Rcpp	80382
2	ggplot2	69504
3	plyr	65837
4	stringr	65371
5	digest	63067
6	RColorBrewer	57606
7	reshape2	57236
8	colorspace	51693
9	labeling	49615
10	scales	47407

107 sites in 49 regions



*CRAN*  
[Mirrors](#)  
[What's new?](#)  
[Task Views](#)  
[Search](#)

*About R*  
[R Homepage](#)  
[The R Journal](#)

*Software*  
[R Sources](#)  
[R Binaries](#)  
[Packages](#)  
[Other](#)

*Documentation*  
[Manuals](#)  
[FAQs](#)  
[Contributed](#)

### CRAN Mirrors

The Comprehensive R Archive Network is available at the following URLs, please choose a location close to you. Some statistics on the status of the mirrors can be found here: [main page](#), [windows release](#), [windows old release](#).

O-Cloud	<a href="http://cran.rstudio.com/">http://cran.rstudio.com/</a>	Rstudio, automatic redirection to servers worldwide	
Algeria	<a href="http://cran.usthb.dz/">http://cran.usthb.dz/</a>	University of Science and Technology Houari Boumediene	
Argentina	<a href="http://mirror.fcaglp.unlp.edu.ar/CRAN/">http://mirror.fcaglp.unlp.edu.ar/CRAN/</a>	Universidad Nacional de La Plata	
Australia	<a href="http://cran.csiro.au/">http://cran.csiro.au/</a> <a href="http://cran.ms.unimelb.edu.au/">http://cran.ms.unimelb.edu.au/</a>	CSIRO University of Melbourne	
Austria	<a href="http://cran.at.r-project.org/">http://cran.at.r-project.org/</a>	Wirtschaftsuniversitaet Wien	
Belgium	<a href="http://www.freestatistics.org/cran/">http://www.freestatistics.org/cran/</a>	K.U.Leuven Association	
Brazil	<a href="http://nbegib.uesc.br/mirrors/cran/">http://nbegib.uesc.br/mirrors/cran/</a> <a href="http://cran-r.cslf.ufpr.br/">http://cran-r.cslf.ufpr.br/</a> <a href="http://cran.fiocruz.br/">http://cran.fiocruz.br/</a> <a href="http://www.vps.fmvz.usp.br/CRAN/">http://www.vps.fmvz.usp.br/CRAN/</a> <a href="http://brieger.esalq.usp.br/CRAN/">http://brieger.esalq.usp.br/CRAN/</a>	Center for Comp. Biol. at Universidade Estadual de Santa Cruz Universidade Federal do Parana Oswaldo Cruz Foundation, Rio de Janeiro University of Sao Paulo, Sao Paulo University of Sao Paulo, Piracicaba	
Canada	<a href="http://cran.stat.sfu.ca/">http://cran.stat.sfu.ca/</a> <a href="http://mirror.its.dal.ca/cran/">http://mirror.its.dal.ca/cran/</a> <a href="http://cran.utsstat.utoronto.ca/">http://cran.utsstat.utoronto.ca/</a> <a href="http://cran.skazkaforyou.com/">http://cran.skazkaforyou.com/</a> <a href="http://cran.parentingamerica.com/">http://cran.parentingamerica.com/</a>	Simon Fraser University, Burnaby Dalhousie University, Halifax University of Toronto iWeb, Montreal iWeb, Montreal	
Chile	<a href="http://dirichlet.mat.puc.cl/">http://dirichlet.mat.puc.cl/</a>	Pontificia Universidad Catolica de Chile, Santiago	
China	<a href="http://ftp.ctex.org/mirrors/CRAN/">http://ftp.ctex.org/mirrors/CRAN/</a> <a href="http://mirror.bjtu.edu.cn/cran/">http://mirror.bjtu.edu.cn/cran/</a> <a href="http://mirrors.opencas.cn/cran/">http://mirrors.opencas.cn/cran/</a>	CTEX.ORG Beijing Jiaotong University, Beijing Chinese Academy of Sciences, Beijing	

# Non CRAN Repositories

<http://www.rdocumentation.org/>

The screenshot shows the RDocumentation website interface. At the top, there is a search bar with the placeholder "Start searching the documentation". Below it is a "TASK VIEWS" sidebar containing a list of R package categories: Bayesian, ChemPhys, ClinicalTrials, Cluster, DifferentialEquations, Distributions, Econometrics, Environmetrics, ExperimentalDesign, Finance, Genetics, gR, Graphics, HighPerformanceComputing, MachineLearning, MedicalImaging, MetaAnalysis, Multivariate, NaturalLanguageProcessing, NumericalMathematics, OfficialStatistics, Optimization, Pharmacokinetics, Phylogenetics, Psychometrics, ReproducibleResearch, Robust, SocialSciences, Spatial, SpatioTemporal, and TimeSeries.

The main content area is titled "Documentation" and features a search bar stating "Search the R documentation of 7393 R packages and 150600 R functions:". Below the search bar is a descriptive text: "Rdocumentation is a tool that helps you easily find and browse the documentation of all current and some past packages on CRAN. Click on the search bar at the top left for instant search or fill out the forms below for advanced search!". There are five input fields for "All Fields", "Package Name", "Function Name", "Title", and "Description", followed by an "Author(s)" field. A large green "Start search" button is located at the bottom of these fields. To the right of the search form, there is a sidebar for DataCamp, which includes a logo, the text "Learn Data Science with R", a price of "\$25/month", a thumbnail of a course featuring a bar chart and a gold medal, and a "Discover All Courses" button. The sidebar also mentions "Data Manipulation, Data Visualization, R Programming, Big Data, and much more".

# github

The screenshot shows the GitHub homepage with the search bar at the top containing "Search GitHub". Below the search bar are navigation links for "Explore", "Gist", "Blog", and "Help". On the right side, there are icons for "decisionstats" and other user profiles. The main content area is titled "Explore GitHub" and features a "Trending" tab selected, indicated by an orange underline. Below this, the heading "Trending repositories" is displayed, followed by the subtext "Find what repositories the GitHub community is most excited about today." A horizontal navigation bar below the heading includes tabs for "Repositories" (selected), "Developers", and a dropdown menu set to "Trending: today". To the right of this bar is a sidebar with language filters: "All languages", "Unknown languages", "C", "C++", "HTML", "Java", "JavaScript", "Python", and "R" (which is highlighted with a blue background). Below the sidebar is a "ProTip!" box with the text "Looking for most forked R repositories? Try this search". The main content area lists four trending repositories:

- rdpeng/ProgrammingAssignment2**  
Repository for Programming Assignment 2 for R Programming on Coursera  
R • Built by 2
- qinwf/awesome-R**  
A curated list of awesome R frameworks, packages and software.  
R • 7 stars today • Built by 2
- berndbischl/mlr**  
mlr: Machine Learning in R  
R • Built by 2
- rstudio/shinyapps**  
R • Built by 2

# bioconductor

<http://www.bioconductor.org/>



Home

Install

Help

Developers

About

Search:

## BioC2015

Join us for morning talks from distinguished speakers and community members, afternoon workshops to hone your skills, and poster sessions and social activities to get to know members of the Bioconductor community at our [Annual Conference](#), July 20 (Developer Day), 21 and 22 in Seattle, WA.

## About Bioconductor

Bioconductor provides tools for the analysis and comprehension of high-throughput genomic data. Bioconductor uses the R statistical programming language, and is open source and open development. It has two releases each year, [1024 software packages](#), and an active user community. Bioconductor is also available as an [AMI](#) (Amazon Machine Image) and a series of [Docker](#) images.

## News

- Bioconductor [3.1](#) is available.
- Orchestrating high-throughput genomic analysis with [Bioconductor](#) ([abstract](#)) and other [recent literature](#).
- Read our latest [newsletter](#) and [course](#)

## Install »

Get started with *Bioconductor*

- [Install Bioconductor](#)
- [Explore packages](#)
- [Get support](#)
- [Latest newsletter](#)
- [Follow us on twitter](#)
- [Install R](#)

## Learn »

Master *Bioconductor* tools

- [Courses](#)
- [Support site](#)
- [Package vignettes](#)
- [Literature citations](#)
- [Common work flows](#)
- [FAQ](#)
- [Community resources](#)
- [Videos](#)

## Use »

Create bioinformatic solutions with *Bioconductor*

- [Software, Annotation, and Experiment packages](#)
- [Amazon Machine Image](#)
- [Latest release announcement](#)
- [Support site](#)

## Develop »

Contribute to *Bioconductor*

- [Use Bioc 'devel'](#)
- ['Devel' Software, Annotation and Experiment packages](#)
- [Package guidelines](#)
- [New package submission](#)
- [Developer resources](#)
- [Build reports](#)

# Pre Requisites

- R Packages

```
install.packages() INSTALLS  
update.packages() UPDATES  
library() LOADS
```

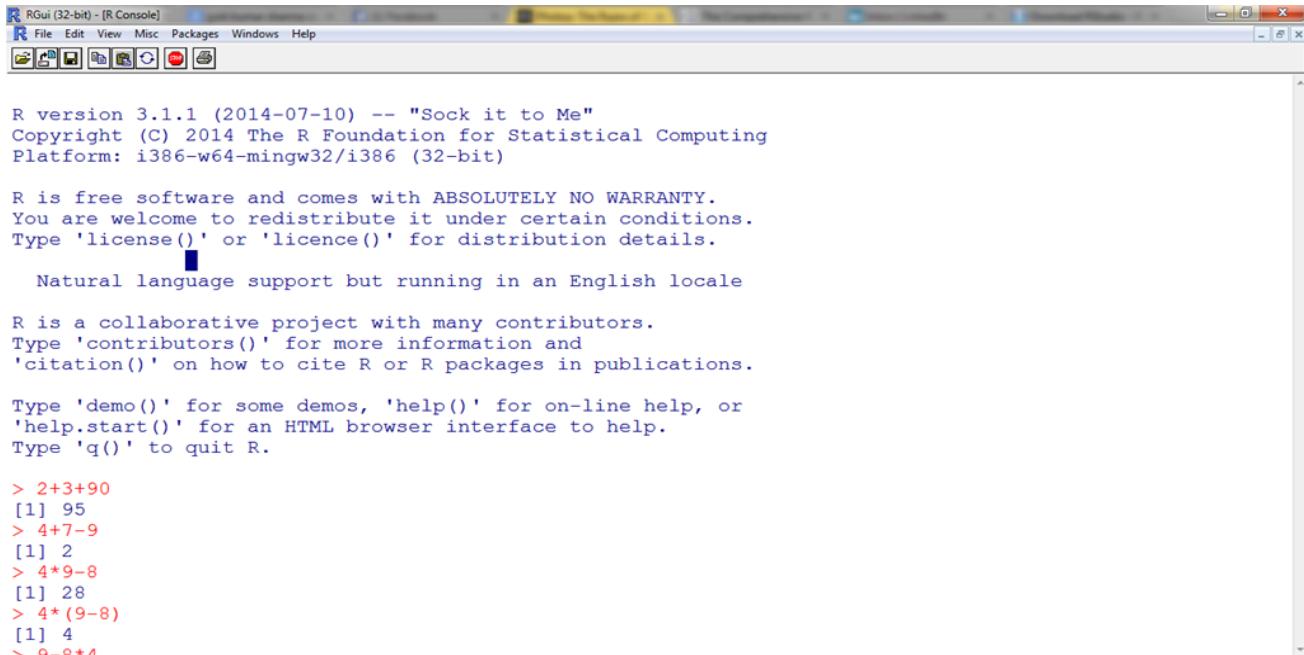
- Packages are **installed** once, updated periodically, but **loaded** every time

# Interfaces to R

- Console

Default

Customization



```
R version 3.1.1 (2014-07-10) -- "Sock it to Me"
Copyright (C) 2014 The R Foundation for Statistical Computing
Platform: i386-w64-mingw32/i386 (32-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

[REDACTED]

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

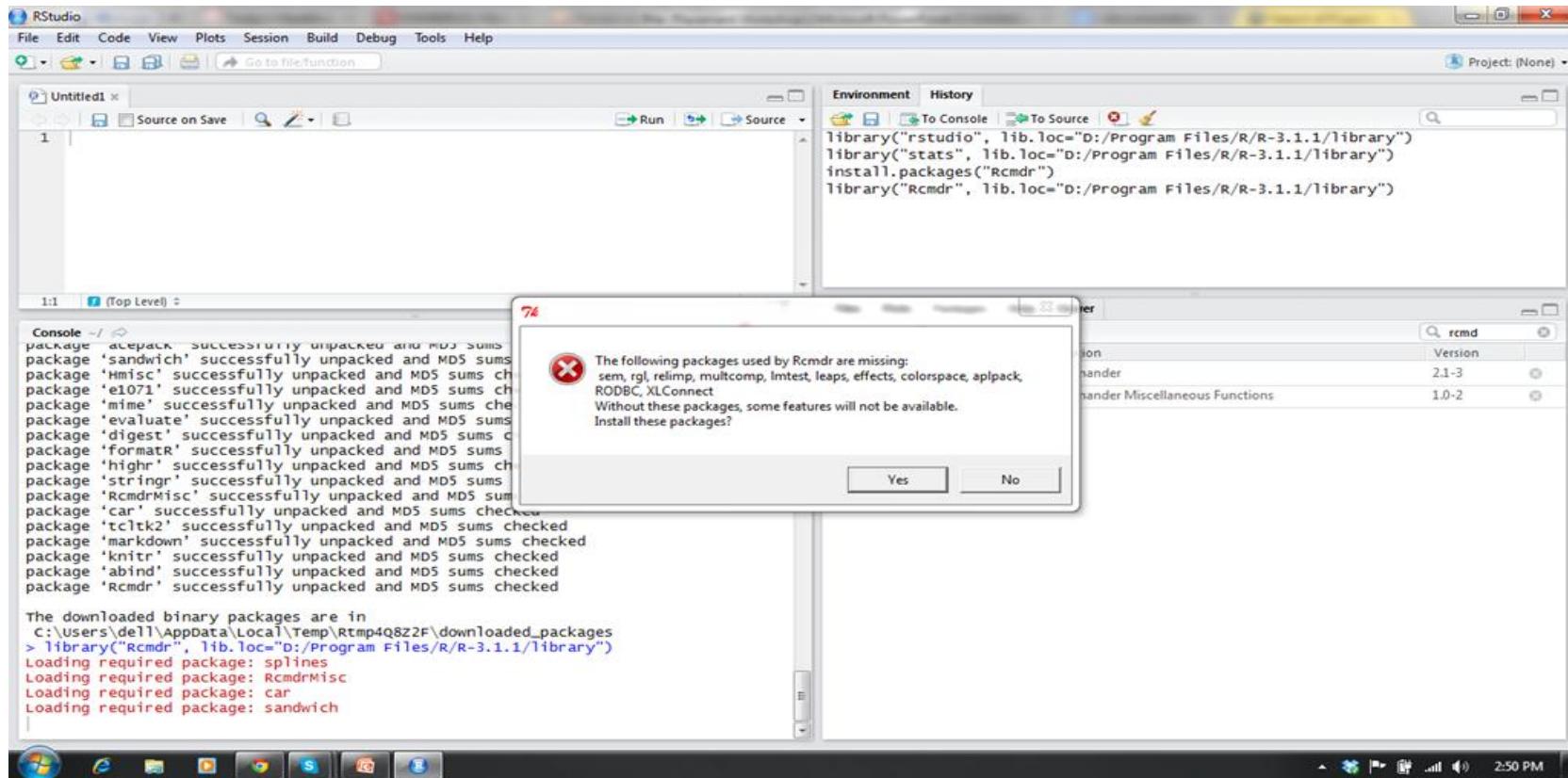
> 2+3+90
[1] 95
> 4+7-9
[1] 2
> 4*9-8
[1] 28
> 4*(9-8)
[1] 4
> q_*8*4
```

- GUI

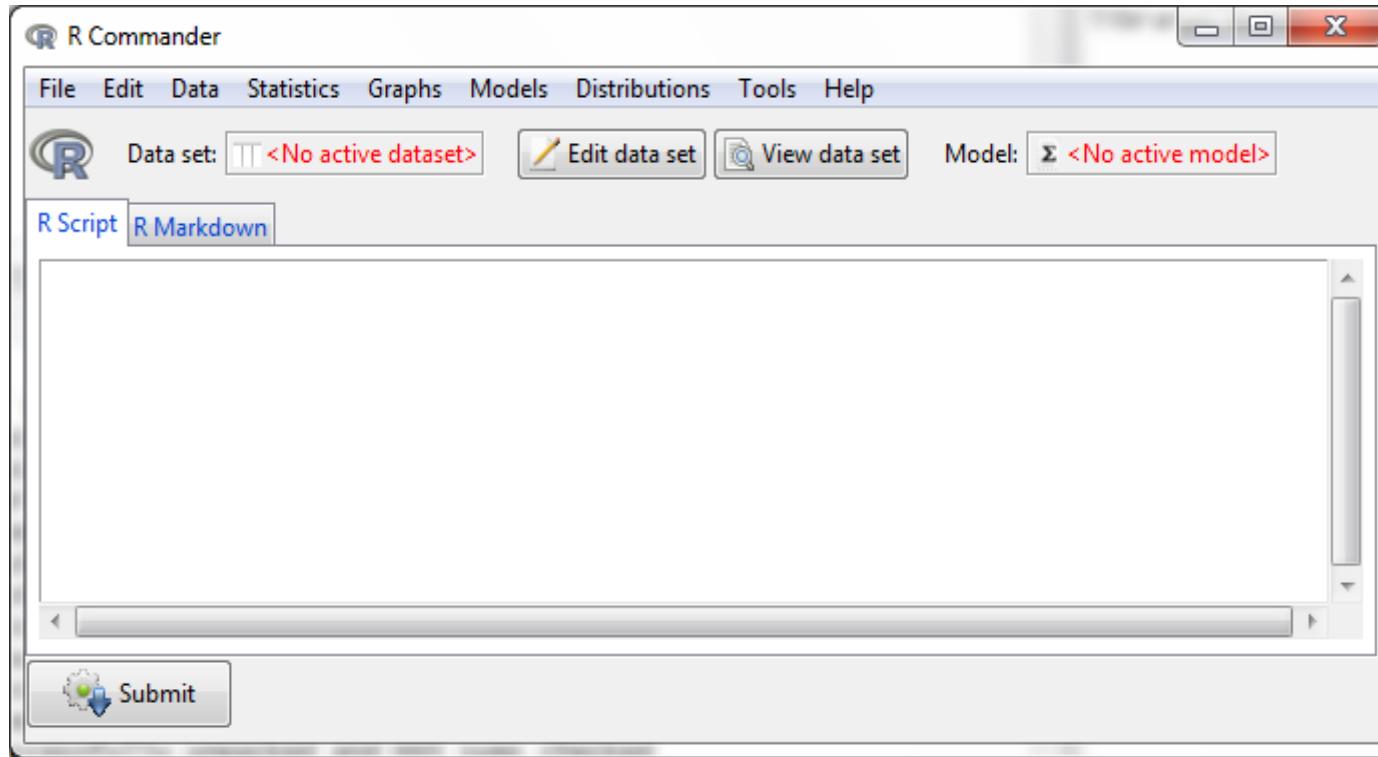
# Graphical Interfaces to R

- R Commander
- Rattle
- Deducer

# Installation of R Commander

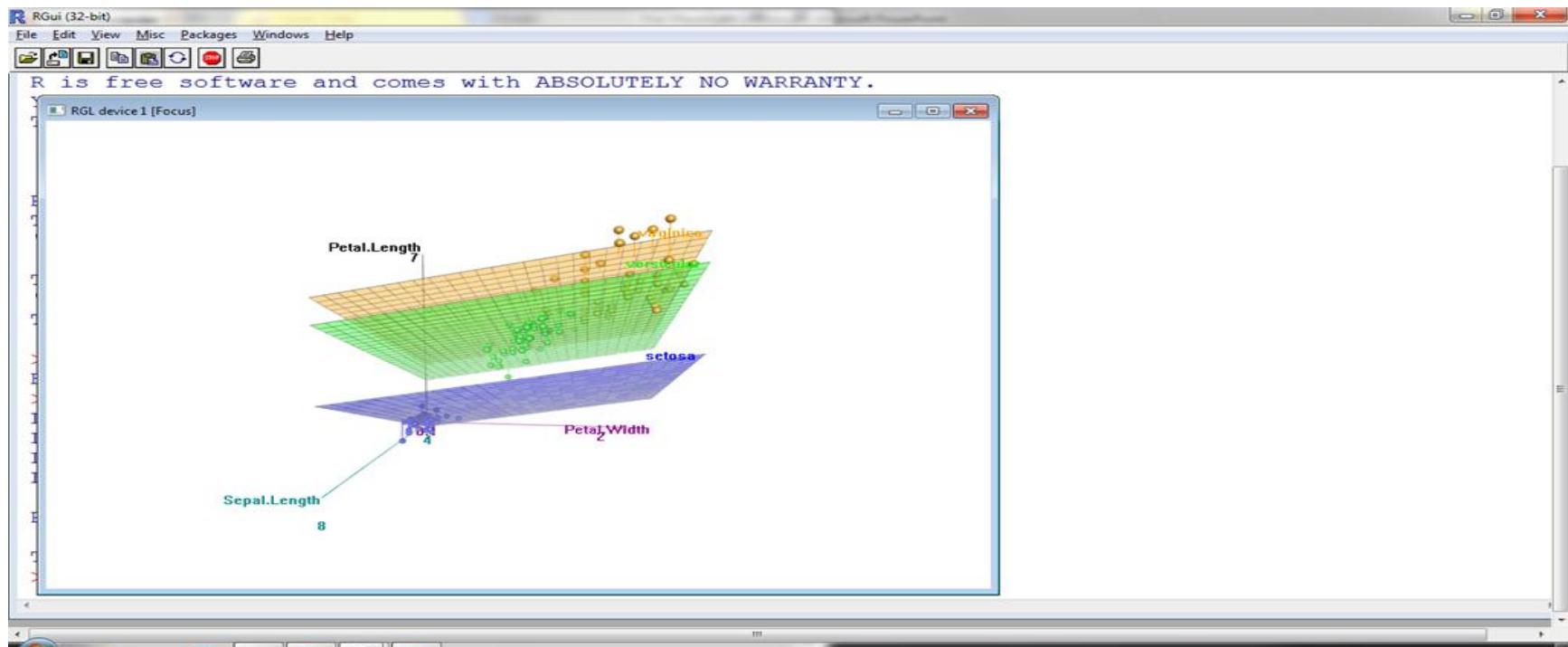


# Overview of R Commander

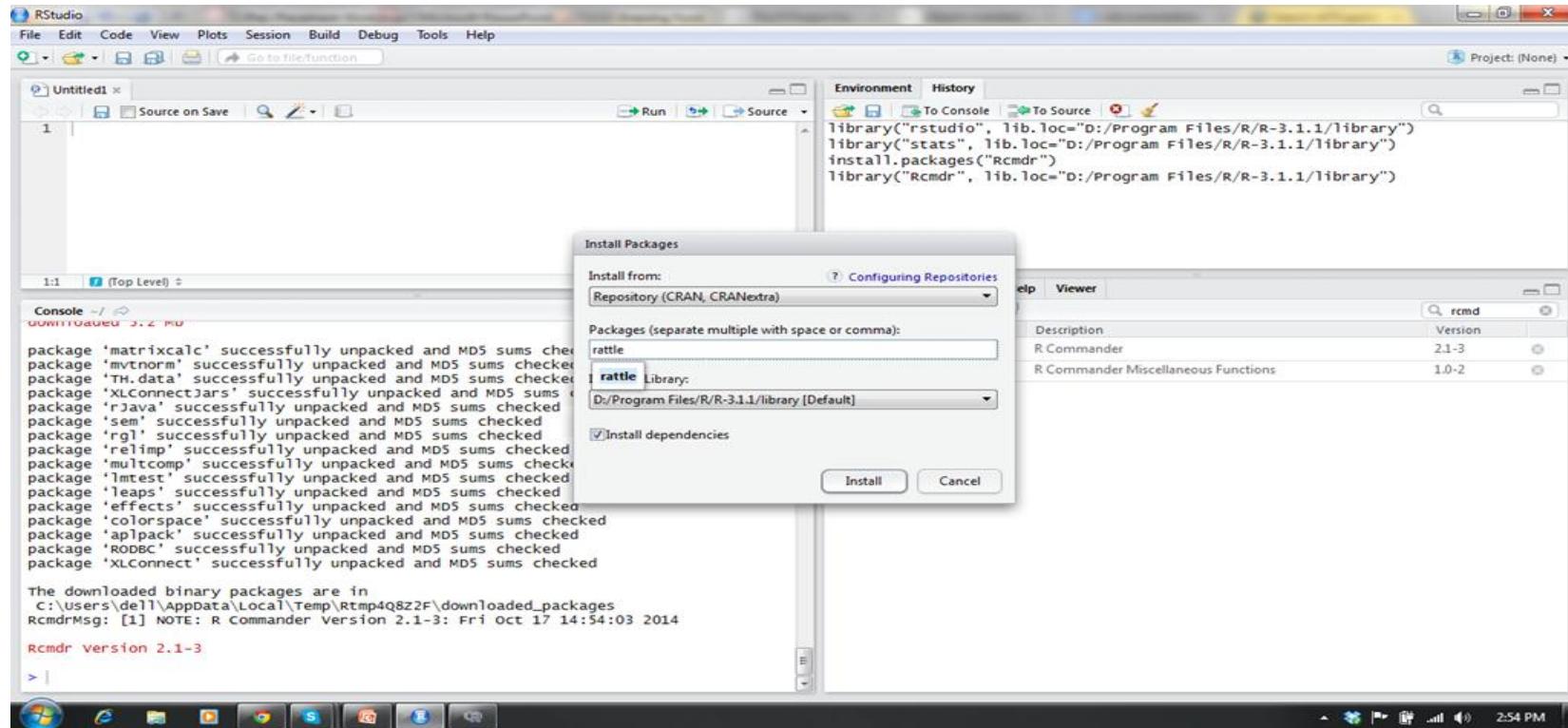


# Demo

## R Commander – 3D Graphs



# Installation of Rattle



# Installation of Rattle

The screenshot shows the RStudio interface with the following components visible:

- File Menu:** File, Edit, Code, View, Plots, Session, Build, Debug, Tools, Help.
- Project Bar:** Untitled1 x, Source on Save, Go to file/function, Project: (None).
- Environment Tab:** Shows library code being run:

```
library("rstudio", lib.loc="D:/Program Files/R/R-3.1.1/library")
library("stats", lib.loc="D:/Program Files/R/R-3.1.1/library")
install.packages("Rcmdr")
library("Rcmdr", lib.loc="D:/Program Files/R/R-3.1.1/library")
install.packages("rattle")
library("rattle", lib.loc="D:/Program Files/R/R-3.1.1/library")
```
- Console Tab:** Displays the output of the R code:

```
package 'leaps' successfully unpacked and MD5 sums checked
package 'leaps' successfully unpacked and MD5 sums checked
package 'effects' successfully unpacked and MD5 sums checked
package 'colorspace' successfully unpacked and MD5 sums checked
package 'alppack' successfully unpacked and MD5 sums checked
package 'RODBC' successfully unpacked and MD5 sums checked
package 'XLConnect' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
C:/Users/dell/AppData/Local/Temp/Rtmp4Q8Z2F/downloaded_packages
RcmdrMsg: [1] NOTE: R Commander version 2.1-3: Fri Oct 17 14:54:03 2014

Rcmdr version 2.1-3

> install.packages("rattle")
trying URL 'http://cran.rstudio.com/bin/windows/contrib/3.1/rattle_3.3.0.zip'
Content type 'application/zip' length 3211375 bytes (3.1 Mb)
opened URL
downloaded 3.1 Mb

package 'rattle' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
C:/Users/dell/AppData/Local/Temp/Rtmp4Q8Z2F/downloaded_packages
> library("rattle", lib.loc="D:/Program Files/R/R-3.1.1/library")
```
- Packages Tab:** Shows a list of installed packages:

Name	Description	Version
multcomp	Simultaneous Inference in General Parametric Models	1.3-7
mvtnorm	Multivariate Normal and t Distributions	1.0-0
nlme	Linear and Nonlinear Mixed Effects Models	3.1-117
nnet	Feed-forward Neural Networks and Multinomial Log-Linear Models	7.3-8
parallel	Support for Parallel computation in R	3.1.1
<b>rattle</b>	Graphical user interface for data mining in R	3.3.0
Rcmdr	R Commander	2.1-3
RcmdrMisc	R Commander Miscellaneous Functions	1.0-2
RColorBrewer	ColorBrewer palettes	1.0-5
relimp	Relative Contribution of Effects in a Regression Model	1.0-3
rgl	3D visualization device system (OpenGL)	0.94.1143
rJava	Low-level R to Java interface	0.9-6
RODBC	ODBC Database Access	1.3-10
rpart	Recursive Partitioning and Regression Trees	4.1-8
rstudio	Tools for RStudio	0.98.1074
sandwich	Robust Standard Error Calculations for Linear and Generalized Linear Models	2.3-2
- Status Bar:** Shows a Dropbox notification: "Screenshot Added A screenshot was added to your Dropbox." and the system time: 2:55 PM.

# Installation of Rattle

The screenshot shows the RStudio interface with the following components:

- Top Bar:** File, Edit, Code, View, Plots, Session, Build, Debug, Tools, Help.
- Project Bar:** Project: (None).
- Code Editor:** Untitled1 (R Script). The code in the editor is:

```
library("rstudio", lib.loc="D:/Program Files/R/R-3.1.1/library")
library("stats", lib.loc="D:/Program Files/R/R-3.1.1/library")
install.packages("Rcmdr")
library("Rcmdr", lib.loc="D:/Program Files/R/R-3.1.1/library")
install.packages("rattle")
library("rattle", lib.loc="D:/Program Files/R/R-3.1.1/library")
rattle()
```
- Console:** The output of the R session is displayed:

```
package 'rattle' successfully unpacked and MD5 sums checked
The downloaded binary packages are in
C:/users/dell/Appdata/Local/Temp/Rtmp4Q8Z2F/downloaded_packages
> library("rattle", lib.loc="D:/Program Files/R/R-3.1.1/library")
Rattle: A free graphical interface for data mining with R.
Version 3.3.0 copyright (c) 2006-2014 Togaware Pty Ltd.
Type 'rattle()' to shake, rattle, and roll your data.
> rattle()
The package 'RGtk2' is required to display the Rattle GUI. It does not
appear to be installed. This package (and its dependencies) can be
installed using the following R command:
install.packages('RGtk2')

This one-time install will allow access to the full functionality of
Rattle.

would you like Rattle to install the package now?
(yes/no) yes
trying URL 'http://cran.rstudio.com/bin/windows/contrib/3.1/RGtk2_2.20.31.zip'
Content type 'application/zip' length 13884133 bytes (13.2 Mb)
opened URL
```
- Environment Tab:** Shows the current environment variables and loaded packages.
- Packages Tab:** Shows a list of available packages in the repository, including rattle, Rcmdr, and RColorBrewer.

# Installation of Rattle

The screenshot shows the RStudio interface with the following components:

- Top Bar:** File, Edit, Code, View, Plots, Session, Build, Debug, Tools, Help.
- Project Bar:** Project: (None).
- Code Editor:** Untitled1 (Source on Save). The code being run is:

```
library("rstudio", lib.loc="D:/Program Files/R/R-3.1.1/library")
library("stats", lib.loc="D:/Program Files/R/R-3.1.1/library")
install.packages("Rcmdr")
library("Rcmdr", lib.loc="D:/Program Files/R/R-3.1.1/library")
install.packages("rattle")
library("rattle", lib.loc="D:/Program Files/R/R-3.1.1/library")
rattle()
```
- Console:** Shows the progress of the download:

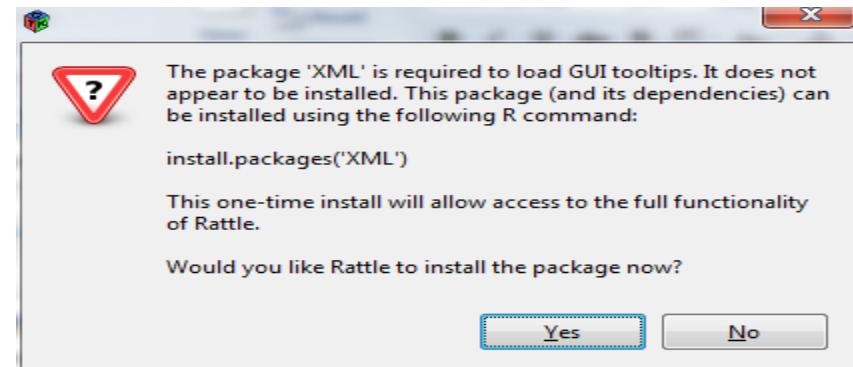
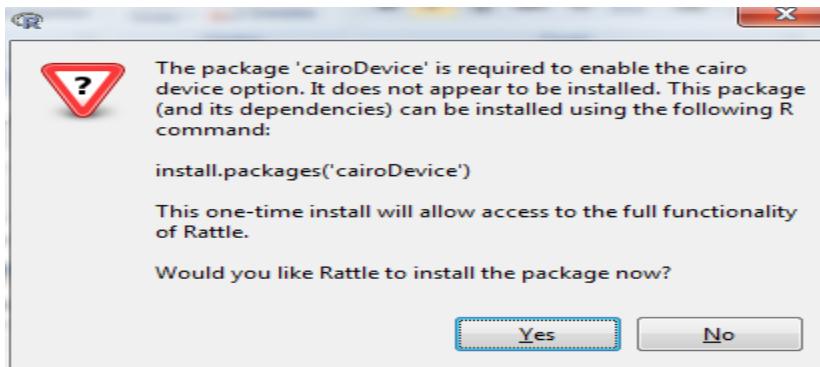
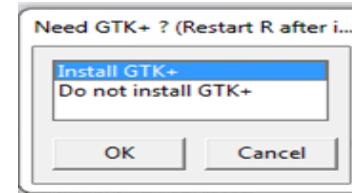
```
90% downloaded
```

URL: ... //cran.rstudio.com/bin/windows/contrib/3.1/RGtk2\_2.20.31.zip

package 'rattle' successfully  
The downloaded binary packages  
C:/Users/dell/AppData/Local/T  
> library("rattle", lib.loc="D:/  
Rattle: A free graphical interface for data mining in R.  
Version 3.3.0 copyright (c) 2006-2014 Togaware Pty Ltd.  
Type 'rattle()' to shake, rattle, and roll your data.  
> rattle()  
The package 'RGtk2' is required to display the Rattle GUI. It does not  
appear to be installed. This package (and its dependencies) can be  
installed using the following R command:  
install.packages('RGtk2')  
This one-time install will allow access to the full functionality of  
Rattle.  
would you like Rattle to install the package now?  
(yes/NO) yes  
trying URL 'http://cran.rstudio.com/bin/windows/contrib/3.1/RGtk2\_2.20.31.zip'  
Content type 'application/zip' length 13884133 bytes (13.2 Mb)  
opened URL
- Environment:** Shows the current environment variables.
- History:** Shows the history of the R session.
- File Explorer:** Shows the file system structure.
- Help:** Shows the help documentation for various packages.
- Task View:** Shows the available tasks and their versions.

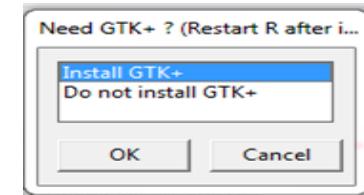
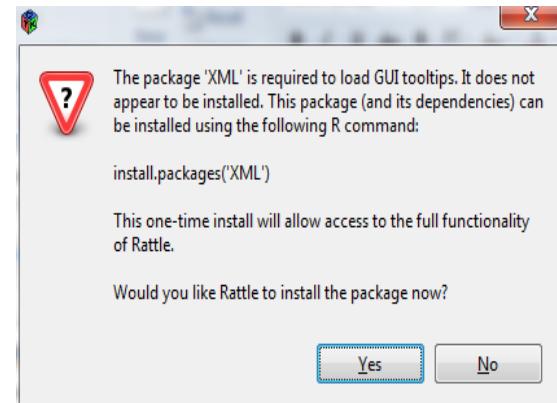
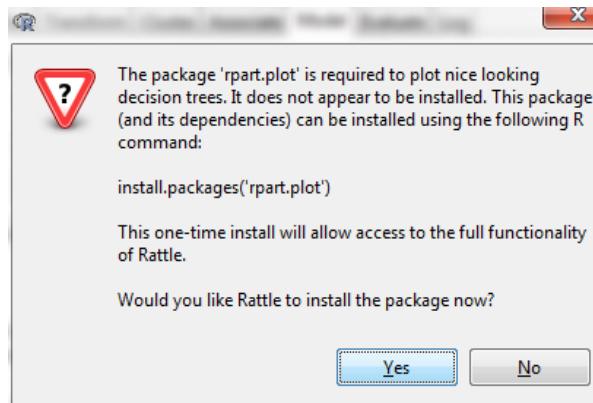
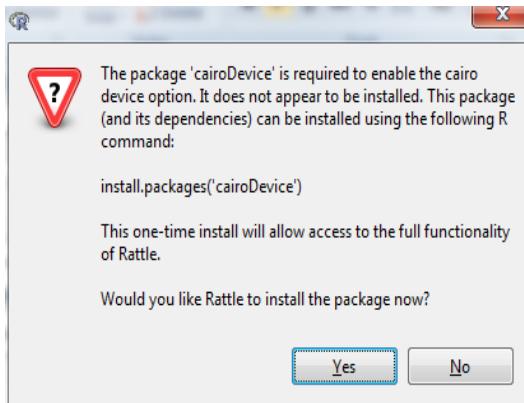
# Installation of Rattle

- GTK+ Installation Necessary
- Install other packages when prompted

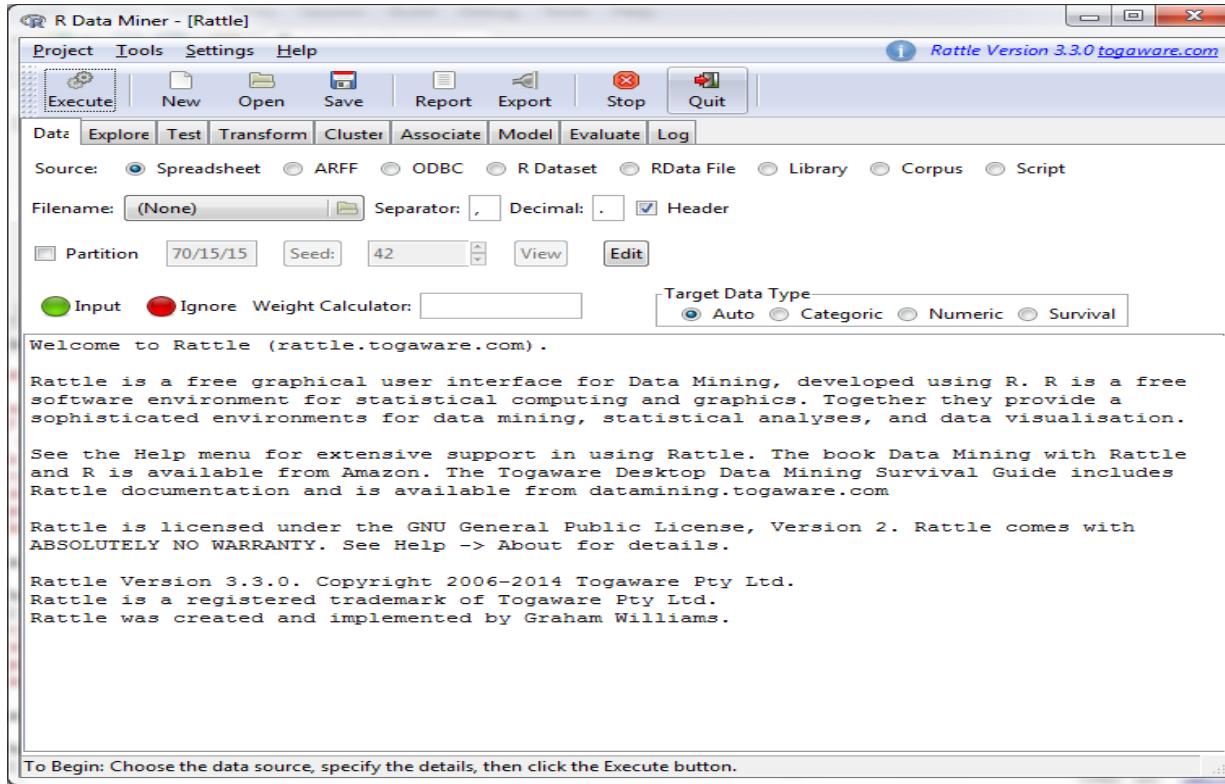


# Installation of Rattle

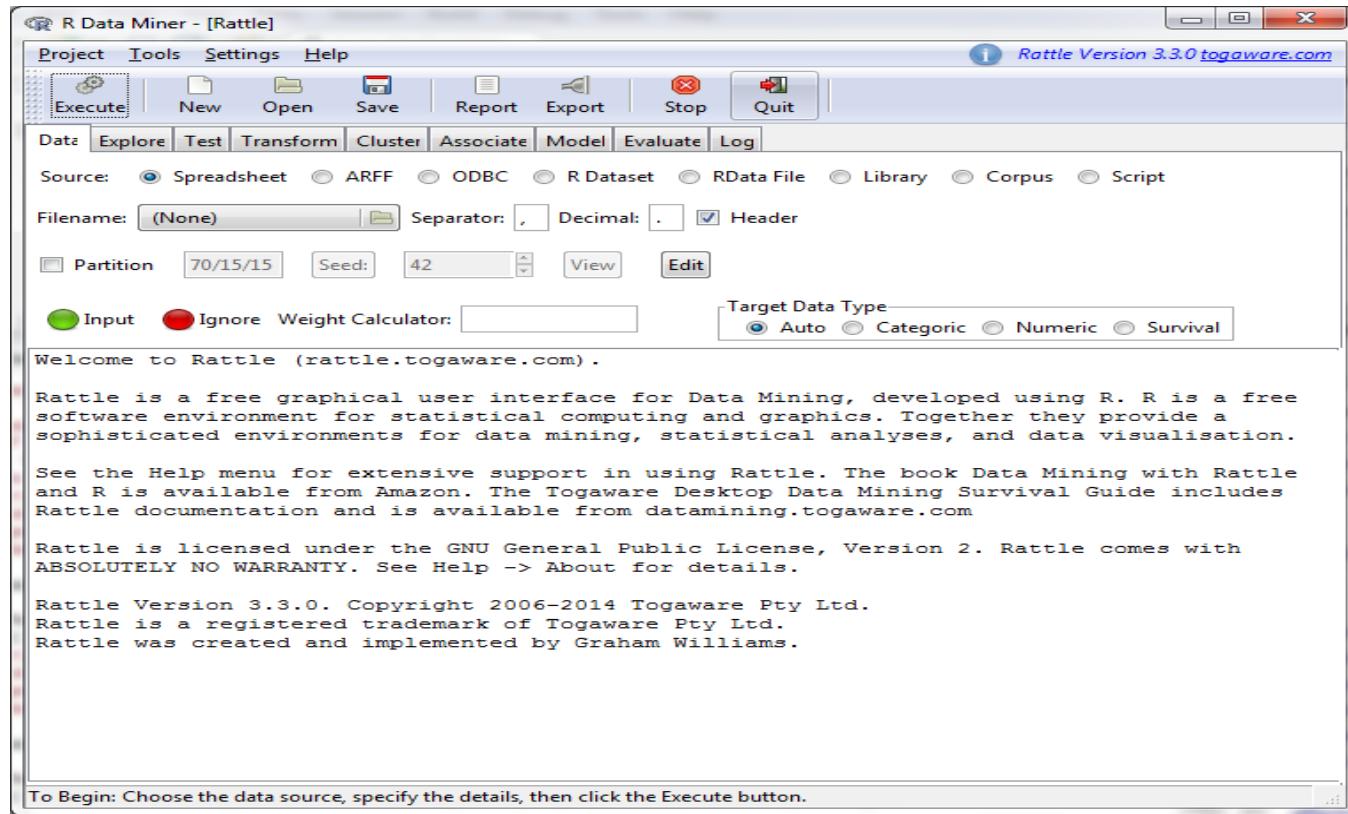
- GTK+ Installation Necessary
- Install other packages when prompted



# Overview of Rattle



# Demo Rattle



# Installation Deducer (with JGR)

The screenshot shows the RStudio interface with the 'Install Packages' dialog box open. The console window displays the process of installing packages from CRAN, including 'Deducer JGR'. The 'Install Packages' dialog has 'Repository (CRAN, CRANextra)' selected and 'Deducer JGR' entered in the 'Packages' field. The 'Install dependencies' checkbox is checked. The background shows the RStudio environment and history panes.

Console output:

```
trying URL 'http://cran.rstudio.com/bin/windows/contrib/3.1/Content type 'application/zip' length 435577 bytes (425 Kb)opened URLdownloaded 425 Kb
```

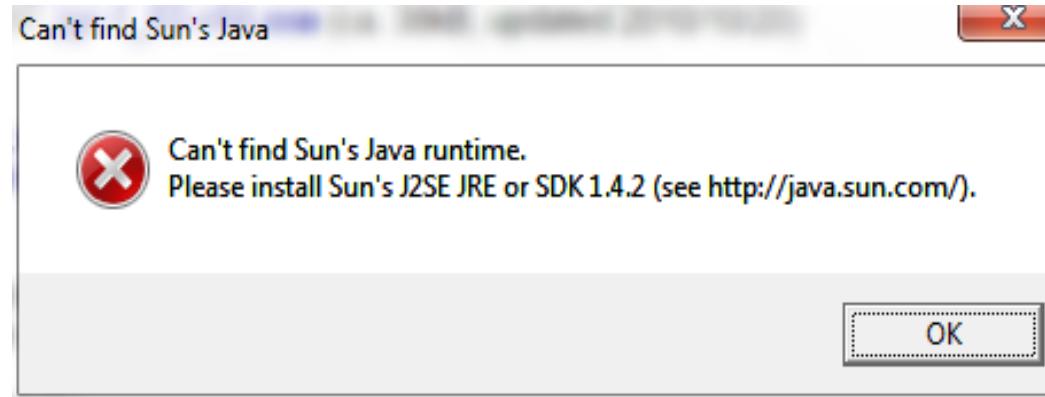
```
package 'modeltools' successfully unpacked and MD5 sums checked
package 'DEoptimR' successfully unpacked and MD5 sums checked
package 'mclust' successfully unpacked and MD5 sums checked
package 'flexmix' successfully unpacked and MD5 sums checked
package 'prabclus' successfully unpacked and MD5 sums checked
package 'diptest' successfully unpacked and MD5 sums checked
package 'robustbase' successfully unpacked and MD5 sums checked
package 'kernlab' successfully unpacked and MD5 sums checked
package 'trimcluster' successfully unpacked and MD5 sums checked
package 'fpc' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
C:/Users/dell/AppData/Local/Temp/Rtmp4Q8Z2F/downloaded_packages
Error in loadNamespace(package, c(which.lib.loc, lib.loc)) :
  cyclic namespace dependency detected when loading 'fpc', already loading 'kernlab'
, 'fpc'
Error in eval(expr, envir, enclos) :
  could not find function "plotcluster"
> |
```

Install Packages dialog:

Description	Version
Combine multi-dimensional arrays	1.4-0
ace() and avas() for selecting regression transformations	1.3-3.3
Another Plot PACKAGE: stem.leaf, bagplot, faces, spin3R, plotsummary, plotl Hulls, and some slider functions	1.3.0
Bootstrap Functions (originally by Angelo Canty for S)	1.3-11
Cairo-based cross-platform antialiased graphics device driver.	2.20
Companion to Applied Regression	2.0-21
Functions for Classification	7.3-10
Cluster Analysis Extended Rousseeuw et al.	1.15.2
Code Analysis Tools for R	0.2-8
Color Space Manipulation	1.2-4
The R Compiler Package	3.1.1
The R Datasets Package	3.1.1
Differential Evolution Optimization in pure R	1.0-1
Create cryptographic hash digests of R objects	0.6.4
Hartigan's dip test statistic for unimodality - corrected code	0.75-5

# Installation Deducer (with JGR)



# Installation Deducer (with JGR)

Java ME  
Java SE Support  
Java SE Advanced & Suite  
Java Embedded  
Java DB  
Web Tier  
Java Card  
Java TV  
New to Java  
Community  
Java Magazine

## Java SE Runtime Environment 7 Downloads

Do you want to run Java™ programs, or do you want to develop Java programs? If you want to run Java programs, but not develop them, download the Java Runtime Environment, or JRE™.

If you want to develop applications for Java, download the Java Development Kit, or JDK™. The JDK includes the JRE, so you do not have to download both separately.

7u71 JRE MD5 Checksum  
7u72 JRE MD5 Checksum

### What is the difference between a Java CPU (7u71) and PSU (7u72) release?

Java SE Critical Patch Updates (CPU) contain fixes to security vulnerabilities and critical bug fixes. Oracle strongly recommends that all Java SE users upgrade to the latest CPU releases as they are made available. Most user should choose this release.

Java SE Patch Set Updates (PSU) contain all of the security fixes in the CPUs released up to that version, as well as additional non-critical fixes. Java PSU releases should only be used if you are being impacted by one of the additional bugs fixed in that version.

Visit Java CPU and PSU Releases Explained for details.

### Java SE Runtime Environment 7u71

You must accept the Oracle Binary Code License Agreement for Java SE to download this software.

Thank you for accepting the Oracle Binary Code License Agreement for Java SE; you may now download this software.

Product / File Description	File Size	Download
Linux x86	31.58 MB	<a href="#">jre-7u71-linux-i586.rpm</a>
Linux x86	46.22 MB	<a href="#">jre-7u71-linux-i586.tar.gz</a>
Linux x64	32.1 MB	<a href="#">jre-7u71-linux-x64.rpm</a>
Linux x64	44.84 MB	<a href="#">jre-7u71-linux-x64.tar.gz</a>
Mac OS X x64	48.57 MB	<a href="#">jre-7u71-macosx-x64.dmg</a>
Mac OS X x64	44.52 MB	<a href="#">jre-7u71-macosx-x64.tar.gz</a>
Solaris x86	52.16 MB	<a href="#">jre-7u71-solaris-i586.tar.gz</a>
Solaris x64	16.14 MB	<a href="#">jre-7u71-solaris-x64.tar.gz</a>
Solaris SPARC	54.95 MB	<a href="#">jre-7u71-solaris-sparc.tar.gz</a>
Solaris SPARC 64-bit	18.1 MB	<a href="#">jre-7u71-solaris-sparcv9.tar.gz</a>
Windows x86 Online	0.89 MB	<a href="#">jre-7u71-windows-i586-iftw.exe</a>
Windows x86 Offline	28.09 MB	<a href="#">jre-7u71-windows-i586.exe</a>
Windows x86	40 MB	<a href="#">jre-7u71-windows-i586.tar.gz</a>
Windows x64	29.59 MB	<a href="#">jre-7u71-windows-x64.exe</a>
Windows x64	41.71 MB	<a href="#">jre-7u71-windows-x64.tar.gz</a>

[Java EE and Glassfish](#)

[Java ME](#)

[Java Card](#)

[NetBeans IDE](#)

[Java Mission Control](#)

### Java Resources

[Java APIs](#)

[Technical Articles](#)

[Demos and Videos](#)

[Forums](#)

[Java Magazine](#)

[Java.net](#)

[Developer Training](#)

[Tutorials](#)

[Java.com](#)

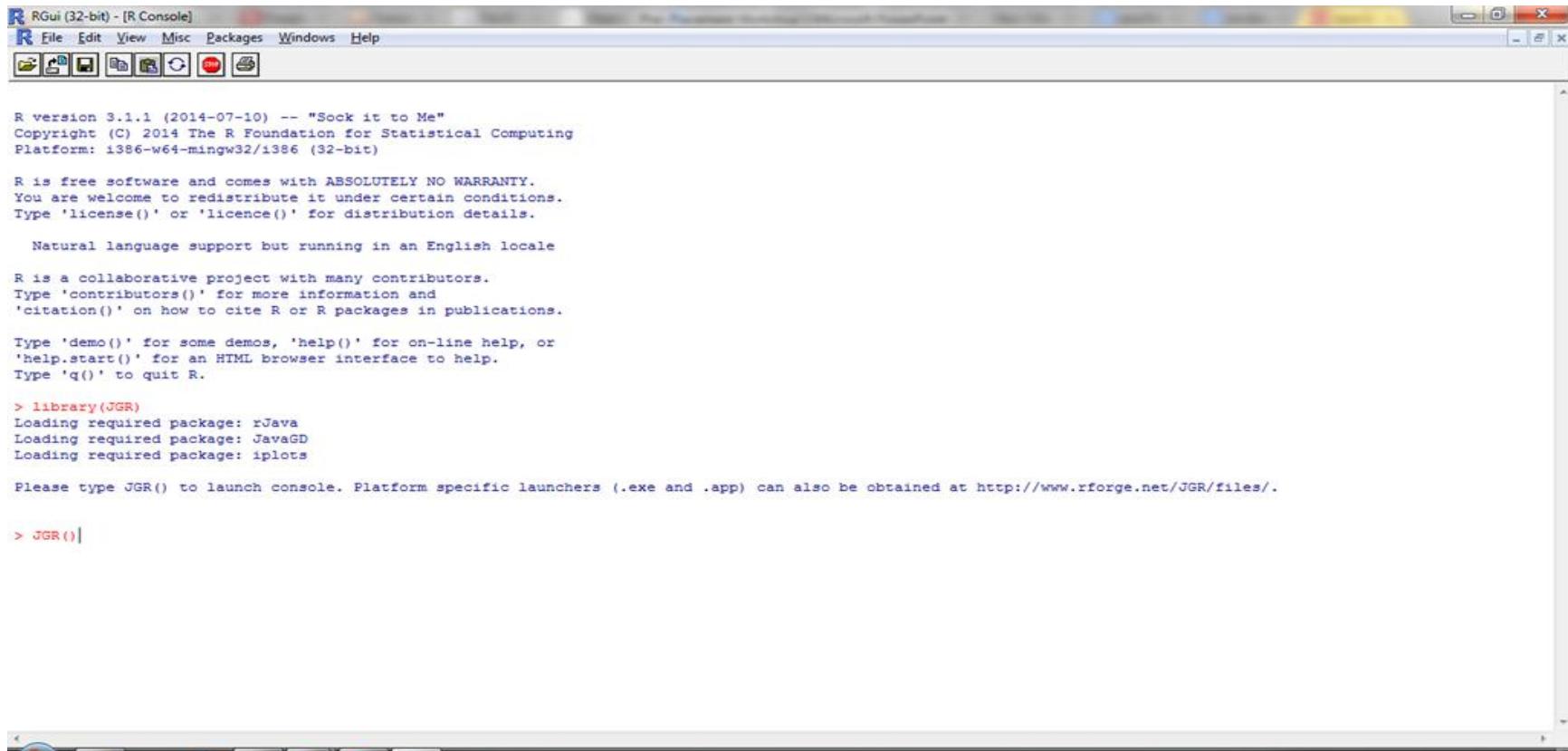


Subscribe Today



Watch Now

# Installation Deducer (with JGR)



R Gui (32-bit) - [R Console]

R File Edit View Misc Packages Windows Help

R version 3.1.1 (2014-07-10) -- "Sock it to Me"  
Copyright (C) 2014 The R Foundation for Statistical Computing  
Platform: i386-w64-mingw32/i386 (32-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.

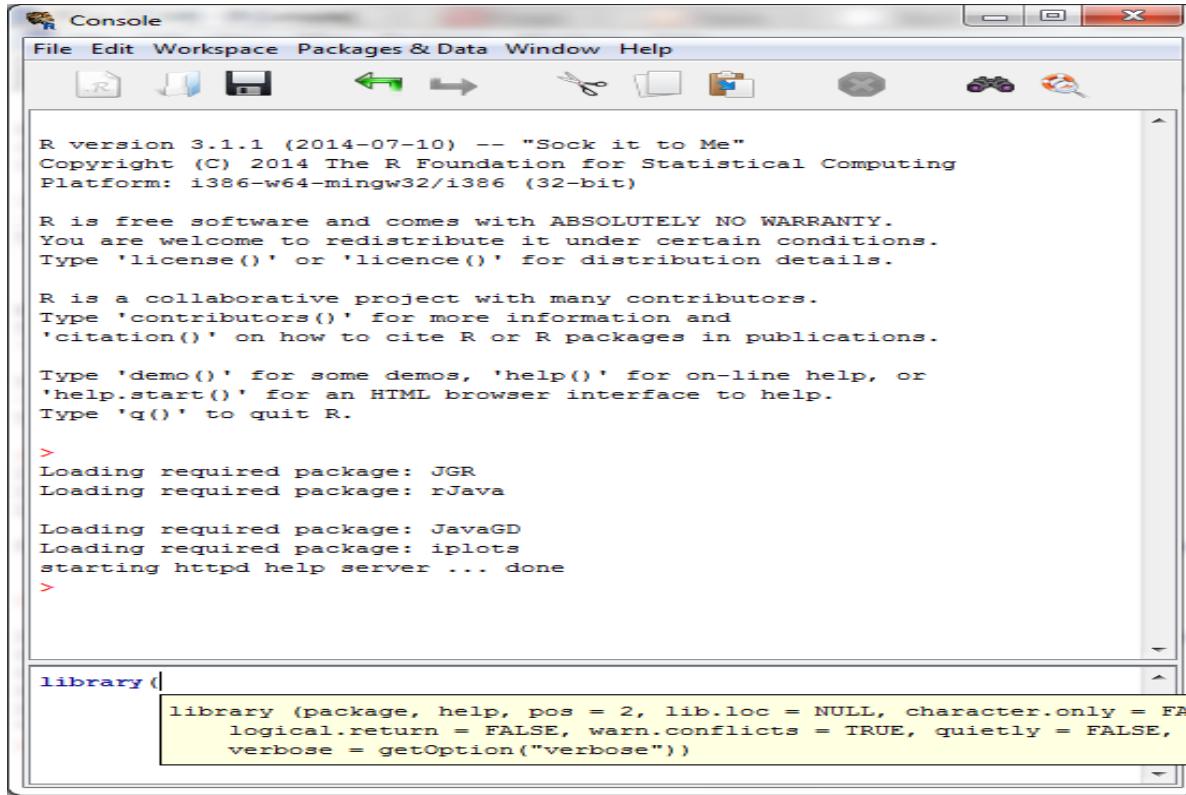
Type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.

```
> library(JGR)
Loading required package: rJava
Loading required package: JavaGD
Loading required package: ipplots

Please type JGR() to launch console. Platform specific launchers (.exe and .app) can also be obtained at http://www.rforge.net/JGR/files/.
```

> JGR()

# Installation Deducer (with JGR)



R version 3.1.1 (2014-07-10) -- "Sock it to Me"  
Copyright (C) 2014 The R Foundation for Statistical Computing  
Platform: i386-w64-mingw32/i386 (32-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.

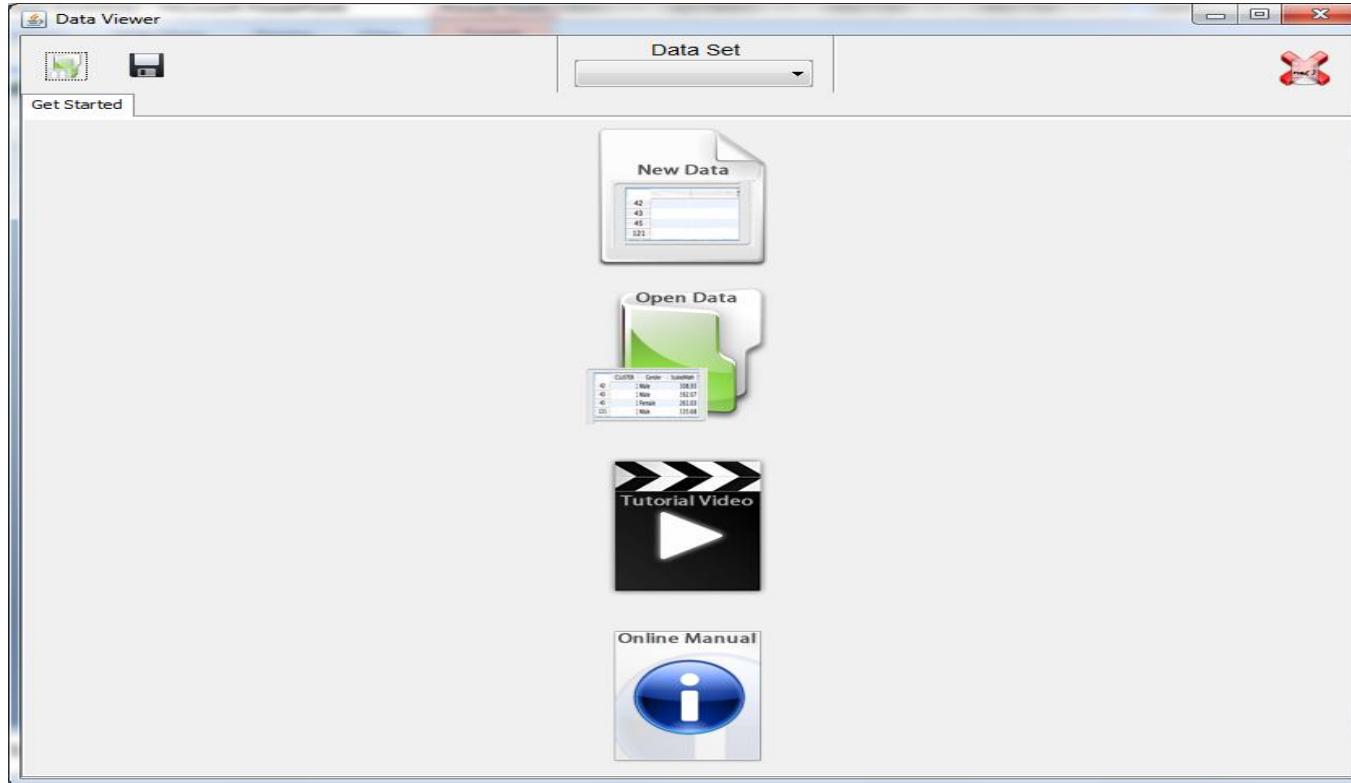
>  
Loading required package: JGR  
Loading required package: rJava

Loading required package: JavaGD  
Loading required package: ipplots  
starting httpd help server ... done  
>

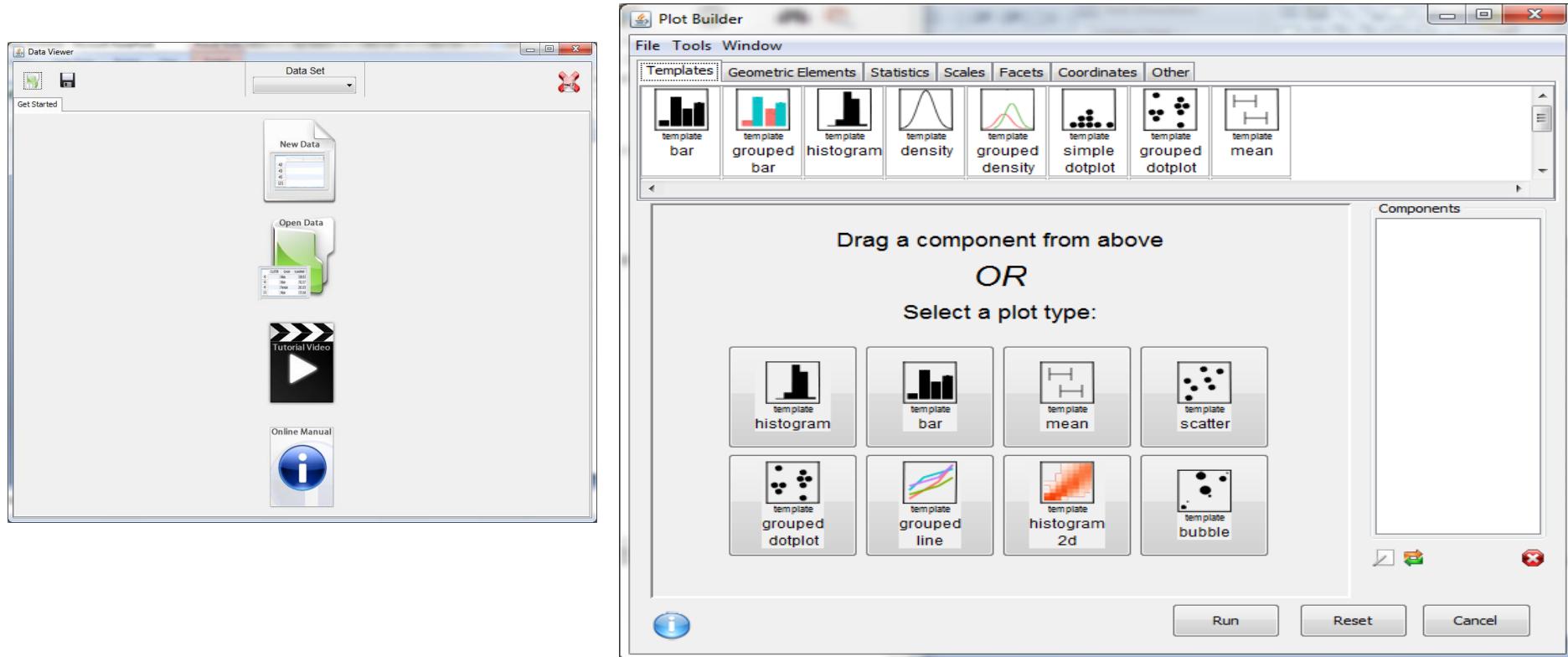
library (

```
library (package, help, pos = 2, lib.loc = NULL, character.only = FALSE,
logical.return = FALSE, warn.conflicts = TRUE, quietly = FALSE,
verbose = getOption("verbose"))
```

# Installation Deducer (with JGR)



# Overview of Deducer (with JGR)



# Demo Deducer

- data()
- data(mtcars)

Data Viewer

File Edit Help

Data Set  
(df) mtcars

Data View Variable View

	mpg	cyl	disp	hp	drat	wt	qsec	vs	am
Mazda RX4	21.0	6	160.0	110	3.90	2.620	16.46	0	
Mazda RX4 Wag	21.0	6	160.0	110	3.90	2.875	17.02	0	
Datsun 710	22.8	4	108.0	93	3.85	2.320	18.61	1	
Hornet 4 Drive	21.4	6	258.0	110	3.08	3.215	19.44	1	
Hornet Sportabout	18.7	8	360.0	175	3.15	3.440	17.02	0	
Valiant	18.1	6	225.0	105	2.76	3.460	20.22	1	
Duster 360	14.3	8	360.0	245	3.21	3.570	15.84	0	
Merc 240D	24.4	4	146.7	62	3.69	3.190	20.00	1	
Merc 230	22.8	4	140.8	95	3.92	3.150	22.90	1	
Merc 280	19.2	6	167.6	123	3.92	3.440	18.30	1	
Merc 280C	17.8	6	167.6	123	3.92	3.440	18.90	1	
Merc 450SE	16.4	8	275.8	180	3.07	4.070	17.40	0	
Merc 450SL	17.3	8	275.8	180	3.07	3.730	17.60	0	
Merc 450SLC	15.2	8	275.8	180	3.07	3.780	18.00	0	
Cadillac Fleetwood	10.4	8	472.0	205	2.93	5.250	17.98	0	
Lincoln Continental	10.4	8	460.0	215	3.00	5.424	17.82	0	
Chrysler Imperial	14.7	8	440.0	230	3.23	5.345	17.42	0	
Fiat 128	32.4	4	78.7	66	4.08	2.200	19.47	1	
Honda Civic	30.4	4	75.7	52	4.93	1.615	18.52	1	
Toyota Corolla	33.9	4	71.1	65	4.22	1.835	19.90	1	
Toyota Corona	21.5	4	120.1	97	3.70	2.465	20.01	1	
Dodge Challenger	15.5	8	318.0	150	2.76	3.520	16.87	0	
AMC Javelin	15.2	8	304.0	150	3.15	3.435	17.30	0	
Camaro Z28	13.3	8	350.0	245	3.73	3.840	15.41	0	
Pontiac Firebird	19.2	8	400.0	175	3.08	3.845	17.05	0	
Fiat X1-9	27.3	4	79.0	66	4.08	1.935	18.90	1	
Porsche 914-2	26.0	4	120.3	91	4.43	2.140	16.70	0	
Lotus Europa	30.4	4	95.1	113	3.77	1.513	16.90	1	
Ford Pantera L	15.8	8	351.0	264	4.22	3.170	14.50	0	
Ferrari Dino	19.7	6	145.0	175	3.62	2.770	15.50	0	
Maserati Bora	15.0	8	301.0	335	3.54	3.570	14.60	0	
Volvo 142E	21.4	4	121.0	109	4.11	2.780	18.60	1	
33									
24									

# RStudio

RStudio Desktop enables you with following advantages of native R console

- Syntax highlighting, code completion, and smart indentation
- Execute R code directly from the source editor
- Quickly jump to function definitions
- Easily manage multiple working directories using projects
- Integrated R help and documentation
- Interactive debugger to diagnose and fix errors quickly
- Extensive package development tools

<http://www.rstudio.com/products/>

# RStudio

RStudio Server enables you to provide a browser based interface (the RStudio IDE) to a version of R running on a remote Linux server. Deploying R and RStudio on a server has a number of benefits, including:

- The ability to access your R workspace from any computer in any location;
- Easy sharing of code, data, and other files with colleagues;
- Allowing multiple users to share access to the more powerful compute resources (memory, processors, etc.) available on a well equipped server; and
- Centralized installation and configuration of R, R packages, TeX, and other supporting libraries.

File Edit Code View Plots Session Build Debug Tools Help

Project: (None)

new2.R x packages.R x chapter1.Rmd x Untitled1\* x

Source on Save Run Source

```

1 library(ggplot2)
2 data(diamonds)
3 barplot(diamonds$price)
4 plot(diamonds$price)
5 plot(diamonds$price,diamonds$carat)
6 pie(table(diamonds$cut))
7 boxplot(diamonds$price)
8 boxplot(diamonds$price-diamonds$cut)
9 boxplot(diamonds$price-diamonds$color)
10 plot(diamonds$cut,diamonds$color)
11 hist(diamonds$price)
12
12:1 (Top Level) R Script

```

Console ~ /

```

> kmeans
function (x, centers, iter.max = 10, nstart = 1, algorithm = c("Hartigan-Wong",
  "Lloyd", "Forgy", "MacQueen"), trace = FALSE)
{
  do_one <- function(nmeth) {
    switch(nmeth, {
      isteps.Qtran <- 50 * m
      iTran <- c(as.integer(isteps.Qtran), integer(max(0,
        k - 1)))
      Z <- .Fortran(C_kmnns, x, m, p, centers = centers,
        as.integer(k), c1 = integer(m), c2 = integer(m),
        nc = integer(k), double(k), ncp = integer(k),
        D = double(m), iTran = iTran, live = integer(k),
        iter = iter.max, wss = double(k), ifault = as.integer(trace))
      switch(Z$ifault, stop("empty cluster: try better set of initial centers",
        call. = FALSE), Z$iter <- max(Z$iter, iter.max +
        1L), stop("number of cluster centres must lie between 1 and nrow(x)",
        call. = FALSE), warning(gettextf("Quick-TRANSFER stage steps exceeded maximum (= %d"),
        isteps.Qtran), call. = FALSE))
    }, {
      Z <- .C(C_kmeans_Lloyd, x, m, p, centers = centers,
        k, c1 = integer(m), iter = iter.max, nc = integer(k),
        wss = double(k)))
    })
  }
}

```

Environment History Import Dataset Clear List

Global Environment

Data

- diamonds 53940 obs. of 10 variables
- iris3 50 obs. of 12 variables

Values

- a NULL (empty)
- i 90L

Files Plots Packages Help Viewer

R: Search Results Find in Topic

## Search Results



The search string was "kmeans"

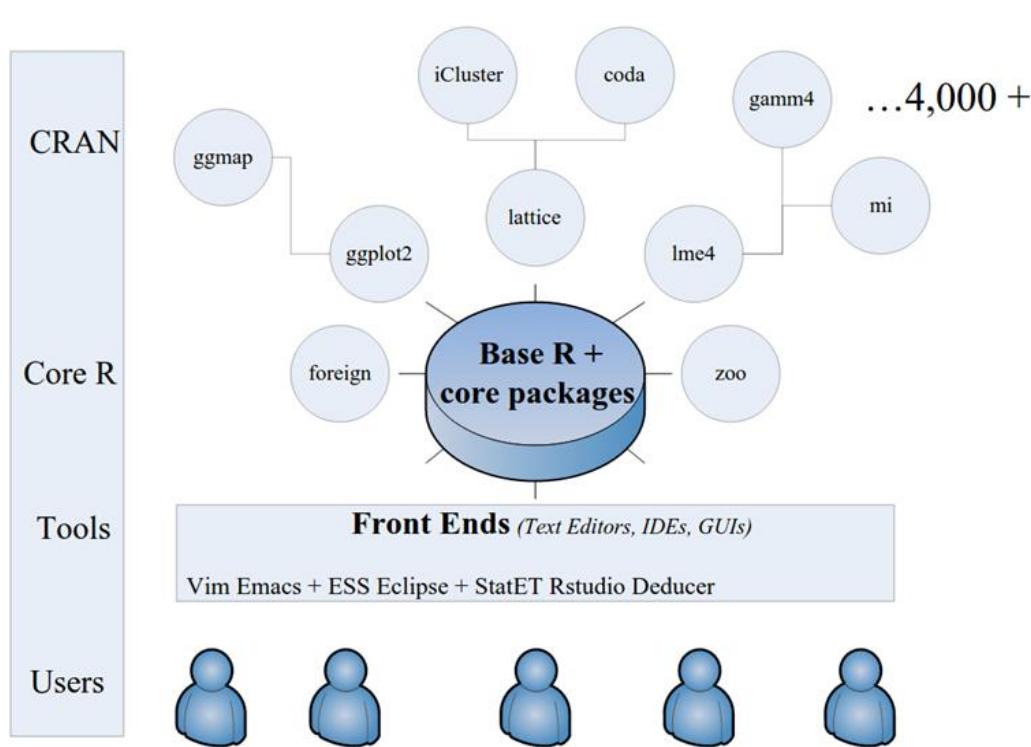
Vignettes:

- [broom::kmeans](#) kmeans with dplyr+broom [HTML](#) [source](#) [R code](#)

Help pages:

- [amap::Kmeans](#) K-Means Clustering
- [broom::augment.kmeans](#) Tidying methods for kmeans objects
- [e1071::cmeans](#) Fuzzy C-Means Clustering

# R Landscape



*edited by the R Development Core Team.*

The following manuals for R were created on Debian Linux and may differ from the manuals for Mac or Windows on platform. Version of the manuals for each platform are part of the respective R installations. The manuals change with R, hence we provide version for the patched release version (R-patched) and finally a version for the forthcoming R version that is still in development.

Here they can be downloaded as PDF files, EPUB files, or directly browsed as HTML:

## Manuals

<http://cran.r-project.org/manuals.html>

Manual	R-release
<b>An Introduction to R</b> is based on the former "Notes on R", gives an introduction to the language and how to use R for doing statistical analysis and graphics.	<a href="#">HTML</a>   <a href="#">PDF</a>   <a href="#">EPUB</a>
<b>R Data Import/Export</b> describes the import and export facilities available either in R itself or via packages which are available from CRAN.	<a href="#">HTML</a>   <a href="#">PDF</a>   <a href="#">EPUB</a>
<b>R Installation and Administration</b>	<a href="#">HTML</a>   <a href="#">PDF</a>   <a href="#">EPUB</a>
<b>Writing R Extensions</b> covers how to create your own packages, write R help files, and the foreign language (C, C++, Fortran, ...) interfaces.	<a href="#">HTML</a>   <a href="#">PDF</a>   <a href="#">EPUB</a>
A draft of <b>The R language definition</b> documents the language <i>per se</i> . That is, the objects that it works on, and the details of the expression evaluation process, which are useful to know when programming R functions.	<a href="#">HTML</a>   <a href="#">PDF</a>   <a href="#">EPUB</a>
<b>R Internals</b> : a guide to the internal structures of R and coding standards for the core team working on R itself.	<a href="#">HTML</a>   <a href="#">PDF</a>   <a href="#">EPUB</a>
<b>The R Reference Index</b> : contains all help files of the R standard and recommended packages in printable form. (9MB, approx. 3500 pages)	<a href="#">PDF</a>

Translations of manuals into other languages than English are available from the [contributed documentation](#) section (only a few are currently available).

The LaTeX or Texinfo sources of the latest version of these documents are contained in every R source distribution (in the `doc/manuals` directory). They can be found in the respective [archives of the R sources](#). The HTML versions of the manuals are also part of most R installations.

Please check the manuals for R-devel before reporting any issues with the released versions.

**R**

# Documentation

## - Vignettes

### ggplot2: An Implementation of the Grammar of Graphics

An implementation of the grammar of graphics in R. It combines the advantages of both base and lattice graphics: by step from multiple data sources. It also implements a sophisticated multidimensional conditioning system and a information, documentation and examples.

Version: 1.0.1  
Depends: R ( $\geq$  2.14), stats, methods  
Imports: [plyr](#) ( $\geq$  1.7.1), [digest](#), grid, [gttable](#) ( $\geq$  0.1.1), [reshape2](#), [scales](#) ( $\geq$  0.2.3), [proto](#), [MASS](#), [quantreg](#), [Hmisc](#), [mapproj](#), [maps](#), [hexbin](#), [maptools](#), [multcomp](#), [nlme](#), [testthat](#), [knitr](#), [mgcv](#)  
Suggests: [sp](#)  
Enhances:  
Published: 2015-03-17  
Author: Hadley Wickham [aut, cre], Winston Chang [aut]  
Maintainer: Hadley Wickham <h.wickham@gmail.com>  
BugReports: <https://github.com/hadley/ggplot2/issues>  
License: [GPL-2](#)  
URL: <http://ggplot2.org>, <https://github.com/hadley/ggplot2>  
NeedsCompilation: no  
Citation: [ggplot2 citation info](#)  
Materials: [README NEWS](#)  
In views: [Graphics](#), [Phylogenetics](#)  
CRAN checks: [ggplot2 results](#)

#### Downloads:

Reference manual: [ggplot2.pdf](#)  
Vignettes: [Contributing to ggplot2 development](#), [ggplot2 release process](#)  
Package source: [ggplot2\\_1.0.1.tar.gz](#)  
Windows binaries: r-devel: [ggplot2\\_1.0.1.zip](#), r-release: [ggplot2\\_1.0.1.zip](#), r-oldrel: [ggplot2\\_1.0.1.zip](#)  
OS X Snow Leopard binaries: r-release: not available, r-oldrel: [ggplot2\\_1.0.1.tgz](#)  
OS X Mavericks binaries: r-release: [ggplot2\\_1.0.1.tgz](#)  
Old sources: [ggplot2 archive](#)

#### Reverse dependencies:

Reverse depends: [alphahull](#), [AmpliconDuo](#), [aoristic](#), [apsimr](#), [bcrm](#), [bde](#), [benchmark](#), [biomod2](#), [bootnet](#), [brms](#), [car](#)

<http://cran.r-project.org/web/views/>

<a href="#">Bayesian</a>	Bayesian Inference
<a href="#">ChemPhys</a>	Chemometrics and Computational Physics
<a href="#">ClinicalTrials</a>	Clinical Trial Design, Monitoring, and Analysis
<a href="#">Cluster</a>	Cluster Analysis & Finite Mixture Models
<a href="#">DifferentialEquations</a>	Differential Equations
<a href="#">Distributions</a>	Probability Distributions
<a href="#">Econometrics</a>	Econometrics
<a href="#">Environmetrics</a>	Analysis of Ecological and Environmental Data
<a href="#">ExperimentalDesign</a>	Design of Experiments (DoE) & Analysis of Experimental Data
<a href="#">Finance</a>	Empirical Finance
<a href="#">Genetics</a>	Statistical Genetics
<a href="#">Graphics</a>	Graphic Displays & Dynamic Graphics & Graphic Devices & Visualization
<a href="#">HighPerformanceComputing</a>	High-Performance and Parallel Computing with R
<a href="#">MachineLearning</a>	Machine Learning & Statistical Learning
<a href="#">MedicalImaging</a>	Medical Image Analysis
<a href="#">MetaAnalysis</a>	Meta-Analysis
<a href="#">Multivariate</a>	Multivariate Statistics
<a href="#">NaturalLanguageProcessing</a>	Natural Language Processing
<a href="#">NumericalMathematics</a>	Numerical Mathematics
<a href="#">OfficialStatistics</a>	Official Statistics & Survey Methodology
<a href="#">Optimization</a>	Optimization and Mathematical Programming
<a href="#">Pharmacokinetics</a>	Analysis of Pharmacokinetic Data
<a href="#">Phylogenetics</a>	Phylogenetics, Especially Comparative Methods
<a href="#">Psychometrics</a>	Psychometric Models and Methods
<a href="#">ReproducibleResearch</a>	Reproducible Research
<a href="#">Robust</a>	Robust Statistical Methods
<a href="#">SocialSciences</a>	Statistics for the Social Sciences
<a href="#">Spatial</a>	Analysis of Spatial Data
<a href="#">SpatioTemporal</a>	Handling and Analyzing Spatio-Temporal Data
<a href="#">Survival</a>	Survival Analysis
<a href="#">TimeSeries</a>	Time Series Analysis
<a href="#">WebTechnologies</a>	Web Technologies and Services
<a href="#">gR</a>	gRaphical Models in R

# R Community

- email groups <http://www.r-project.org/mail.html>

**R-announce**

**R-help**

**R-package-devel**

**R-devel**

**R-packages**

**Special Interest Groups**

- Stack Overflow [r]
- Twitter #rstats
- Blogs at <http://www.r-bloggers.com/> (573 blogs)

# Stack Overflow

StackExchange 85 9 help [r]

stackoverflow Questions Tags Users Badges Unanswered Ask Question

Tagged Questions newest featured frequent votes active unanswered

90,861 questions tagged

about »

Featured on Meta

April 2015 Community Moderator Election Results

Hot Meta Posts

Failed edit to a question says: "Your answer couldn't be submitted"

The Font Awesome child tags are too specific - are they even necessary?

Flagging questions with details only in comments

R Count number of rows in one column of a data frame? asked 2 mins ago

I just want to know how to get r to list the number of occupied rows of a specific column of a data frame. My guess was nrow(dataframe\$column) though that didn't work.

r RyanMe321 6 2

Create interactive webmap with markers in R using Shiny, Leaflet and rCharts

I am trying to create an interactive webmap in R to display storms using Shiny, Leaflet and rCharts (the structure is loosely based on the <http://ramnathv.github.io/bikeshare> app). The idea is that ...

r dictionary leaflet shiny rcharts Louise 18 5

asked 5 mins ago

R - gsub a specific character of a specific position

I would like to delete the last character of a variable. I was wondering if it is possible to select the position with gsub and delete the character of this particular position. In this example, I ...

regex r position gsub giacomoV 1 1

asked 15 mins ago

Looking for a job?

Chief Software Architect - Java + \$100K

Crossover Bengaluru, India / remote

airbnb instacart

http://stackoverflow.com/questions/tagged/

Collabera Value. Accelerated.

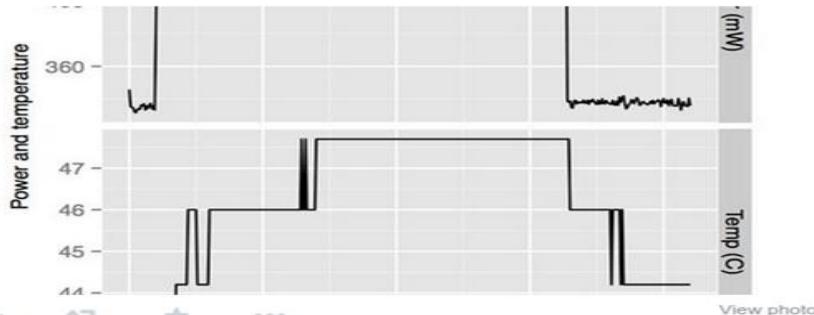
Visit us at [collaberatact.com](http://collaberatact.com)

153

...

## Results for #rstats

Top / All

**Mark Benson** @markbenson · 5mPower and heat are related. Here's an R plot I did that proves it on the Kindle Fire. [#rstats](#) [vanilladraft.com/stmes/](#)[View photo](#)**Stéphane Fréchette** @sfrechette · 8mHow to get your very own RStudio Server and Shiny Server with DigitalOcean [r-bloggers.com/how-to-get-you...](#) #datascience #feedly #rstats #shiny**Ankit kansal** @sinisterinankit · 9m

Interesting post on configuring parallel computing on #r #rstudio #rstats #dataprocessing #data

**Learn R** @R\_ProgrammingHow to do parallel computing with R? [rstatistics.net/parallel-compu...](#) #rstats #datascience

<https://twitter.com/search?q=rstats&src=sprv>

# Help within R

? "keyword"

?? "keyword"

Example-

```
> ?kmeans
```

```
> ??kmeans
```

# Introductory R

```
> Sys.Date()
```

```
[1] "2015-05-10,,
```

```
> Sys.time()
```

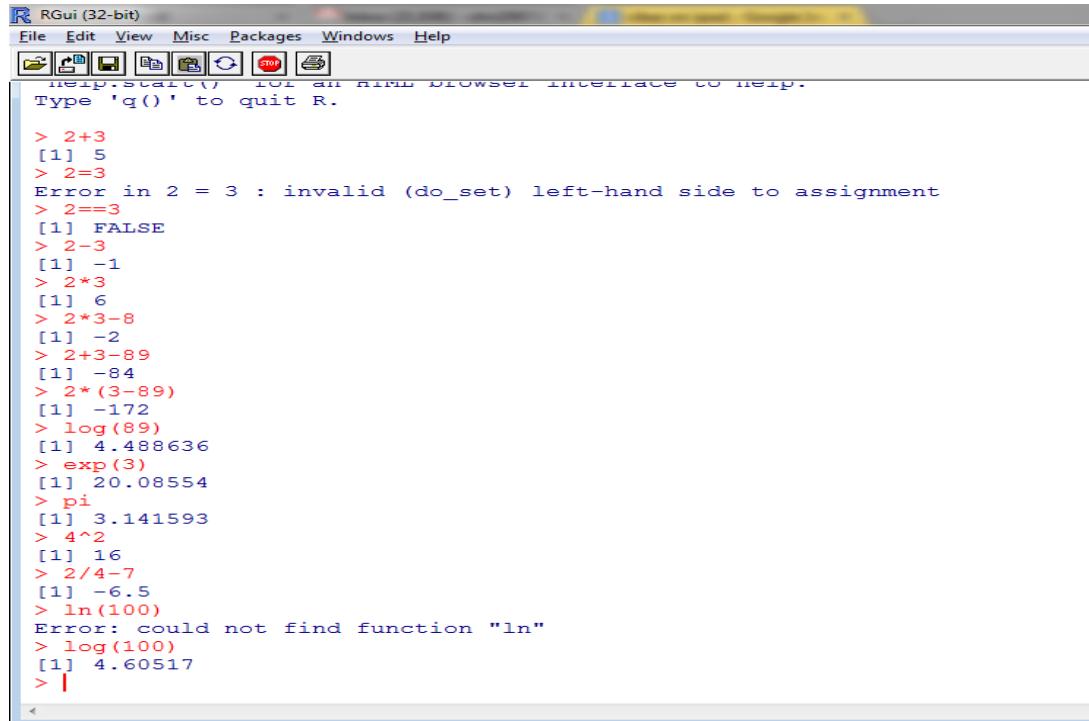
```
[1] "2015-05-10 18:28:32 IST"
```

# R as a Calculator

## Basic Math on R Console

- +
- -
- Log
- Exp
- \*
- /
- ()
- mean
- sum
- sd
- log
- median
- exp

# Demo- Basic Math on R Console



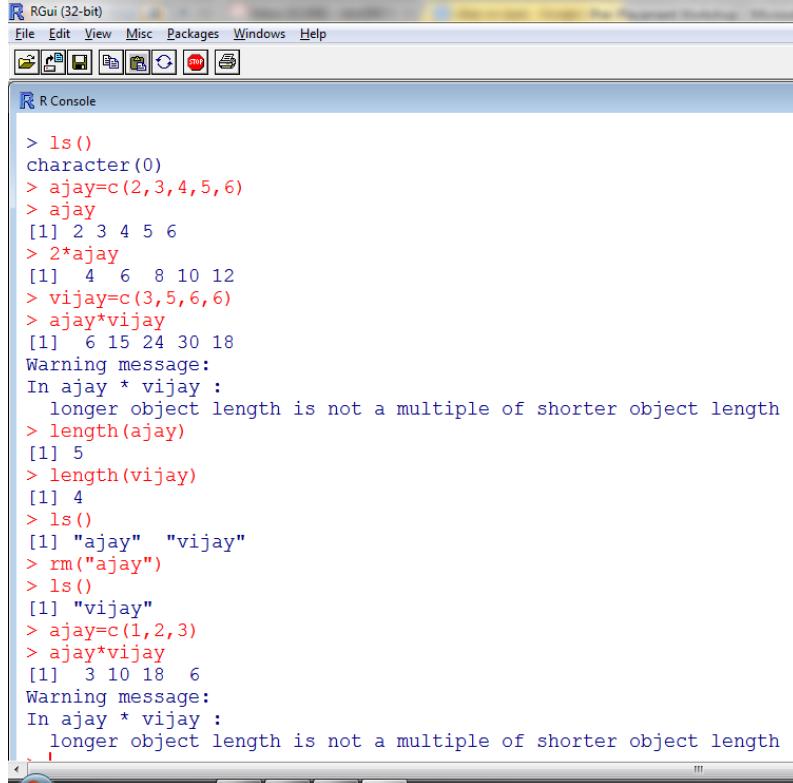
The screenshot shows the RGui (32-bit) application window. The title bar reads "RGui (32-bit)". The menu bar includes File, Edit, View, Misc, Packages, Windows, and Help. Below the menu is a toolbar with various icons. The main window displays the R console output. It starts with a help message: "help.start() for an HTML browser interface to help. Type 'q()' to quit R." followed by several basic arithmetic operations and function calls:

```
help.start() for an HTML browser interface to help.
Type 'q()' to quit R.

> 2+3
[1] 5
> 2=3
Error in 2 = 3 : invalid (do_set) left-hand side to assignment
> 2==3
[1] FALSE
> 2-3
[1] -1
> 2*3
[1] 6
> 2*3-8
[1] -2
> 2+3-89
[1] -84
> 2*(3-89)
[1] -172
> log(89)
[1] 4.488636
> exp(3)
[1] 20.08554
> pi
[1] 3.141593
> 4^2
[1] 16
> 2/4-7
[1] -6.5
> ln(100)
Error: could not find function "ln"
> log(100)
[1] 4.60517
> |
```

Hint- Ctrl +L clears screen

# Demo- Basic Objects on R Console



R Gui (32-bit)

File Edit View Misc Packages Windows Help

R Console

```
> ls()
character(0)
> ajay=c(2,3,4,5,6)
> ajay
[1] 2 3 4 5 6
> 2*ajay
[1] 4 6 8 10 12
> vijay=c(3,5,6,6)
> ajay*vijay
[1] 6 15 24 30 18
Warning message:
In ajay * vijay :
  longer object length is not a multiple of shorter object length
> length(ajay)
[1] 5
> length(vijay)
[1] 4
> ls()
[1] "ajay" "vijay"
> rm("ajay")
> ls()
[1] "vijay"
> ajay=c(1,2,3)
> ajay*vijay
[1] 3 10 18 6
Warning message:
In ajay * vijay :
  longer object length is not a multiple of shorter object length
```

## Functions-

ls() – what objects are here

rm("foo") removes object named foo

## Assignment

Using = or -> assigns object names to values

Hint- Up arrow gives you last typed command ↑

# Functions and Loops

- Loops

```
for (number in 1:5){ print (number) }
```

```
> for (number in 1:5){ print (number) }
[1] 1
[1] 2
[1] 3
[1] 4
[1] 5
> for (i in 1:5){ print (i) }
[1] 1
[1] 2
[1] 3
[1] 4
[1] 5
> for (i in 1:5){ rnorm(i,10,10) }
> for (i in 1:5){ print(rnorm(i,10,10)) }
[1] 1.090406
[1] 8.611727 16.670168
[1] 10.84623 13.13938 11.56230
[1] 6.068250 -18.723389 33.174107 -1.320091
[1] 13.939702 -9.037375 13.755986 9.459680 9.625309
> |
```

# Functions and Loops

- Function

```
functionajay=function(a)(a^2+2*a+1)
```

```
> functionajay=function(a) (a^2+2*a+1)
> functionajay(1)
[1] 4
> functionajay(2)
[1] 9
> for (i in 1:5){ print(rnorm(i) )
Error: unexpected ')' in "for (i in 1:5){ print(rnorm(i) )"
>
> for (i in 1:5){ print(functionajay(i)) }
[1] 4
[1] 9
[1] 16
[1] 25
[1] 36
> |
```

Hint: Always match brackets

Each ( deserves a )

Each { deserves a }

Each [ deserves a ]

# Other sources to learn R

swirlstats

<http://swirlstats.com/>

datacamp

<https://www.datacamp.com/>

codeschool

<http://tryr.codeschool.com/>

coursera

<https://www.coursera.org/course/compdata>

<https://www.coursera.org/course/rprog>

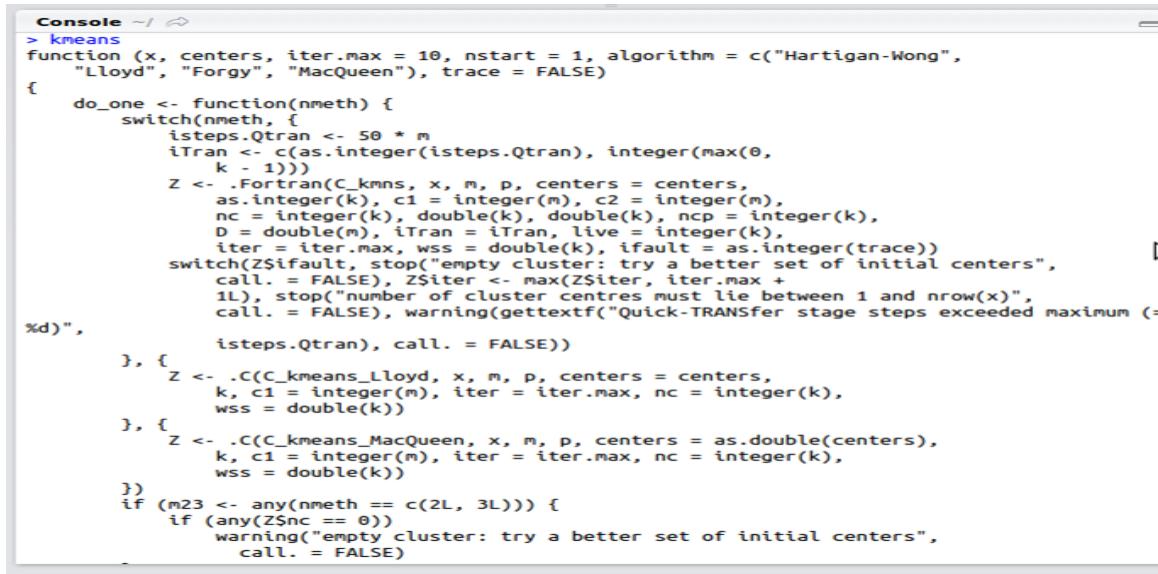


# Good Coding Practices

- Use # for comment
- Use git for version control
- Use Rstudio for multiple lines of code

# Functions in R

- custom functions
- source code for a function
- Understanding help ?, ??



The screenshot shows the R console window with the title "Console". The code displayed is the source for the kmeans function. It defines a function that takes parameters x, centers, iter.max, nstart, algorithm, and trace. The function uses a switch statement to handle different algorithms: Hartigan-Wong, Lloyd, Forgy, and MacQueen. For each algorithm, it calls a C function (C\_kmeans, C\_kmeans\_Lloyd, or C\_kmeans\_MacQueen) with arguments x, m, p, and centers. The C functions return a list Z containing cluster assignments, number of clusters k, and total within-cluster sum of squares wss.

```
> kmeans
function (x, centers, iter.max = 10, nstart = 1, algorithm = c("Hartigan-Wong",
    "Lloyd", "Forgy", "MacQueen"), trace = FALSE)
{
  do_one <- function(nmeth) {
    switch(nmeth, {
      isteps.Qtran <- 50 * m
      iTran <- c(as.integer(isteps.Qtran), integer(max(0,
        k - 1)))
      Z <- .Fortran(C_kmeans, x, m, p, centers = centers,
        as.integer(k), c1 = integer(m), c2 = integer(m),
        nc = integer(k), double(k), ncp = integer(k),
        D = double(m), iTran = iTran, live = integer(k),
        iter = iter.max, wss = double(k), ifault = as.integer(trace))
      switch(Z$efault, stop("empty cluster: try a better set of initial centers",
        call. = FALSE), Z$iter <- max(Z$iter, iter.max +
          1L), stop("number of cluster centres must lie between 1 and nrow(x)",
        call. = FALSE), warning(gettextf("Quick-TRANSfer stage steps exceeded maximum (%d)",
          isteps.Qtran), call. = FALSE))
    }, {
      Z <- .C(C_kmeans_Lloyd, x, m, p, centers = centers,
        k, c1 = integer(m), iter = iter.max, nc = integer(k),
        wss = double(k))
    }, {
      Z <- .C(C_kmeans_MacQueen, x, m, p, centers = as.double(centers),
        k, c1 = integer(m), iter = iter.max, nc = integer(k),
        wss = double(k))
    })
  if (m23 <- any(nmeth == c(2L, 3L))) {
    if (any(Z$nc == 0))
      warning("empty cluster: try a better set of initial centers",
        call. = FALSE)
  }
}
```

# HOMEWORK TIME !



# Functions Used in this Lesson

- function(x)
- for
- library
- install.packages
- update.packages
- ls
- rm
- print

# Citations and References

> citation()

To cite R in publications use:

R Core Team (2015). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <http://www.R-project.org/>.

**Collabera®**  
*Value. Accelerated.*

Collabera