# Data Science Using R

Lesson06–Introduction to R Markdown and Rattle

# Objective

After completing this lesson you will be able to:

- Describe R Markdown and Rattle
- Build a basic R Markdown document
- Explain the various features of Rattle
- Run a dataset in Rattle through a set of commonly used techniques of data analysis.

# R Markdown–An Introduction

- R Markdown is an authoring format that enables easy creation of dynamic documents, presentations, and reports from R.

- R markdown can be used to create reports in the following format:

| Report Format | Output Format |
|---|---|
| Document | HTML, PDF, WORD |
| Presentation | HTML(ioslides), HTML(Slidy), PDF(Beamer) |
| Interactive Shiny Report | Shiny Document, Shiny Presentation |

- R Markdown documents can be automatically regenerated whenever underlying R code or data changes.

# R Markdown–Install Package

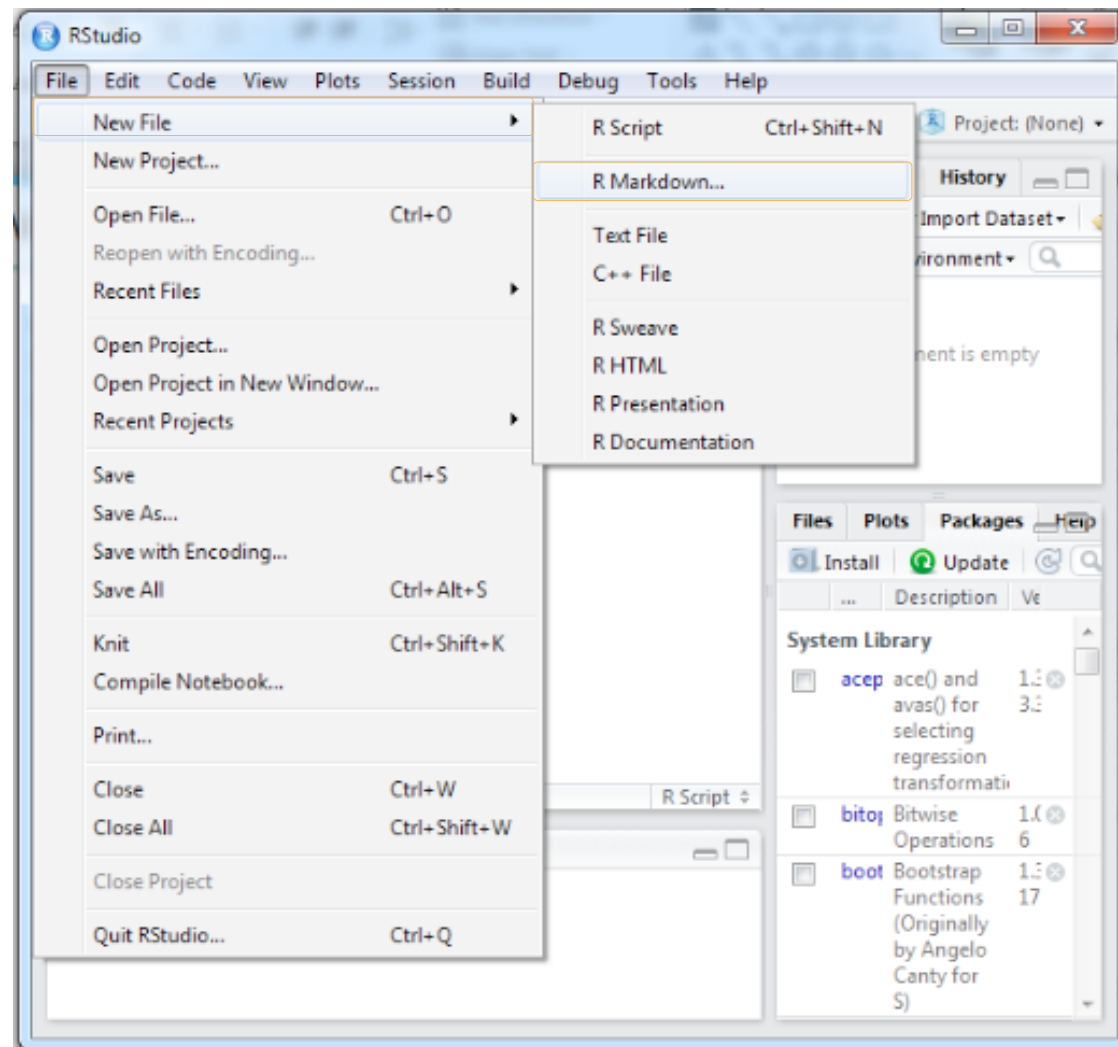- The first step to use R markdown is to install the package.

> ***On R Studio Console***:
> `>Install.packages("rmarkdown")`
>
> *Or install using the Rstudio Install packages options*
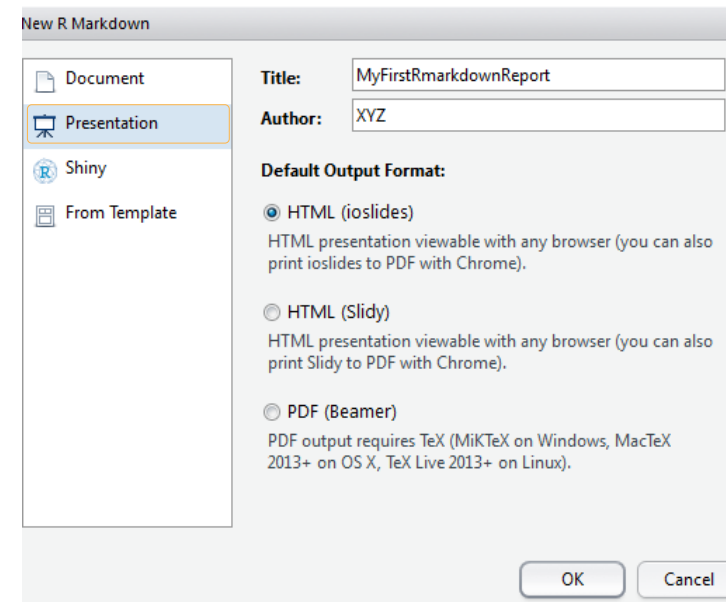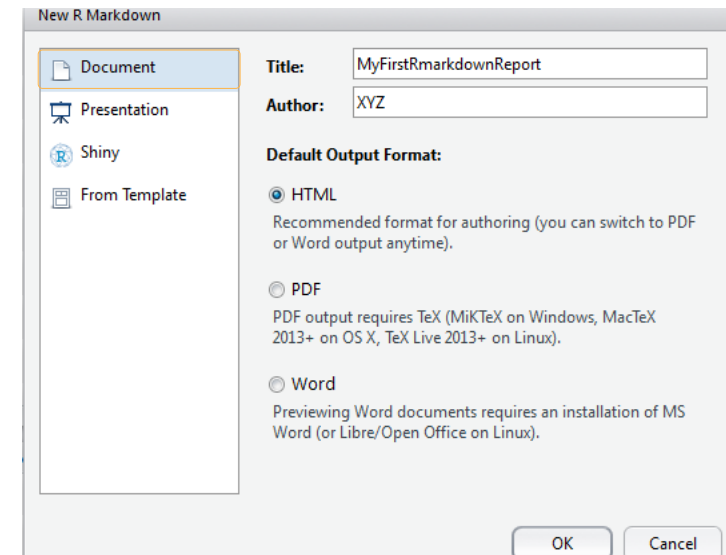
# Working with R Markdown

- Open a new R markdown file from the R Studio file option.

# Working with R Markdown

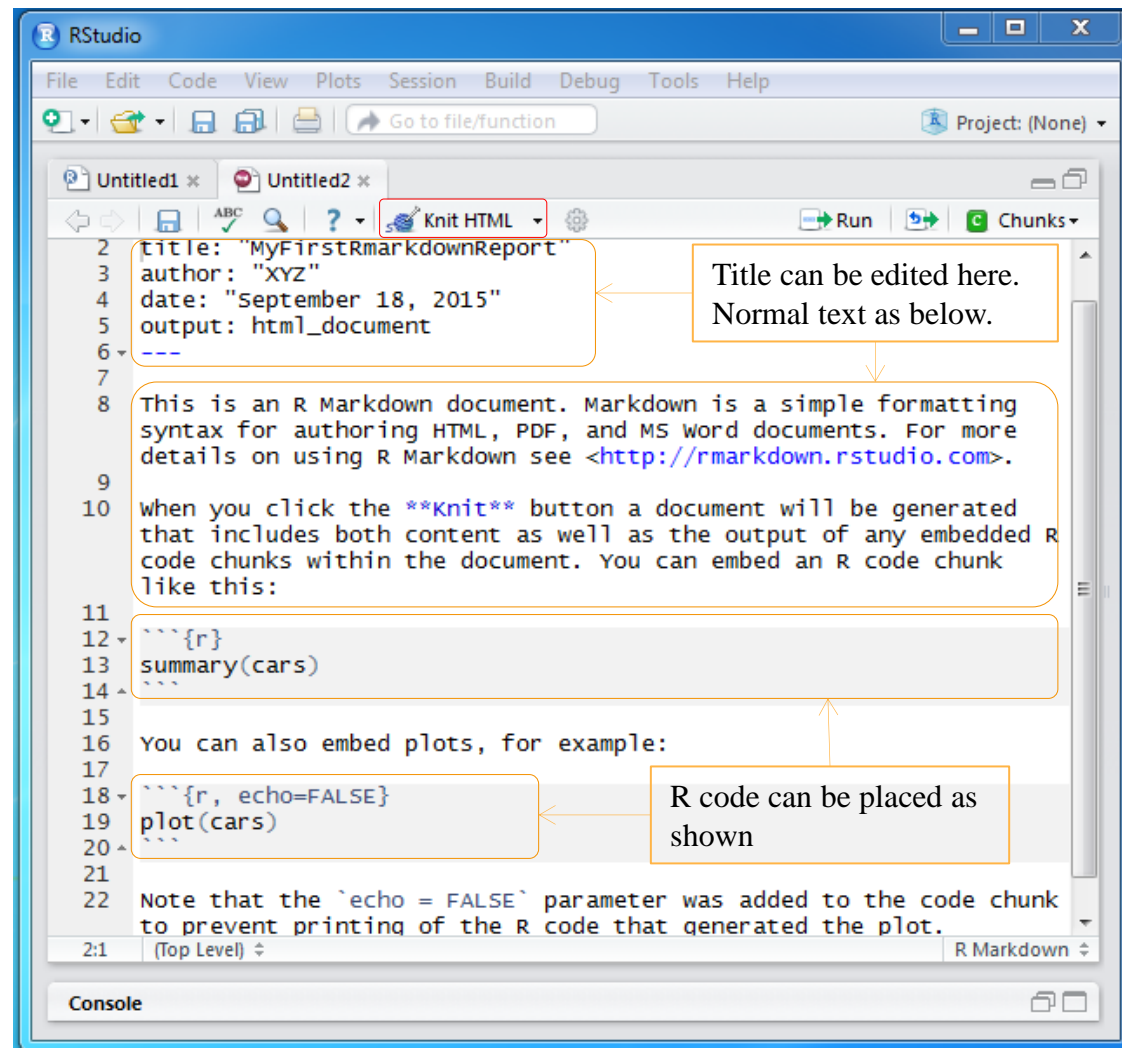- Select the type of report from the window that follows.

> - *Select 'Document' as the report type if creating an HTML, PDF or Word document.*
> - *Select 'Presentation' as the report type if creating HTML or PDF presentation.*
> - *Select 'Shiny' as the report type if creating an interactive shiny report.*
> - *There are specific templates which can be picked up to create report.*
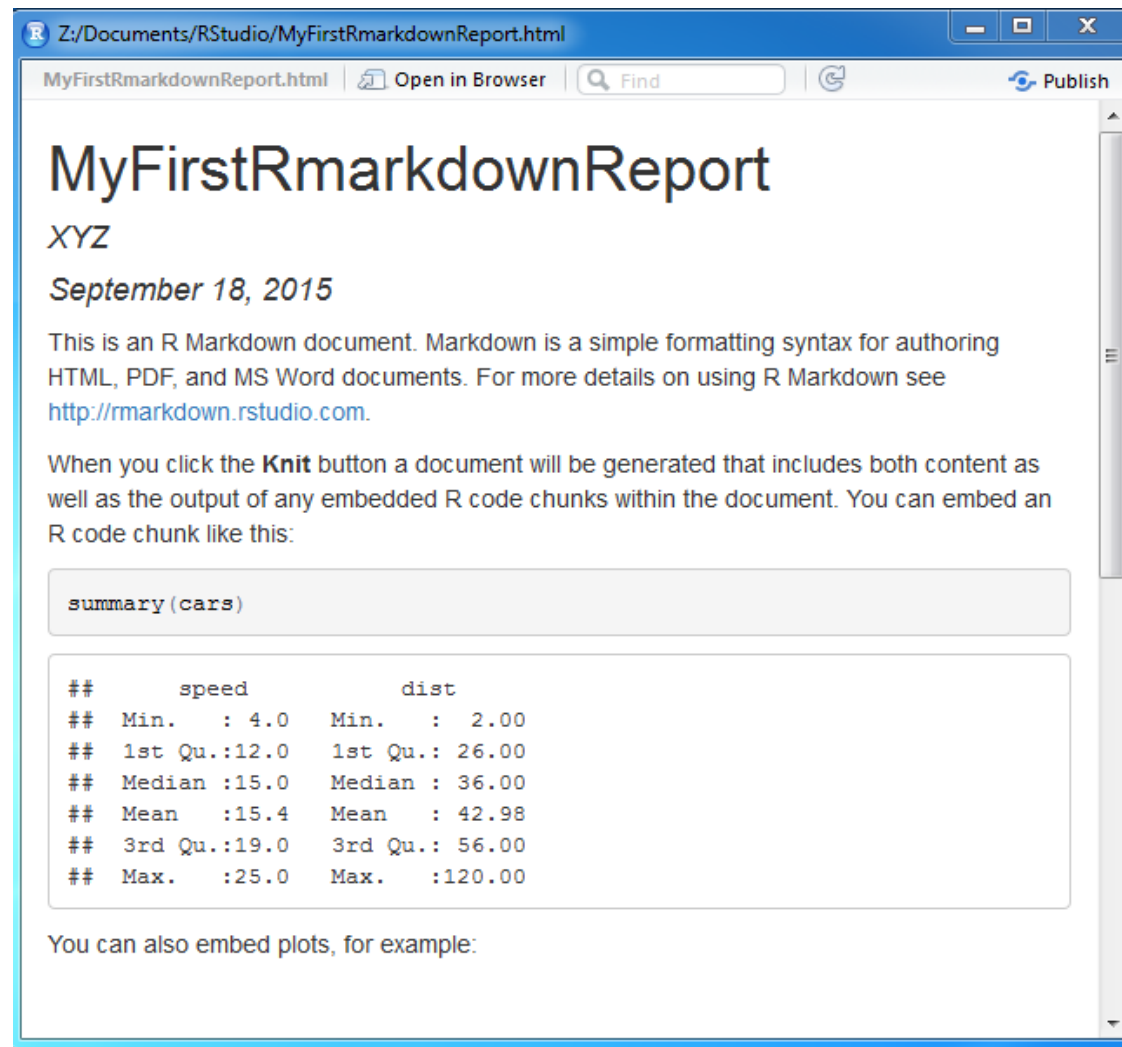
# My First R Markdown Code

- The R markdown code structure is simple to follow.
- Click on the Knit HTML icon to save the file.

- File gets saved with '.Rmd' extension in the current working directory.
- Report can be opened up in a separate window or inside the R Studio viewer.
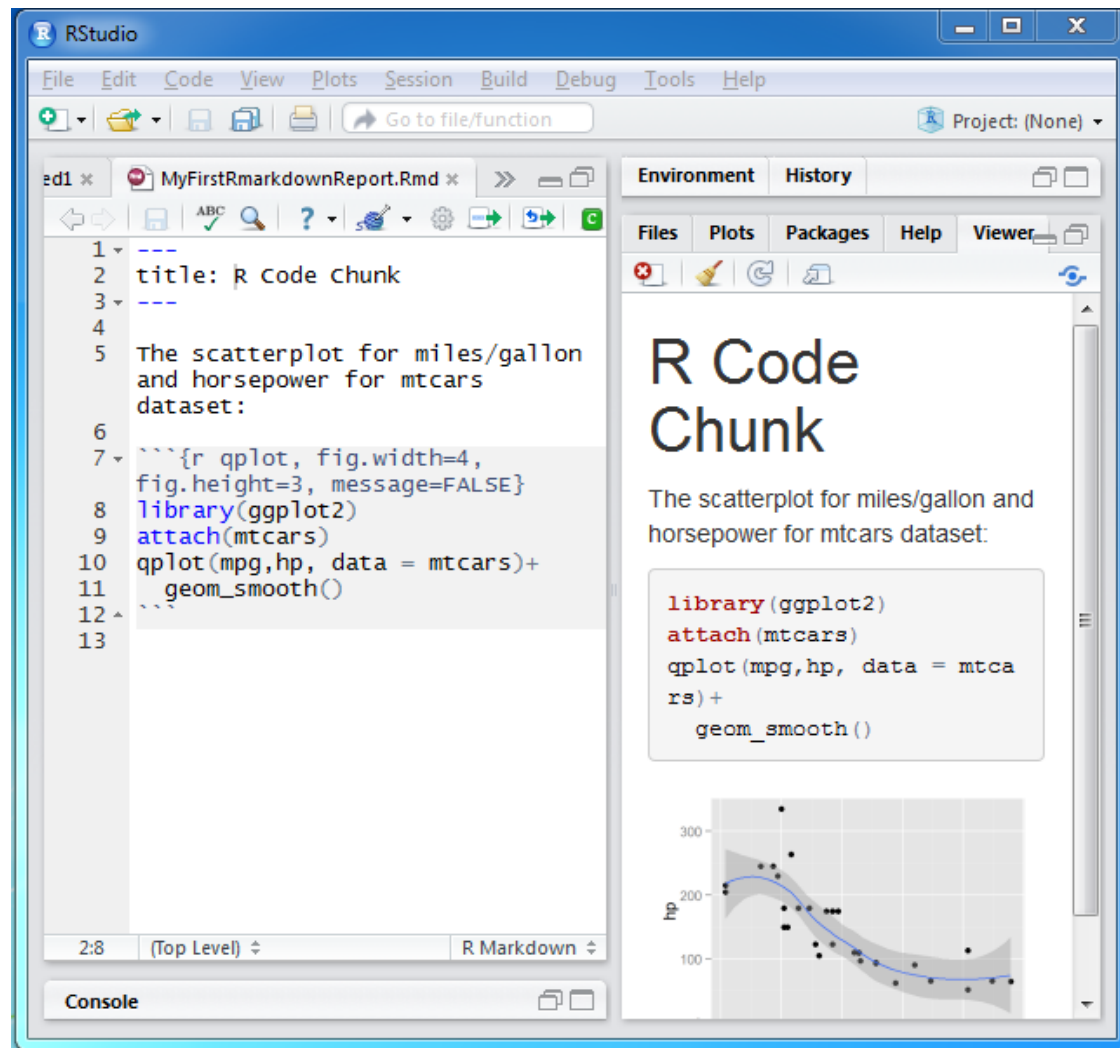
# My First R Markdown Report

- The report will look like a formatted report.

- Very sophisticated formatting can be applied on the text including writing equations, hyperlinks, appending images etc.

# R Markdown Code and Viewer

- The R code and viewer can be used side by side as a regular R scripting tool.
- The code for scatter plot and resulting output in the viewer is depicted here.

More on Rmarkdown at:
http://rmarkdown.rstudio.com/

# Demo of the RMarkdown using an example dataset.

# Rattle–An Introduction

- R Analytical Tool to Learn Easily (Rattle) is a user interface based data mining tool built on top of R.

> *On R Studio Console*:
> `>Install.packages("rattle")`
>
> *To force the installation of all dependency:*
> `>install.packages("rattle", dep=c("Suggests"))`
>
>
> *Or install using the Rstudio Install packages options*

- Rattle relies on extensive collection of R packages which powers the Rattle UI.

Dependent packages for Rattle are RGtk2, cairoDevice and XML. Troubleshooting at
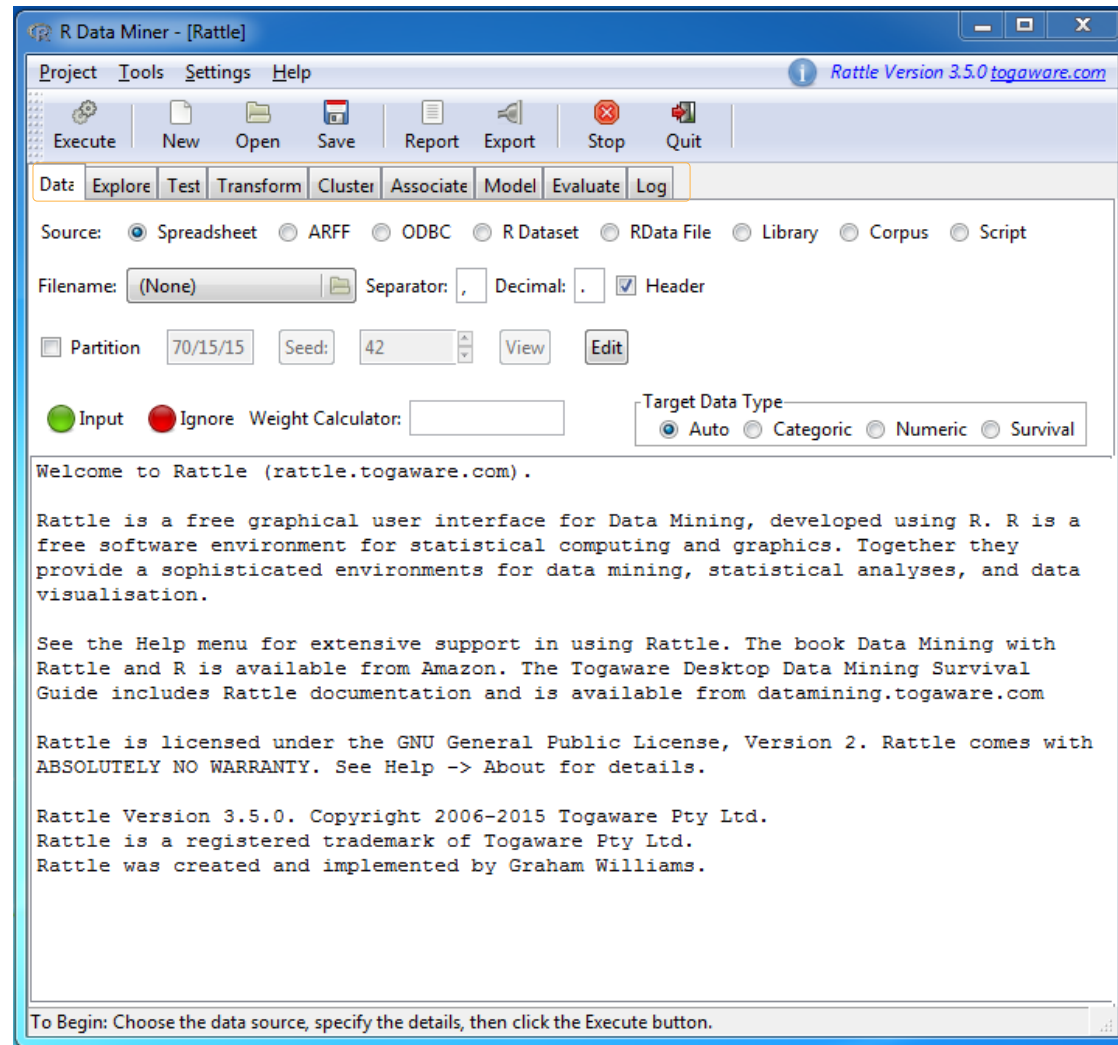http://rattle.togaware.com/rattle-install-troubleshooting.html

# Rattle User Interface

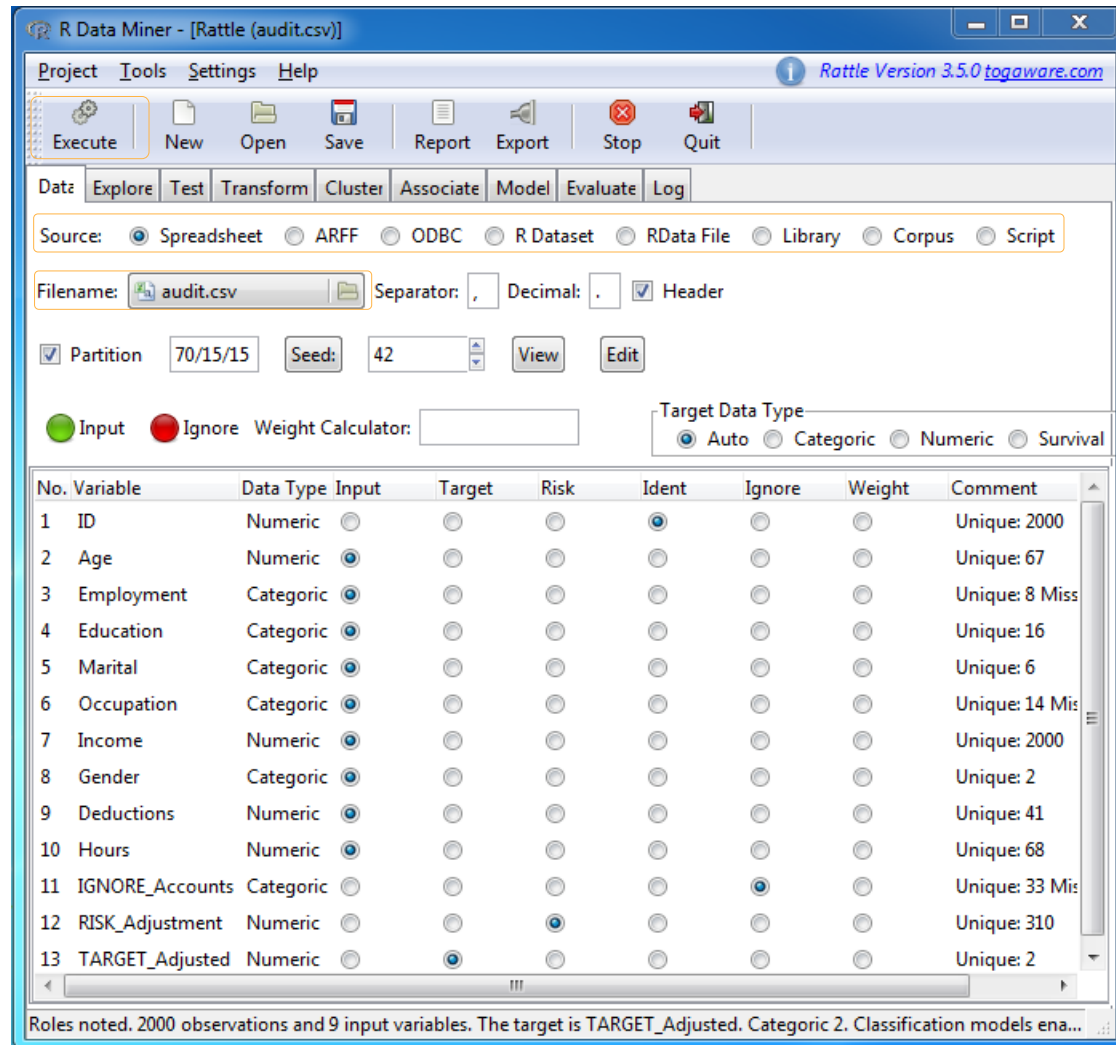- The user interface can be invoked as follows:

*On R Studio Console:*
>`library(rattle)`
>`rattle()`

- Tab based view with options to:
  - Load dataset
  - Explore dataset
  - Test distributions
  - Transform data
  - Clustering and association
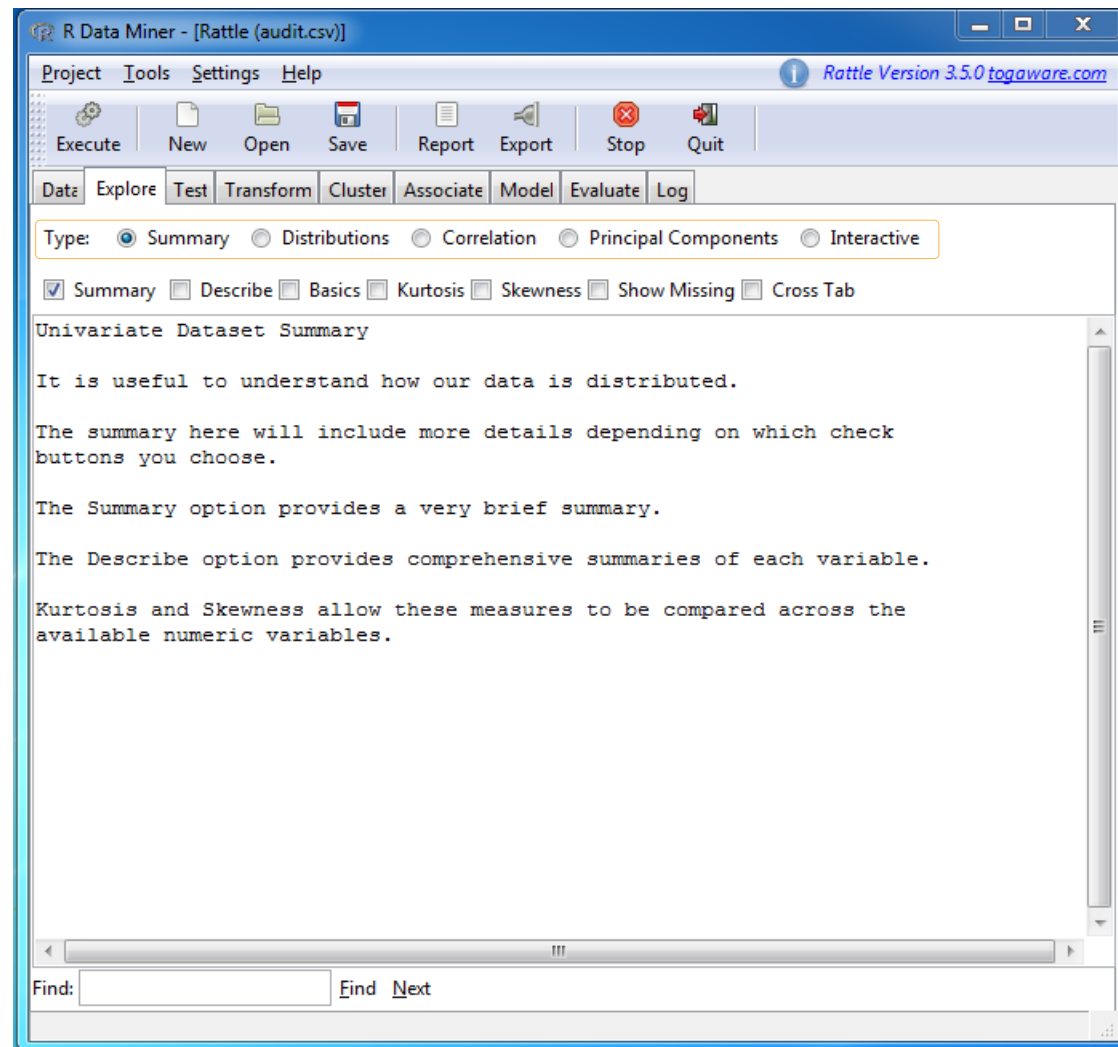  - Build models
  - Evaluate models
  - Code log

# Rattle–Load Dataset

- A dataset is executed by the execute command.

  - *If execute is clicked without any dataset, Rattle gives an option to load example dataset.*

  - *Rattle recognizes special pre-fixes for default variable role*
    - *'ID_'*
    - *'IGNORE_'*
    - *'RISK_' (measure of size of the target)*
    - *'IMP_'*
    - *'TARGET_'*

# Rattle–Explore Dataset

- Explore tab provides various options for exploratory data analysis

  - **Summary**: *Provides basic univariate summary and extended summary.*
  - **Distributions**: *Provide various plots for numeric as well as categorical data*
  - **Correlation**: *provides insights into the independence of the numeric input variables.*
  - **Principal component**: *Provides insight into the importance of variables in explaining the variation.*
  - **Interactive**: *Provides option for Interactive data exploration.*

# Rattle–Test Dataset

- Provides access to number of statistical tests of distributions.

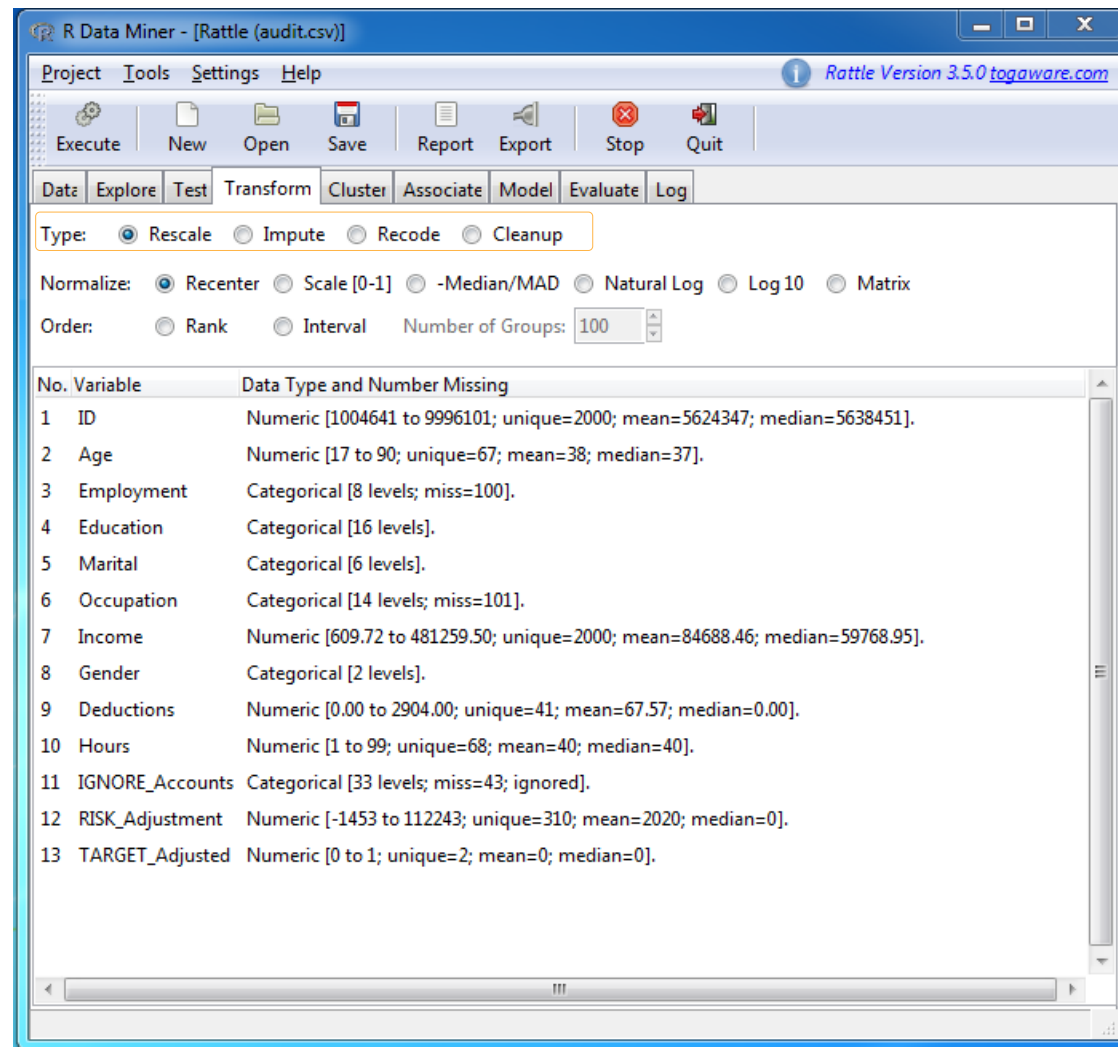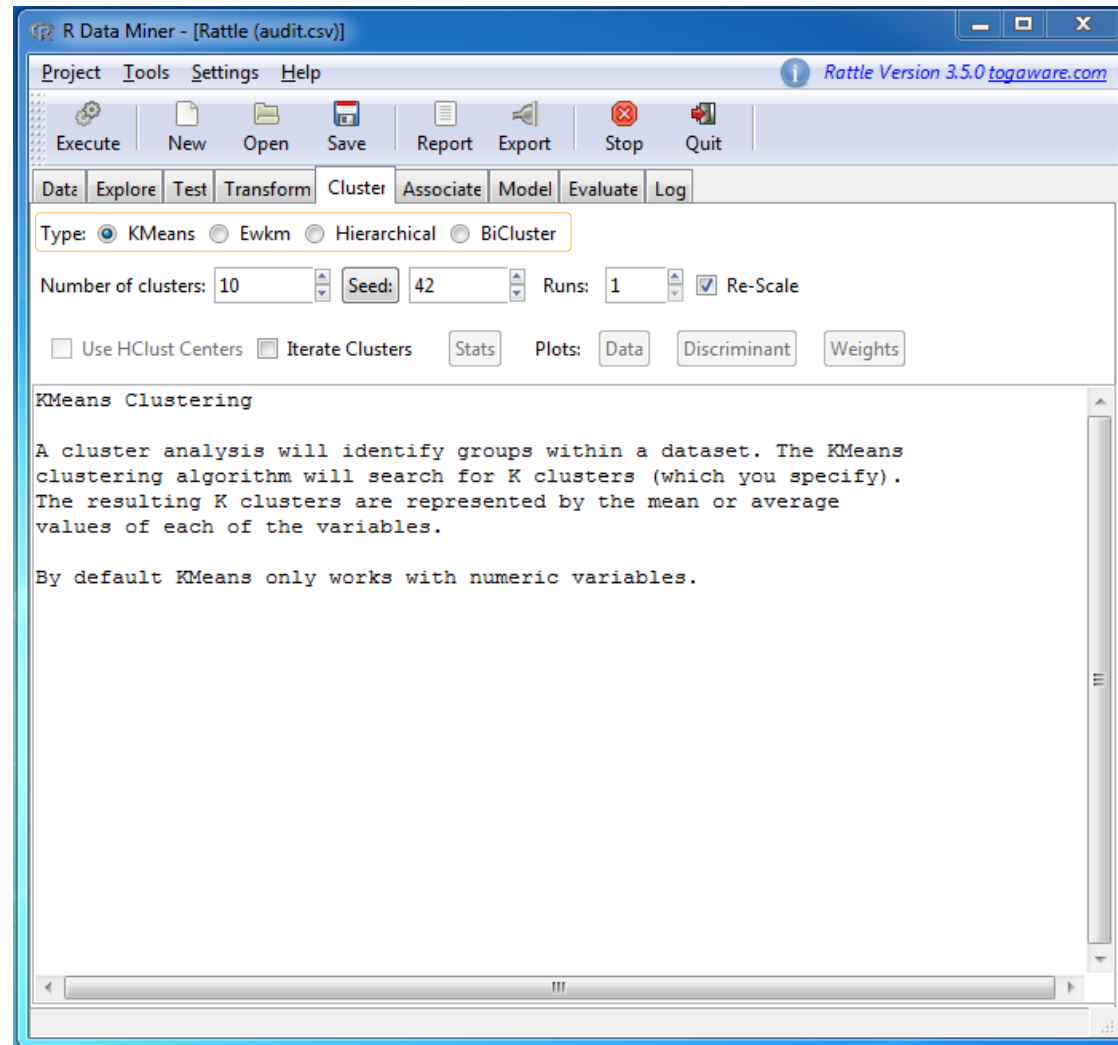# Rattle–Transform Dataset

- Cleaning data and creating new features (derived variables) takes significant time in data analysis.

  - ***Rescale****: Provides options for re-centering and scaling around zero.*
  - ***Impute****: Provides basic imputation of missing values using mean, median and mode.*
  - ***Recode****: Provides options for recoding/binning the variables with a default of 4 bins.*
  - ***Cleanup****: Provides option to treat the missing values after having tried imputation etc.*

R Data Miner - [Rattle (audit.csv)]

Rattle Version 3.5.0 togaware.com

Project   Tools   Settings   Help

Execute   New   Open   Save   Report   Export   Stop   Quit

Data   Explore   Test   Transform   Cluster   Associate   Model   Evaluate   Log

Type:   ● Rescale   ○ Impute   ○ Recode   ○ Cleanup

Normalize:   ● Recenter   ○ Scale [0-1]   ○ -Median/MAD   ○ Natural Log   ○ Log 10   ○ Matrix

Order:   ○ Rank   ○ Interval   Number of Groups: 100

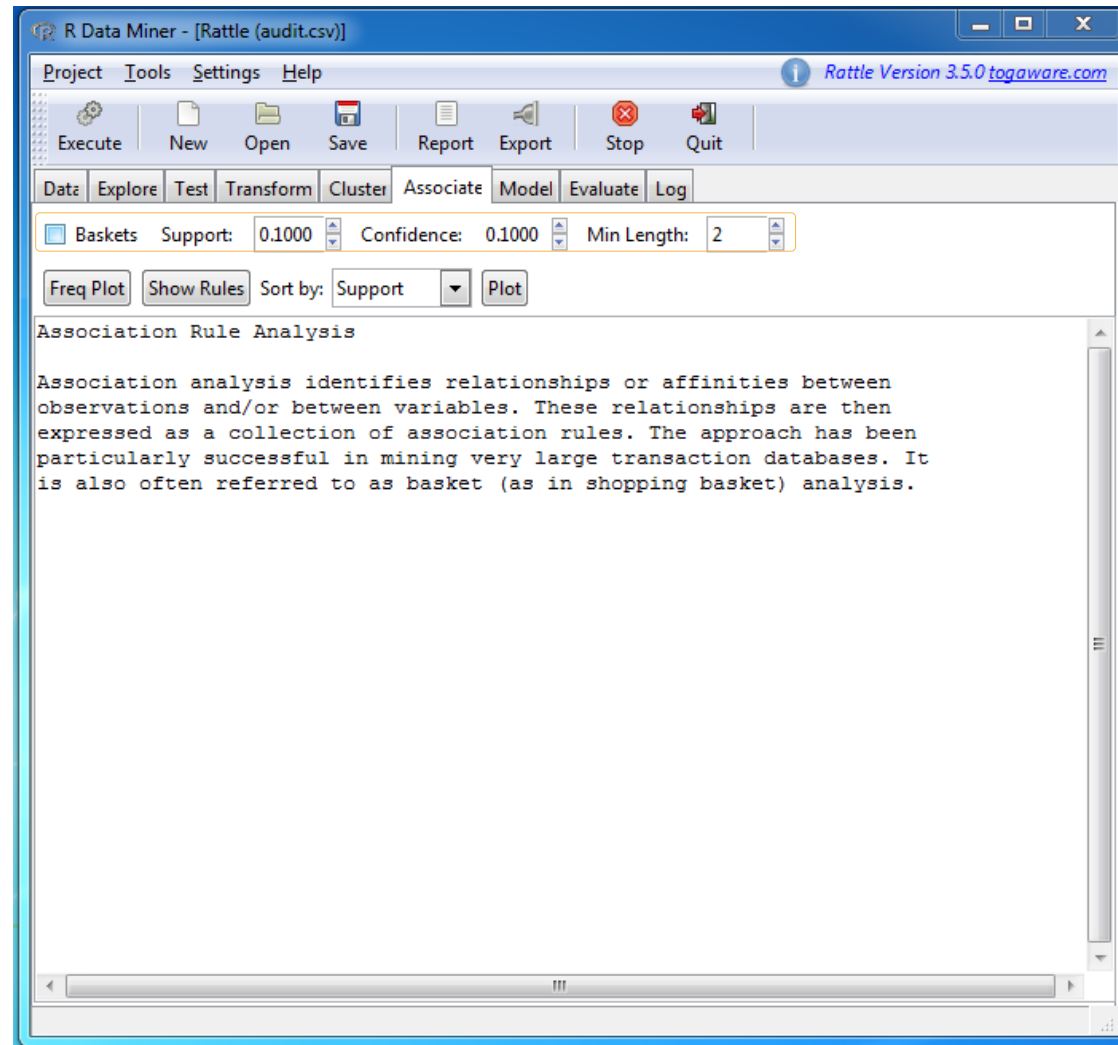| No. | Variable | Data Type and Number Missing |
|-----|----------|------------------------------|
| 1 | ID | Numeric [1004641 to 9996101; unique=2000; mean=5624347; median=5638451]. |
| 2 | Age | Numeric [17 to 90; unique=67; mean=38; median=37]. |
| 3 | Employment | Categorical [8 levels; miss=100]. |
| 4 | Education | Categorical [16 levels]. |
| 5 | Marital | Categorical [6 levels]. |
| 6 | Occupation | Categorical [14 levels; miss=101]. |
| 7 | Income | Numeric [609.72 to 481259.50; unique=2000; mean=84688.46; median=59768.95]. |
| 8 | Gender | Categorical [2 levels]. |
| 9 | Deductions | Numeric [0.00 to 2904.00; unique=41; mean=67.57; median=0.00]. |
| 10 | Hours | Numeric [1 to 99; unique=68; mean=40; median=40]. |
| 11 | IGNORE_Accounts | Categorical [33 levels; miss=43; ignored]. |
| 12 | RISK_Adjustment | Numeric [-1453 to 112243; unique=310; mean=2020; median=0]. |
| 13 | TARGET_Adjusted | Numeric [0 to 1; unique=2; mean=0; median=0]. |

# Rattle–Cluster Analysis

- Cluster tab provides option to build descriptive or unsupervised model.

- Several clustering algorithm available as options to identify groups within the dataset.
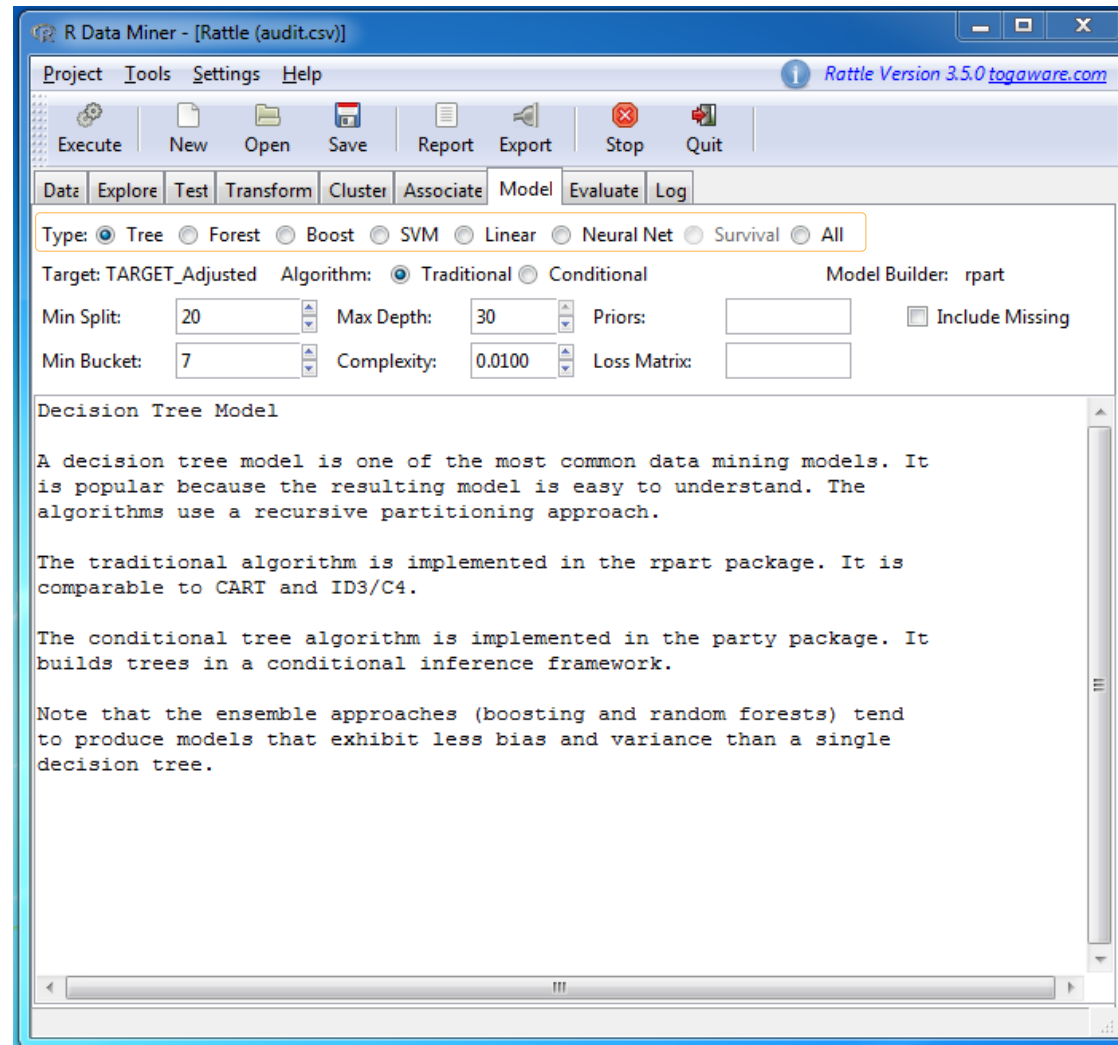
# Rattle–Basket Analysis

- Associate tab gives another option to build descriptive or unsupervised model.

- Option available for market basket analysis to identify affinities between observations and/or between variables.



R Data Miner - [Rattle (audit.csv)]

Rattle Version 3.5.0 togaware.com

Project   Tools   Settings   Help

Execute   New   Open   Save   Report   Export   Stop   Quit

Data   Explore   Test   Transform   Cluster   Associate   Model   Evaluate   Log

☐ Baskets   Support: 0.1000   Confidence: 0.1000   Min Length: 2

Freq Plot   Show Rules   Sort by: Support   Plot

Association Rule Analysis

Association analysis identifies relationships or affinities between observations and/or between variables. These relationships are then expressed as a collection of association rules. The approach has been particularly successful in mining very large transaction databases. It is also often referred to as basket (as in shopping basket) analysis.

# Rattle–Model Dataset

- Model tab provides a comprehensive list of techniques to build predictive models.

  - *Provides an option to use all the model building techniques over the same dataset.*
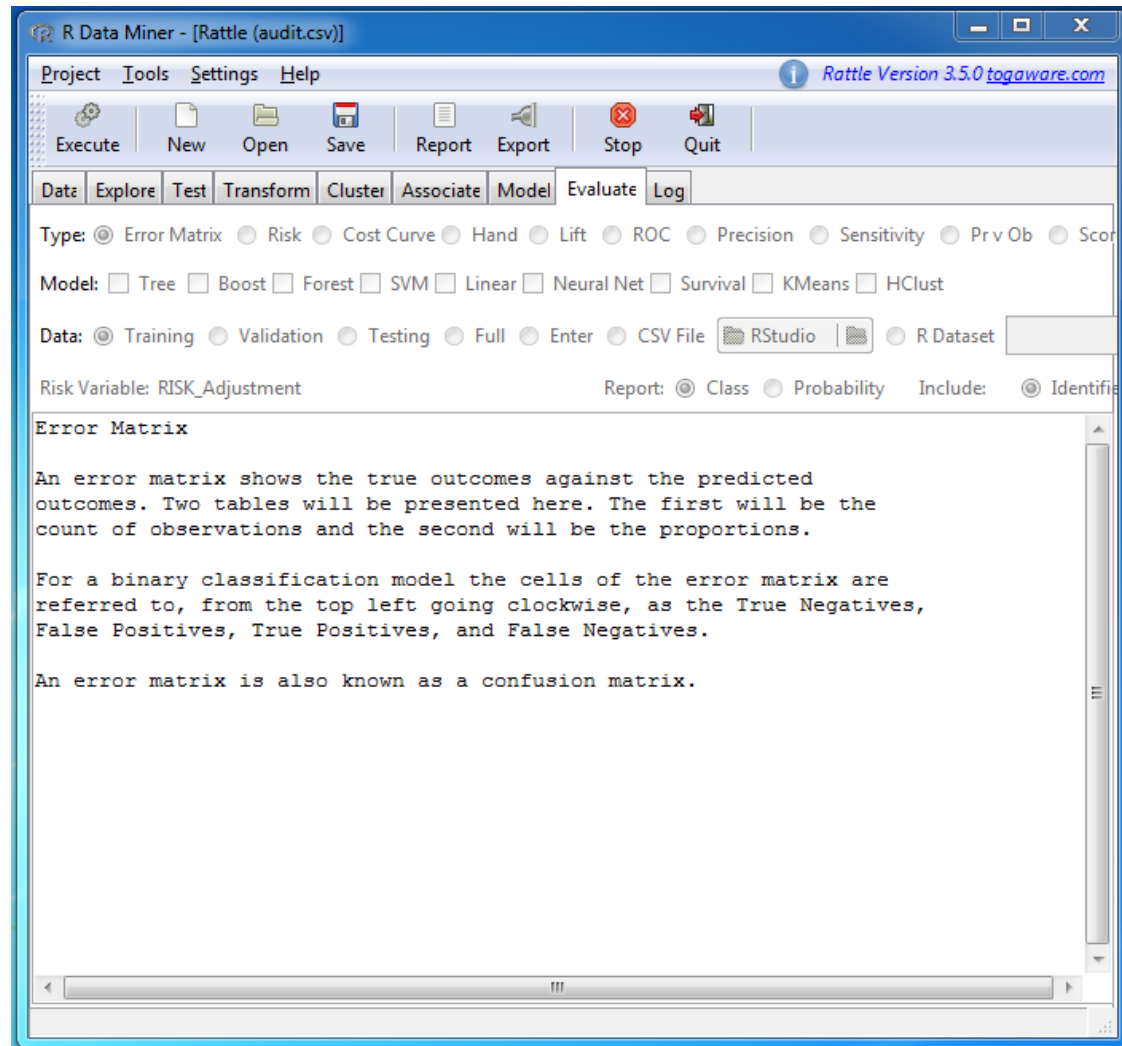  - *The models can be evaluated for performance and the best model can be selected.*

# Rattle–Evaluate Model

- Evaluate tab provides a collection of techniques for evaluating the performance of models

  - *Some of the commonly used techniques for model comparison can be seen as options:*
    - *Error matrix*
    - *ROC curve*
    - *Lift Chart*

  - *Rattle supports deployment of the model through the 'Score' option.*
    - *The complete model can be saved as a Rattle project and can later be used on the new dataset to score the*

# Rattle–Log Generation

- Log tab records the process of building the model.
- The recorded script gives the flexibility to fine tune the analysis using R directly.

- *The log can be used for deployment to score a new dataset.*

Demo of the Rattle tool using an example dataset.

# Summary

Summary of the topics covered in this lesson:

- R Analytical Tool to Learn Easily (Rattle) is a user interface based data mining tool built on top of R.

- Rattle provides a tab based options to load, explore, test, transform a dataset; followed by building and evaluating models.
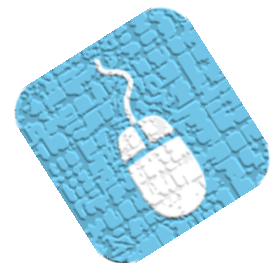
# QUIZ TIME

# Quiz Question 1

<table>
<tr><td>Quiz 1</td><td>What is the command line syntax to install rattle?<br><em>Select all that apply.</em></td></tr>
</table>

a.     *install.packages("rattle", dep=c("Suggests"))*

b.     *install.packages("rattle")*

c.     *install.package("rattle")*

d.     *install.package("rattle", dep=c("Suggests"))*

| | |
|---|---|
| **Quiz 1** | What is the command line syntax to install rattle? *Select all that apply.* |

a.     *install.packages("rattle", dep=c("Suggests"))*

b.     *install.packages("rattle")*

c.     *install.package("rattle")*

d.     *install.package("rattle", dep=c("Suggests"))*

**Correct answer is:**

*a & b*

Both a and b has the correct syntax. Option a has an optional argument of forcing the dependent packages to be installed.

# End of Lesson06–Introduction to R Markdown and Rattle