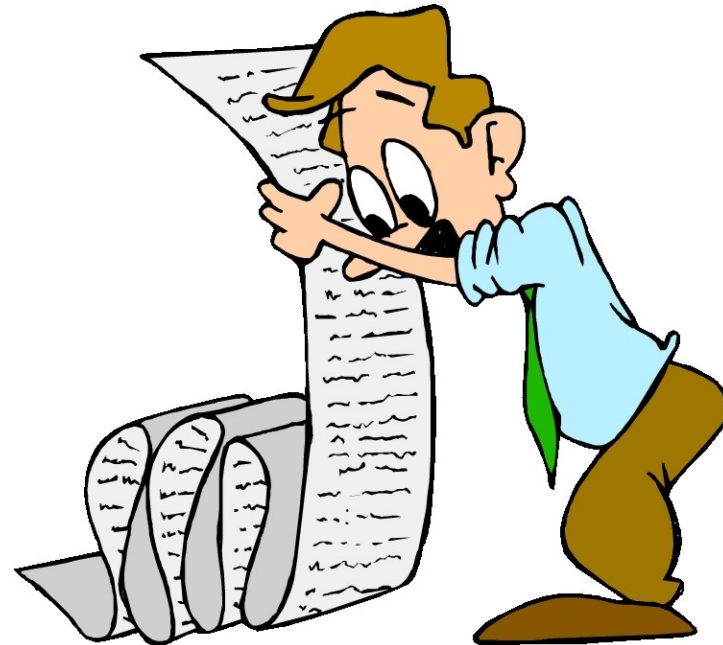The Data Co.™

# "Data Mining with Rattle and R"

**Colman McMahon**
colman@thedata.co

Linked in ®
www.linkedin.com/in/colmanmcmahon
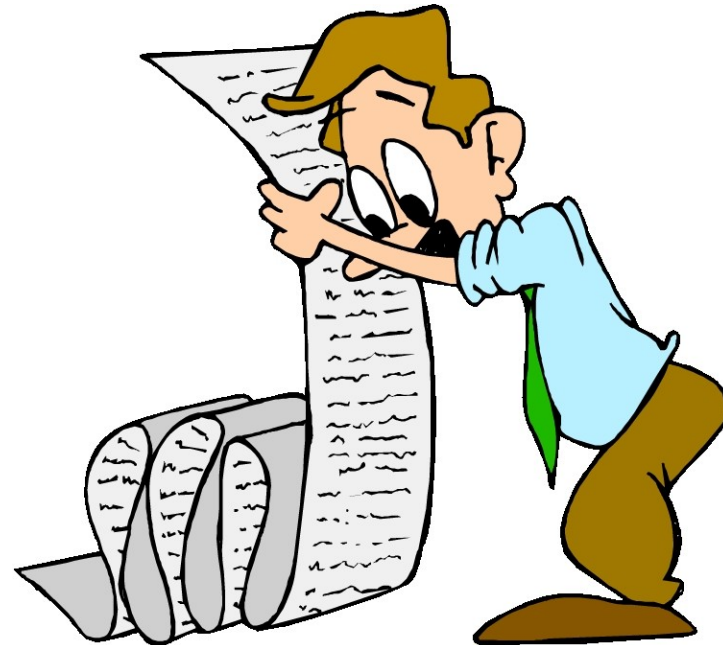
# Agenda

- ▶ Intro

- ▶ Data mining

- ▶ Rattle installation

- ▶ Rattle workflow

- ▶ Appendix

# Agenda

- ▶ **Intro**
- ▶ Data mining
- ▶ Rattle installation
- ▶ Rattle workflow
- ▶ Appendix

# Bio

▶ Professional

» Background: Film, VFX, Digital Media
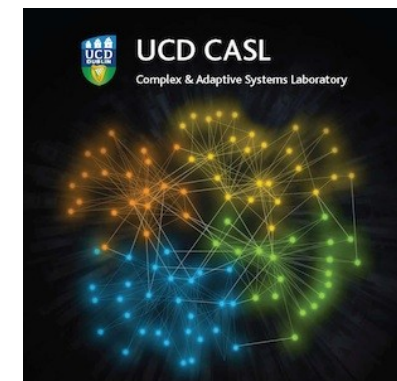
» Currently: PhD Fellow, UCD Dynamics Lab

▶ PhD:

» Policy network analysis

» Creative Industries

www.dl.ucd.ie

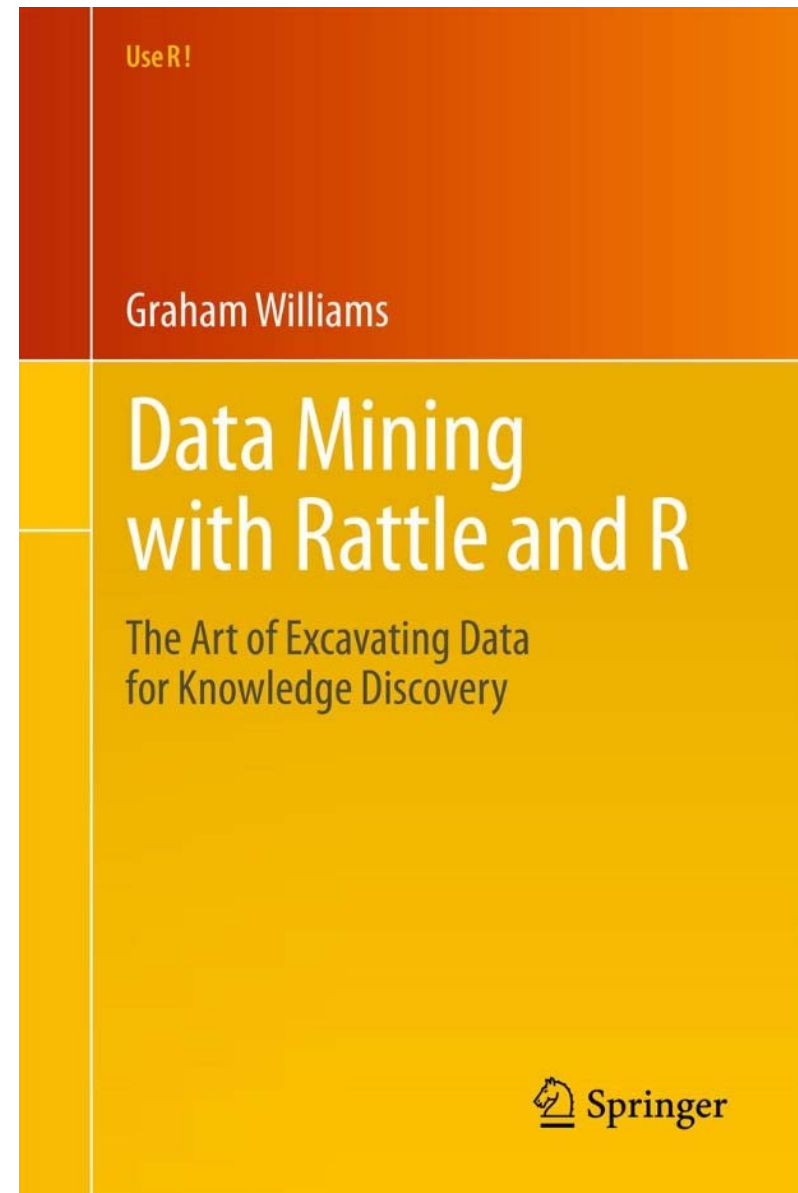▶ Technologies:

» Social Network Analysis

» Agent Based Modelling

» Data/statistical analysis

www.ucd.ie/casl/

▶ Data Mining with Rattle and R

» The Art of Excavating Data

for Knowledge Discovery

▶ Graham Williams

▶ Springer, 2011

▶ ISBN 978-1-4419-9889-7

# Rattle

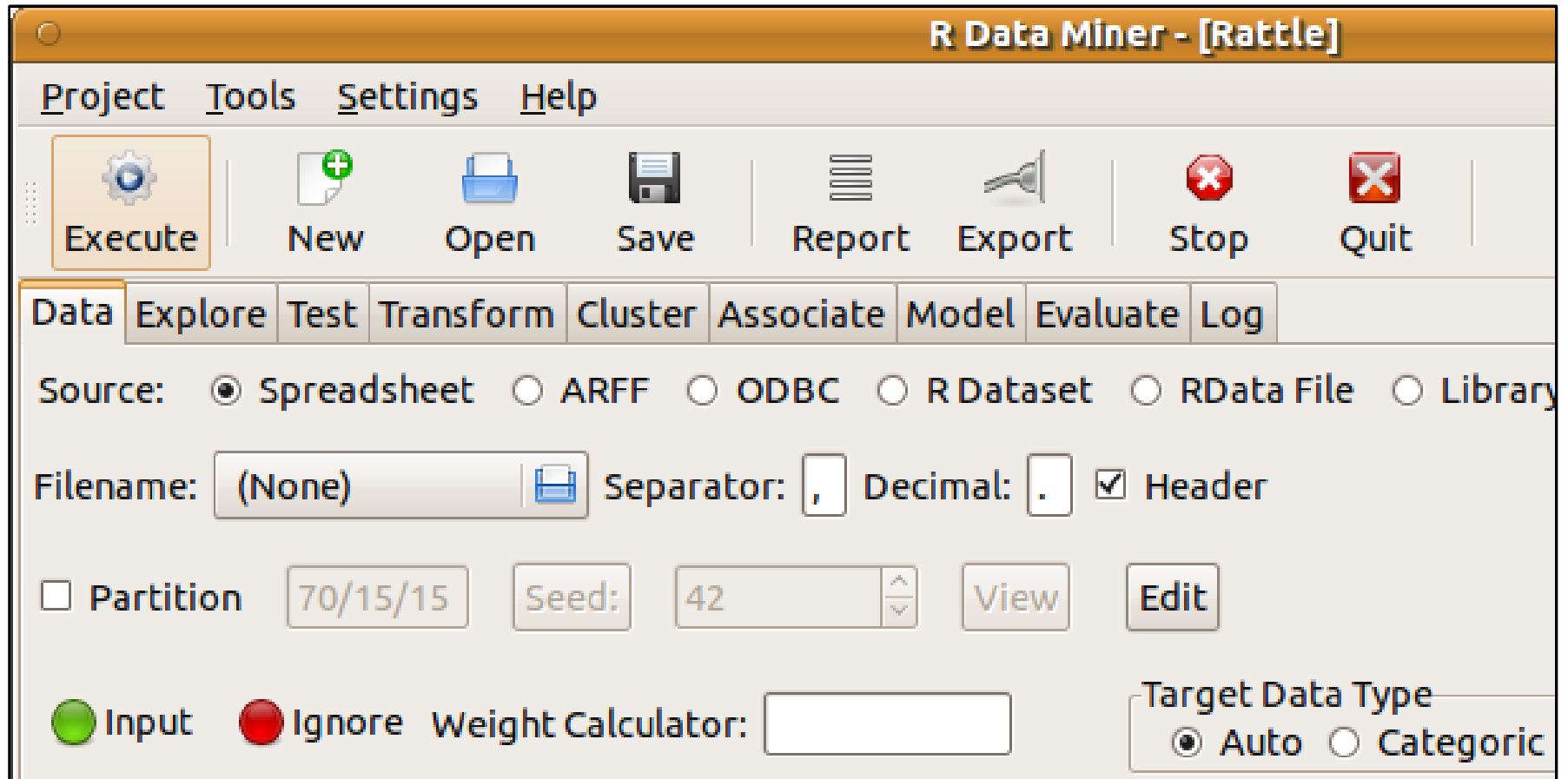**togaware**

## Rattle – "the R Analytical Tool To Learn Easily"

▶ Presents statistical & visual summaries of data

▶ Transforms data into forms that can be readily modelled

▶ Builds models (unsupervised and supervised) from the data

▶ Graphically presents the performance of models

▶ Scores new datasets

# Rattle

- Built on the statistical language R

  » an understanding of R is not required in order to use it

- Simple to use, quick to deploy, and allows us to rapidly work through the data processing, modelling, and evaluation phases of a data mining project

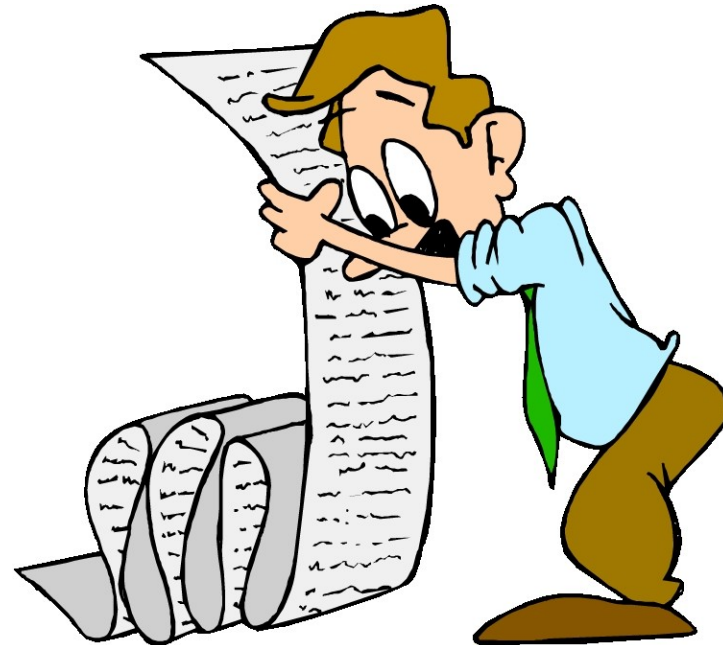- Can migrate from Rattle to R when we need to fine-tune and further develop our data mining projects

# Rattle GUI

# Agenda

- ▶ Intro
- ▶ Data mining
- ▶ Rattle stack
- ▶ Rattle workflow
- ▶ Appendix

# CRISP-DM

► Cross Industry Process for Data Mining (CRISP-DM, 1996)

» framework for delivering data mining projects.

1) **Problem Understanding**

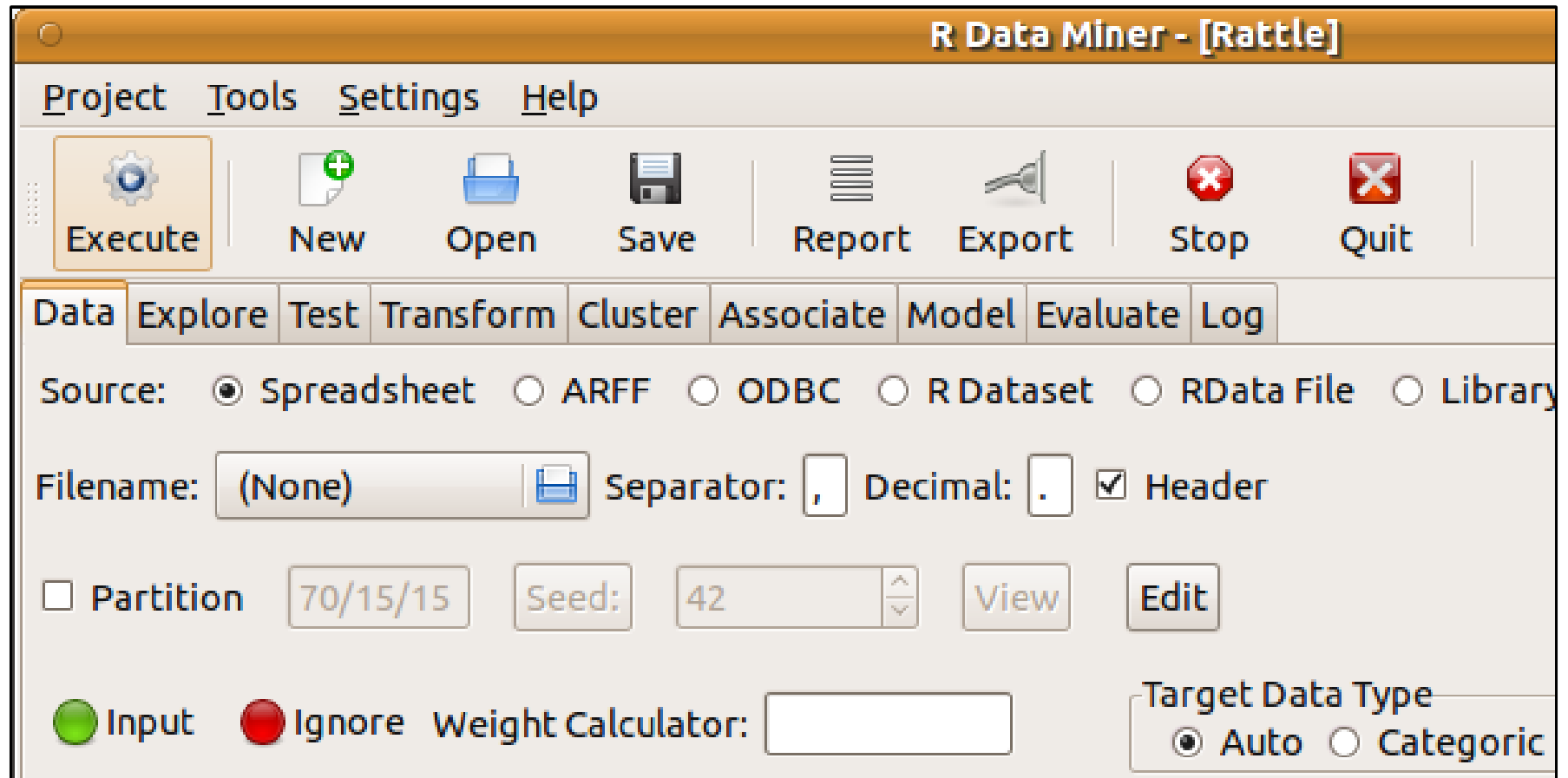2) **Data Understanding**

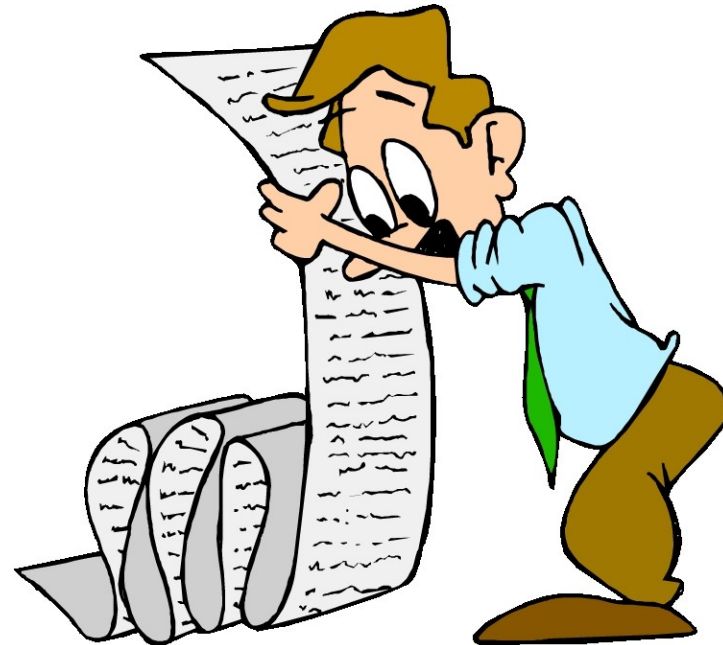3) **Data Preparation**

4) **Modelling**

5) **Evaluation**

6) **Deployment**

# Rattle GUI

# Agenda

- ▶ Intro

- ▶ Data mining

- ▶ Rattle installation

- ▶ Rattle workflow

- ▶ Appendix

# Rattle stack

| Rattle |
|:------:|
| **R** |

# R (www.r-project.org)

- R - a sophisticated statistical software package

  » easily installed, instructional, state-of-the-art, and it is free and open source

- The basic *modus operandi* - write scripts using the R language

- Steeper learning curve than using GUI based systems, but once over the hurdle, becomes relatively easy

# R Project (www.r-project.org)

# Rattle (rattle.togware.com)



www.thedata.co

# Installation

Install `R`

`www.r-project.org`


Start `R`


Install Rattle

`> install.packages("rattle")`


Load rattle into the `R` library

`> library(rattle)`

`> rattle()`

# Agenda

- ▶ Intro

- ▶ Data mining
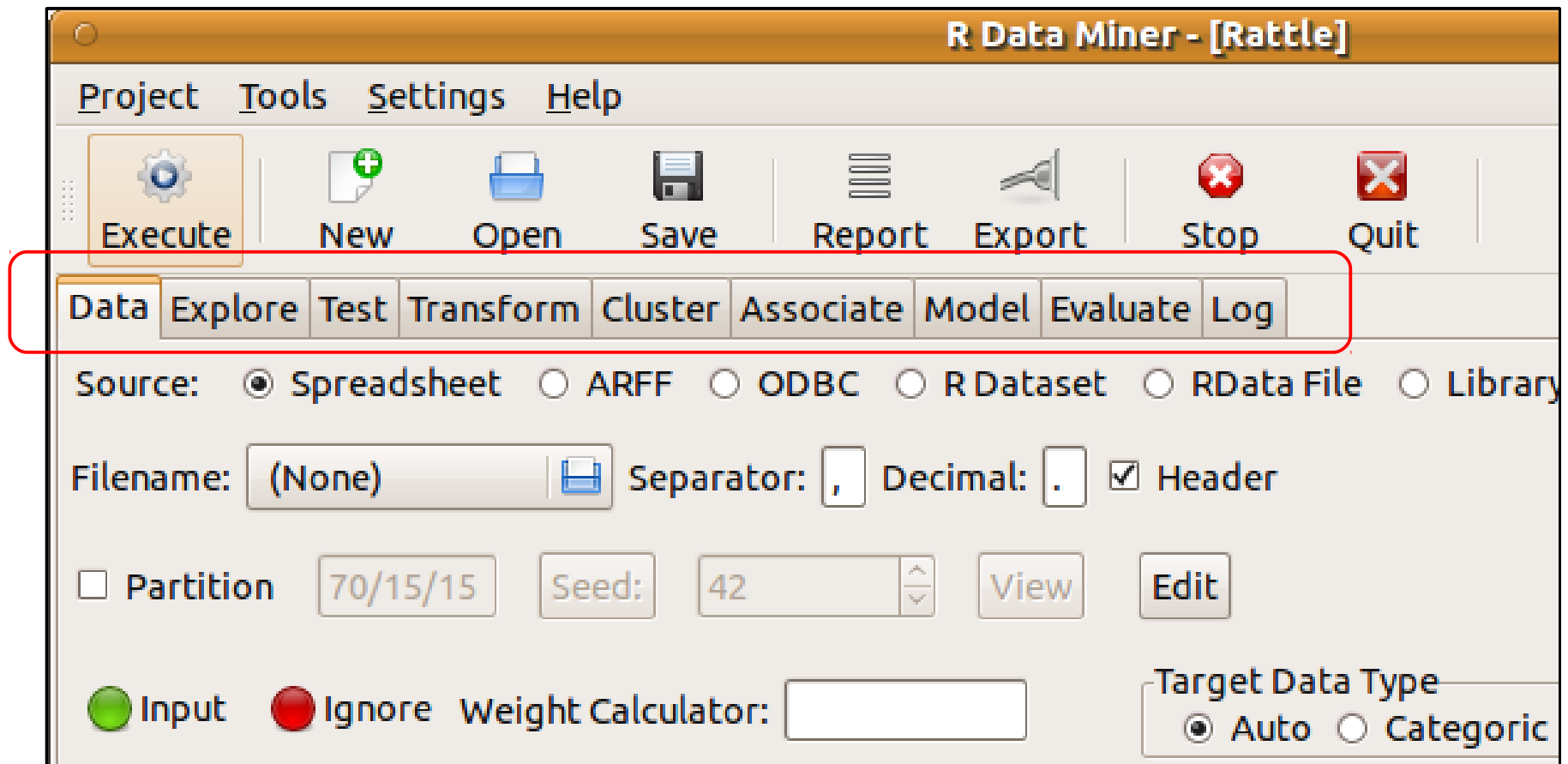
- ▶ Rattle installation

- ▶ Rattle workflow

- ▶ Appendix

# Rattle workflow

1) Load a dataset

2) Select variables and entities for exploring and mining

3) Explore the data to understand how it is distributed or spread

4) Transform the data to suit our data mining purposes

5) Build our models

6) Evaluate the models on other datasets

7) Export the models for deployment

# Workflow - tabs

www.thedata.co

www.thedata.co

# Explore



```
> summary(weather[7:9])

     Sunshine           WindGustDir      WindGustSpeed
 Min.    : 0.00     NW        : 73     Min.    :13.0
 1st Qu.: 5.95     NNW       : 44     1st Qu.:31.0
 Median : 8.60     E         : 37     Median :39.0
 Mean    : 7.91     WNW       : 35     Mean    :39.8
 3rd Qu.:10.50     ENE       : 30     3rd Qu.:46.0
 Max.    :13.60     (Other):144     Max.    :98.0
 NA's    : 3.00     NA's      :  3     NA's    : 2.0
```

▶ Provides a textual overview of the data

# Explore – detailed

```
> library(fBasics)
> basicStats(weather$Sunshine)

            X..weather.Sunshine
nobs                 366.0000
NAs                    3.0000
Minimum                0.0000
Maximum               13.6000
1. Quartile            5.9500
3. Quartile           10.5000
Mean                   7.9094
Median                 8.6000
Sum                 2871.1000
SE Mean                0.1827
LCL Mean               7.5500
UCL Mean               8.2687
Variance              12.1210
Stdev                  3.4815
Skewness              -0.7235
Kurtosis              -0.2706
```

# Explore (graphically)

# Explore (graphically)

Figure 6.3: The parallel coordinates plot from latticist

Figure 6.9: Colourful brushing of multiple scatterplots

Figure 6.10: GGobi's scatter plot matrix.

Figure 6.11: GGobi's parallel coordinates plot.

# Test



► various statistical tests, e.g. the T-test and F-test

# Transform



- ▶ Normalising

- ▶ Filling in missing values

- ▶ Turning numeric variables into categoric variables (and vice versa)

- ▶ Dealing with outliers

- ▶ Removing variables or observations with missing values

# Cluster



Cluster sizes:

[1] "23 17 22 22 17 36 23 34 22 32"

Data means:

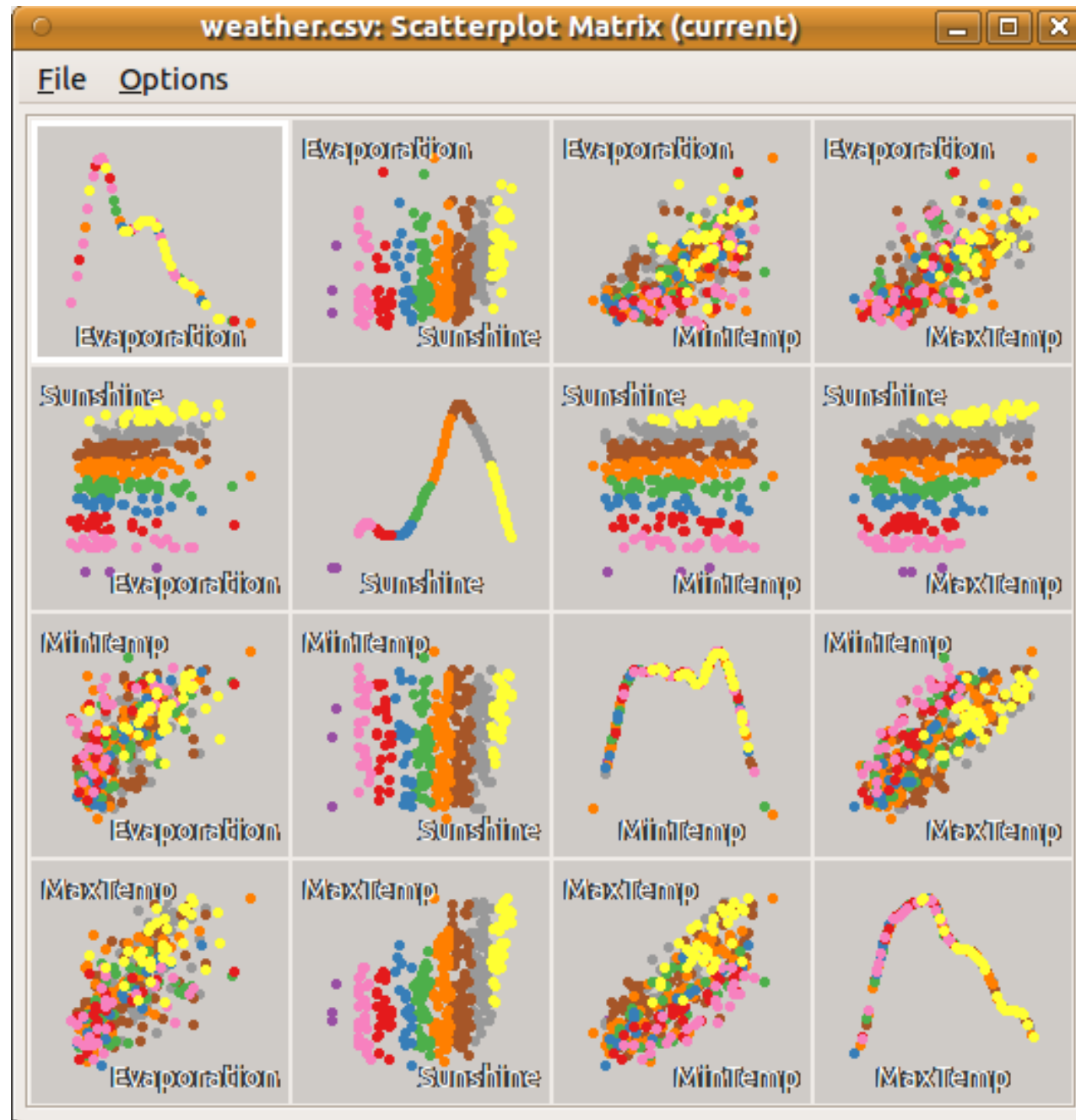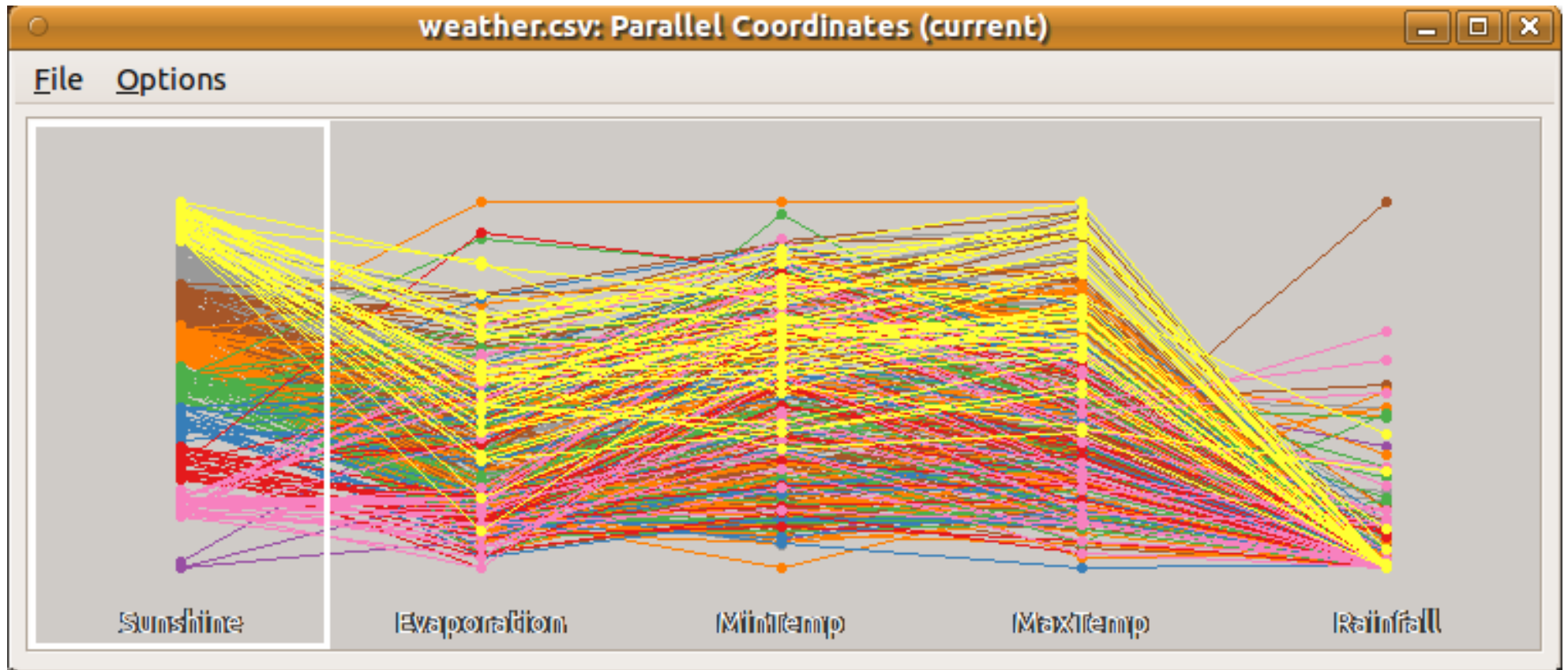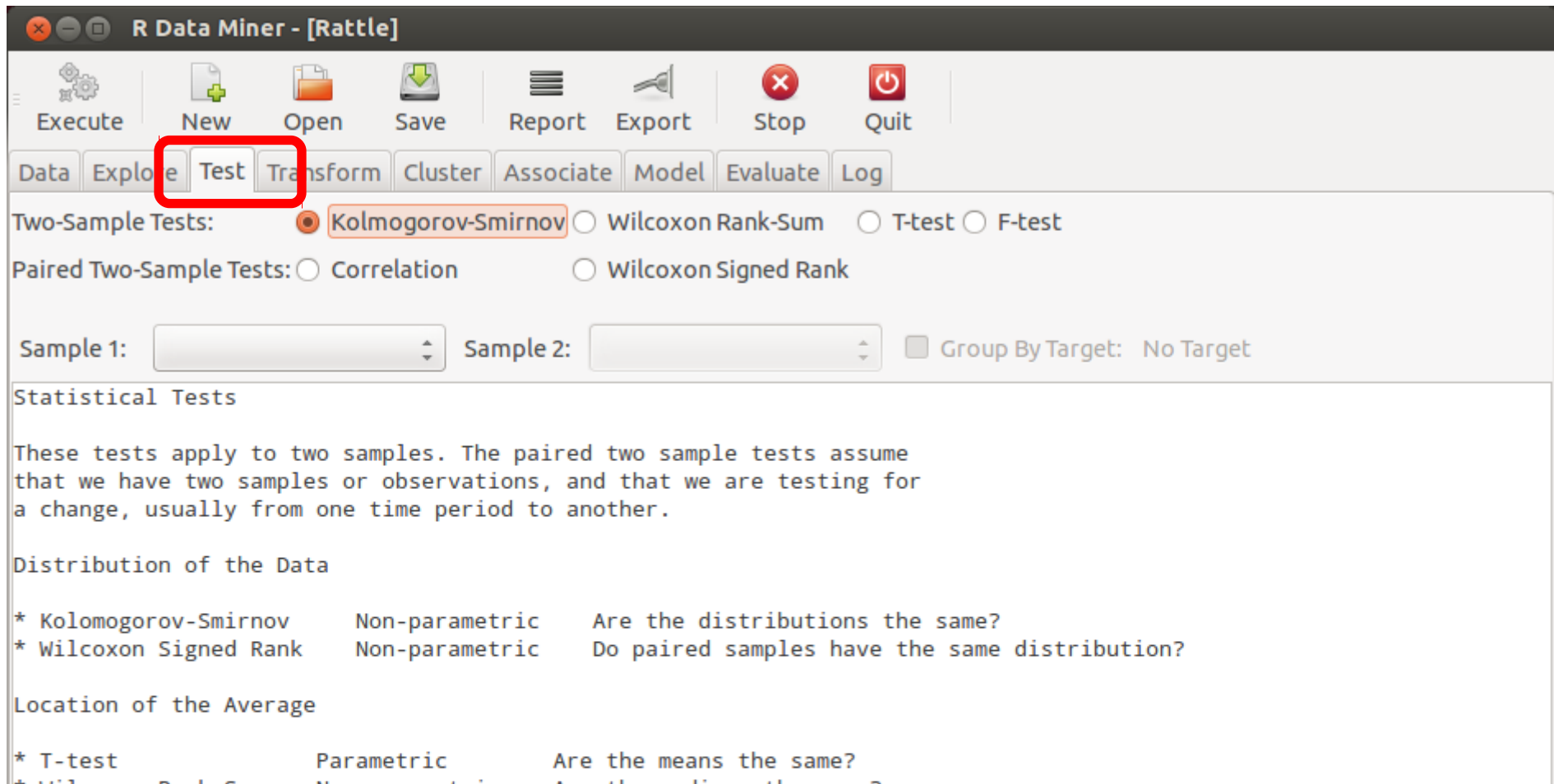|  |  |  |  |
|---|---|---|---|
| MinTemp | MaxTemp | Rainfall | Evaporation |
| 7.146 | 20.372 | 1.377 | 4.544 |
| Sunshine | WindGustSpeed | WindSpeed9am | WindSpeed3pm |
| 8.083 | 39.944 | 9.819 | 18.056 |
| Humidity9am | Humidity3pm | Pressure9am | Pressure3pm |
| 71.472 | 43.859 | 1019.748 | 1016.979 |
| Cloud9am | Cloud3pm | Temp9am | Temp3pm |
| 3.690 | 3.851 | 12.269 | 19.081 |

Cluster centers:

| | MinTemp | MaxTemp | Rainfall | Evaporation | Sunshine |
|---|---|---|---|---|---|
| 1 | 8.5000 | 21.05 | 1.27826 | 6.330 | 10.496 |
| 2 | 11.6059 | 30.95 | 0.11765 | 7.647 | 11.276 |
| 3 | 13.4136 | 28.77 | 1.02727 | 6.200 | 9.464 |
| 4 | 9.1818 | 16.90 | 4.94545 | 3.800 | 2.191 |
| 5 | 7.7412 | 15.19 | 3.58824 | 3.306 | 5.659 |

• Allows data miners to break data into more meaningful groups and then contrast the different clusters against each other

# Associate

# Model

# Decision Trees



Rattle: Plot 2

**Decision Tree weather.csv $ RainTomorrow**

Pressure3pm >= < 1011.9

Cloud3pm < => 7.5

Sunshine >= < 8.85

| 4 | 5 | 6 | 7 |
| --- | --- | --- | --- |
| No | Yes | No | Yes |
| 195 obs | 9 obs | 25 obs | 27 obs |
| 94.9% | 66.7% | 80% | 74.1% |

Save  Print  Close

# Boosting



The Ada Boost model has been built. Time taken: 1.62 secs

# Boosting



Figure 13.3: The variable importance plot for a boosted model.

# Evaluate



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

Project  Tools  Settings  Help

Rattle Version 2.6.7 togaware.com

Execute  New  Open  Save  Report  Export  Stop  Quit

Data  Explore  Test  Transform  Cluster  Associate  Model  **Evaluate**  Log

Type: ● **Error Matrix**  ○ Risk  ○ Cost Curve  ○ Hand  ○ Lift  ○ ROC  ○ Precision  ○ Sensitivity  ○ Pr v Ob  ○ Score

Model: ☑ **Tree**  ☐ Boost  ☐ Forest  ☐ SVM  ☐ Linear  ☐ Neural Net  ☐ Survival  ☐ KMeans  ☐ HClust

Data: ○ Training  ○ Validation  ● **Testing**  ○ Full  ○ Enter  ○ CSV File  ☐ dm...  ☐  ○ R Dataset

Risk Variable: RISK_MM

Report: ● Class  ○ Probability  Include:  ● Identifiers  ○ All

```
Error matrix for the Decision Tree model on weather.csv [test] (counts):

       Predicted
Actual No Yes
   No  35   6
   Yes  5  10

Error matrix for the Decision Tree model on weather.csv [test] (%):

       Predicted
Actual No Yes
   No  62  11
   Yes  9  18

Overall error: 0.1964

===================================================================
```

Generated confusion matrix.

## Evaluation Using the Training Dataset:

| Count | | Predict No | Predict Yes | | Percentage | | Predict No | Predict Yes |
|---|---|---|---|---|---|---|---|---|
| Actual | No | 205 | 10 | | Actual | No | 80 | 4 |
| | Yes | 15 | 26 | | | Yes | 6 | 10 |

## Evaluation Using the Validation Dataset:

| Count | | Predict No | Predict Yes | | Percentage | | Predict No | Predict Yes |
|---|---|---|---|---|---|---|---|---|
| Actual | No | 39 | 5 | | Actual | No | 72 | 9 |
| | Yes | 5 | 5 | | | Yes | 9 | 9 |

## Evaluation Using the Testing Dataset:

| Count | | Predict No | Predict Yes | | Percentage | | Predict No | Predict Yes |
|---|---|---|---|---|---|---|---|---|
| Actual | No | 35 | 6 | | Actual | No | 62 | 11 |
| | Yes | 5 | 10 | | | Yes | 9 | 18 |

## Evaluation Using the Full Dataset:

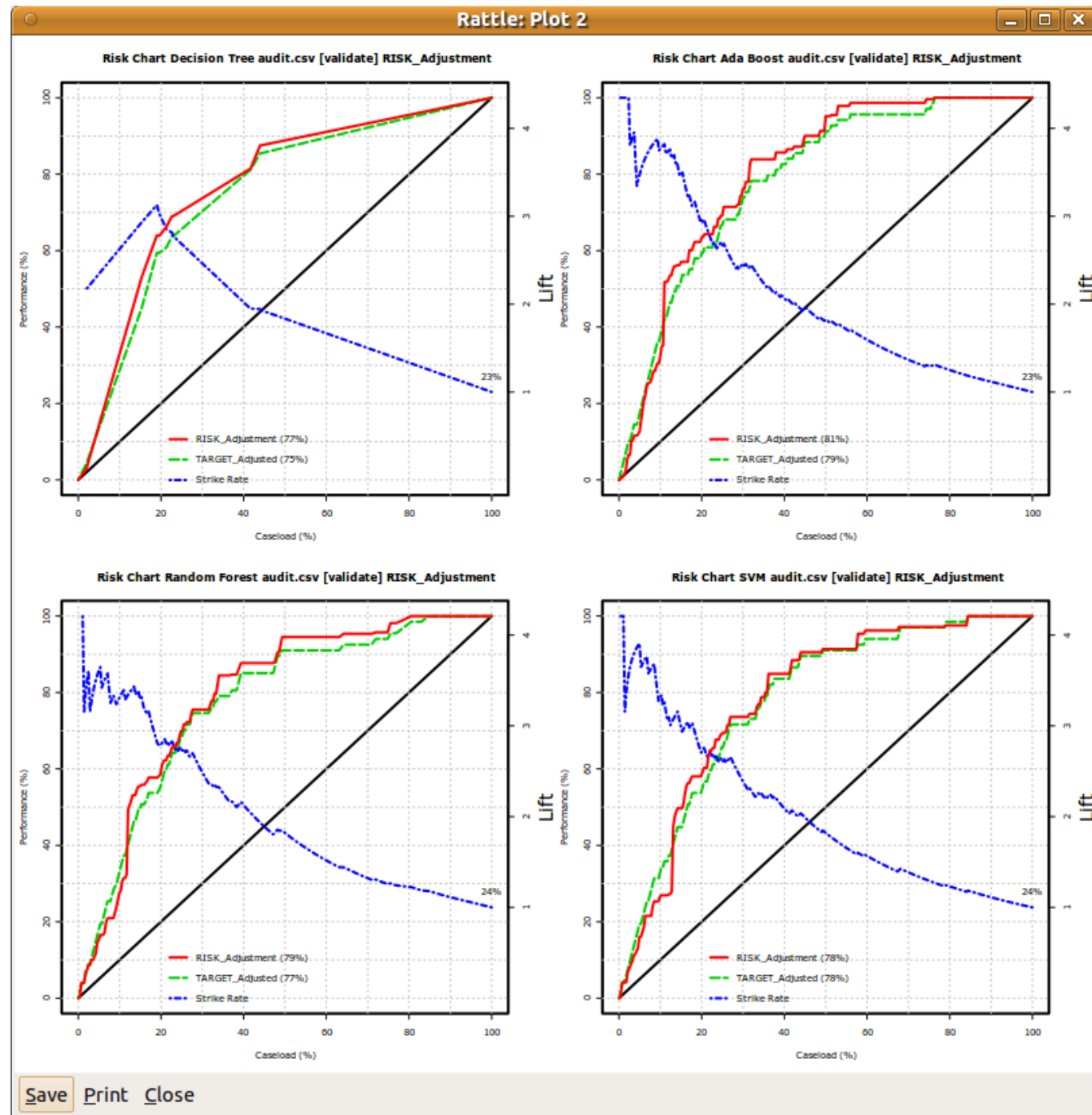| Count | | Predict No | Predict Yes | | Percentage | | Predict No | Predict Yes |
|---|---|---|---|---|---|---|---|---|
| Actual | No | 279 | 21 | | Actual | No | 76 | 6 |
| | Yes | 25 | 41 | | | Yes | 7 | 11 |

Figure 15.5: Four risk charts displayed to compare performances of multiple model builders on the audit dataset.

# First Model
## (example)

# First model

▶ Once we have processed our data, we can build a model

1) Click on the Execute button

- Rattle will notice that no dataset has been identified

2) The sample "weather" dataset will be offered

- Click "Yes"

3) Click on the Model tab

- This is where we tell Rattle what kind of model we want to build

4) Click on the Execute button.

# Building a Model

Figure 2.5

Figure 2.6: The target variable, RainTomorrow, is skewed, with Yes being quite under-represented

# Appendix

# Installation

Install `R`

```
www.r-project.org
```

```
Start R
```

Install Rattle

```
> install.packages("rattle")
```

Load rattle into the `R` library

```
> library(rattle)
```

```
> rattle()
```

# Articles on Rattle

- **Rattle: A Data Mining GUI for R**

- http://journal.r-project.org/archive/2009-2/RJournal_2009-2_Williams.pdf

- **Getting started with data mining in R using Rattle**

- http://techpad.co.uk/content.php?sid=240

# R Resources (sample)

► **Data Manipulation with R** (Spector, 2008) – covers basic data structures, reading and writing data, subscripting, manipulating, aggregating, and reshaping data

► **Introductory Statistics with R** (Dalgaard, 2008) - good introduction to statistics using R.

► **Modern Applied Statistics with S** (Venables and Ripley, 2002 - an extensive introduction to statistics using R.

► **Data Analysis and Graphics Using R** (Maindonald and Braun, 2007) -  excellent practical coverage of many aspects of exploring and modelling data using R

► **The Elements of Statistical Learning** (Hastie et al., 2009) is a more mathematical treatise, covering all of the machine learning techniques discussed in this book in quite some mathematical depth.
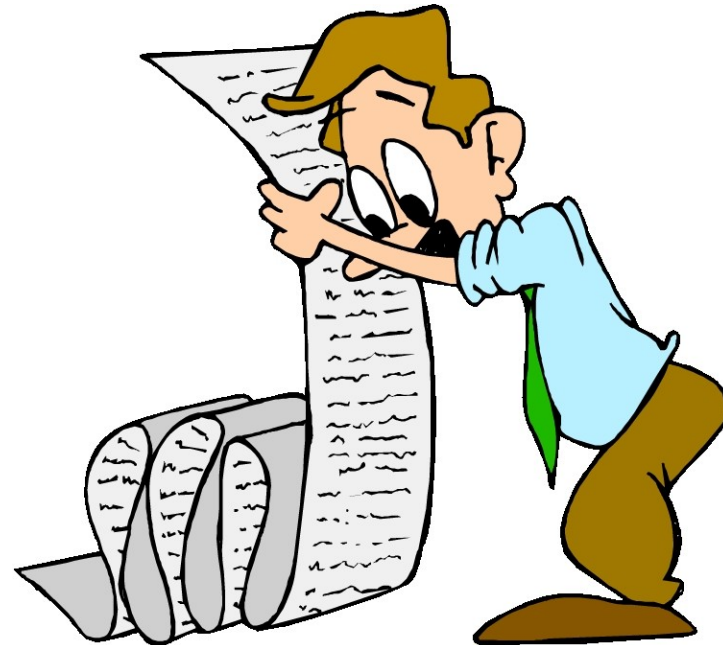
► **R for SAS and SPSS Users** (Muenchen, 2008) is an excellent choice

► **Lattice: Multivariate Data Visualization with R** (Sarkar, 2008) - covers

► the extensive capabilities of one of the graphics/plotting packages available for R.

► **ggplot2: Elegant Graphics for Data Analysis** (Wickham, 2009) - newer graphics framework is detailed

► **Bivand et al.** (2008) cover applied spatial data analysis,

► **Kleiber and Zeileis** (2008) cover applied econometrics

► **Cowpertwait and Metcalfe** (2009) cover time series

# Agenda

► Intro

► Data mining

► Rattle installation

► Rattle workflow

► Appendix

**The Data Co.™**

# "Data Mining with Rattle and R"

DIT Analytics Club
07 March, 2013
"Peadar Kearney's", Dame Street

**Colman McMahon**
colman@thedata.co

Linked in ®
www.linkedin.com/in/colmanmcmahon