## An Introductory Tutorial

### Data Analytics using R

Graham.Williams@togaware.com

Senior Director and Data Scientist, Analytics
Australian Taxation Office

Adjunct Professor, Australian National University
Fellow, Institute of Analytics Professionals of Australia

Graham.Williams@togaware.com
http://datamining.togaware.com

Visit: `http://onepager.togaware.com` for Tutorial Notes

---

## Tutorial Overview

1. Motivating R – A Language for Data Mining

2. Data Mining in R – Hands-on Rattle GUI

3. Programming Data in R – Scripting our Analyses

4. Disseminate Research in R – Ensembles and wsrf

---

## Tutorial Overview

1. Motivating R – A Language for Data Mining

2. Data Mining in R – Hands-on Rattle GUI

3. Programming Data in R – Scripting our Analyses

4. Disseminate Research in R – Ensembles and wsrf

---

## Installing R

- Instructions on Togaware: `http://rattle.togaware.com`
- Visit CRAN: `http://mirrors.ustc.edu.cn/CRAN/`
- Linux: Install package for your distribution
  $ `wajig install r-recommended` (Debian/Ubuntu)
- Windows: Download and install from CRAN
- MacOSX: Download and install from CRAN

---

## Why Big Data and Ensembles with R?

- Most widely used Data Mining and Machine Learning Package
  - Machine Learning
  - Statistics
  - Software Engineering and Programming with Data
  - Not the nicest of languages for a computer scientist
- Free (Libre) Open Source Statistical Software
  - . . . all modern statistical approaches
  - . . . many/most machine learning algorithms
  - . . . opportunity to readily add new algorithms
- That is important for us in the research community
  Get our algorithms out there and being used—impact!!!

---

## Why Big Data and Ensembles with R?

- Most widely used Data Mining and Machine Learning Package
  - Machine Learning
  - Statistics
  - Software Engineering and Programming with Data
  - Not the nicest of languages for a computer scientist
- Free (Libre) Open Source Statistical Software
  - . . . all modern statistical approaches
  - . . . many/most machine learning algorithms
  - . . . opportunity to readily add new algorithms
- That is important for us in the research community
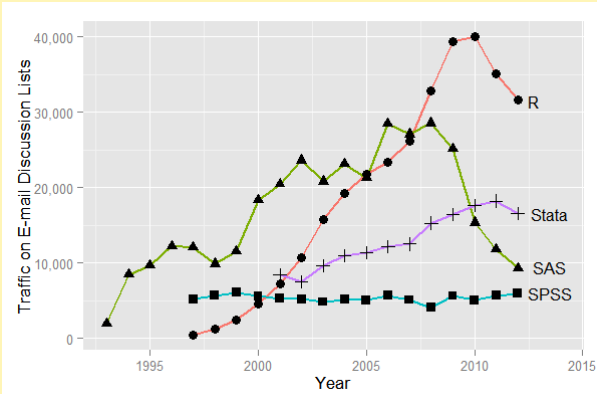  Get our algorithms out there and being used—impact!!!

## Why Big Data and Ensembles with R?

- Most widely used Data Mining and Machine Learning Package
  - Machine Learning
  - Statistics
  - Software Engineering and Programming with Data
  - Not the nicest of languages for a computer scientist
- Free (Libre) Open Source Statistical Software
  - ...all modern statistical approaches
  - ...many/most machine learning algorithms
  - ...opportunity to readily add new algorithms
- That is important for us in the research community
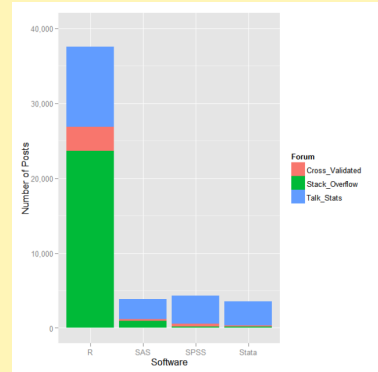  Get our algorithms out there and being used—impact!!!

## Why Big Data and Ensembles with R?

- Most widely used Data Mining and Machine Learning Package
  - Machine Learning
  - Statistics
  - Software Engineering and Programming with Data
  - Not the nicest of languages for a computer scientist
- Free (Libre) Open Source Statistical Software
  - ...all modern statistical approaches
  - ...many/most machine learning algorithms
  - ...opportunity to readily add new algorithms
- That is important for us in the research community
  Get our algorithms out there and being used—impact!!!

## How Popular is R? Discussion List Traffic

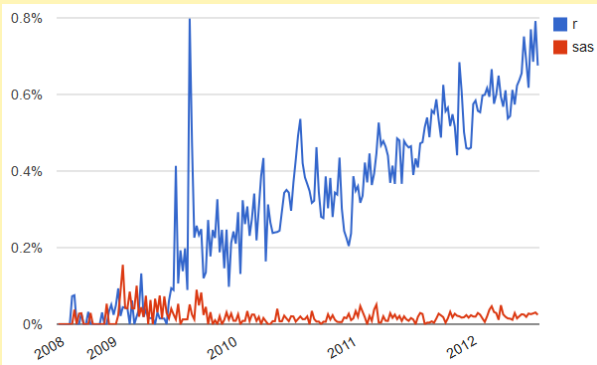Monthly email traffic on software's main discussion list.

## How Popular is R? Discussion Topics

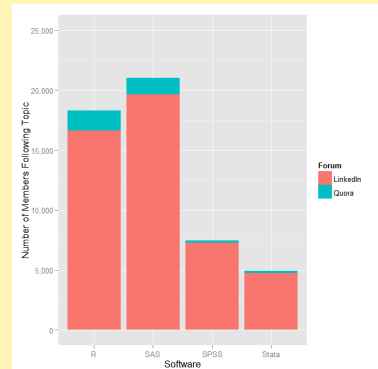Number of discussions on popular QandA forums 2013.

## How Popular is R? R versus SAS

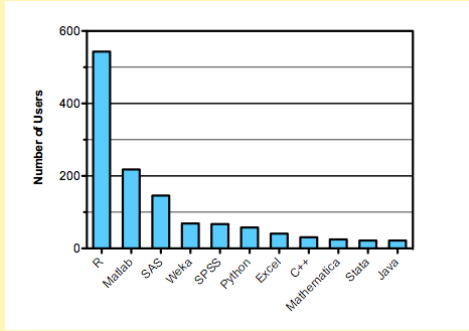Number of R/SAS related posts to Stack Overflow by week.

## How Popular is R? Professional Forums

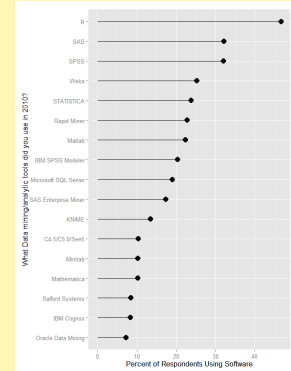Registered for the main discussion group for each software.

# HOW POPULAR IS R? USED IN ANALYTICS COMPETITIONS

Software used in data analysis competitions in 2011.

---

# HOW POPULAR IS R? USER SURVEY

Rexer Analytics Survey 2010 results for data mining/analytic tools.

---

# WHAT IS R?

Video from Revolution Analytics - 90 seconds

http://www.revolutionanalytics.com/what-is-open-source-r/

---

# TUTORIAL OVERVIEW

1. MOTIVATING R – A LANGUAGE FOR DATA MINING

2. DATA MINING IN R – HANDS-ON RATTLE GUI

3. PROGRAMMING DATA IN R – SCRIPTING OUR ANALYSES

4. DISSEMINATE RESEARCH IN R – ENSEMBLES AND WSRF

---

# OVERVIEW

1. AN INTRODUCTION TO DATA MINING

2. THE RATTLE PACKAGE FOR DATA MINING

3. MOVING INTO R

---

# OVERVIEW

1. AN INTRODUCTION TO DATA MINING

2. THE RATTLE PACKAGE FOR DATA MINING

3. MOVING INTO R

# Data Mining and Big Data

- Application of
  - Machine Learning
  - Statistics
  - Software Engineering and Programming with Data
  - Intuition

- To Big Data — Volume, Velocity, Variety, Value, Veracity

- . . . to discover new knowledge
- . . . to improve business outcomes
- . . . to deliver better tailored services

# The Business of Data Mining

- Australian Taxation Office
  - Lodgment ($110M)
  - Tax Havens ($150M)
  - Tax Fraud ($250M)

- Department of Immigration

- IBM Buys SPSS for $1.2B in 2009
- SAS has annual revenue approaching $3B
- Analytics is >$100B business and >$320B by 2020 (McKinsey)
- Amazon, eBay/PayPal, Google . . .

# Basic Tools: Data Mining Algorithms

- Linear Discriminant Analysis (lda)
- Logistic Regression (glm)
- Decision Trees (rpart, wsrpart)
- Random Forests (randomForest, wsrf)
- Boosted Stumps (ada)
- Neural Networks (nnet)
- Support Vector Machines (kernlab)
- . . .

*That's a lot of tools to learn in R!*
*Many with different interfaces and options.*

# Overview

1. **An Introduction to Data Mining**

2. **The Rattle Package for Data Mining**

3. **Moving Into R**

# Why a GUI?

- Statistics can be complex and traps await
- **So many** tools in R to deliver insights
- Effective analyses should be scripted
- Scripting also required for repeatability
- R is a language for **programming** with data

How to remember how to do all of this in R?
How to skill up 150 data analysts with Data Mining?

# Users of Rattle

Today, Rattle is used world wide in many industries
- Health analytics
- Customer segmentation and marketing
- Fraud detection
- Government

It is used by
- Consultants and Analytics Teams across business
- Universities to teach Data Mining

It is and will remain freely available.

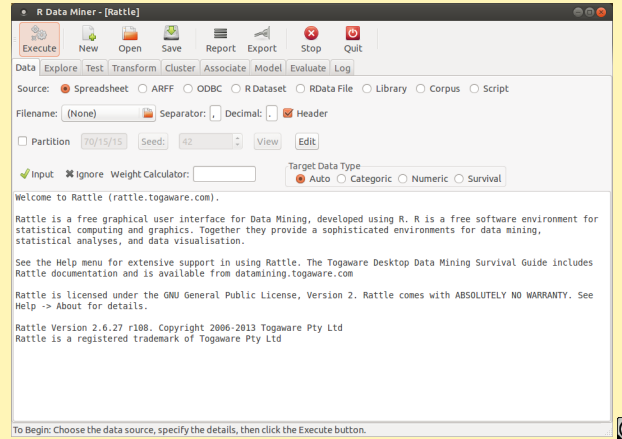CRAN and `http://rattle.togaware.com`

# Installation

- Rattle is built using R
- Need to download and install R from cran.r-project.org
- Recommend also install RStudio from www.rstudio.org

- Then start up RStudio and install Rattle:
  `install.packages("rattle")`
- Then we can start up Rattle:
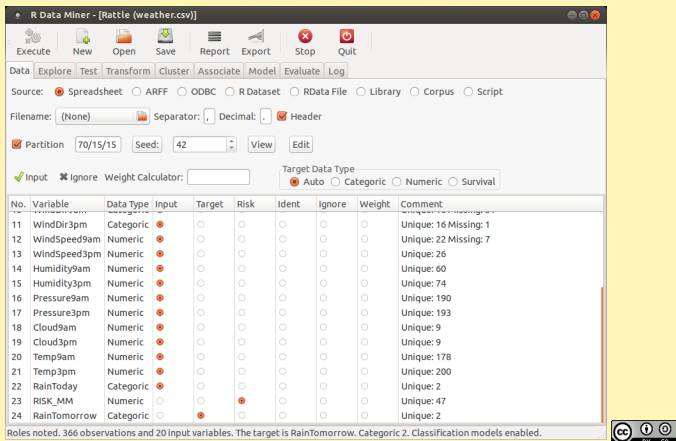  `rattle()`

- Required packages are loaded as needed.

---

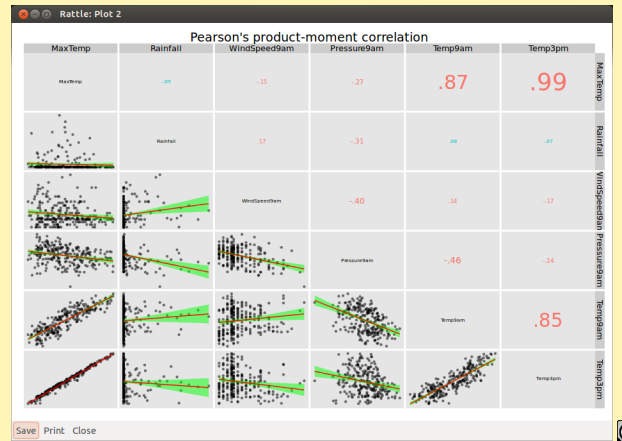# A Tour Thru Rattle: Startup
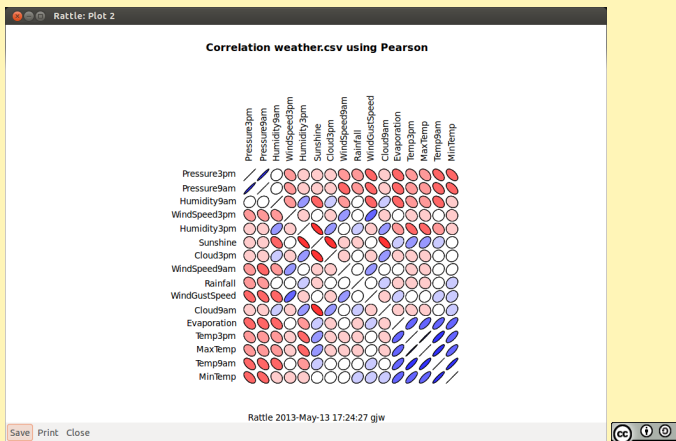
R Data Miner - [Rattle]

Execute  New  Open  Save  Report  Export  Stop  Quit

Data  Explore  Test  Transform  Cluster  Associate  Model  Evaluate  Log

Source: ● Spreadsheet ○ ARFF ○ ODBC ○ R Dataset ○ RData File ○ Library ○ Corpus ○ Script

Filename: (None)  Separator: ,  Decimal: .  ☑ Header

☐ Partition  70/15/15  Seed: 42  View  Edit

✓ Input  ✖ Ignore  Weight Calculator:

Target Data Type
● Auto ○ Categoric ○ Numeric ○ Survival

Welcome to Rattle (rattle.togaware.com).

Rattle is a free graphical user interface for Data Mining, developed using R. R is a free software environment for statistical computing and graphics. Together they provide a sophisticated environments for data mining, statistical analyses, and data visualisation.

See the Help menu for extensive support in using Rattle. The Togaware Desktop Data Mining Survival Guide includes Rattle documentation and is available from datamining.togaware.com

Rattle is licensed under the GNU General Public License, Version 2. Rattle comes with ABSOLUTELY NO WARRANTY. See Help -> About for details.

Rattle Version 2.6.27 r108. Copyright 2006-2013 Togaware Pty Ltd
Rattle is a registered trademark of Togaware Pty Ltd

To Begin: Choose the data source, specify the details, then click the Execute button.

---

# A Tour Thru Rattle: Loading Data

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

Execute  New  Open  Save  Report  Export  Stop  Quit

Data  Explore  Test  Transform  Cluster  Associate  Model  Evaluate  Log

Source: ● Spreadsheet ○ ARFF ○ ODBC ○ R Dataset ○ RData File ○ Library ○ Corpus ○ Script

Filename: (None)  Separator: ,  Decimal: .  ☑ Header

☑ Partition  70/15/15  Seed: 42  View  Edit

✓ Input  ✖ Ignore  Weight Calculator:

Target Data Type
● Auto ○ Categoric ○ Numeric ○ Survival

| No. | Variable | Data Type | Input | Target | Risk | Ident | Ignore | Weight | Comment |
|-----|----------|-----------|-------|--------|------|-------|--------|--------|---------|
| 11 | WindDir3pm | Categoric | ● | ○ | ○ | ○ | ○ | | Unique: 16 Missing: 1 |
| 12 | WindSpeed9am | Numeric | ● | ○ | ○ | ○ | ○ | | Unique: 22 Missing: 7 |
| 13 | WindSpeed3pm | Numeric | ● | ○ | ○ | ○ | ○ | | Unique: 26 |
| 14 | Humidity9am | Numeric | ● | ○ | ○ | ○ | ○ | | Unique: 60 |
| 15 | Humidity3pm | Numeric | ● | ○ | ○ | ○ | ○ | | Unique: 74 |
| 16 | Pressure9am | Numeric | ● | ○ | ○ | ○ | ○ | | Unique: 190 |
| 17 | Pressure3pm | Numeric | ● | ○ | ○ | ○ | ○ | | Unique: 193 |
| 18 | Cloud9am | Numeric | ● | ○ | ○ | ○ | ○ | | Unique: 9 |
| 19 | Cloud3pm | Numeric | ● | ○ | ○ | ○ | ○ | | Unique: 9 |
| 20 | Temp9am | Numeric | ● | ○ | ○ | ○ | ○ | | Unique: 178 |
| 21 | Temp3pm | Numeric | ● | ○ | ○ | ○ | ○ | | Unique: 200 |
| 22 | RainToday | Categoric | ● | ○ | ○ | ○ | ○ | | Unique: 2 |
| 23 | RISK_MM | Numeric | ○ | ○ | ● | ○ | ○ | | Unique: 47 |
| 24 | RainTomorrow | Categoric | ○ | ● | ○ | ○ | ○ | | Unique: 2 |

Roles noted. 366 observations and 20 input variables. The target is RainTomorrow. Categoric 2. Classification models enabled.

---

# A Tour Thru Rattle: Explore Distribution

Rattle: Plot 2

**Pearson's product-moment correlation**

Save  Print  Close

---

# A Tour Thru Rattle: Explore Correlations

Rattle: Plot 2

**Correlation weather.csv using Pearson**

Rattle 2013-May-13 17:24:27 gjw

Save  Print  Close

---

# A Tour Thru Rattle: Hierarchical Cluster

Rattle: Plot 2

**Cluster Dendrogram weather.csv**

Height

Rattle 2013-May-14 14:42:30 gjw

Save  Print  Close

# A Tour Thru Rattle: Decision Tree

# A Tour Thru Rattle: Decision Tree Plot

# A Tour Thru Rattle: Random Forest

# A Tour Thru Rattle: Risk Chart

# Overview

1. An Introduction to Data Mining

2. The Rattle Package for Data Mining

3. Moving Into R

# Data Miners are Programmers of Data

- Data miners are programmers of data
- A GUI can only do so much
- R is a powerful statistical language

- Professional data mining
  - Scripting
  - Transparency
  - Repeatability

# From GUI to CLI — Rattle's Log Tab

# From GUI to CLI — Rattle's Log Tab

# Step 1: Load the Dataset

```
dsname <- "weather"
ds     <- get(dsname)
dim(ds)

## [1] 366  24

names(ds)

##  [1] "Date"         "Location"       "MinTemp"       "...
##  [5] "Rainfall"     "Evaporation"    "Sunshine"      "...
##  [9] "WindGustSpeed" "WindDir9am"     "WindDir3pm"    "...
## [13] "WindSpeed3pm"  "Humidity9am"    "Humidity3pm"   "...
....
```

# Step 2: Observe the Data — Observations

```
head(ds)

##          Date Location MinTemp MaxTemp Rainfall Evapora...
## 1 2007-11-01 Canberra     8.0    24.3      0.0          ...
## 2 2007-11-02 Canberra    14.0    26.9      3.6          ...
## 3 2007-11-03 Canberra    13.7    23.4      3.6          ...
....

tail(ds)

##            Date Location MinTemp MaxTemp Rainfall Evapo...
## 361 2008-10-26 Canberra     7.9    26.1        0        ...
## 362 2008-10-27 Canberra     9.0    30.7        0        ...
## 363 2008-10-28 Canberra     7.1    28.4        0        ...
....
```

# Step 2: Observe the Data — Structure

```
str(ds)

## 'data.frame': 366 obs. of  24 variables:
##  $ Date         : Date, format: "2007-11-01" "2007-11-...
##  $ Location     : Factor w/ 46 levels "Adelaide","Alba...
##  $ MinTemp      : num  8 14 13.7 13.3 7.6 6.2 6.1 8.3 ...
##  $ MaxTemp      : num  24.3 26.9 23.4 15.5 16.1 16.9 1...
##  $ Rainfall     : num  0 3.6 3.6 39.8 2.8 0 0.2 0 0 16...
##  $ Evaporation  : num  3.4 4.4 5.8 7.2 5.6 5.8 4.2 5.6...
##  $ Sunshine     : num  6.3 9.7 3.3 9.1 10.6 8.2 8.4 4....
##  $ WindGustDir  : Ord.factor w/ 16 levels "N"<"NNE"<"N...
##  $ WindGustSpeed: num  30 39 85 54 50 44 43 41 48 31 ...
##  $ WindDir9am   : Ord.factor w/ 16 levels "N"<"NNE"<"N...
##  $ WindDir3pm   : Ord.factor w/ 16 levels "N"<"NNE"<"N...
....
```

# Step 2: Observe the Data — Summary

```
summary(ds)

##       Date                    Location        MinTemp  ...
##  Min.   :2007-11-01   Canberra    :366   Min.   :-5.3...
##  1st Qu.:2008-01-31   Adelaide    :  0   1st Qu.: 2.3...
##  Median :2008-05-01   Albany      :  0   Median : 7.4...
##  Mean   :2008-05-01   Albury      :  0   Mean   : 7.2...
##  3rd Qu.:2008-07-31   AliceSprings:  0   3rd Qu.:12.5...
##  Max.   :2008-10-31   BadgerysCreek:  0  Max.   :20.9...
##                       (Other)     :  0                ...
##     Rainfall        Evaporation       Sunshine       Wind...
##  Min.   : 0.00   Min.   : 0.20   Min.   : 0.00   NW   ...
##  1st Qu.: 0.00   1st Qu.: 2.20   1st Qu.: 5.95   NNW  ...
##  Median : 0.00   Median : 4.20   Median : 8.60   E    ...
....
```

## Step 2: Observe the Data — Variables

```
id      <- c("Date", "Location")
target <- "RainTomorrow"
risk   <- "RISK_MM"
(ignore <- union(id, risk))

## [1] "Date"     "Location" "RISK_MM"

(vars   <- setdiff(names(ds), ignore))

## [1] "MinTemp"      "MaxTemp"      "Rainfall"     "...
## [5] "Sunshine"     "WindGustDir"  "WindGustSpeed" "...
## [9] "WindDir3pm"   "WindSpeed9am" "WindSpeed3pm"  "...
## [13] "Humidity3pm"  "Pressure9am"  "Pressure3pm"   "...
....
```

## Step 3: Clean the Data — Remove Missing

```
dim(ds)

## [1] 366  24

sum(is.na(ds[vars]))

## [1] 47

ds <- ds[-attr(na.omit(ds[vars]), "na.action"),]
```

## Step 3: Clean the Data — Remove Missing

```
dim(ds)

## [1] 328  24

sum(is.na(ds[vars]))

## [1] 0
```

## Step 3: Clean the Data—Target as Categoric

```
summary(ds[target])

##   RainTomorrow
## Min.   :0.000
## 1st Qu.:0.000
## Median :0.000
## Mean   :0.183
## 3rd Qu.:0.000
## Max.   :1.000
....


ds[target] <- as.factor(ds[[target]])
levels(ds[target]) <- c("No", "Yes")
```
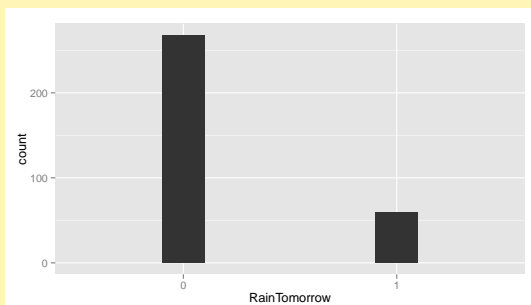
## Step 3: Clean the Data—Target as Categoric

```
summary(ds[target])

## RainTomorrow
## 0:268
## 1: 60
```

## Step 4: Prepare for Modelling

```
(form <- formula(paste(target, "~ .")))

## RainTomorrow ~ .

(nobs <- nrow(ds))

## [1] 328

train <- sample(nobs, 0.70*nobs)
length(train)

## [1] 229

test  <- setdiff(1:nobs, train)
length(test)

## [1] 99
```
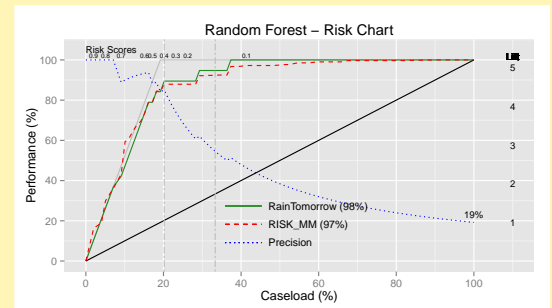
# Step 5: Build the Model—Random Forest

```
library(randomForest)
model <- randomForest(form, ds[train, vars], na.action=na.omit)
model

##
## Call:
##  randomForest(formula=form, data=ds[train, vars], ...
##                Type of random forest: classification
##                      Number of trees: 500
## No. of variables tried at each split: 4
....
```

---

# Step 6: Evaluate the Model—Risk Chart

```
pr <- predict(model, ds[test,], type="prob")[,2]
riskchart(pr, ds[test, target], ds[test, risk],
          title="Random Forest - Risk Chart",
          risk=risk, recall=target, thresholds=c(0.35, 0.15))
```

---

# Tutorial Overview

1. Motivating R – A Language for Data Mining

2. Data Mining in R – Hands-on Rattle GUI

3. Programming Data in R – Scripting our Analyses

4. Disseminate Research in R – Ensembles and wsrf

---

# Overview

1. R Tool Suite

2. RStudio

3. Introduction to R

4. Knitting

---

# Overview

1. R Tool Suite

2. RStudio

3. Introduction to R

4. Knitting

---

# Tools

- Ubuntu GNU/Linux operating system
  - Feature rich toolkit, up-to-date, easy to install, FLOSS

- RStudio
  - Easy to use integrated development environment, FLOSS

- R Statistical Software Language
  - Extensive, powerful, thousands of contributors, FLOSS

- KnitR
  - Produce beautiful documents, easily reproducible, FLOSS

## Using Ubuntu

- Desktop Ubuntu

- Connecting to Analytics Servers
  - Using XWin
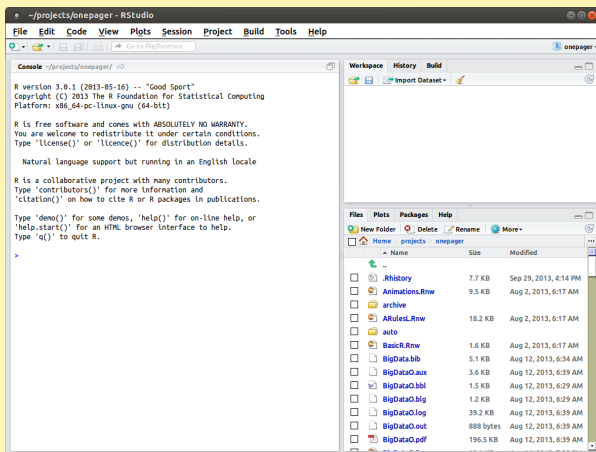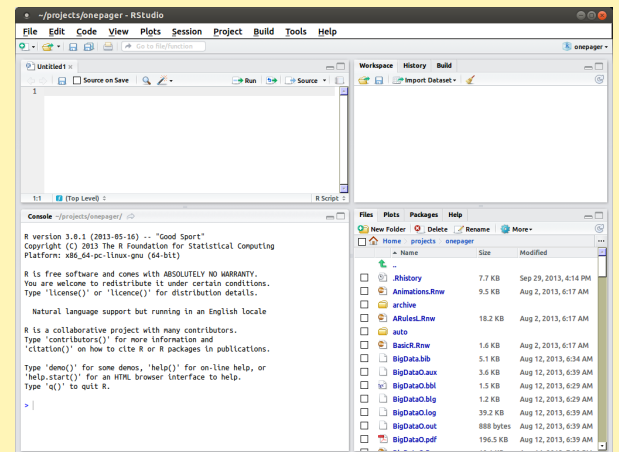  - Using VNC

- Start up RStudio from the Dash

---

## Overview

① R Tool Suite

② RStudio

③ Introduction to R

④ Knitting

---

## RStudio—The Default Three Panels

---

## RStudio—With R Script File—Editor Panel

---

## Overview

① R Tool Suite

② RStudio

③ Introduction to R

④ Knitting

---

## Scatterplot—R Code

Our first little bit of R code:

- Load a couple of *packages* into the R *library*

```
library(rattle)  # Provides the weather dataset
library(ggplot2) # Provides the qplot() function
```
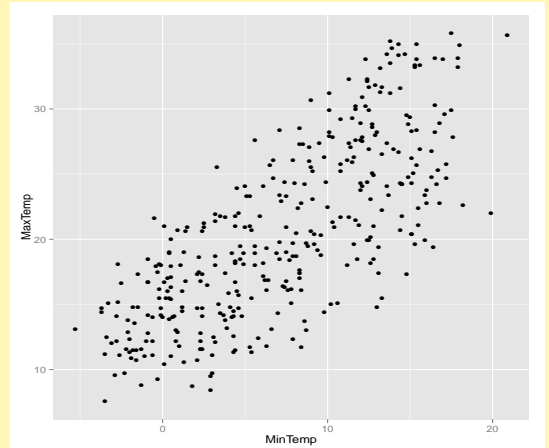
- Then produce a quick plot using `qplot()`

```
ds <- weather
qplot(MinTemp, MaxTemp, data=ds)
```
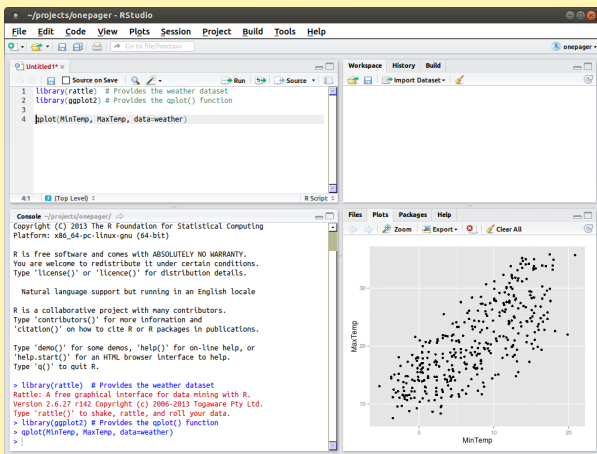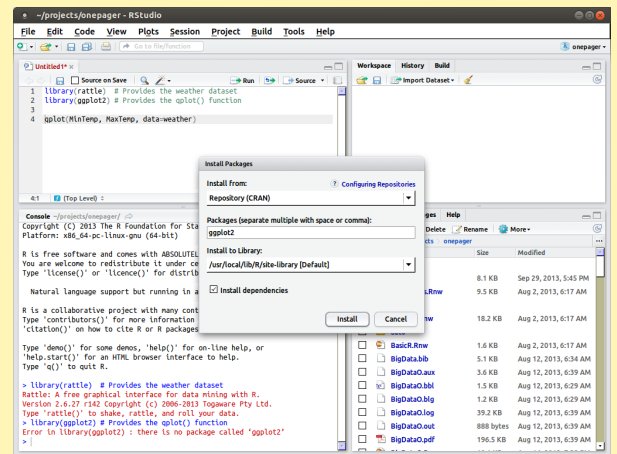
- Your turn: give it a go.

# Scatterplot—R Code

Our first little bit of R code:

- Load a couple of *packages* into the R *library*

```r
library(rattle)  # Provides the weather dataset
library(ggplot2) # Provides the qplot() function
```

- Then produce a quick plot using `qplot()`

```r
ds <- weather
qplot(MinTemp, MaxTemp, data=ds)
```

- Your turn: give it a go.

---

# Scatterplot—Plot

---

# Scatterplot—RStudio

---

# Missing Packages–Tools→Install Packages. . .

---

# RStudio—Installing ggplot2

---

# RStudio—Keyboard Shortcuts

These will become very useful!

- Editor:
  - Ctrl-Enter will send the line of code to the R console
  - Ctrl-2 will move the cursor to the Console

- Console:
  - UpArrow will cycle through previous commands
  - Ctrl-UpArrow will search previous commands
  - Tab will complete function names and list the arguments
  - Ctrl-1 will move the cursor to the Editor

Your turn: try them out.

# RStudio—Keyboard Shortcuts

These will become very useful!

- Editor:
  - Ctrl-Enter will send the line of code to the R console
  - Ctrl-2 will move the cursor to the Console

- Console:
  - UpArrow will cycle through previous commands
  - Ctrl-UpArrow will search previous commands
  - Tab will complete function names and list the arguments
  - Ctrl-1 will move the cursor to the Editor

Your turn: try them out.

---

# Basic R

```
library(rattle)    # Load the weather dataset.
head(weather)      # First 6 observations of the dataset.

##         Date Location MinTemp MaxTemp Rainfall Evapora...
## 1 2007-11-01 Canberra     8.0    24.3      0.0       ...
## 2 2007-11-02 Canberra    14.0    26.9      3.6       ...
## 3 2007-11-03 Canberra    13.7    23.4      3.6       ...
....

str(weather)       # Structure of the variables in the dataset.

## 'data.frame': 366 obs. of  24 variables:
## $ Date       : Date, format: "2007-11-01" "2007-11-...
## $ Location   : Factor w/ 46 levels "Adelaide","Alba...
## $ MinTemp    : num  8 14 13.7 13.3 7.6 6.2 6.1 8.3 ...
....
```

---

# Basic R

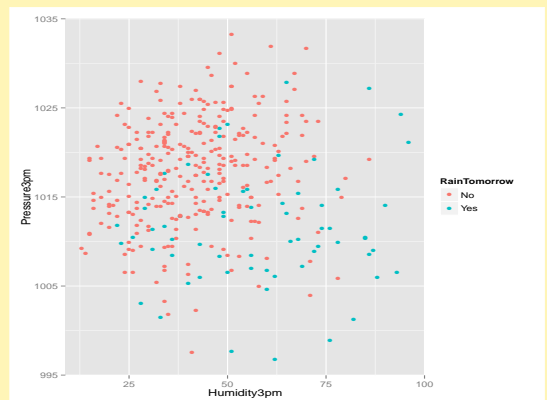```
summary(weather)   # Univariate summary of the variables.

##      Date                   Location      MinTemp      ...
## Min.   :2007-11-01   Canberra     :366   Min.   :-5.30   ...
## 1st Qu.:2008-01-31   Adelaide     :  0   1st Qu.: 2.30   ...
## Median :2008-05-01   Albany       :  0   Median : 7.45   ...
## Mean   :2008-05-01   Albury       :  0   Mean   : 7.27   ...
## 3rd Qu.:2008-07-31   AliceSprings :  0   3rd Qu.:12.50   ...
## Max.   :2008-10-31   BadgerysCreek:  0   Max.   :20.90   ...
##                      (Other)      :  0                   ...
##    Rainfall       Evaporation       Sunshine      WindGust...
## Min.   : 0.00   Min.   : 0.20   Min.   : 0.00   NW    : ...
## 1st Qu.: 0.00   1st Qu.: 2.20   1st Qu.: 5.95   NNW   : ...
## Median : 0.00   Median : 4.20   Median : 8.60   E     : ...
## Mean   : 1.43   Mean   : 4.52   Mean   : 7.91   WNW   : ...
## 3rd Qu.: 0.20   3rd Qu.: 6.40   3rd Qu.:10.50   ENE   : ...
....
```
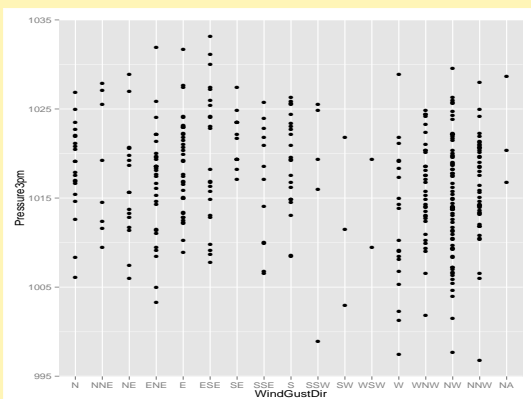
---

# Visual Summaries—Add A Little Colour

```
qplot(Humidity3pm, Pressure3pm, colour=RainTomorrow, data=ds)
```
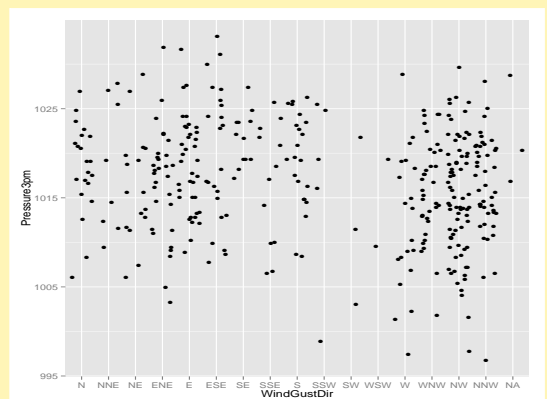
---

# Visual Summaries—Careful with Categorics

```
qplot(WindGustDir, Pressure3pm, data=ds)
```
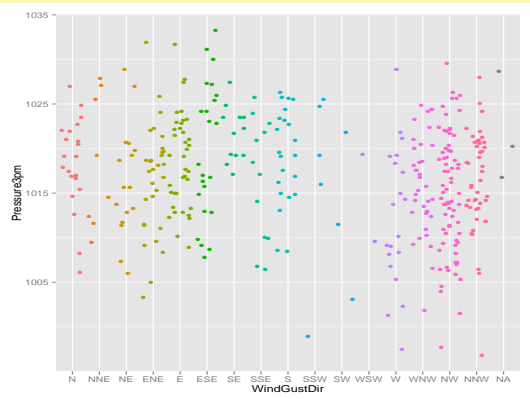
---

# Visual Summaries—Add A Little Jitter

```
qplot(WindGustDir, Pressure3pm, data=ds, geom="jitter")
```
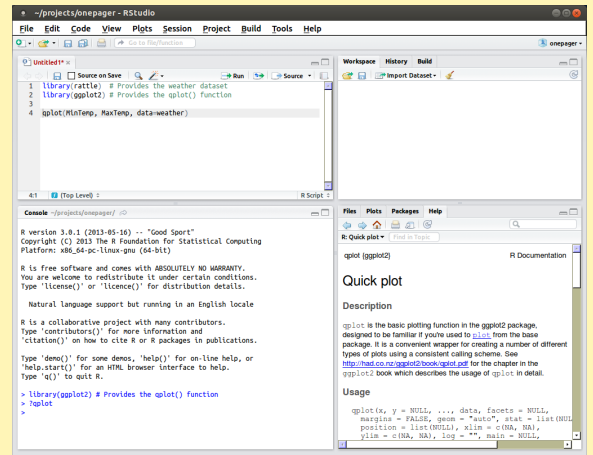
# Visual Summaries—And Some Colour

```
qplot(WindGustDir, Pressure3pm, data=ds, colour=WindGustDir, geom="jitter")
```
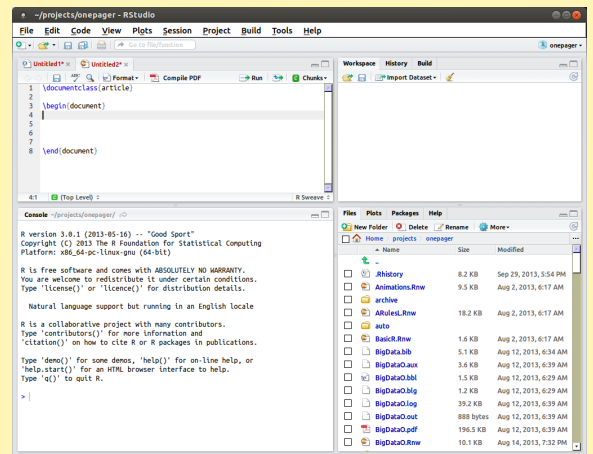
# Getting Help—Precede Command with ?

# Overview

1. R Tool Suite

2. RStudio

3. Introduction to R

4. Knitting

# Create a KnitR Document: New→R Sweave

# Setup KnitR

We wish to use KnitR rather than the older Sweave processor

In RStudio we can configure the options to use knitr:
- Select Tools→Options
- Choose the Sweave group
- Choose **knitr** for *Weave Rnw files using:*
- The remaining defaults should be okay
- Click **Apply** and then **OK**

# Simple KnitR Document

Insert the following into your new KnitR document:

```
\title{Sample KnitR Document}
\author{Graham Williams}
\maketitle

\section*{My First Section}

This is some text that is automatically typeset
by the LaTeX processor to produce well formatted
quality output as PDF.
```
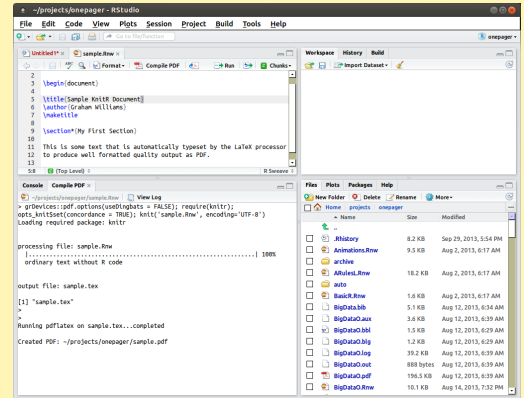
Your turn—Click **Compile PDF** to view the result.

# SIMPLE KNITR DOCUMENT

Insert the following into your new KnitR document:

```
\title{Sample KnitR Document}
\author{Graham Williams}
\maketitle

\section*{My First Section}

This is some text that is automatically typeset
by the LaTeX processor to produce well formatted
quality output as PDF.
```

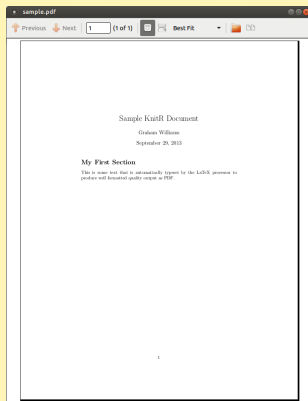Your turn—Click **Compile PDF** to view the result.

# SIMPLE KNITR DOCUMENT

# SIMPLE KNITR DOCUMENT—RESULTING PDF

Result of **Compile PDF**

# KNITR: ADD R COMMANDS

R code can be used to generate results into the document:

```
<<echo=FALSE, message=FALSE>>=
library(rattle)  # Provides the weather dataset
library(ggplot2) # Provides the qplot() function

ds <- weather
qplot(MinTemp, MaxTemp, data=ds)
@
```

Your turn—Click **Compile PDF** to view the result.

# KNITR: ADD R COMMANDS

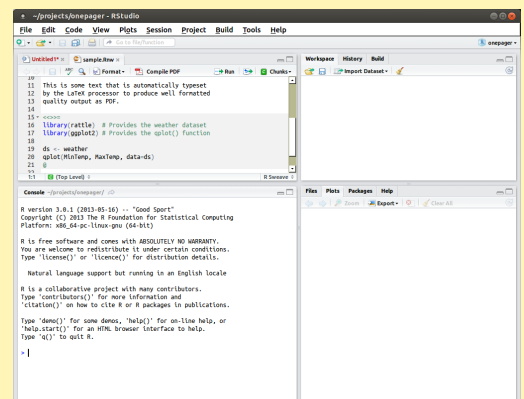R code can be used to generate results into the document:

```
<<echo=FALSE, message=FALSE>>=
library(rattle)  # Provides the weather dataset
library(ggplot2) # Provides the qplot() function

ds <- weather
qplot(MinTemp, MaxTemp, data=ds)
@
```
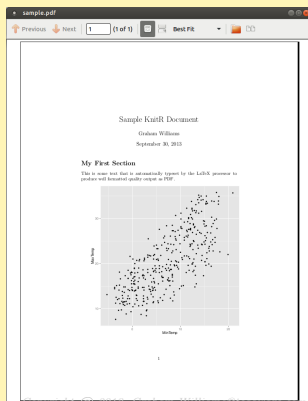
Your turn—Click **Compile PDF** to view the result.

# KNITR DOCUMENT WITH R CODE

# Simple KnitR Document—Resulting PDF with Plot

Result of **Compile PDF**

---

# LaTeX Basics

```
\subsection*{...}        % Introduce a Sub Section

\subsubsection*{...}     % Introduce a Sub Sub Section

\textbf{...}             % Bold font
\textit{...}             % Italic font

\begin{itemize}          % A bullet list
  \item ...
  \item ...
\end{itemize}
```

Plus an extensive collection of other markup and capabilities.

---

# KnitR Basics

```
echo=FALSE       # Do not display the R code
eval=TRUE        # Evaluate the R code

results="hide"   # Hide the results of the R commands

fig.width=10     # Extend figure width from 7 to 10 inches
fig.height=8     # Extend figure height from 7 to 8 inches

out.width="0.8\\textwidth"    # Fit figure 80% page width
out.height="0.5\\textheight"  # Fit figure 50% page height
```

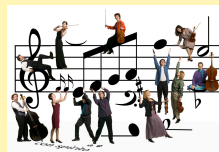Plus an extensive collection of other options.

---

# Tutorial Overview

1. Motivating R – A Language for Data Mining

2. Data Mining in R – Hands-on Rattle GUI

3. Programming Data in R – Scripting our Analyses

4. Disseminate Research in R – Ensembles and wsrf

---

# Case Study – Ensembles in R
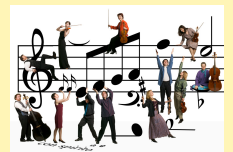
Major advances in Data Mining

- The best off-the-shelf technology includes random forests, boosting and support vector machines?

- Available for investigation now through open source solutions, with closed source tools catching up.

---

# Case Study – Ensembles in R

Major advances in Data Mining

- The best off-the-shelf technology includes random forests, boosting and support vector machines?

- Available for investigation now through open source solutions, with closed source tools catching up.

# Case Study – Ensembles in R

Major advances in Data Mining

- The best off-the-shelf technology includes random forests, boosting and support vector machines?

- Available for investigation now through open source solutions, with closed source tools catching up.

# Introducing Random Forests

Research with Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences

- Random forests are a popular classification method building an ensemble of a single type of decision tree.

- It is unsurpassed in accuracy among current algorithms.

- Algorithmically intuitive and simple.

- It is used widely in numerous research domains including bioinformatics, image classification, text classification.

# Random Forests Algorithm

- Build many decision trees (e.g., 500).

- For each tree:
  - Select a random subset of the training set ($N$);
  - Choose different subsets of features for each node of the decision tree ($m << M$);
  - Build the tree without pruning (i.e., overfit)

- Classify a new entity using every decision tree:
  - Each tree "votes" for the entity.
  - The decision with the largest number of votes wins!
  - The proportion of votes is the resulting score.

# Using Weighted Variable Subspaces

- Performance of a random forest is improved by
  - **Strengthening** each tree
  - Reducing **correlation** between each tree

- Problem of large number of variables:
  - Random selection means too many irrelevant variables

- Introduce the concept of weighted subspace random forests
  - Bias the selection of variables toward most important variables
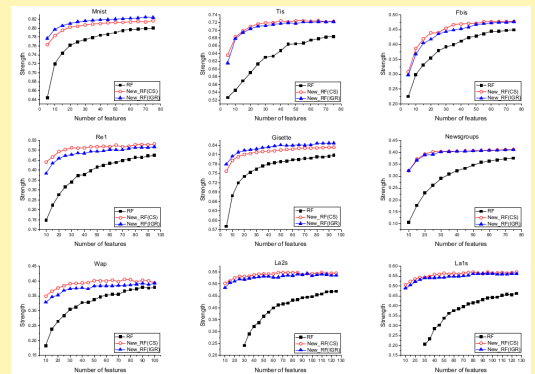
# Datasets

| Name | # Features | # Train Set | # Test Set | # Classes |
|---|---|---|---|---|
| Mnist | 780 | 60,000 | 10,000 | 2 |
| Tis | 927 | 5200 | 6875 | 2 |
| Fbis | 2000 | 1711 | 752 | 17 |
| Re1 | 3758 | 1147 | 510 | 25 |
| Gisette | 5000 | 5000 | 1000 | 2 |
| Newsgroups | 5000 | 11,268 | 7504 | 20 |
| Wap | 8460 | 1104 | 456 | 20 |
| La2s | 12,432 | 1855 | 845 | 6 |
| La1s | 13,195 | 1963 | 887 | 6 |

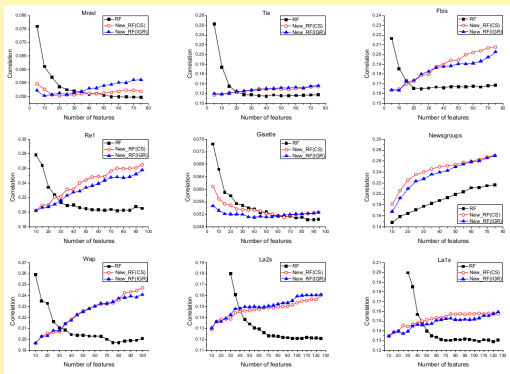*(From International Journal of Data Warehousing and Mining 2012)*
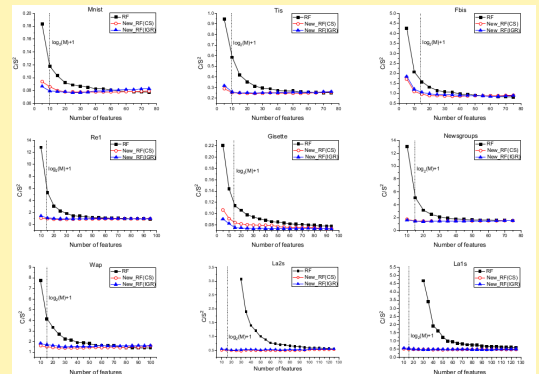
# Comparison of Strength vs Features



*(From International Journal of Data Warehousing and Mining 2012)*

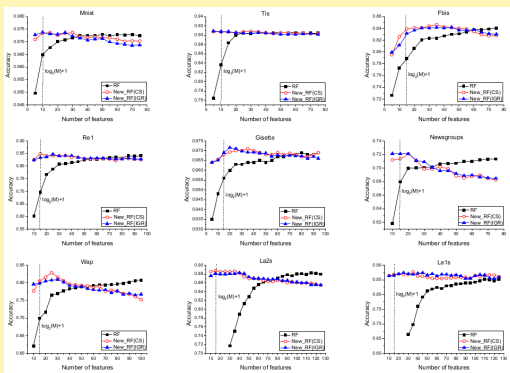## Comparison of Correlation vs Features



*(From International Journal of Data Warehousing and Mining 2012)*

## $c/s^2$ vs Features



*(From International Journal of Data Warehousing and Mining 2012)*

## Accuracy vs Features



*(From International Journal of Data Warehousing and Mining 2012)*

## Repeatable and Transparent Research

- SIAT have implemented the research as C++ code
- Enhanced for parallel environment - multi-core and multi-node
- At least as good as random forest but always very much quicker.
- Integrated into R using the Rcpp package of R
- Now openly available for use and peer review:
  `install.packages("wsrf", repos="http://rattle.togaware.com")`
- Similarly wskm, wsrpart, eqrf.
- Will be published to CRAN shortly.
- **Publish-by-Example** — use wsrf as template
- C, C++, Fortran, Java (RWeka)

## Using the Package

```
install.packages("wsrf", repos="http://rattle.togaware.com")
library(help=wsrf)
library(wsrf)
model <- wsrf(form, ds[train, vars])
pr    <- predict(model, na.omit(ds[test, vars]))
```
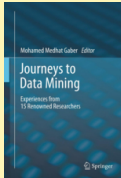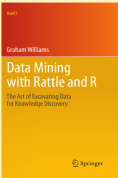
## Tutorial Overview

1. Motivating R – A Language for Data Mining

2. Data Mining in R – Hands-on Rattle GUI

3. Programming Data in R – Scripting our Analyses

4. Disseminate Research in R – Ensembles and wsrf

# Resources and References

- **OnePageR**: http://onepager.togaware.com – Tutorial Notes
- Rattle: http://rattle.togaware.com
- Guides: http://datamining.togaware.com
- Practise: http://analystfirst.com

- Book: Data Mining using Rattle/R
- Chapter: Rattle and Other Tales
- Paper: A Data Mining GUI for R — R Journal, Volume 1(2)

---

# Thank You

## Question Time