# Data Science with R

## Introducing Data Mining with Rattle and 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

# Overview

# 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

# 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

# A Tour Thru Rattle: Loading Data

# A TOUR THRU RATTLE: EXPLORE DISTRIBUTION

# A Tour Thru Rattle: Explore Correlations

# A TOUR THRU RATTLE: HIERARCHICAL CLUSTER

# A Tour Thru Rattle: Decision Tree

# A Tour Thru Rattle: Decision Tree Plot

# A TOUR THRU RATTLE: RANDOM FOREST

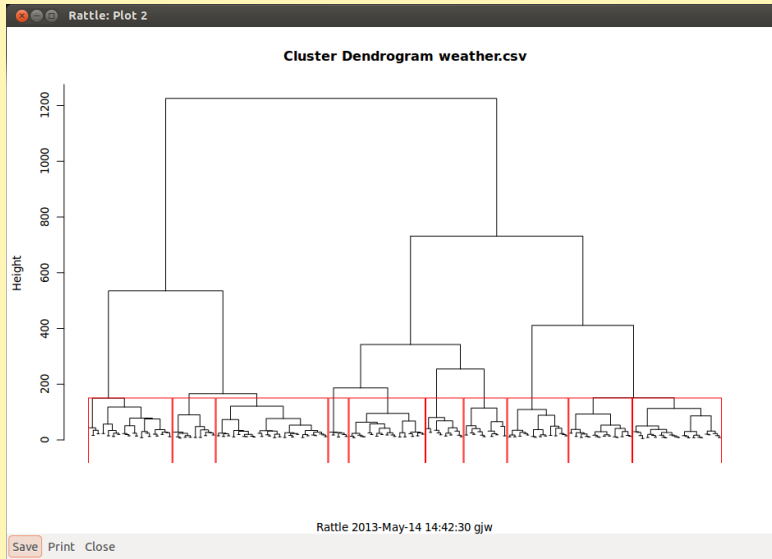# A Tour Thru Rattle: Risk Chart



Risk Chart Random Forest weather.csv [test] RainTomorrow

# OVERVIEW

# 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
....
```
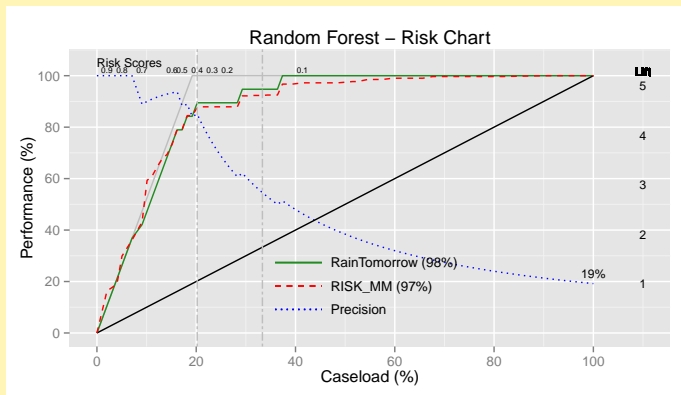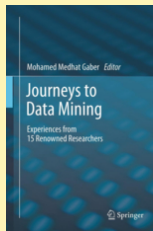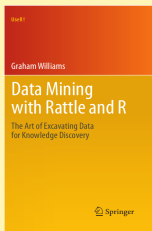
# Step 6: Evaluate the Model—Risk Chart

```
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))
```

# 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

*This document, sourced from StartL.Rnw revision 282, was processed by KnitR version 1.5 of 2013-09-28 and took 4 seconds to process. It was generated by gjw on nyx running Ubuntu 13.10 with Intel(R) Xeon(R) CPU W3520 @ 2.67GHz having 4 cores and 12.3GB of RAM. It completed the processing 2014-02-14 06:20:04.*