

# DATA ANALYTICS AND BUSINESS INTELLIGENCE (8696/8697)

## INTRODUCING DATA SCIENCE WITH RATTLE AND R

Graham.Williams@togaware.com

Chief Data Scientist  
Australian Taxation Office

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

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



# OVERVIEW

- 1 AN INTRODUCTION TO DATA MINING
- 2 WHY CHOOSE R FOR DATA SCIENCE?
- 3 THE RATTLE PACKAGE FOR DATA MINING
- 4 MOVING INTO R
- 5 KNITTING



# OVERVIEW

- 1 AN INTRODUCTION TO DATA MINING
- 2 WHY CHOOSE R FOR DATA SCIENCE?
- 3 THE RATTLE PACKAGE FOR DATA MINING
- 4 MOVING INTO R
- 5 KNITTING



# DATA MINING

A data driven analysis to uncover otherwise unknown but useful patterns in large datasets, to discover new knowledge and to develop predictive models, turning data and information into knowledge and (one day perhaps) wisdom, in a timely manner.



# DATA MINING

- Application of
  - Machine Learning
  - Statistics
  - Software Engineering and Programming with Data
  - Effective Communications and Intuition
- ...to Datasets that vary by  
Volume, Velocity, Variety, Value, Veracity
- ...to discover new knowledge
- ...to improve business outcomes
- ...to deliver better tailored services



# DATA MINING IN RESEARCH

- **Health Research**

Adverse reactions using linked Pharmaceutical, General Practitioner, Hospital, Pathology datasets.

- **Astronomy**

Microlensing events in the Large Magellanic Cloud of several million observed stars (out of 10 billion).

- **Psychology**

Investigation of age-of-onset for Alzheimer's disease from 75 variables for 800 people.

- **Social Sciences**

Survey evaluation. Social network analysis - identifying key influencers.



# DATA MINING IN GOVERNMENT

- **Australian Taxation Office**

- Lodgment (\$110M)
- Tax Havens (\$150M)
- Tax Fraud (\$250M)

- **Immigration and Border Control**

- Check passengers before boarding

- **Health and Human Services**

- Doctor shoppers
- Over servicing



# THE BUSINESS OF DATA MINING

- SAS has annual revenues of \$3B (2013)
- IBM bought SPSS for \$1.2B (2009)
- Analytics is >\$100B business and >\$320B by 2020
- Amazon, eBay/PayPal, Google, Facebook, LinkedIn, ...
- Shortage of 180,000 data scientists in US in 2018 (McKinsey) ...





# BASIC DATA MINING ALGORITHMS

- Cluster Analysis (kmeans, wskm)
- Association Analysis (arules)
- Linear Discriminant Analysis (lda)
- Logistic Regression (glm)
- Decision Trees (rpart, wsrpart)
- Random Forests (randomForest, wsrfr)
- 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 WHY CHOOSE R FOR DATA SCIENCE?**
- 3 THE RATTLE PACKAGE FOR DATA MINING
- 4 MOVING INTO R
- 5 KNITTING

# INSTALLING R AND RATTLE

- **First task is to install R**

As free/libre open source software (FLOSS or FOSS), R and Rattle are available to all, with no limitations on our freedom to use and share the software, except to share and share alike.

- Visit CRAN at <http://cran.rstudio.com>
- Visit Rattle at <http://rattle.togaware.com>
- Linux: Install packages (Ubuntu is recommended)  
`$ wajig install r-recommended r-cran-rattle`
- Windows: Download and install from CRAN
- MacOSX: Download and install from CRAN



# WHY DO DATA SCIENCE WITH R?

- Most widely used Data Mining and Machine Learning Package
  - Machine Learning
  - Statistics
  - Software Engineering and Programming with Data
  - But 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 DO DATA SCIENCE WITH R?

- Most widely used Data Mining and Machine Learning Package
  - Machine Learning
  - Statistics
  - Software Engineering and Programming with Data
  - But 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 DO DATA SCIENCE WITH R?

- Most widely used Data Mining and Machine Learning Package
  - Machine Learning
  - Statistics
  - Software Engineering and Programming with Data
  - But 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 DO DATA SCIENCE WITH R?

- Most widely used Data Mining and Machine Learning Package
  - Machine Learning
  - Statistics
  - Software Engineering and Programming with Data
  - But 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!!!



## R: A DANGEROUS TOOL?

“I think it addresses a niche market for high-end data analysts that want free, readily available code. We have customers who build engines for aircraft. I am happy they are not using freeware when I get on a jet.” Anne H. Milley, director of technology product marketing at SAS (New York Times, 7 January 2009).

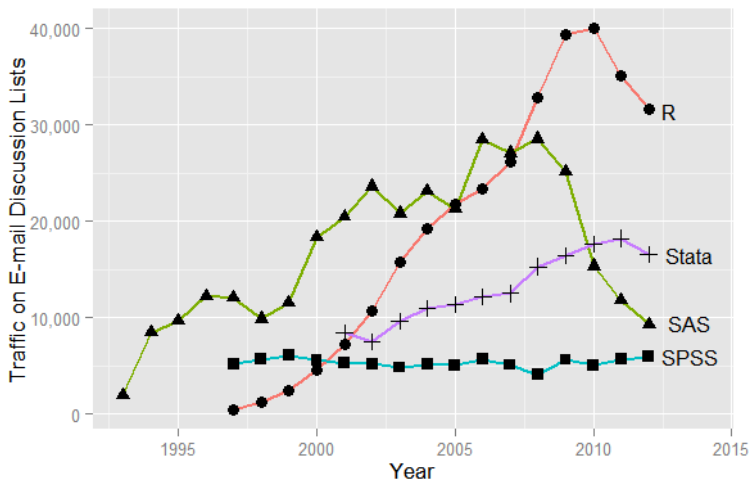
It's interesting that SAS Institute feels that non-peer-reviewed software with hidden implementations of analytic methods that cannot be reproduced by others should be trusted when building aircraft engines. (Frank Harrell)





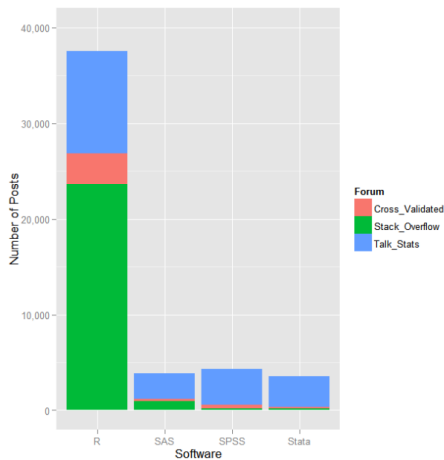
# HOW POPULAR IS R? DISCUSSION LIST TRAFFIC

Monthly email traffic on software's main discussion list.



# HOW POPULAR IS R? DISCUSSION TOPICS

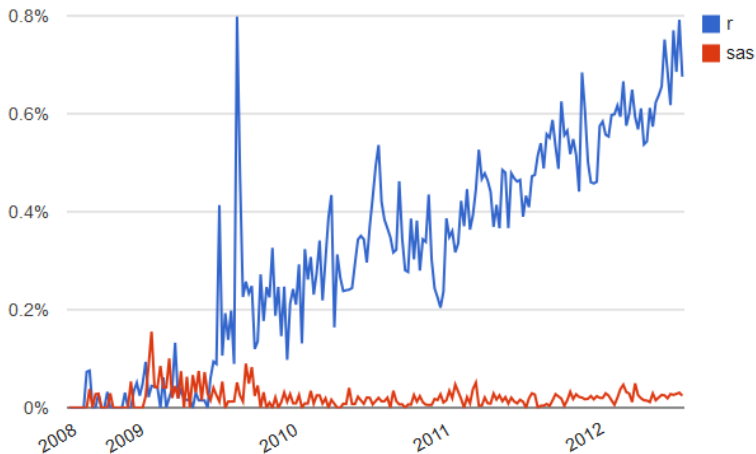
Number of discussions on popular QandA forums 2013.



Source: <http://r4stats.com/articles/popularity/>

# HOW POPULAR IS R? R VERSUS SAS

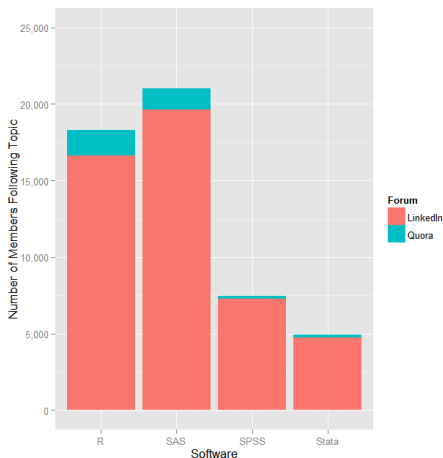
Number of R/SAS related posts to Stack Overflow by week.



Source: <http://r4stats.com/articles/popularity/>

# HOW POPULAR IS R? PROFESSIONAL FORUMS

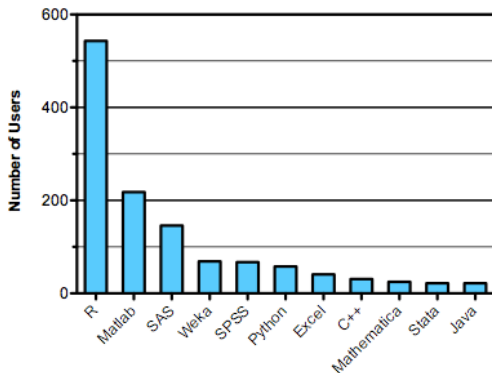
Registered for the main discussion group for each software.



Source: <http://r4stats.com/articles/popularity/>

# HOW POPULAR IS R? USED IN ANALYTICS COMPETITIONS

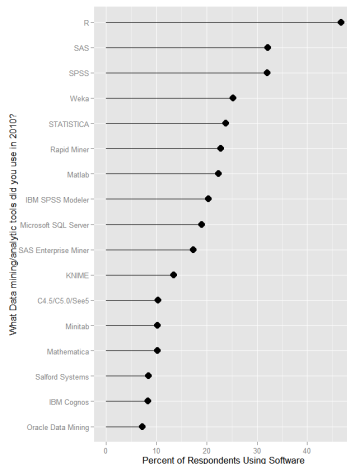
Software used in data analysis competitions in 2011.



Source: <http://r4stats.com/articles/popularity/>

# HOW POPULAR IS R? USER SURVEY

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



Source: <http://r4stats.com/articles/popularity/>

# R SKILLS ATTRACT GOOD SALARIES

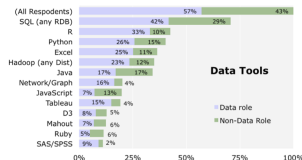
- 2014 survey of average US tech salaries by Dice Tech puts R at the top of the list at \$115,531

ComputerWorld

- 2013 O'Reilly Strata Conference: Data Scientists use R over other data programming languages and Data Scientists using Open Source earn \$130,000 on average.

Revolution Analytics

AVERAGE SALARY FOR High Paying Skills and Experience		
SKILL	2013	YR/YR CHANGE
R	\$ 115,531	n/a
NoSQL	\$ 114,796	1.6%
MapReduce	\$ 114,396	n/a
PMBok	\$ 112,382	1.3%
Cassandra	\$ 112,382	n/a
Omnigraffle	\$ 111,039	0.3%
Pig	\$ 109,561	n/a
SOA (Service Oriented Architecture)	\$ 108,997	-0.5%
Hadoop	\$ 108,669	-5.6%
Mongo DB	\$ 107,825	-0.4%



# WHAT IS R?

## R — The Video

### A 90 Second Promo from Revolution Analytics

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





# CHOOSING R FOR DATA SCIENCE

- Data Science is about Analysing Data;
- R is freely available to all to analyse data;
- R has the most extensive suite of functionality available;
- Nothing else is any longer even close.

*This document, sourced from StartL.Rnw revision 436, was processed by KnitR version 1.6 of 2014-05-24 and took 1 seconds to process. It was generated by gjw on nyx running Ubuntu 14.04 LTS with Intel(R) Xeon(R) CPU W3520 @ 2.67GHz having 4 cores and 12.3GB of RAM. It completed the processing 2014-06-21 20:26:13.*



# OVERVIEW

- 1 AN INTRODUCTION TO DATA MINING
- 2 WHY CHOOSE R FOR DATA SCIENCE?
- 3 THE RATTLE PACKAGE FOR DATA MINING**
- 4 MOVING INTO R
- 5 KNITTING

# WHY A GUI FOR DATA SCIENCE IN R?

- 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

- Universities to teach Data Mining
- Within research projects for basic analyses
- Consultants and Analytics Teams across business

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](http://cran.r-project.org)
- Recommend also install RStudio from [www.rstudio.org](http://www.rstudio.org)

- Then start up RStudio and install Rattle:

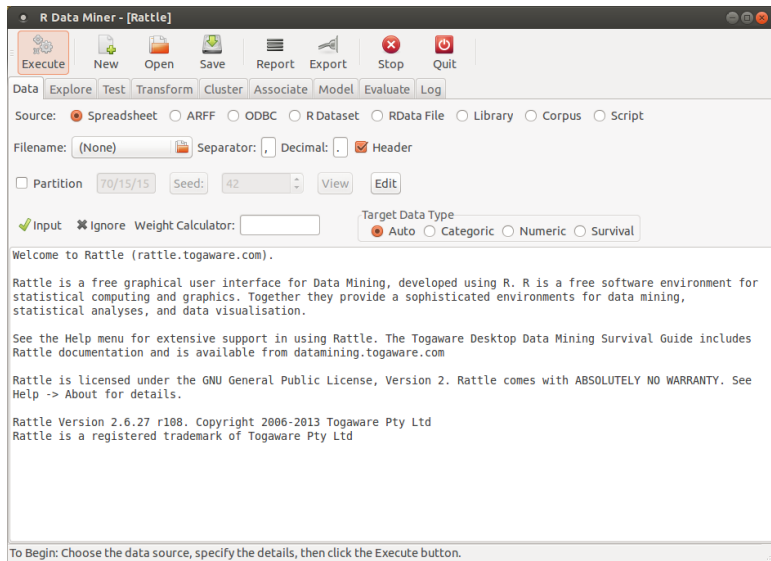
```
install.packages("rattle")
```

- Then we can start up Rattle:

```
rattle()
```

- Required packages are loaded as needed.

# A TOUR THRU RATTLE: STARTUP



# A TOUR THRU RATTLE: LOADING DATA

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

Project Tools Settings Help Rattle Version 3.0.2 [togaware.com](http://togaware.com)

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:  Separator:  Decimal:  ☒ Header

☒ Partition  Seed:  View Edit

☒ Input ☐ Ignore Weight Calculator:

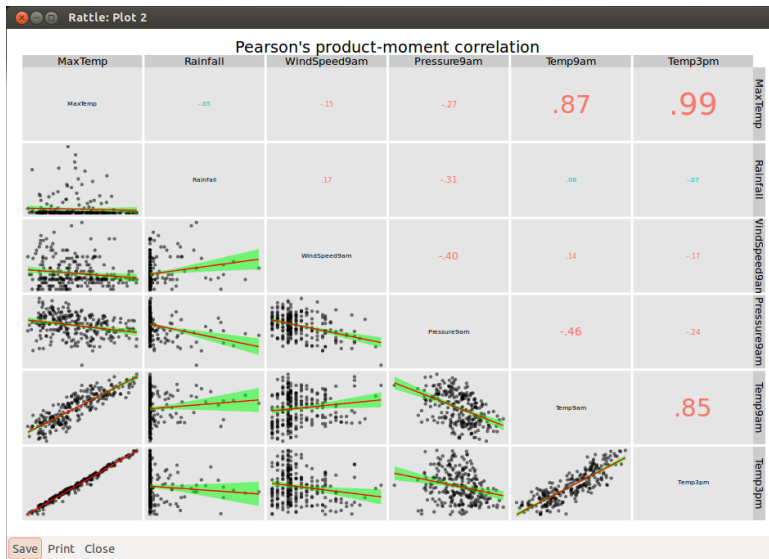
Target Data Type  
☒ Auto ☐ Categorical ☐ Numeric ☐ Survival

No.	Variable	Data Type	Input	Target	Risk	Ident	Ignore	Weight	Comment
16	Pressure9am	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 190
17	Pressure3pm	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 193
18	Cloud9am	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 9
19	Cloud3pm	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 9
20	Temp9am	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 178
21	Temp3pm	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 200
22	RainToday	Categorical	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2
23	RISK_MM	Numeric	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 47
24	RainTomorrow	Categorical	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2

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

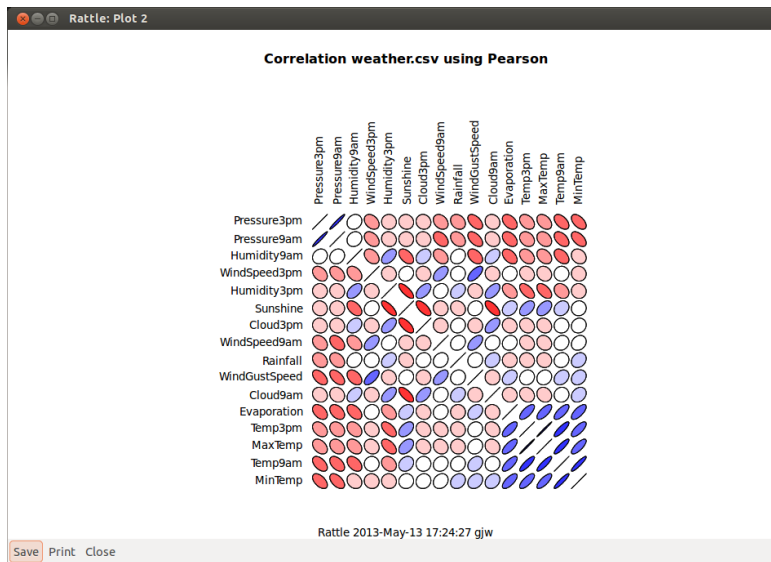


# A TOUR THRU RATTLE: EXPLORE DISTRIBUTION

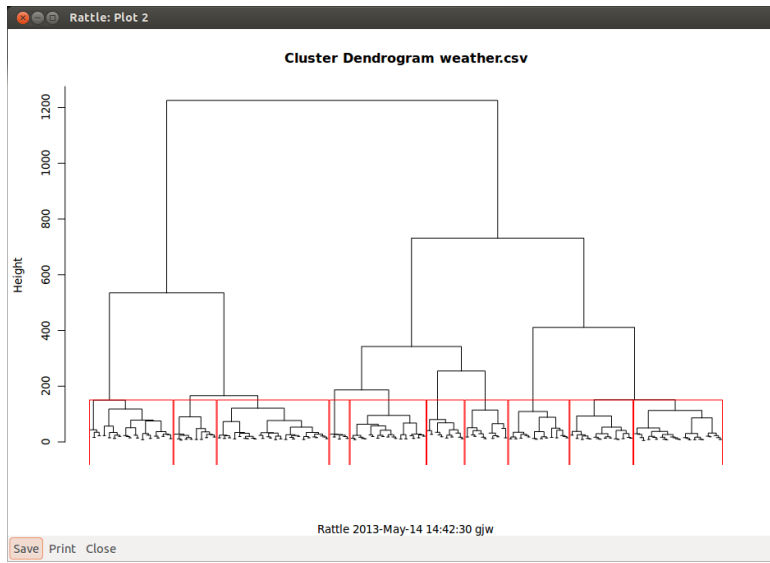




# A TOUR THRU RATTLE: EXPLORE CORRELATIONS



# A TOUR THRU RATTLE: HIERARCHICAL CLUSTER



# A TOUR THRU RATTLE: DECISION TREE

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

Execute New Open Save Report Export Stop Quit

Data Explore Test Transform Cluster Associate Model Evaluate Log

Type: ☒ Tree ☐ Forest ☐ Boost ☐ SVM ☐ Linear ☐ Neural Net ☐ Survival ☐ All

Target: RainTomorrow Algorithm: ☒ Traditional ☐ Conditional Model Builder: rpart

Min Split: 20 Max Depth: 30 Priors:  ☐ Include Missing

Min Bucket: 7 Complexity: 0.0100 Loss Matrix:

Summary of the Decision Tree model for Classification (built using 'rpart'):

n= 256

node), split, n, loss, yval, (yprob)  
\* denotes terminal node

- 1) root 256 41 No (0.83984375 0.16015625)
- 2) Pressure3pm>=1011.9 204 16 No (0.92156863 0.07843137)
- 4) Cloud3pm< 7.5 195 10 No (0.94871795 0.05128205) \*
- 5) Cloud3pm>=7.5 9 3 Yes (0.33333333 0.66666667) \*
- 3) Pressure3pm< 1011.9 52 25 No (0.51923077 0.48076923)
- 6) Sunshine>=8.85 25 5 No (0.80000000 0.20000000) \*
- 7) Sunshine< 8.85 27 7 Yes (0.25925926 0.74074074) \*

Classification tree:  
rpart(formula = RainTomorrow ~ ., data = crs\$dataset[crs\$strain,  
c(crs\$input, crs\$target)], method = "class", parms = list(split = "information"),  
control = rpart.control(usesurrogate = 0, maxsurrogate = 0))

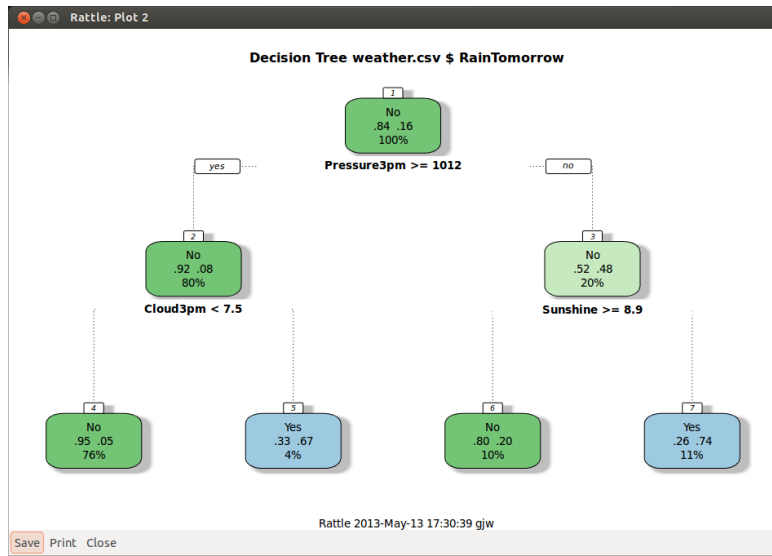
Variables actually used in tree construction:  
[1] Cloud3pm Pressure3pm Sunshine

Root node error: 41/256 = 0.16016

The Decision Tree model has been built. Time taken: 0.09 secs



## A TOUR THRU RATTLE: DECISION TREE PLOT



# A TOUR THRU RATTLE: RANDOM FOREST

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

Execute New Open Save Report Export Stop Quit

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

Type: ☐ Tree ☒ Forest ☐ Boost ☐ SVM ☐ Linear ☐ Neural Net ☐ Survival ☐ All

Target: RainTomorrow Algorithm: ☒ Traditional ☐ Conditional Model Builder: randomForest

Number of Trees: 500 Sample Size: Importance Rules 1

Number of Variables: 4 ☒ Impute Errors OOB ROC

**Summary of the Random Forest Model**

=====

Number of observations used to build the model: 256  
Missing value imputation is active.

Call:  
randomForest(formula = RainTomorrow ~ .,  
data = crs\$dataset[crs\$sample, c(crs\$input, crs\$target)],  
ntree = 500, mtry = 4, importance = TRUE, replace = FALSE, na.action = na.roughfix)

Type of random forest: classification  
Number of trees: 500  
No. of variables tried at each split: 4

OOB estimate of error rate: 13.28%

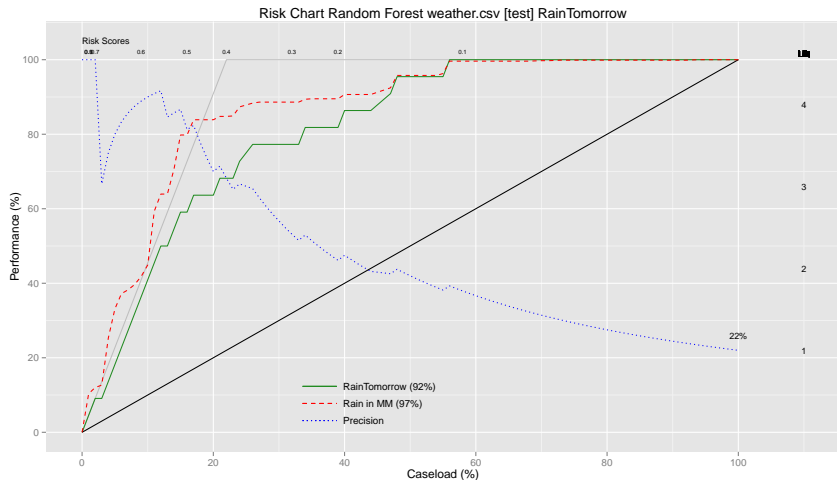
Confusion matrix:  
No Yes class.error  
No 207 8 0.0372093  
Yes 26 15 0.6341463

Analysis of the Area Under the Curve (AUC)  
=====

The Random Forest model has been built. Time taken: 0.87 secs



# A TOUR THRU RATTLE: RISK CHART



# RATTLE INTERFACE NOTES

- Work through the tabs from left to right
- After setting up a tab we need to Execute it
- Projects save the current Rattle state
- Projects can be restored at a later time



# OVERVIEW

- 1 AN INTRODUCTION TO DATA MINING
- 2 WHY CHOOSE R FOR DATA SCIENCE?
- 3 THE RATTLE PACKAGE FOR DATA MINING
- 4 MOVING INTO R**
- 5 KNITTING



# DATA SCIENTISTS ARE PROGRAMMERS OF DATA

## But...

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

## Data Scientists Desire...

- Scripting
- Transparency
- Repeatability
- Sharing



# FROM GUI TO CLI — RATTLE'S LOG TAB

```

# Rattle is Copyright (c) 2006-2013 Togaware Pty Ltd.

#=====
# Rattle timestamp: 2013-05-13 16:49:53 x86_64-pc-linux-gnu

# Rattle version 2.6.27 user 'gjlw'

# Export this log textview to a file using the Export button or the Tools
# menu to save a log of all activity. This facilitates repeatability. Exporting
# to file 'myrf01.R', for example, allows us to type in the R Console
# the command source('myrf01.R') to repeat the process automatically.
# Generally, we may want to edit the file to suit our needs. We can also directly
# edit this current log textview to record additional information before exporting.

# Saving and loading projects also retains this log.

library(rattle)

# This log generally records the process of building a model. However, with very
# little effort the log can be used to score a new dataset. The logical variable
# 'building' is used to toggle between generating transformations, as when building
# a model, and simply using the transformations, as when scoring a dataset.

building <- TRUE
scoring <- ! building

# The colorspace package is used to generate the colours used in plots, if available.

library(colorspace)
  
```



# FROM GUI TO CLI — RATTLE'S LOG TAB

```
# Export Comments ☒ Rename Internal Variables: From crs$ to MY ☐
# Rattle timestamp: 2013-05-13 17:35:07 x86_64-pc-linux-gnu
# Random Forest
# The 'randomForest' package provides the 'randomForest' function.
require(randomForest, quietly=TRUE)
# Build the Random Forest model.
set.seed(crv$seed)
crs$rfr <- randomForest(RainTomorrow ~ .,
  data=crs$dataset[crs$sample,c(crs$input, crs$target)],
  ntree=500,
  mtry=4,
  importance=TRUE,
  na.action=na.roughfix,
  replace=FALSE)
# Generate textual output of 'Random Forest' model.
crs$rfr
# The 'pROC' package implements various AUC functions.
require(pROC, quietly=TRUE)
# Calculate the Area Under the Curve (AUC).|
```



# 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/ 49 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])
```

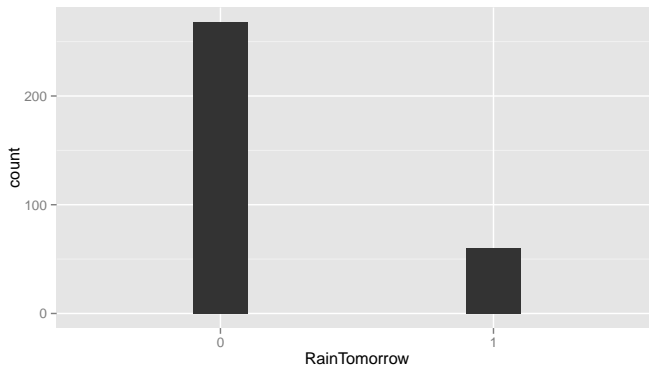
```
##    RainTomorrow  
##  Min.      :0.000  
## 1st Qu.:0.000  
##  Median :0.000  
##   Mean  :0.183  
## 3rd Qu.:0.000  
##   Max.   :1.000  
....
```

```
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
```



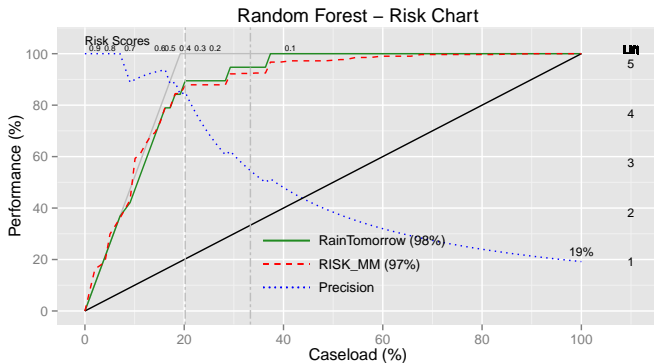
## STEP 5: BUILD THE MODEL—RANDOM FOREST

```
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))
```



# TOOLS

- Ubuntu GNU/Linux operating system
  - Feature rich toolkit, up-to-date, easy to install, FLOSS
- RStudio
  - Easy to use integrated development environment, FLOSS
  - Powerful alternative is Emacs (Speaks Statistics), FLOSS
- R Statistical Software Language
  - Extensive, powerful, thousands of contributors, FLOSS
- KnitR and  $\text{\LaTeX}$ 
  - Produce beautiful documents, easily reproducible, FLOSS



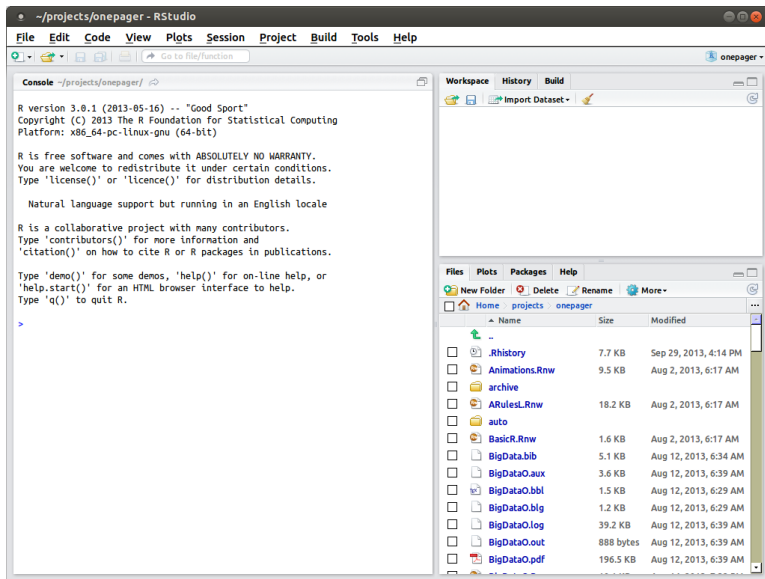


# USING UBUNTU

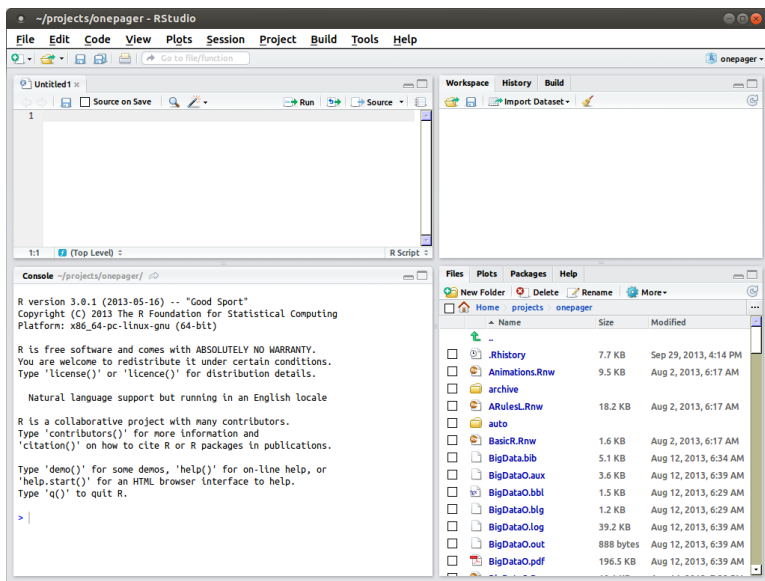
- Desktop Operating System (GNU/Linux)
- Replacing Windows and OSX
- The GNU Tool Suite based on Unix — significant heritage
- Multiple specialised single task tools, working well together
- Compared to single application trying to do it all
- Powerful data processing from the command line:  
grep, awk, head, tail, wc, sed, perl, python, most, diff, make,  
paste, join, patch, ...
- For interacting with R — start up RStudio from the Dash



# RSTUDIO—THE DEFAULT THREE PANELS



# RSTUDIO—WITH R SCRIPT FILE—EDITOR PANEL



# 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*

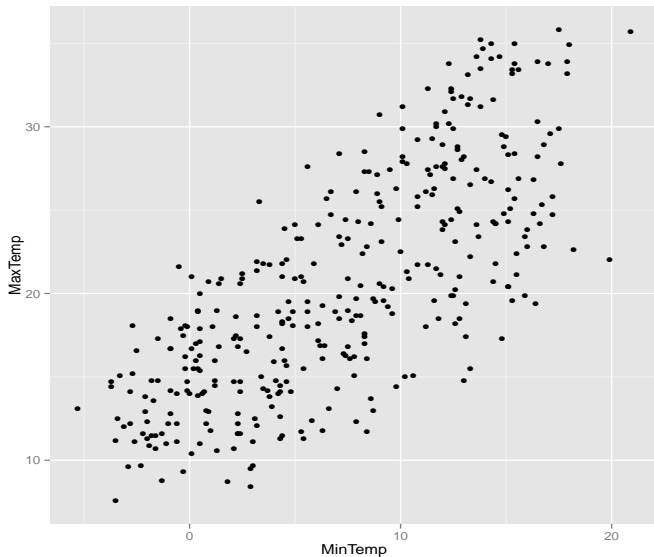
```
library(rattle) # Provides the weather dataset  
library(ggplot2) # Provides the qplot() function
```

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

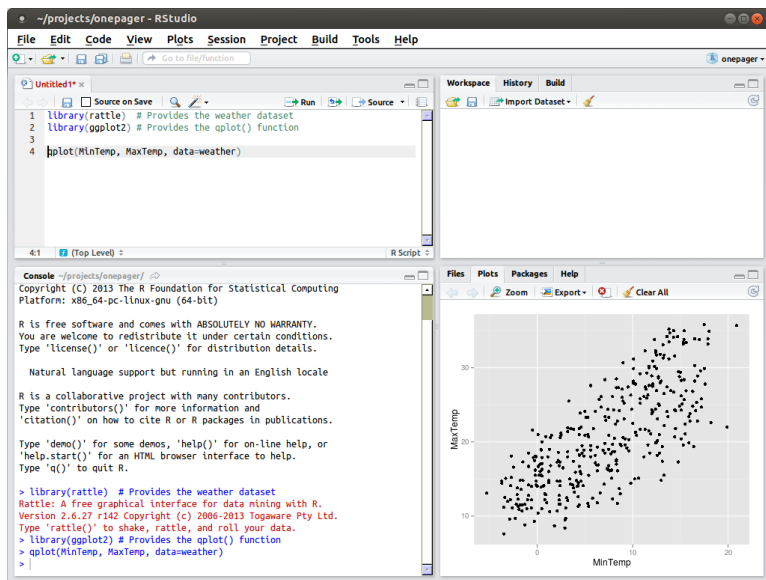
```
ds <- weather  
qplot(MinTemp, MaxTemp, data=ds)
```

- Your turn: give it a go.

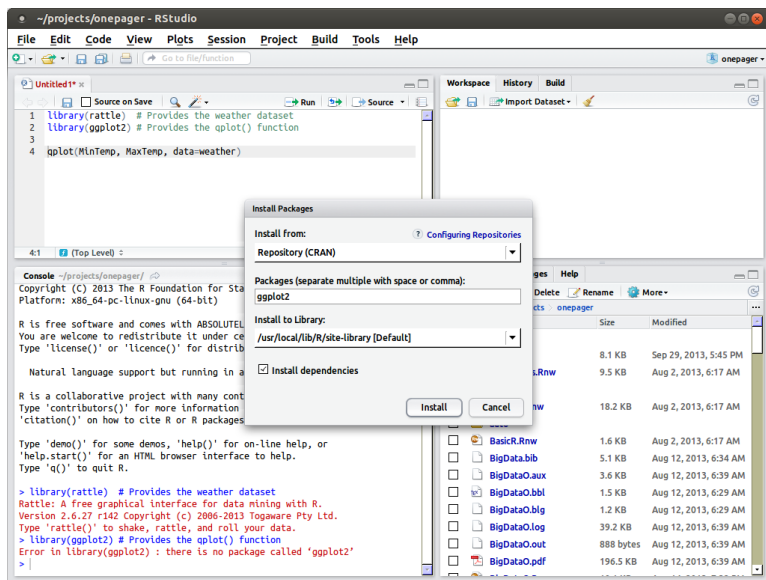
# SCATTERPLOT—PLOT



# SCATTERPLOT—RSTUDIO



# MISSING PACKAGES→TOOLS→INSTALL PACKAGES...





# RSTUDIO—INSTALLING GGPLOT2

The screenshot shows the RStudio interface with the following components:

- Source Editor:** Contains a script with the following code:
 

```
1 library(rattle) # Provides the weather dataset
2 library(ggplot2) # Provides the qplot() function
3
4 qplot(MinTemp, MaxTemp, data=weather)
```
- Console:** Shows the output of the command `install.packages("ggplot2")`. The output indicates that the package was successfully installed from CRAN.
 

```
> install.packages("ggplot2")
Installing package into '/usr/local/lib/R/site-library'
(as 'lib' is unspecified)
trying URL 'http://cran.at.r-project.org/src/contrib/ggplot2_0.9.3.1.tar.gz'
Content type 'application/x-gzip' length 2330942 bytes (2.2 Mb)
opened URL
=====
downloaded 2.2 Mb

* installing *source* package 'ggplot2' ...
** package 'ggplot2' successfully unpacked and MD5 sums checked
** R
** data
*** moving datasets to lazyload DB
** inst
** preparing package for lazy loading
** help
*** installing help indices
** building package indices
** testing if installed package can be loaded
* DONE (ggplot2)

The downloaded source packages are in
'/tmp/Rtmprk1onur/downloaded_packages'
```
- Files Panel:** Displays the file structure of the 'projects/onepager' directory. The files listed are:
 

Name	Size	Modified
..		
.Rhistory	8.1 KB	Sep 29, 2013, 5:45 PM
Animations.Rnw	9.5 KB	Aug 2, 2013, 6:17 AM
archive		
ARulesL.Rnw	18.2 KB	Aug 2, 2013, 6:17 AM
auto		
BasicR.Rnw	1.6 KB	Aug 2, 2013, 6:17 AM
BigData.bib	5.1 KB	Aug 12, 2013, 6:34 AM
BigDataO.aux	3.6 KB	Aug 12, 2013, 6:39 AM
BigDataO.bbl	1.5 KB	Aug 12, 2013, 6:29 AM
BigDataO.bib	1.2 KB	Aug 12, 2013, 6:29 AM
BigDataO.log	39.2 KB	Aug 12, 2013, 6:39 AM
BigDataO.out	888 bytes	Aug 12, 2013, 6:39 AM
BigDataO.pdf	196.5 KB	Aug 12, 2013, 6:39 AM
BigDataO.Rnw	10.1 KB	Aug 14, 2013, 7:32 PM

# 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/ 49 levels "Adelaide","Alba..."
## $ MinTemp        : num  8 14 13.7 13.3 7.6 6.2 6.1 8.3 ...
## ...
```



# BASIC R

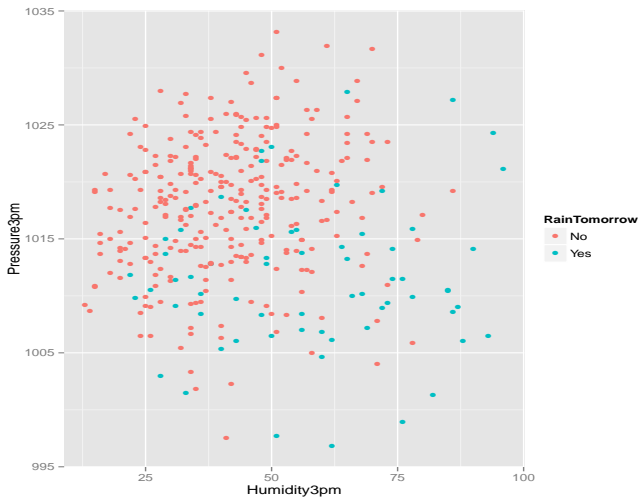
```
summary(weather) # Univariate summary of the variables.
```

```
##           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     : ...
## .....
```



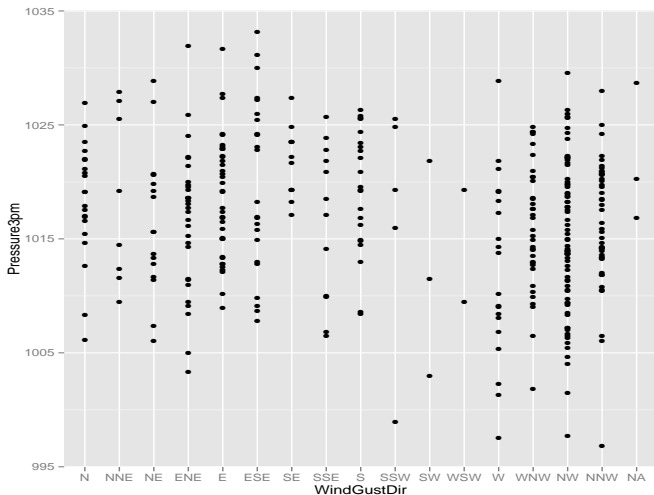
# VISUAL SUMMARIES—ADD A LITTLE COLOUR

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



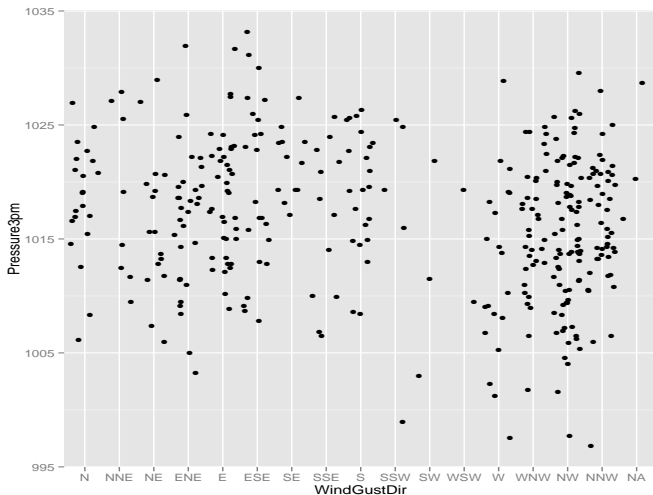
# VISUAL SUMMARIES—CAREFUL WITH CATEGORIES

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



# VISUAL SUMMARIES—ADD A LITTLE JITTER

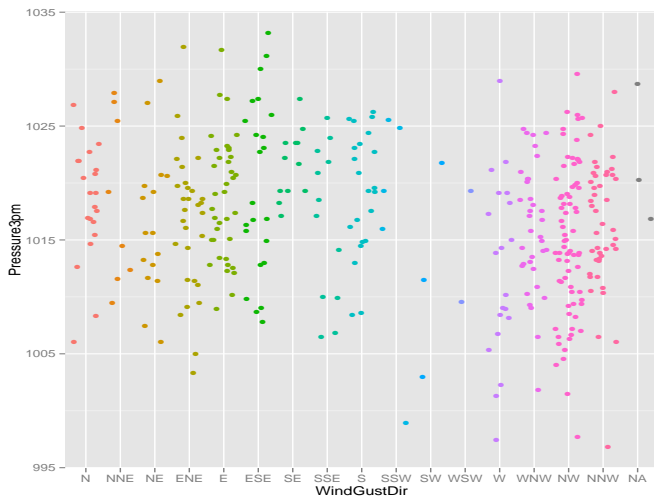
```
qplot(WindGustDir, Pressure3pm, data=ds, geom="jitter")
```





# VISUAL SUMMARIES—AND SOME COLOUR

```
qplot(WindGustDir, Pressure3pm, data=ds, colour=WindGustDir, geom="jitter")
```



# GETTING HELP—PRECEDE COMMAND WITH ?

The screenshot shows the RStudio interface with the following components:

- Source Editor:** Contains an R script with the following code:
 

```
1 library(rattle) # Provides the weather dataset
2 library(ggplot2) # Provides the qplot() function
3
4 qplot(MinTemp, MaxTemp, data=weather)
```
- Console:** Displays the R version (3.0.1), copyright information, and a list of help topics including 'demo()', 'help()', 'help.start()', and 'q()'. It also shows the command to load the ggplot2 package and the qplot function.
 

```
R version 3.0.1 (2013-05-16) -- "Good Sport"
Copyright (C) 2013 The R Foundation for Statistical Computing
Platform: x86_64-pc-linux-gnu (64-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(ggplot2) # Provides the qplot() function
> ?qplot
>
```
- Help Pane:** Displays the documentation for the `qplot` function from the `ggplot2` package. It includes a description of the function, its usage, and a link to the `ggplot2` book for more details.
 

```
R Documentation

Quick plot

Description

qplot is the basic plotting function in the ggplot2 package,
designed to be familiar if you're used to plot from the base
package. It is a convenient wrapper for creating a number of different
types of plots using a consistent calling scheme. See
http://had.co.nz/ggplot2/book/qplot.pdf for the chapter in the
ggplot2 book which describes the usage of qplot in detail.

Usage

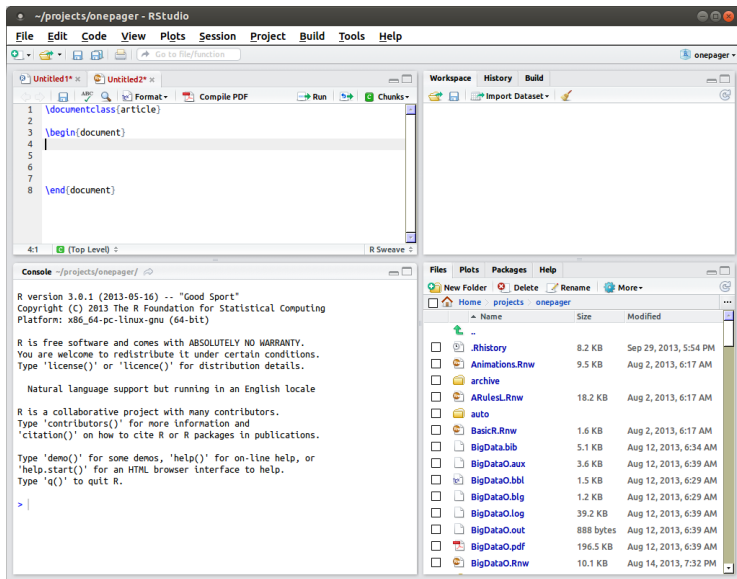
qplot(x, y = NULL, ..., data, facets = NULL,
      margins = FALSE, geom = "auto", stat = list(NULL,
      position = list(NULL), xlim = c(NA, NA),
      ylim = c(NA, NA), log = "", main = NULL,
```



# OVERVIEW

- 1 AN INTRODUCTION TO DATA MINING
- 2 WHY CHOOSE R FOR DATA SCIENCE?
- 3 THE RATTLE PACKAGE FOR DATA MINING
- 4 MOVING INTO R
- 5 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

The screenshot displays the RStudio interface with the following components:

- Editor:** Shows the `sample.Rnw` file with the following content:
 

```

1  \begin{document}
2
3
4
5  \title{Sample Knitr Document}
6  \author{Graham Williams}
7  \maketitle
8
9  \section{My First Section}
10
11 This is some text that is automatically typeset by the LaTeX processor
12 to produce well formatted quality output as PDF.
13
14 
```
- Console:** Shows the execution of the Knitr command:
 

```

> grDevices::pdf(options(useingbats = FALSE); require(knitr);
opts_knit$set(concordance = TRUE); knitr('sample.Rnw', encoding='UTF-8')
Loading required package: knitr

processing file: sample.Rnw
|.....| 100%
ordinary text without R code

output file: sample.tex
[1] "sample.tex"
>
Running pdflatex on sample.tex...completed
Created PDF: ~/projects/onepager/sample.pdf
      
```
- File Explorer:** Shows the project structure for `onepager`:
 

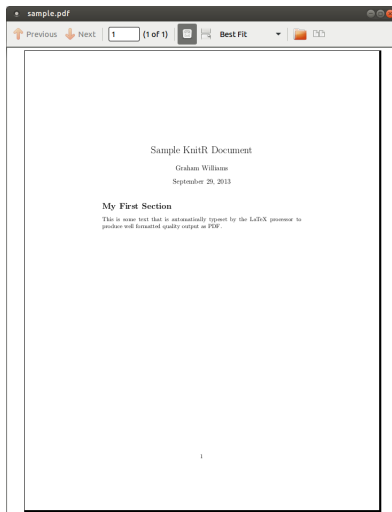
Name	Size	Modified
..		
.Rhistory	8.2 KB	Sep 29, 2013, 5:54 PM
Animations.Rnw	9.5 KB	Aug 2, 2013, 6:17 AM
archive		
ARules.Rnw	18.2 KB	Aug 2, 2013, 6:17 AM
auto		
Basic.R.Rnw	1.6 KB	Aug 2, 2013, 6:17 AM
BigData.bib	5.1 KB	Aug 12, 2013, 6:34 AM
BigDataO.aux	3.6 KB	Aug 12, 2013, 6:39 AM
BigDataO.bbl	1.5 KB	Aug 12, 2013, 6:29 AM
BigDataO.blg	1.2 KB	Aug 12, 2013, 6:29 AM
BigDataO.log	39.2 KB	Aug 12, 2013, 6:39 AM
BigDataO.out	888 bytes	Aug 12, 2013, 6:39 AM
BigDataO.pdf	196.5 KB	Aug 12, 2013, 6:39 AM
BigDataO.Rnw	10.1 KB	Aug 14, 2013, 7:32 PM





# 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.

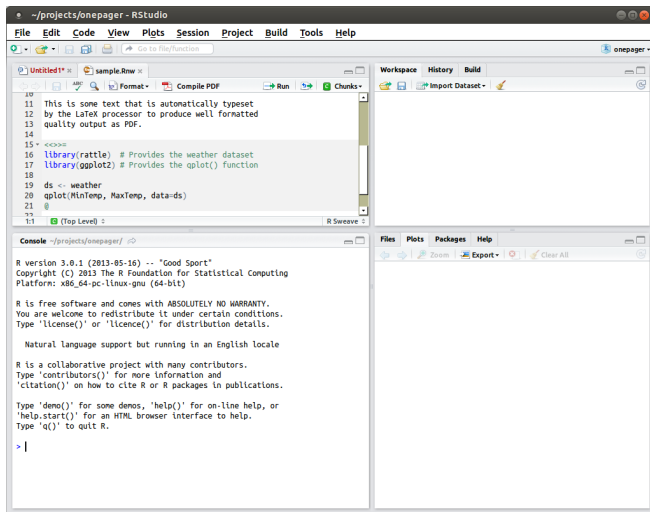
# KNITR: ADD R COMMANDS

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)  
@
```

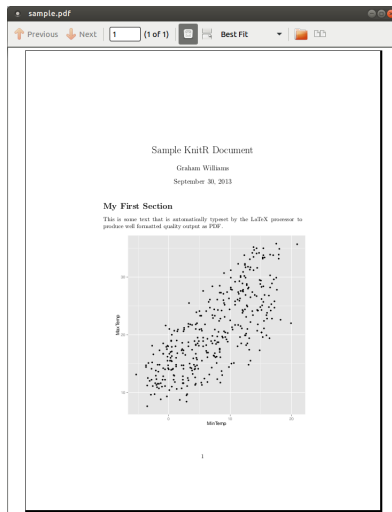
Your turn—Click **Compile PDF** to view the result.

# KNITR DOCUMENT WITH R CODE



# SIMPLE KNITR DOCUMENT—PDF WITH PLOT

## Result of **Compile PDF**



# L<sup>A</sup>T<sub>E</sub>X BASICS

<code>\subsection*{...}</code>	% Introduce a Sub Section
<code>\subsubsection*{...}</code>	% Introduce a Sub Sub Section
<code>\textbf{...}</code>	% Bold font
<code>\textit{...}</code>	% Italic font
<code>\begin{itemize}</code>	% A bullet list
<code>\item ...</code>	
<code>\item ...</code>	
<code>\end{itemize}</code>	

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

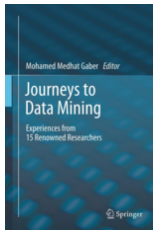
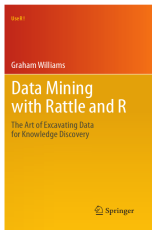
out.width="0.8\\textwidth"    # Fit figure 80% page width
out.height="0.5\\textheight"  # Fit figure 50% page height
```

Plus an extensive collection of other options.



# 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)





# LECTURE SUMMARY

- Data Science—Analytics—Data Mining;
- Rattle as a GUI for Quickly Analysing Data;
- The Power is with R.

*This document, sourced from StartL.Rnw revision 436, was processed by KnitR version 1.6 of 2014-05-24 and took 4.9 seconds to process. It was generated by gjw on nyx running Ubuntu 14.04 LTS with Intel(R) Xeon(R) CPU W3520 @ 2.67GHz having 4 cores and 12.3GB of RAM. It completed the processing 2014-06-21 20:26:17.*

