# Data Science with R

## Introducing and Interacting with 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

# Overview

# 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
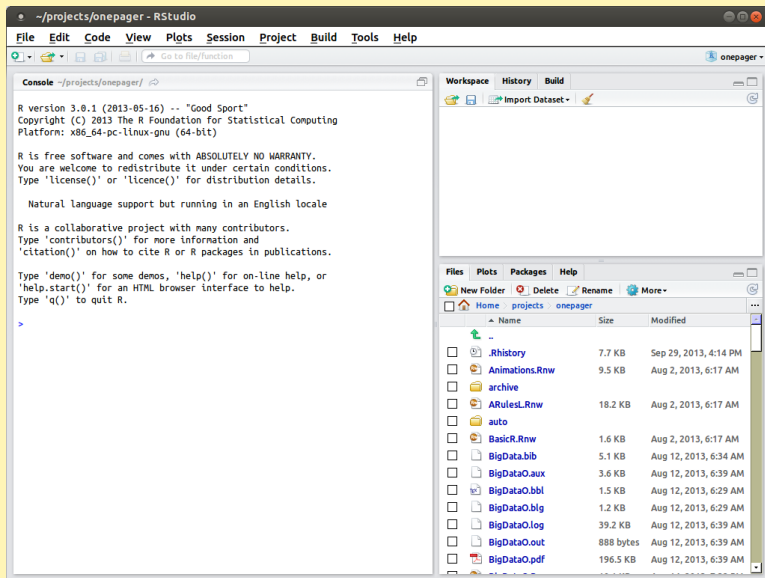  - Produce beautiful documents, easily reproducible, FLOSS

# Using Ubuntu

- Desktop Ubuntu

- Connecting to Analytics Servers
  - Using XWin
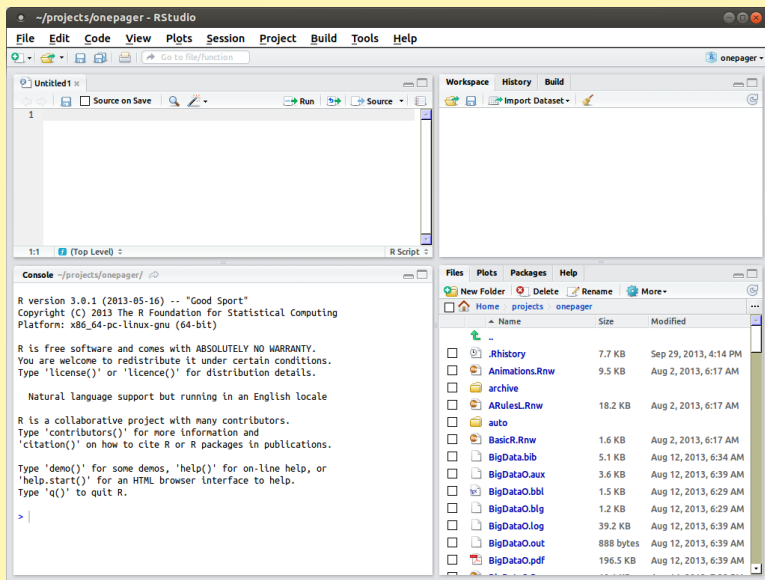  - Using VNC

- Start up RStudio from the Dash

# Overview

# RStudio—The Default Three Panels

# RStudio—With R Script File—Editor Panel

# Overview

# Scatterplot—R Code

Our first little bit of R code:

- Load a couple of *packages* into the R *library*
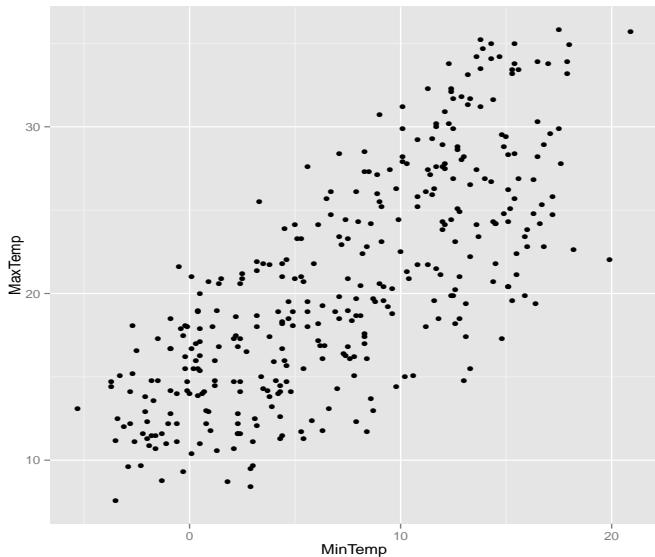
```r
library(rattle)   # Provides the weather dataset
library(ggplot2) # Provides the qplot() function
```

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

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

- Your turn: give it a go.

# SCATTERPLOT—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
```

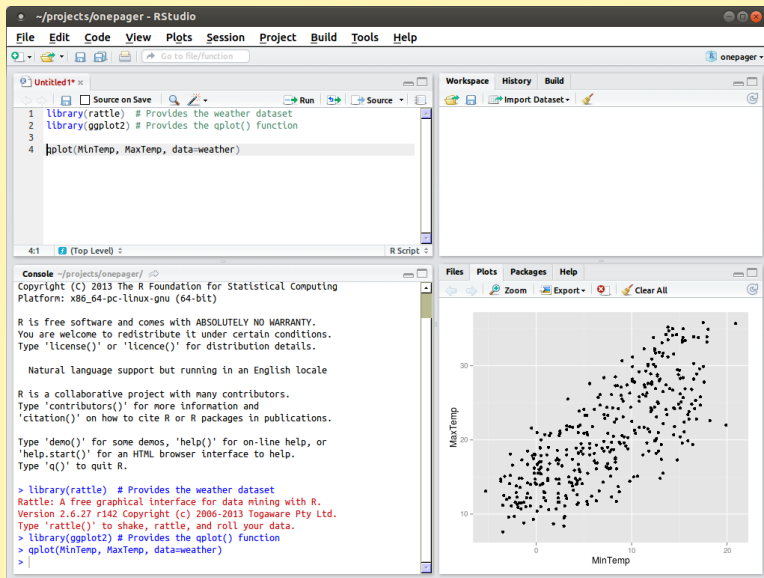- Then produce a quick plot using `qplot()`

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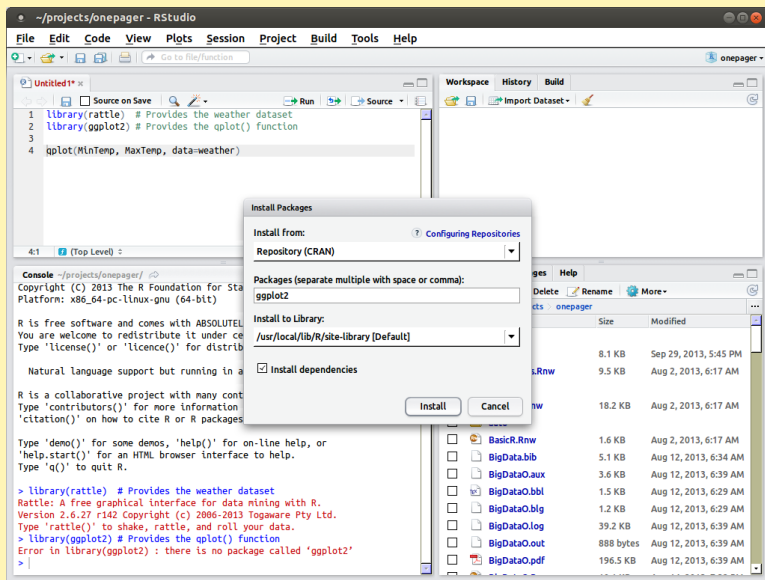
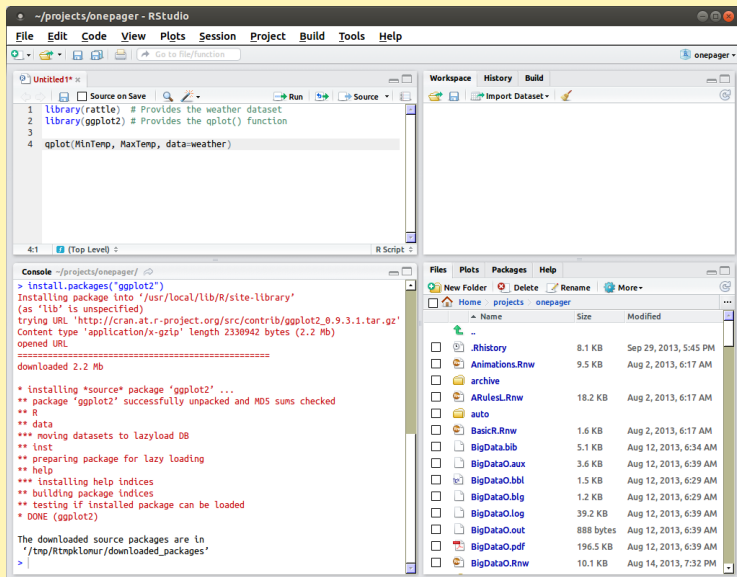- Your turn: give it a go.

# SCATTERPLOT—PLOT

# Scatterplot—RStudio

# Missing Packages–Tools→Install Packages. . .

# RStudio—Installing ggplot2

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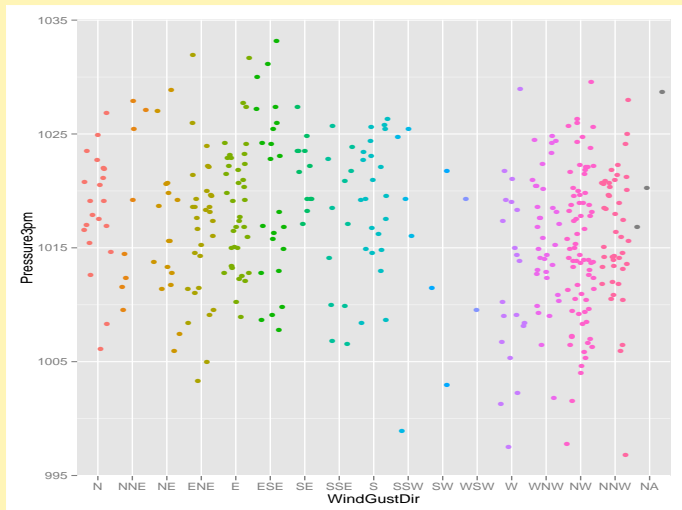
`qplot(WindGustDir, Pressure3pm, data=ds)`

# 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

# 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)
@
```

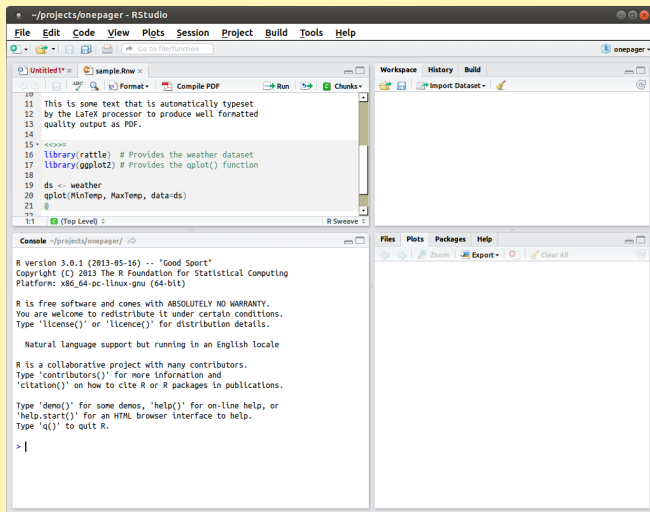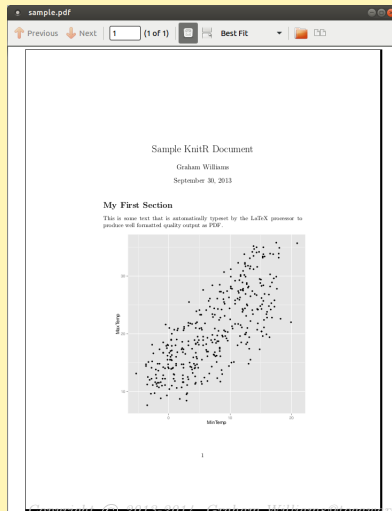Your turn—Click **Compile PDF** to view the result.

# 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.

 32/

# 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.

# Thank You

Question Time

This document, sourced from IntroRL.Rnw revision 282, was processed by KnitR version 1.5 of 2013-09-28 and took 2.4 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 06:19:56.