

DATA SCIENCE WITH R

DOCUMENTING PROJECTS WITH KNITR

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

- 1 MOTIVATION
- 2 USING RSTUDIO
- 3 BASIC L^AT_EX MARKUP
- 4 INCORPORATING R CODE
- 5 FORMATTING TABLES AND PLOTS
- 6 GETTING SOPHISTICATED
- 7 SUMMARY



OVERVIEW

- 1 MOTIVATION
- 2 USING RSTUDIO
- 3 BASIC L^AT_EX MARKUP
- 4 INCORPORATING R CODE
- 5 FORMATTING TABLES AND PLOTS
- 6 GETTING SOPHISTICATED
- 7 SUMMARY



WHY IS REPRODUCIBILITY IMPORTANT?

Your SES drops by and asks:

- “Remember that analysis you did last year? I’ve heard there is an update on the data that you used. Can you add the new data in and repeat the same analysis?”
- “Jo Bloggs did a great analysis of the company returns data just before she left. Can you get someone else to analyse the new data set using the same methods, and so produce an updated report that will be understandable to the Exec?”
- “That case you provided an analysis of last year has finally reached the courts. We need to ensure we have a clear trail of the data sources, the analyses performed, and the results obtained, to stand up in court. Could you document these please.”



LITERATE DATA MINING OVERVIEW

- One document to intermix the analysis, code, and results
- Authors productive with narrative and code in one document
- Sweave (Leisch 2002) and now KnitR (Yihui 2011)
- Embed R code into \LaTeX documents for typesetting
- KnitR also supports publishing to the web



WHY REPRODUCIBLE DATA MINING?

- **Automatically** regenerate documents when code, data, or assumptions change.
- Eliminate errors that occur when transcribing results into documents.
- Record the context for the analysis and decisions made about the type of analysis to perform in the one place.
- Document the processes to provide integrity for the conclusions of the analysis.
- Share approach with others for peer review and for learning from each other—engender a continuous learning environment.



PRIME OBJECTIVE: TRUSTWORTHY SOFTWARE

*Those who receive the results of modern data analysis have limited opportunity to **verify the results** by direct observation. Users of the analysis have no option but to **trust the analysis**, and by extension the software that produced it. This places an **obligation** on all creators of software to program in such a way that the **computations can be understood and trusted**. This obligation I label the Prime Directive.*

John Chambers, *Software for Data Analysis: Programming with R*



BEAUTIFUL OUTPUT BY DEFAULT

The reader wants to read the document and easily do so!

- Code highlighting is done automatically
- Default theme is carefully designed
- Many other themes are available
- R Code is “properly” reformatted

OVERVIEW

- 1 MOTIVATION
- 2 USING RSTUDIO
- 3 BASIC L^AT_EX MARKUP
- 4 INCORPORATING R CODE
- 5 FORMATTING TABLES AND PLOTS
- 6 GETTING SOPHISTICATED
- 7 SUMMARY



SUPPORTING TECHNOLOGY

A suite of Free and Open Source Software — FLOSS

- RStudio — Creating, managing, compiling documents
- L^AT_EX — Markup language for typesetting a document
- R — Statistical analysis language
- KnitR — Integrator of typesetting and analysis



USING RSTUDIO

- Simplified interaction with R, \LaTeX , and KnitR
- Executes R code one line at a time
- Formats \LaTeX documents and provides spell checking
- A single click compile to PDF and synchronised views

Demonstrate: Startup and explore RStudio.

OVERVIEW

- 1 MOTIVATION
- 2 USING RSTUDIO
- 3 BASIC L^AT_EX MARKUP
- 4 INCORPORATING R CODE
- 5 FORMATTING TABLES AND PLOTS
- 6 GETTING SOPHISTICATED
- 7 SUMMARY



INTRODUCING L^AT_EX

- A text markup language rather than a WYSIWYG.
- Based on T_EX from 1977 — very stable and powerful.
- L^AT_EX is easier to use macro package built on T_EX.
- Ensures consistent style (layout, fonts, tables, maths, etc.)
- Automatic indexes, footnotes and references.
- Documents are well structured and are clear text.
- Has a learning curve.



BASIC L^AT_EX USAGE

```
\documentclass{article}
```

```
\begin{document}
```

```
\end{document}
```

Demonstrate Create a new Sweave document in RStudio

STRUCTURES

```
\documentclass{article}
```

```
\begin{document}
```

```
\section{Introduction}
```

```
...
```

```
\subsection{Concepts}
```

```
...
```

```
\end{document}
```



FORMATS

```
\documentclass{article}
```

```
\begin{document}
```

```
\begin{itemize}
```

```
  \item ABC
```

```
  \item DEF
```

```
\end{itemize}
```

```
This if \textbf{bold} text or \textbf{italic} text, ...
```

```
\end{document}
```


RSTUDIO SUPPORT FOR L^AT_EX

RStudio provides excellent support for working with L^AT_EX documents

Helps to avoid having to know too much about L^AT_EX

Best illustrated through a demonstration

- Format menu
 - Section commands
 - Font commands
 - List commands
 - Verbatim/Block commands
- Spell Checker
- Compile PDF

Demonstrate: Start a new document, add contents, format to PDF.

OVERVIEW

- 1 MOTIVATION
- 2 USING RSTUDIO
- 3 BASIC L^AT_EX MARKUP
- 4 INCORPORATING R CODE
- 5 FORMATTING TABLES AND PLOTS
- 6 GETTING SOPHISTICATED
- 7 SUMMARY



INCORPORATING R CODE

- We insert R code in a *Chunk* starting with `<< >>=`
- We terminate the Chunk with `@`
- Save \LaTeX with extension `Rnw`

This Chunk

```
<<simple_example>>=
x <- sum(1:10)
x
@
```

Produces

```
x <- sum(1:10)
x

## [1] 55
```

- *Demonstrate:* Do this in RStudio

MAKING YOU LOOK GOOD

```
<<format_example>>=
for(i in 1:5){j<-cos(sin(i)*i^2)+3;print(j-5)}
@

for(i in 1:5)
{
  j <- cos(sin(i)*i^2)+3
  print(j-5)
}

## [1] -1.334
## [1] -2.88
## [1] -1.704
## [1] -1.103
....
```

R WITHIN THE TEXT

- Include information about data within the narrative.
- We can do that with `\Sexpr{...}`.

Our dataset has `\Sexpr{nrow(ds)}` observations of `\Sexpr{ncol(ds)}` variables.

Becomes

Our dataset has 82169 observations of 24 variables.

Better Still: `\Sexpr{format(nrow(ds), big.mark=",")}`

Our dataset has 82,169 observations of 24 variables.



OVERVIEW

- 1 MOTIVATION
- 2 USING RSTUDIO
- 3 BASIC L^AT_EX MARKUP
- 4 INCORPORATING R CODE
- 5 FORMATTING TABLES AND PLOTS**
- 6 GETTING SOPHISTICATED
- 7 SUMMARY



A SIMPLE TABLE

```
library(xtable)
obs <- sample(1:nrow(weatherAUS), 8)
vars <- 2:6
xtable(weatherAUS[obs, vars])
```

	Location	MinTemp	MaxTemp	Rainfall	Evaporation
73129	Walpole	12.90	29.00	0.00	
34446	Bendigo	12.00	24.30	0.00	3.00
19601	SydneyAirport	20.00	28.20	0.00	5.80
61331	Woomera	7.30	18.80	0.00	6.40
738	Albury	20.30	29.70	3.20	
70231	Perth	11.00	26.90	0.00	5.40
12376	NorfolkIsland	14.60	20.60	0.00	3.20
48795	Brisbane	13.30	22.50	0.00	



TABLE: EXCLUDE ROW NAMES

```
print(xtable(weatherAUS[obs, vars]),
      include.rownames=FALSE)
```

Location	MinTemp	MaxTemp	Rainfall	Evaporation
Walpole	12.90	29.00	0.00	
Bendigo	12.00	24.30	0.00	3.00
SydneyAirport	20.00	28.20	0.00	5.80
Woomera	7.30	18.80	0.00	6.40
Albury	20.30	29.70	3.20	
Perth	11.00	26.90	0.00	5.40
NorfolkIsland	14.60	20.60	0.00	3.20
Brisbane	13.30	22.50	0.00	2.20



TABLE: LIMIT NUMBER OF DIGITS

```
print(xtable(weatherAUS[obs, vars],
             digits=1),
      include.rownames=FALSE)
```

Location	MinTemp	MaxTemp	Rainfall	Evaporation
Walpole	12.9	29.0	0.0	
Bendigo	12.0	24.3	0.0	3.0
SydneyAirport	20.0	28.2	0.0	5.8
Woomera	7.3	18.8	0.0	6.4
Albury	20.3	29.7	3.2	
Perth	11.0	26.9	0.0	5.4
NorfolkIsland	14.6	20.6	0.0	3.2
Brisbane	13.3	22.5	0.0	2.2



TABLE: TINY FONT

```
vars <- 2:8
print(xtable(weatherAUS[obs, vars],
             digits=0),
      size="tiny",
      include.rownames=FALSE)
```

Location	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir
Walpole	13	29	0			NE
Bendigo	12	24	0	3		NE
SydneyAirport	20	28	0	6	11	NE
Woomera	7	19	0	6	12	SSE
Albury	20	30	3			NNW
Perth	11	27	0	5	11	SSW
NorfolkIsland	15	21	0	3	4	WNW
Brisbane	13	22	0	2	3	W

TABLE: COLUMN ALIGNMENT

```
vars <- 2:8
print(xtable(weatherAUS[obs, vars],
             digits=0,
             align="rlrrrrrr"),
      size="tiny")
```

	Location	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir
73129	Walpole	13	29	0			NE
34446	Bendigo	12	24	0	3		NE
19601	SydneyAirport	20	28	0	6	11	NE
61331	Woomera	7	19	0	6	12	SSE
738	Albury	20	30	3			NNW
70231	Perth	11	27	0	5	11	SSW
12376	NorfolkIsland	15	21	0	3	4	WNW
48795	Brisbane	13	22	0	2	3	W

TABLE: CAPTION

```
print(xtable(weatherAUS[obs, vars],
            digits=1,
            caption="This is the table caption."),
      size="tiny")
```

	Location	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir
73129	Walpole	12.9	29.0	0.0			NE
34446	Bendigo	12.0	24.3	0.0	3.0		NE
19601	SydneyAirport	20.0	28.2	0.0	5.8	10.8	NE
61331	Woomera	7.3	18.8	0.0	6.4	11.6	SSE
738	Albury	20.3	29.7	3.2			NNW
70231	Perth	11.0	26.9	0.0	5.4	11.1	SSW
12376	NorfolkIsland	14.6	20.6	0.0	3.2	3.9	WNW
48795	Brisbane	13.3	22.5	0.0	2.2	2.9	W

TABLE : This is the table caption.

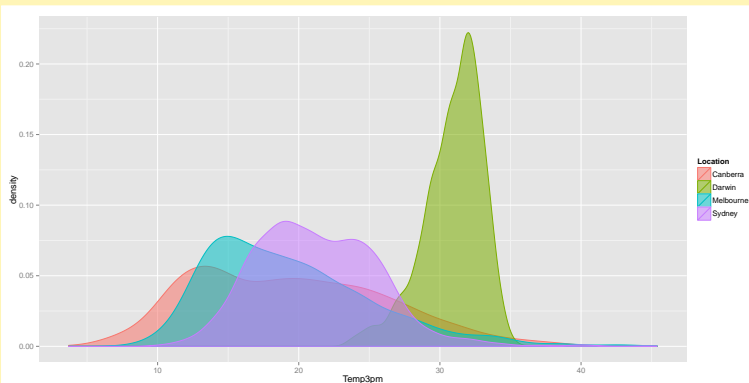
OVERVIEW

- 1 MOTIVATION
- 2 USING RSTUDIO
- 3 BASIC L^AT_EX MARKUP
- 4 INCORPORATING R CODE
- 5 FORMATTING TABLES AND PLOTS**
- 6 GETTING SOPHISTICATED
- 7 SUMMARY



PLOTS

```
library(ggplot2)
cities <- c("Canberra", "Darwin", "Melbourne", "Sydney")
ds <- subset(weatherAUS, Location %in% cities & ! is.na(Temp3pm))
g <- ggplot(ds, aes(Temp3pm, colour=Location, fill=Location))
g <- g + geom_density(alpha = 0.55)
print(g)
```



OVERVIEW

- 1 MOTIVATION
- 2 USING RSTUDIO
- 3 BASIC L^AT_EX MARKUP
- 4 INCORPORATING R CODE
- 5 FORMATTING TABLES AND PLOTS
- 6 GETTING SOPHISTICATED
- 7 SUMMARY



ADVANCED TOPIC—KNITR AND ESS AND EMACS

Demonstration



ACTUAL EXAMPLES

- Linked Risk Visualisations
- Visualising Clusters
- Siebel Case Profile Attachments
- Specifications of Rule-Based Model Logic



OVERVIEW

- 1 MOTIVATION
- 2 USING RSTUDIO
- 3 BASIC L^AT_EX MARKUP
- 4 INCORPORATING R CODE
- 5 FORMATTING TABLES AND PLOTS
- 6 GETTING SOPHISTICATED
- 7 SUMMARY



SUMMARY

- Document as we go to record all modelling activity
- Ensure transparency, repeatability, sharing
- Mature technology: \LaTeX and R
- Modern support: KnitR and RStudio



FURTHER READING

- <http://onepager.togaware.com/>
- <http://yihui.name/knitr/>
- <http://www.rstudio.org/>
- <http://yihui.name/slides/2012-knitr-RStudio.html>
- <http://bcb.dfci.harvard.edu/~aedin/courses/ReproducibleResearch/ReproducibleResearch.pdf>
- https://dl.dropbox.com/u/233041/Bios301/lecture2_knitr.html
- <http://biostat.mc.vanderbilt.edu/wiki/pub/Main/ReproducibleResearchTutorial/HarrellScottTutorial-useR2012.pdf>



FURTHER READING

- <http://onepager.togaware.com/>
- <http://yihui.name/knitr/>
- <http://www.rstudio.org/>
- <http://yihui.name/slides/2012-knitr-RStudio.html>
- <http://bcb.dfci.harvard.edu/~aedin/courses/ReproducibleResearch/ReproducibleResearch.pdf>
- https://dl.dropbox.com/u/233041/Bios301/lecture2_knitr.html
- <http://biostat.mc.vanderbilt.edu/wiki/pub/Main/ReproducibleResearchTutorial/HarrellScottTutorial-useR2012.pdf>

