

DATA MINING IN GOVERNMENT

FROM R TO RATTLE AND BACK AGAIN

Graham Williams

Chief Data Miner and Director of Analytics
Office of the Chief Knowledge Officer
Australian Taxation Office

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

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

Download instructions: <http://rattle.togaware.com>

Copyright © 2008 Graham.Williams@togaware.com

OVERVIEW

ANALYTICS CASE STUDY

Data Mining
Australian Taxation Office

CHALLENGE — DELIVER DATA MINING CAPABILITY

Competing Technologies
Commodity and Open Source

RATTLE IN ACTION

Rattle Screenshots
Moving Forward with Rattle

Copyright © 2008 Graham.Williams@togaware.com

DATA MINING — IN A NUTSHELL

- The **non-trivial** extraction of **novel, implicit, and actionable** knowledge from **large** databases **and in a timely manner**.
- Build **models** (regression, neural networks, decision trees, random forests, boosted trees, support vector machines, association and correlation analysis, cluster analysis, ...) **from data** that represent observations of the world, and often whole populations.
- **Knowledge from data any way we can**.
- Technology to enable data processing, data understanding and visualisation, and data analysis of very large databases for insight and hypothesis generation.

Copyright © 2008 Graham.Williams@togaware.com

AUSTRALIAN TAXATION OFFICE — CASE STUDY

- Employs 22,000 staff Australia wide
- Revenue Collection and Refund Management
- Compliance and Risk Modelling

- 12M Individuals, \$450B Income, \$100B Tax
- 2M Companies..., \$1800B Income, \$40B Tax
- Deductions \$B

- Tax payer's charter:
Fair but firm; Protect privacy; Assume honest
- Service standards — turn around refunds
- **Whilst protecting integrity of revenue collection**

Copyright © 2008 Graham.Williams@togaware.com

ATO ANALYTICS — DEPLOYING DATA MINING

- Established as an organisational capability in 2003
Add a scientific basis to managing Risk
- Team of 16 data mining specialists
Statistics and Machine Learning
- Support 120 data analysts
- Spread data mining throughout the organisation through a central centre of excellence
- Provide a unified framework for Risk Management across the organisation

Copyright © 2008 Graham.Williams@togaware.com

ANALYTICS CASE STUDY — HIGH RISK REFUNDS

- High Risk Refunds (HRR) identified prior to payment of refunds.
- Current business rules identify too many "high risk" refunds.
 - We might identify 100,000 cases each year.
 - Sometimes as few as 5% require adjustment.
 - But revenue at risk is significant (from \$10m to \$1b).

Current practise of simple business rules based on experience:

- Total claimed investment deductions > \$N
- Ratio of self education deductions to total income > N
- Total international transfers > N times taxable income
- Luxury vehicle purchase \$M > N times taxable income

Copyright © 2008 Graham.Williams@togaware.com

ANALYTICS CASE STUDY — HIGH RISK REFUNDS

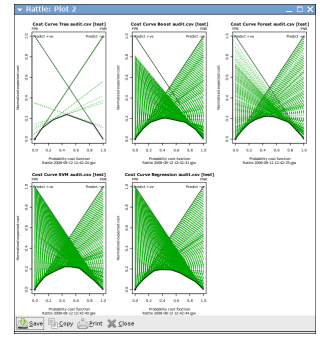
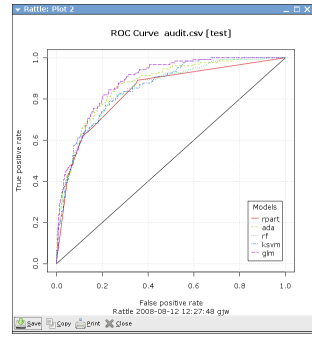
Almost, any kind of modelling will help ... data mining for HRR

- Major task is all about the data:
 - data understanding/preparation, feature generation/selection
 - 100,000 cases by 1,000 variables
- Stock and trade: glm, rpart, ada, randomForest, kernlab
- Simple binary classification and \$ regression
- Identify new characteristics to target high risk (5%);
- Focus resources on productive cases - \$ and tax payer benefit;
- Decision trees and ensembles (random forests) are often effective.

Copyright © 2008 Graham.Williams@togaware.com

ANALYTICS CASE STUDY — HIGH RISK REFUNDS

“Out of the box” model performance:



Copyright © 2008 Graham.Williams@togaware.com

OTHER AREAS OF MODELLING

- High Risk Refunds
- Required to Lodge (\$110M)
- Assessing Levels of Debt
 - Propensity to Pay
 - Capacity to Pay
- Determining Optimal Treatment Strategies
- Identity Theft — eTax and International
- Project Wickenby Text Mining
- Tax Havens — AUSTRAC

Copyright © 2008 Graham.Williams@togaware.com

OVERVIEW

ANALYTICS CASE STUDY
Data Mining
Australian Taxation Office

CHALLENGE — DELIVER DATA MINING CAPABILITY
Competing Technologies
Commodity and Open Source

RATTLE IN ACTION
Rattle Screenshots
Moving Forward with Rattle

Copyright © 2008 Graham.Williams@togaware.com

PACKAGING ANALYTICS

- Simple enough for the R jockey
- But how to roll this capability out to the point-n-click generation
- SAS/EM and SPSS/Clementine and KXEN and Salford Systems ...
- But, why not invest in people, not vendors?
- R provides quite a capable data mining environment, and has much more to offer than the competitors.
- Challenge: How to roll R out to a large community of data analysts?
- Answer: “Package” (or template) analyses tuned for the kinds of analyses required, reviewed and vetted by our Statisticians and Data Miners.

Copyright © 2008 Graham.Williams@togaware.com

BACKGROUND — TECHNOLOGIES

- Originally purchased data mining tools (SAS/EM, Teradata Warehouse Miner) and hardware (Big Iron MS/Windows 32 bit).
- **But** data mining needs skilled people, not off the shelf solutions (yet).
- **Also** data mining technology is rapidly developing, and commercial tools not always up to date.



Copyright © 2008 Graham.Williams@togaware.com

NEW APPROACHES ENSEMBLES

Commercial software is lagging behind advances in Data Mining (PhD 1989 combining multiple decision trees.... not new technology)

- Current best off the shelf technology includes random forests, boosting and support vector machines - SAS/EM, SPSS?
- Open source solutions allow organisational investment in people.
- Decide in 2005 to move to open source platform — 3 years to deploy!!



Copyright © 2008 Graham.Williams@togaware.com

ANALYTICSNET

Build a network of DataMining Nodes:

- 1 CPU (2 Cores), AMD64, 16GB RAM, 300GB Disk
- 4 CPU (8 Cores), AMD64, 32GB RAM, 1TB Disk (Good Price Point)
- 8 CPU (16 Cores), AMD64, 128GB RAM, 10TB Disk (Near Term)



- Best of class open source operating system (Debian GNU/Linux)
- Open Source data mining tools R, Rattle, Weka, AlphaMiner
- Open Source **does** deliver quality software
A difficult sell (commercial interests).

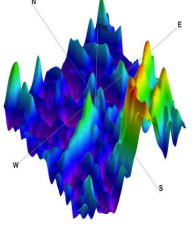
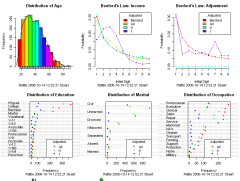
Data source: Data Warehouse (Teradata or Netezza or SQLite) as the workhorse data server

Copyright © 2008 Graham.Williams@togaware.com

THE CHALLENGE REMAINS...

But how to roll out R across an organisation?

- Shortage of skilled data miners
- Skilled data analysts
- What do you mean "a CLI?"
- Development of training - using a GUI
- Integrity of statistical analyses
- Community of Practise for peer review



Copyright © 2008 Graham.Williams@togaware.com

OVERVIEW

ANALYTICS CASE STUDY
Data Mining
Australian Taxation Office

CHALLENGE — DELIVER DATA MINING CAPABILITY
Competing Technologies
Commodity and Open Source

RATTLE IN ACTION
Rattle Screenshots
Moving Forward with Rattle

Copyright © 2008 Graham.Williams@togaware.com

INTRODUCING RATTLE

A stepping stone into the full power of R (quick ref) — original plan or

A self contained tool for data mining — sufficient but limiting

- 1 Load Data
CSV, ARFF, ODBC
- 2 Explore
Plots, GGobi
- 3 Transform
Rescale, Remap
- 4 Cluster, Associate
- 5 Predictive Models
Tree, Forest, NNet
- 6 Evaluate and Score
ROC, Risk, Cost
- 7 Log

```

Rattle: R Console
> packageDescription("rattle")
Package: rattle
Type: Package
Title: A graphical user interface for data mining in R using GTK
Version: 2.3.69
Date: 2008-08-10
Author: Graham Williams <Graham.Williams@togaware.com>
Maintainer: Graham Williams <Graham.Williams@togaware.com>
Depends: R (>= 2.2.0)
Suggests: RGtk2, ada, amap, arules, bitops, cairoDevice, car, cba,
  combinat, doby, e1071, ellipse, fEcofin, fCalendar, fBasics,
  foreign, fpc, gdata, gtools, gplots, Hmisc, kernlab, MASS,
  Matrix, mice, network, odfeave, playwith, pml, randomForest,
  reshape, rggobi, ROCR, RODBC, rpart, RSvgDevice, XML
Description: Rattle provides a Gnome (RGtk2) based interface to R
  functionality for data mining. The aim is to provide a simple
  and intuitive interface that allows a user to quickly load data
  from a CSV file (or via ODBC), transform and explore the data,
  and build and evaluate models, and export models as PMML
  (predictive modelling markup language). All of this with
  knowing little about R. All R commands are logged and available
  for the user, as a tool to then begin interacting directly with
  R itself, if so desired. Rattle also exports a number of
  utility functions and the graphical user interface does not
  need to be run to deploy these.
License: GPL (>= 2)
URL: http://rattle.togaware.com/
Packaged: Sun Aug 10 02:00:29 2008; gjw
Built: R 2.7.1; 2008-08-10 02:02:40; unix
  
```

Copyright © 2008 Graham.Williams@togaware.com

OVERVIEW

ANALYTICS CASE STUDY
Data Mining
Australian Taxation Office

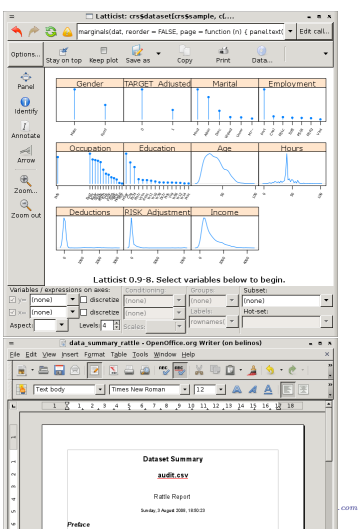
CHALLENGE — DELIVER DATA MINING CAPABILITY
Competing Technologies
Commodity and Open Source

RATTLE IN ACTION
Rattle Screenshots
Moving Forward with Rattle

Copyright © 2008 Graham.Williams@togaware.com

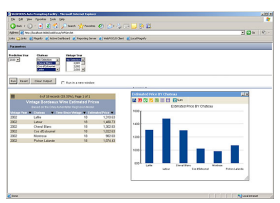
FUTURE DIRECTIONS

- latticist and playwith
- odfWeave report generator
- tm for text mining
- pmml consumer
- dealing with large data: wide and deep



BUSINESS INTELLIGENCE AND DATA MINING

- Information Builders BI Tool — WebFOCUS (Announced 2 Jun)
- Incorporate Open Source Rattle Badged as RStat Open Source Contributions
- SPSS can now create an R data frame object and (somehow — but not in the same memory space?) shares the object between R and SPSS.
- Allows customers to drive the SPSS technology directions!



INTEROPERABILITLY USING PMML

- Share models: import and export
- Commercial (... extensions are bad) and Open Source vendor support
- Uptake of the technology?
- Rattle's PMML now exports PMML for:
 - lm; glm; rpart; kmeans; randomForest; randomSurvivalForest; arules, nnet
- Import into Zementis' ADAPA (Amazon Clouds)
- Integrated into WebFOCUS

RESOURCES

- Data Mining
 - datamining.togaware.com
- Tools:
 - rattle.togaware.com (Rattle and PMML beta versions)
 - www.cs.waikato.ac.nz/ml/weka/
 - www.knime.org
 - rapid-i.com

Thank You