

DELIVERY OF AN ANALYTICS CAPABILITY

ENSEMBLES AND MODEL DELIVERY FOR TAX COMPLIANCE

Graham Williams

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

Adjunct Professor, University of Canberra
Adjunct Professor, Australian National University
Visiting Professor, Chinese Academy of Sciences
Fellow, Institute of Analytics Professionals of Australia

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

OVERVIEW

REVENUE COLLECTION

ESTABLISHING CAPABILITY

DELIVERING OUTCOMES

ONGOING CHALLENGES

Model Deployment

Big Data: Many Variables

“nothing ... [is] certain, except death and taxes.”

Benjamin Franklin, 1789

“The hardest thing in the world to understand is the income tax.”

attributed to Albert Einstein

“The avoidance of taxes is the only intellectual pursuit that carries any reward.”

John Maynard Keynes

“nothing ... [is] certain, except death and taxes.”

Benjamin Franklin, 1789

“The hardest thing in the world to understand is the income tax.”

attributed to Albert Einstein

“The avoidance of taxes is the only intellectual pursuit that carries any reward.”

John Maynard Keynes

“nothing ... [is] certain, except death and taxes.”

Benjamin Franklin, 1789

“The hardest thing in the world to understand is the income tax.”

attributed to Albert Einstein

“The avoidance of taxes is the only intellectual pursuit that carries any reward.”

John Maynard Keynes

TAXATION RECORDS - HISTORIC DATA COLLECTION

- **Taxation and Governance**

Revenue is fundamental to government.

- **Data Collection**

2500 B.C. Sumerian (Iraq) and Elam (Iran) peoples marked their tax records onto dried mud tablets.

- We still inscribe tax records today — now electronically

- **Data Mining Opportunities**

Analyse very large collections (20M by 5K) to check tax payer identity, facilitate reporting (lodgement), ensure compliance, pay refunds or collect debts: the four pillars.



Sources <http://www.upenn.edu/almanac/v48/n28/AncientTaxes.html> <http://www.crystalinks.com/cuneiformtablets.html>

TAXATION RECORDS - HISTORIC DATA COLLECTION

- **Taxation and Governance**

Revenue is fundamental to government.

- **Data Collection**

2500 B.C. Sumerian (Iraq) and Elam (Iran) peoples marked their tax records onto dried mud tablets.

- We still inscribe tax records today — now electronically

- **Data Mining Opportunities**

Analyse very large collections (20M by 5K) to check tax payer identity, facilitate reporting (lodgement), ensure compliance, pay refunds or collect debts: the four pillars.



Sources <http://www.upenn.edu/almanac/v48/n28/AncientTaxes.html> <http://www.crystalinks.com/cuneiformtablets.html>

TAXATION RECORDS - HISTORIC DATA COLLECTION

- **Taxation and Governance**

Revenue is fundamental to government.

- **Data Collection**

2500 B.C. Sumerian (Iraq) and Elam (Iran) peoples marked their tax records onto dried mud tablets.

- **We still inscribe tax records today — now electronically**

- **Data Mining Opportunities**

Analyse very large collections (20M by 5K) to check tax payer identity, facilitate reporting (lodgement), ensure compliance, pay refunds or collect debts: the four pillars.



Sources <http://www.upenn.edu/almanac/v48/n28/AncientTaxes.html> <http://www.crystalinks.com/cuneiformtablets.html>

AUSTRALIAN TAXATION OFFICE

...continues this “fine tradition”

- Employs 22,000 staff Australia wide
- Role to collect revenue and process refunds

- 12M Individuals, \$600B Income, \$120B Tax
- 2M Companies..., \$2,200B Income, \$50B Tax (after expenses)
- GST \$46B, Excise \$26B, FBT \$3B, ... \approx \$350B total income

- Tax payer's charter:
Fair but firm; Assume honesty; Protect privacy
- Service standards — refunds within days ... hours ... seconds?
- Whilst protecting the integrity of the revenue collection

AUSTRALIAN TAXATION OFFICE

...continues this “fine tradition”

- Employs 22,000 staff Australia wide
- Role to collect revenue and process refunds

- 12M Individuals, \$600B Income, \$120B Tax
- 2M Companies..., \$2,200B Income, \$50B Tax (after expenses)
- GST \$46B, Excise \$26B, FBT \$3B, ... \approx \$350B total income

- Tax payer's charter:
Fair but firm; Assume honesty; Protect privacy
- Service standards — refunds within days ... hours ... seconds?
- Whilst protecting the integrity of the revenue collection

AUSTRALIAN TAXATION OFFICE

...continues this “fine tradition”

- Employs 22,000 staff Australia wide
- Role to collect revenue and process refunds

- 12M Individuals, \$600B Income, \$120B Tax
- 2M Companies..., \$2,200B Income, \$50B Tax (after expenses)
- GST \$46B, Excise \$26B, FBT \$3B, ... \approx \$350B total income

- Tax payer's charter:
Fair but firm; Assume honesty; Protect privacy
- Service standards — refunds within days ... hours ... seconds?
- Whilst protecting the integrity of the revenue collection

AUSTRALIAN TAXATION OFFICE

...continues this “fine tradition”

- Employs 22,000 staff Australia wide
- Role to collect revenue and process refunds

- 12M Individuals, \$600B Income, \$120B Tax
- 2M Companies..., \$2,200B Income, \$50B Tax (after expenses)
- GST \$46B, Excise \$26B, FBT \$3B, ... \approx \$350B total income

- Tax payer's charter:
Fair but firm; Assume honesty; Protect privacy
- Service standards — refunds within days ... hours ... seconds?
- **Whilst protecting the integrity of the revenue collection**

OVERVIEW

REVENUE COLLECTION

ESTABLISHING CAPABILITY

DELIVERING OUTCOMES

ONGOING CHALLENGES

Model Deployment

Big Data: Many Variables

CREATING AN ANALYTICS CAPABILITY

- Established as a corporate capability in 2004
- Strong support by visionary CEO and senior executives
- Core team of 15 data miners then and now 30 data miners
- Wider team of 150 data analysts throughout the organisation
- Shared technology throughout organisation
- Provide framework for whole of ATO risk management

Every tax return lodged in Australia today is risk assessed by at least one data mining model.

- Models delivering benefit:
 - Revenue assurance: impact in \$ millions
 - More efficient targeting: resource and tax payer annoyance
 - Better tax payer experience: briefer involvement

CREATING AN ANALYTICS CAPABILITY

- Established as a corporate capability in 2004
- Strong support by visionary CEO and senior executives
- Core team of 15 data miners then and now 30 data miners
- Wider team of 150 data analysts throughout the organisation
- Shared technology throughout organisation
- Provide framework for whole of ATO risk management

Every tax return lodged in Australia today is risk assessed by at least one data mining model.

- Models delivering benefit:
 - Revenue assurance: impact in \$ millions
 - More efficient targeting: resource and tax payer annoyance
 - Better tax payer experience: briefer involvement

CREATING AN ANALYTICS CAPABILITY

- Established as a corporate capability in 2004
- Strong support by visionary CEO and senior executives
- Core team of 15 data miners then and now 30 data miners
- Wider team of 150 data analysts throughout the organisation
- Shared technology throughout organisation
- Provide framework for whole of ATO risk management

Every tax return lodged in Australia today is risk assessed by at least one data mining model.

- Models delivering benefit:
 - Revenue assurance: impact in \$ millions
 - More efficient targeting: resource and tax payer annoyance
 - Better tax payer experience: briefer involvement

CHALLENGE 1: ANALYTICS AS IT

- Traditional IT: Buy software first then buy expertise
c.f. accounting software and accountants

- c.f. data miners as innovative programmers of data
- A culture of multiple tools and dynamic environments
- A culture of sharing algorithms and experiences
- Flexibility in using any technology as and when required
- Not constrained by traditional corporate IT SOE

LESSON: ANALYST FIRST

“Analytics is a non-repetitive, exploratory and creative process where the outcome is not known at the start, and only a fraction of efforts are expected to result in success.”

analystfirst.com

- Analytics is not IT and process.
- Focus is not the tools nor the algorithms!
- Focus needs to be with the skills of the people.
- The people who perform, manage, request and envision analytics.
- Required skill set: Programming, Statistics, Algorithms, Intuition

CHALLENGE 2: MISSING TECHNOLOGY ENSEMBLES

Ensemble concept developed in the 80's

- Multiple Decision Trees (PhD 1987)
- Remains one of the best off-the-shelf technologies: random forests and boosting.
- Not available in the closed source offerings until recently.
- Yet readily available in the open source community.



LESSON: COMMODITY PLATFORM

A data mining capability need not be expensive.

Build a network of workstations:

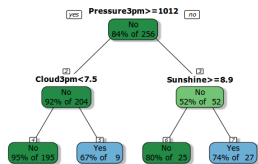
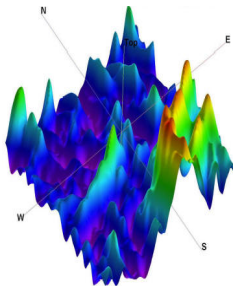
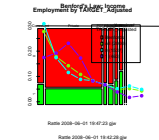
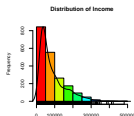
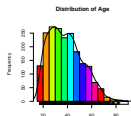
- 16 Cores, 64 bit, top CPU speeds
512GB RAM, 10TB Disk
- Flexible and open source OS
Ubuntu GNU/Linux
- Open Source data mining tools
R, **Rattle**, Weka + SAS, SPSS as reqd
- Open Source **does** deliver quality



Advert: Strive for open scientific transparency and repeatability – release an R package (or Weka) then it is also available in SAS, SPSS, Netezza, Micro Strategy. . .

CHALLENGE 3: PROGRAMMING DATA ANALYSIS

- Exploratory Data Analysis + Mining: R is second to none
- Analytics is about programming with data, and today's culture is all about GUI - simplified interfaces
- Skills Shortage → Train 150 Data Analysts
- Rattle developed as a stepping stone to transparent and repeatable analytics in R.



GRATUITOUS RATTLE SCREENSHOT

R Data Miner - [Rattle]
Project Tools Settings Help Rattle Version 2.6.20 togaware.com

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

No Target Algorithm: Traditional Conditional Model Builder: rpart

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

Min Bucket: 7 Complexity: 0.0100 Loss Matrix:

Decision Tree Model

A decision tree model is one of the most common data mining models. It is popular because the resulting model is easy to understand. The algorithms use a recursive partitioning approach.

The traditional algorithm is implemented in the rpart package. It is comparable to CART and ID3/C4.

The conditional tree algorithm is implemented in the party package. It builds trees in a conditional inference framework.

Note that the ensemble approaches (boosting and random forests) tend to produce models that exhibit less bias and variance than a single

OVERVIEW

REVENUE COLLECTION

ESTABLISHING CAPABILITY

DELIVERING OUTCOMES

ONGOING CHALLENGES

Model Deployment

Big Data: Many Variables

ANALYTICS IN ACTION

High Risk Refunds (HRR) identified prior to issuing of refunds.

- Traditional rules identify too many “high risk” refunds.
- Some tests might identify 100,000 cases each year.
- Sometimes as few as 5% are found to require adjustment.
- Revenue at risk can be very significant (from \$10m to \$1b).

Data Mining modelling for HRR.

- Has identified numerous characteristics to better target risk (5%)
- More effectively deploy resources on productive cases.
- Avoid non-productive audits.
- Uses decision trees and ensembles (random forests).

ANALYTICS IN ACTION

High Risk Refunds (HRR) identified prior to issuing of refunds.

- Traditional rules identify too many “high risk” refunds.
- Some tests might identify 100,000 cases each year.
- Sometimes as few as 5% are found to require adjustment.
- Revenue at risk can be very significant (from \$10m to \$1b).

Data Mining modelling for HRR.

- Has identified numerous characteristics to better target risk (5%)
- More effectively deploy resources on productive cases.
- Avoid non-productive audits.
- Uses decision trees and ensembles (random forests).

OTHER SUCCESSFUL PROJECTS

- High Risk Refunds – ensembles of trees
- Required to Lodge (\$110M) – ensembles of trees
- Assessing Levels of Debt – Propensity/Capacity to Pay
- Optimal Treatment Strategies
- Identity Theft – Outliers, Unusual, Out of Pattern
- International and Tax Havens – Text Mining
- Complex Structures – Network and Link Analysis

OVERVIEW

REVENUE COLLECTION

ESTABLISHING CAPABILITY

DELIVERING OUTCOMES

ONGOING CHALLENGES

Model Deployment

Big Data: Many Variables

CHALLENGE 4: MODEL DEPLOYMENT IS IT

- Deliver PMML from SAS into Teradata open standards and interoperability (KDD98)
- Much data mining is **not** “deployed” models are run ad-hoc as required.
- Models “deployed” by conversion
 - SQL: 2 million lines for RF — 20x200x500
 - C: Netezza to score 15M entities in 90 seconds
 - PMML: deployment in e.g. Zementis’ Adapa
- Challenge is that many (hundreds of) models have been developed and many need to be “deployed”

MODEL MANAGER - MODELS IN PRODUCTION

- Running automatically on demand and event driven.
- Requires a professional production environment:
 - models now have life cycles: traditional dev/test/deploy
 - monitoring of model performance:
 - alerts for out-of-spec models
 - daily dashboards
 - 24/7 support of models
- SAS and SPSS provide model management solutions
- Our developing solution using open source tools:
 - R, Python, Shell, Make (scripting)
 - Bazaar (dev/prod version control)
 - Jenkins (continuous integration framework web interface)
 - ...

CHALLENGE 5: BIG DATA — MANY VARIABLES

- Not really too “big” in the ATO
 - 100M transactions
 - 20M entities
 - 5K variables
- Big enough to present some challenges to traditional tools

- The variables are the “big” issue
- Random Forests are good, but could be better!
- Issues when there are many irrelevant variables.

CHALLENGE 5: BIG DATA — MANY VARIABLES

- Not really too “big” in the ATO
 - 100M transactions
 - 20M entities
 - 5K variables
- Big enough to present some challenges to traditional tools

- The variables are the “big” issue
- Random Forests are good, but could be better!
- Issues when there are many irrelevant variables.

SUBSPACED RANDOM FORESTS

Research with Chinese Academy of Sciences
Shenzhen Institutes of Advanced Technology.

- Random forests are a popular classification method building an ensemble of a single type of decision tree.
- Algorithmically intuitive and simple.
- Aim is for an ensemble of very different decision trees, with each decision tree by itself being a good model of the data it is based on.
- How to increase diversity and individual accuracy?

SUBSPACED RANDOM FORESTS

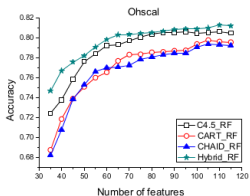
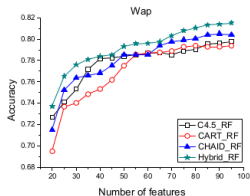
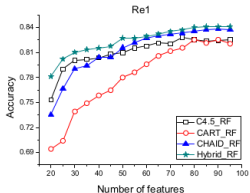
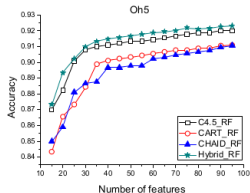
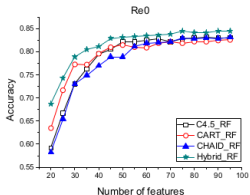
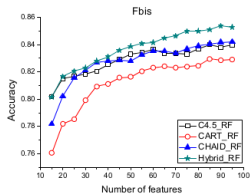
Research with Chinese Academy of Sciences
Shenzhen Institutes of Advanced Technology.

- Random forests are a popular classification method building an ensemble of a single type of decision tree.
- Algorithmically intuitive and simple.
- Aim is for an ensemble of very different decision trees, with each decision tree by itself being a good model of the data it is based on.
- How to increase diversity and individual accuracy?

HYBRID RANDOM FORESTS

- Simple Idea: *Build different types of decision trees—different algorithms—for each randomly sampled training dataset and choose the best tree model for each.*
- Tree types considered: C4.5, CART, CHAID
- Further opportunities:
 - HD, EDD, VED (Decision Master)
 - ctree (R/party)
- A series of experiments suggest that the hybrid approach **always** delivers the best model.

EXPERIMENTAL RESULTS



The hybrid random forest always performs best for this collection of datasets!!!

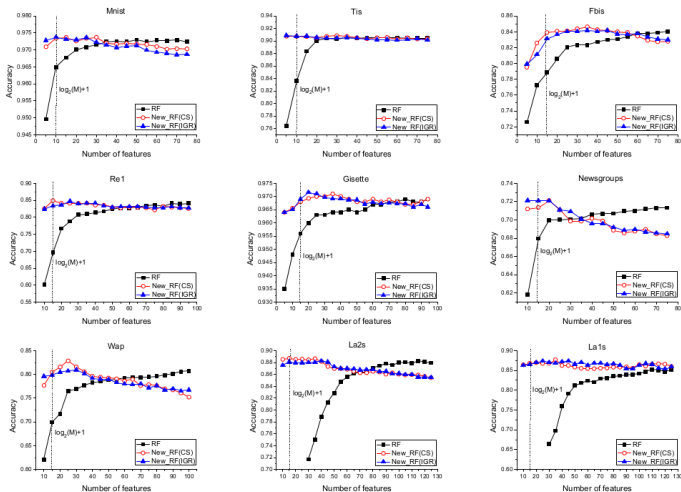
ANOTHER APPROACH: WEIGHTED SUBSPACES

- Performance of a random forest is improved by
 - **Strengthening** each tree
 - Reducing **correlation** between each tree

- Problem of large number of variables:
 - Random selection means too many irrelevant variables

- Introduce the concept of weighted subspace random forests
 - Bias the selection of variables toward most important variables

EXPERIMENTAL RESULTS



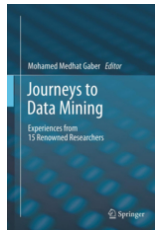
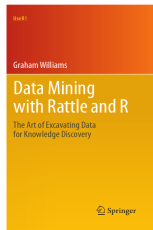
For this collection of experiments the weighted subspace random forest always performs better with many fewer features.

SUMMARY

- Focus on the “Analyst First”
- Commodity hardware and software provide excellent capability
- Sharing algorithms through R packages (rattle, wsrpart, wsrfr)
- Deployment continues to challenge

RESOURCES AND REFERENCES

- Rattle: rattle.togaware.com
- Guides: datamining.togaware.com
- Practise: analystfirst.com
- Book: Data Mining using Rattle/R
- Chapter: Rattle and Other Tales



Thank You