

[www.ala.org.au](http://www.ala.org.au)

# Data Quality

## ALA-ERIN Collaboration

**Simon Bennett**

ALA / ERIN DSEWPAC, 14 Oct 2011



**Australian Government**  
**Department of the Environment,  
Water, Heritage and the Arts**

The Atlas is funded by the Australian Government under the National Collaborative Research Infrastructure Strategy and further supported by the Super Science Initiative of the Education Investment Fund

- ALA Data Quality Wiki
- Spatial Precision Issues
- Spatial Outlier Detection -
  - Jackknife
  - Convex Hull Trimming
  - Alpha Hulls
  - Environmental Surface Issues
  - Maxent Pair wise Analysis
- Python Tool

# ALA Data Quality Wiki



ala-dataquality - Atlas of Living Australia - Data Quality - Google Project Hosting - Windows Internet Explorer

http://code.google.com/p/ala-dataquality/

File Edit View Favorites Tools Help

ausavim@gmail.com | My favorites | Profile | Sign out

ala-dataquality  
Atlas of Living Australia - Data Quality

Project Home Wiki Administer

Summary Updates People

Tip: Project owners, see our [Getting Started](#) guide for steps to configure your project.

## Project Information

Stared by 1 user  
Activity ▲▲ High  
[Project feeds](#)

Code license  
[Mozilla Public License 1.1](#)

Labels  
ALA, biodiversity, Australia, TWDG

Members  
[ausavim](#), [brvnc.ala](#),  
[milo.nicholls@hotmail.com](#),  
[movesyside](#), [leebelein](#), [ghobern](#)  
14 committers

Your role  
Owner

## Atlas of Living Australia - Data Quality Portal

Welcome to the Atlas of Living Australia Data quality checks and issues pages.

- The [Data Quality Checks Spreadsheet](#) is a summary of the data quality checks implemented and planned for the Atlas of Living Australia.
- We have also set up a [Wiki](#) page for each check to document more detail where needed and for you to provide comments on each check.
- Please post your comments of a broader nature, along with any suggestion for content inclusion to the [general](#) comments page.

Links to key ALA data quality documents and resources on the internet are provided below.

### Key ALA Data Quality Documents

- [Data Quality Checks Spreadsheet](#)
- [Data Quality Wiki](#)
- [Completeness model](#) (previously named 'quality model')
- [Darwin Core Key Spatial Concepts](#)
- [ALA continues to focus on data quality](#)

### Key Data Quality References and Tools

- [Darwin Core: Quick Reference Guide](#)
- [Tann, J. & P. Flemons 2008. ALA Review of Online and Desktop Tools.](#)
- [Chapman, A. D. 2005. Principles of Data Quality.](#)
- [Chapman, A. D. 2005. Principles and Methods of Data Cleaning – Primary Species and Species-Occurrence Data](#)
- [Chapman, A. D. 2005. Uses of Primary Species-Occurrence Data](#)
- [Diva GIS](#)
- [Georeferencing Calculator: Online Tool](#)
- [Georeferencing Calculator: Online manual](#)
- [Georeferencing Calculator: MaNIS/HerpNet/ORNIS Georeferencing Guidelines](#)
- [Georeferencing Calculator: Georef-calculator on google code. See Wiki for User Manual](#)
- [BioGeomancer WorkBench](#)
- [BioGeomancer-core - Google code](#)
- [speciesLink Data Cleaning](#)
- [speciesLink spOutlier](#)
- [Maps and the GDA: Geocentric Datum of Australia: Info Sheet](#)

## Data Quality One Stop Shop

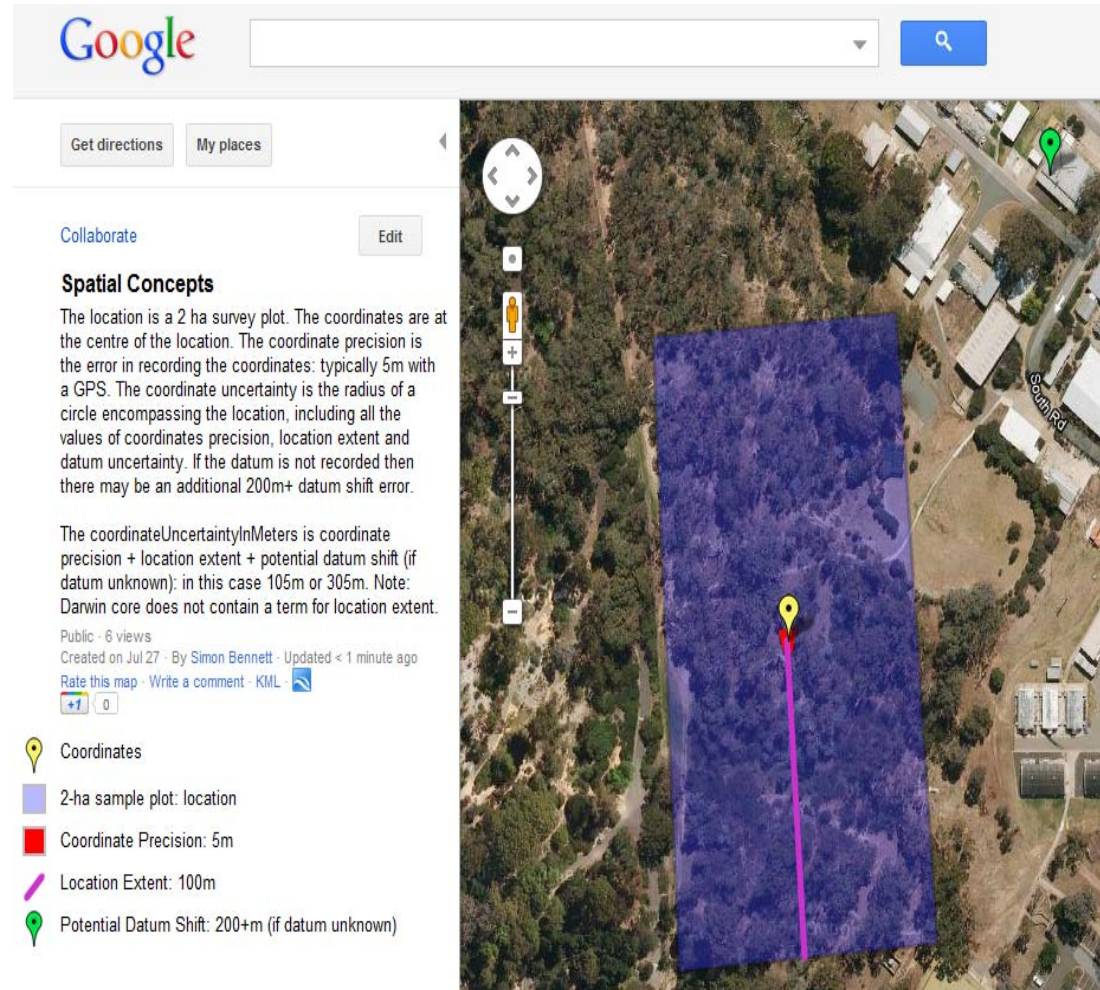
- Links to key documents
- List of data quality checks
- Methodologies
- Discussion

<http://code.google.com/p/ala-dataquality/>

©2011 Google - [Terms](#) - [Privacy](#) - [Project Hosting Help](#)  
Powered by [Google Project Hosting](#)

# Spatial Precision Issues

- Spatial precision and uncertainty, extent – What does it all mean?
- coordinateUncertaintyInMeters is overall error not just extent of location.
- Darwin core has not have an extent of location term – even though recommended georeferencing tools use it.
- All is now clear – but re-education needed.



Google

Get directions My places

Collaborate Edit

### Spatial Concepts

The location is a 2 ha survey plot. The coordinates are at the centre of the location. The coordinate precision is the error in recording the coordinates: typically 5m with a GPS. The coordinate uncertainty is the radius of a circle encompassing the location, including all the values of coordinates precision, location extent and datum uncertainty. If the datum is not recorded then there may be an additional 200m+ datum shift error.

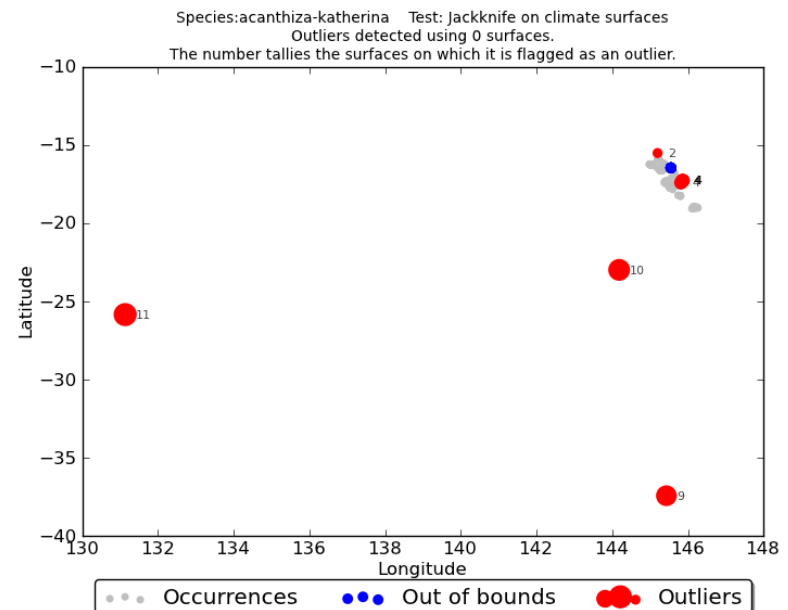
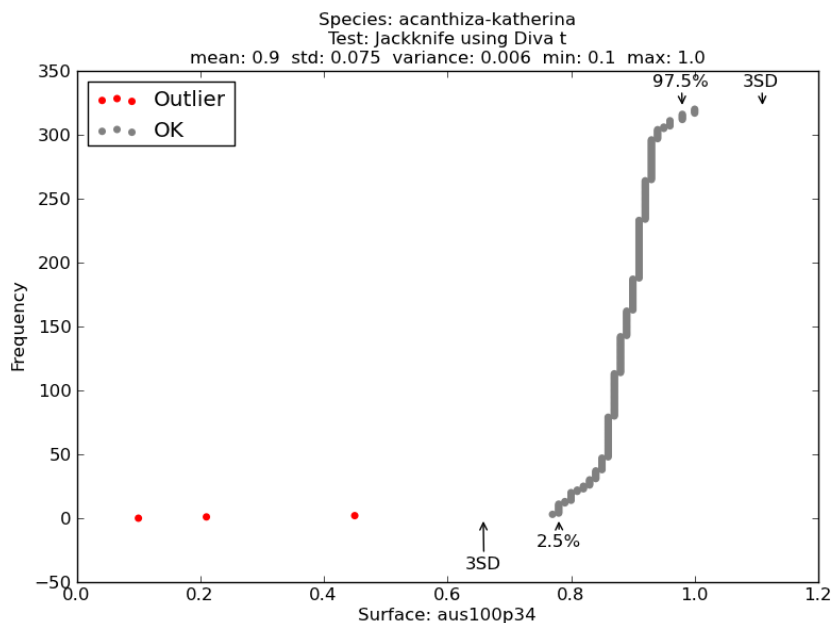
The coordinateUncertaintyInMeters is coordinate precision + location extent + potential datum shift (if datum unknown): in this case 105m or 305m. Note: Darwin core does not contain a term for location extent.

Public - 6 views  
Created on Jul 27 · By Simon Bennett · Updated < 1 minute ago  
[Rate this map](#) · [Write a comment](#) · [KML](#)

- Coordinates
- 2-ha sample plot: location
- Coordinate Precision: 5m
- Location Extent: 100m
- Potential Datum Shift: 200+m (if datum unknown)

# Outlier - Jackknife

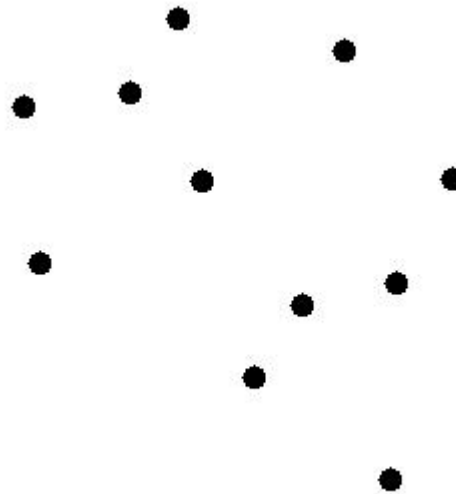
- [http://code.google.com/p/ala-dataquality/wiki/DETECTED\\_OUTLIER\\_JACKKNIFE](http://code.google.com/p/ala-dataquality/wiki/DETECTED_OUTLIER_JACKKNIFE)



- Presently being implemented as an ALA Data quality check
- Prototyped in Python Data Quality Workbench – provided to ERIN

# Outlier – Convex Hull

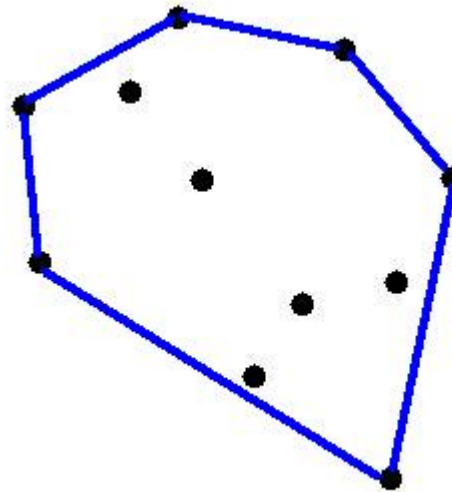
- With Ross Cunningham, ANU Fenner School



Convex hull peeling

*From: Greg Aloupis - Université Libre de Bruxelles*

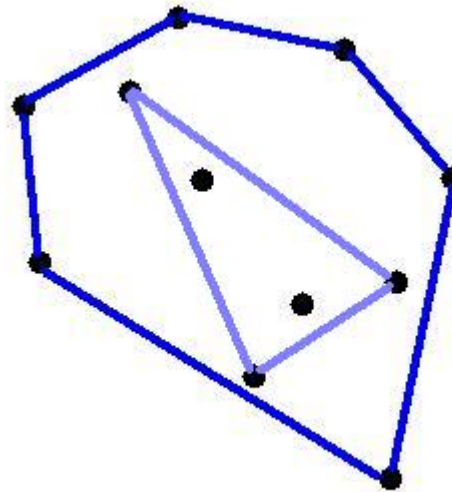
# Outlier – Convex Hull



Convex hull peeling

Outliers tend to occur on first few hulls

# Outlier – Convex Hull

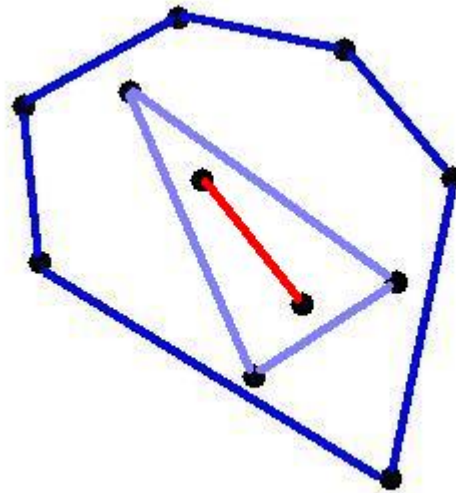


Convex hull peeling

Hull with 50% points inside is inter-quartile range



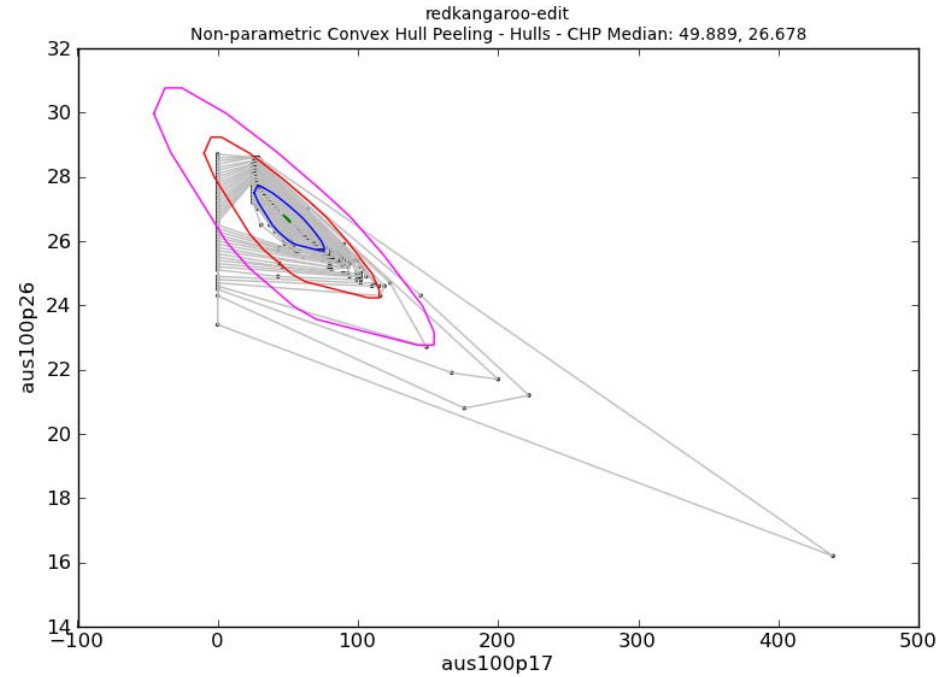
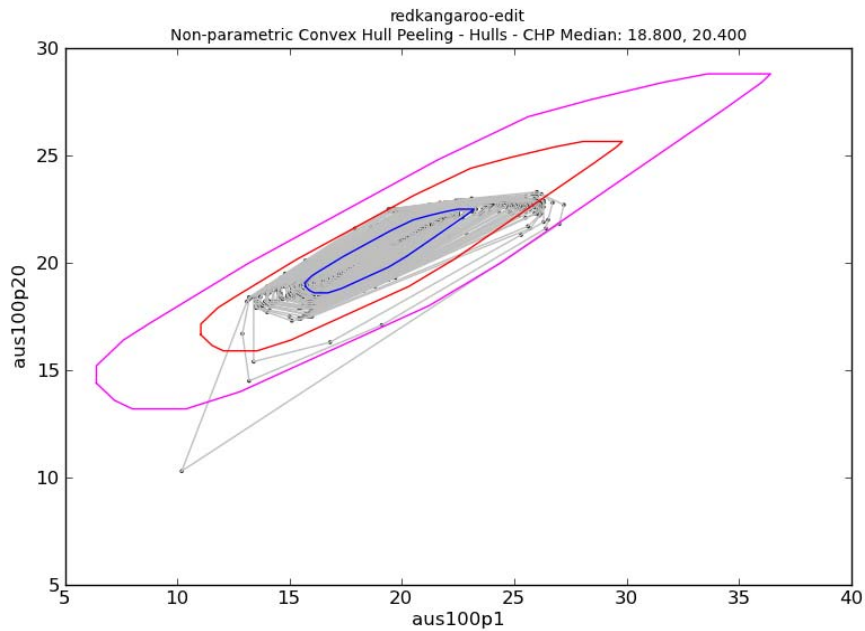
# Outlier - Convex hulls

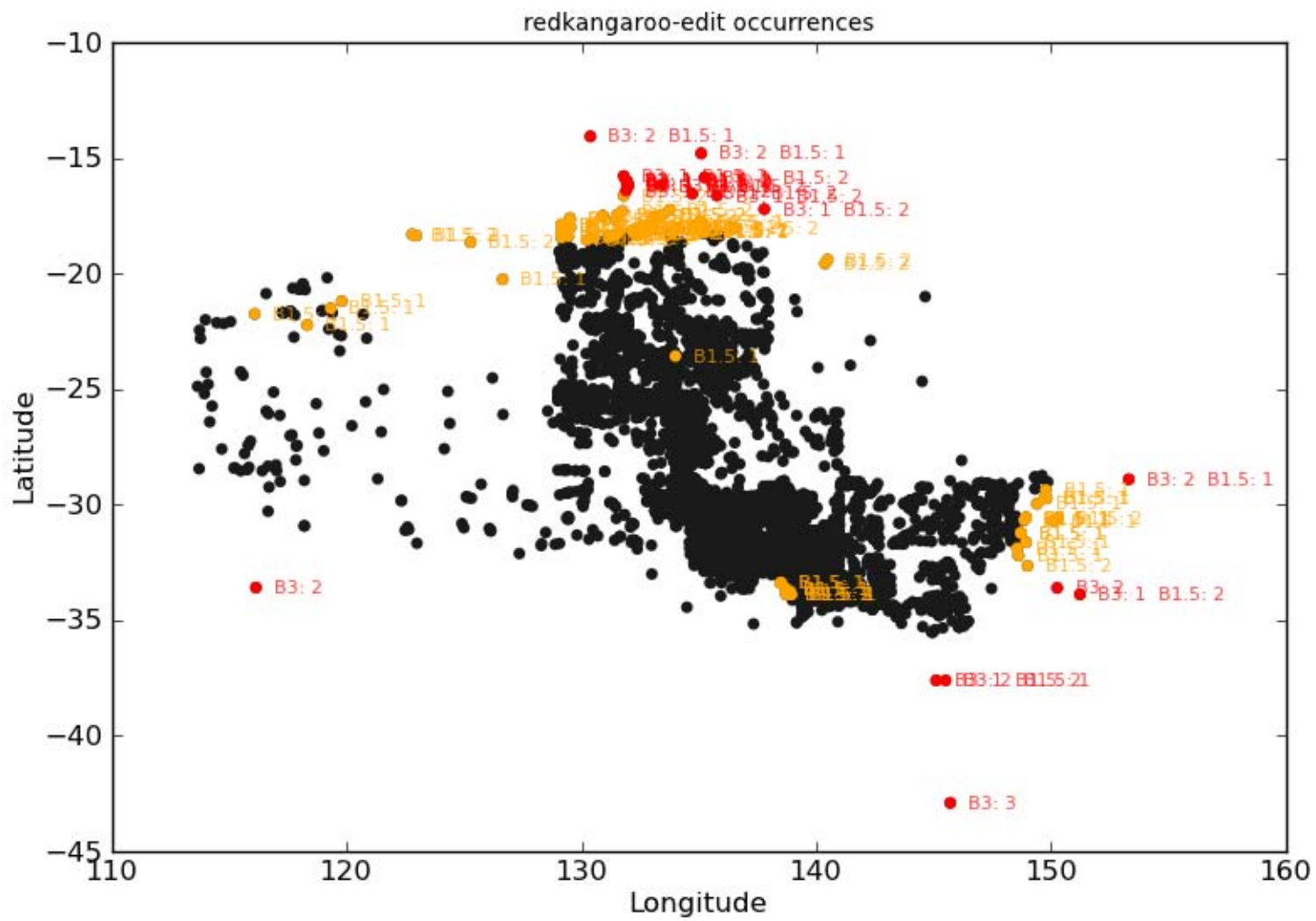


Convex hull peeling

Inner hull is median

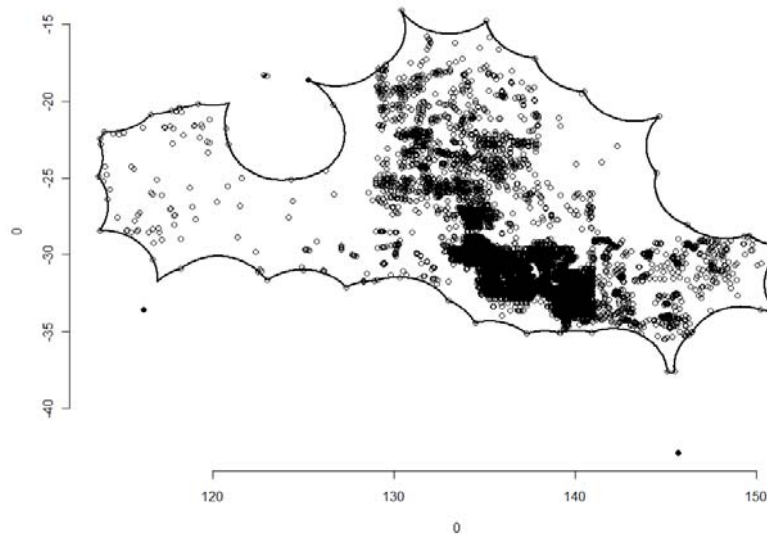
# Outlier – Convex Hull



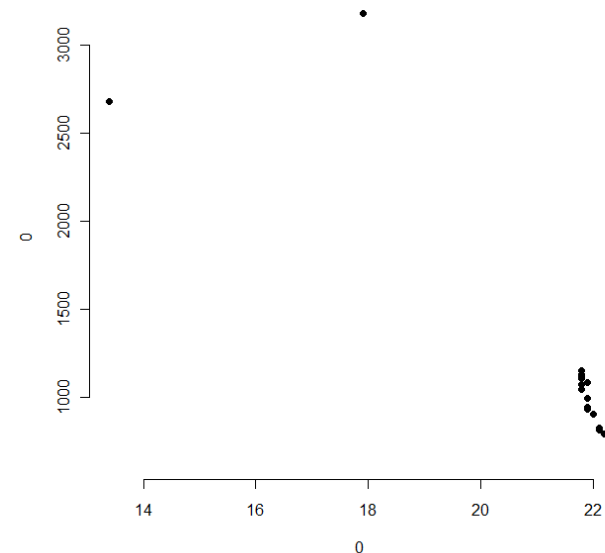


# Outlier – Alpha Hull

- From Burgman and Fox 2003
- Being used by Birds Australia for mass mapping exercise.



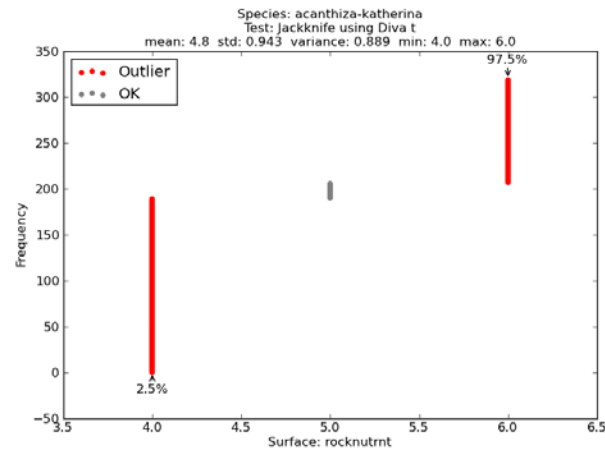
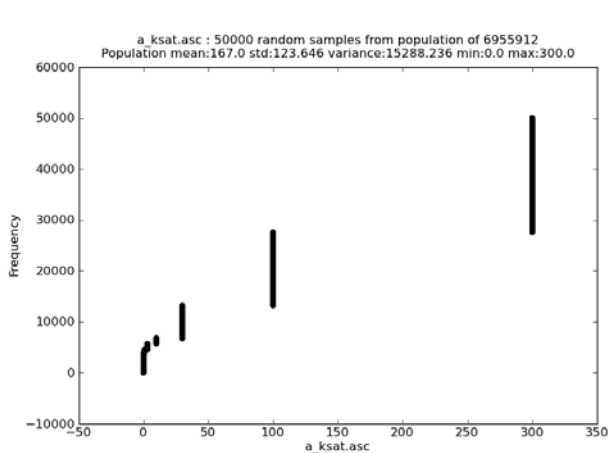
Red Kangaroo: lat/long



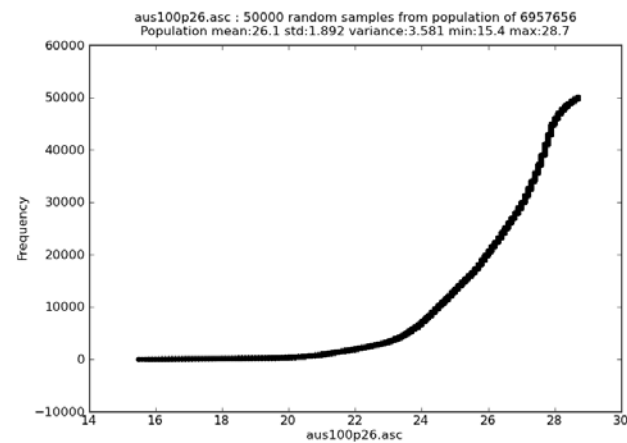
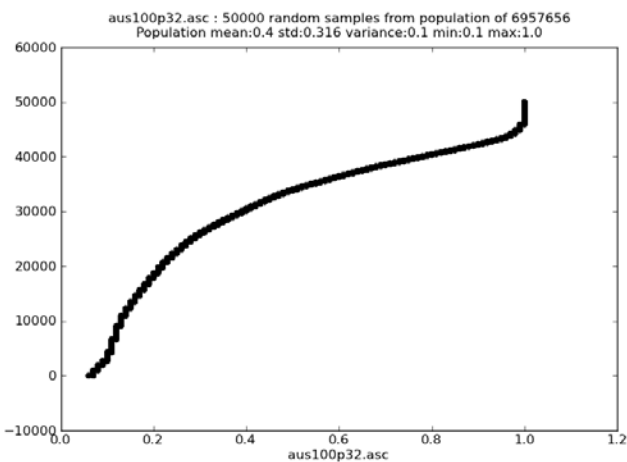
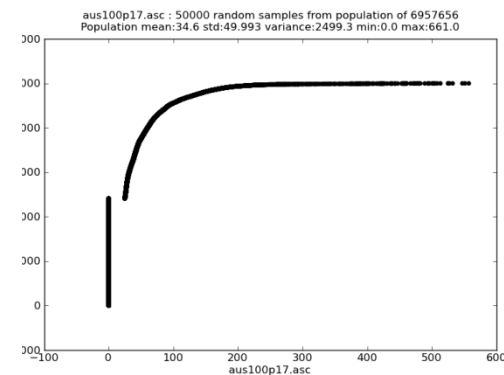
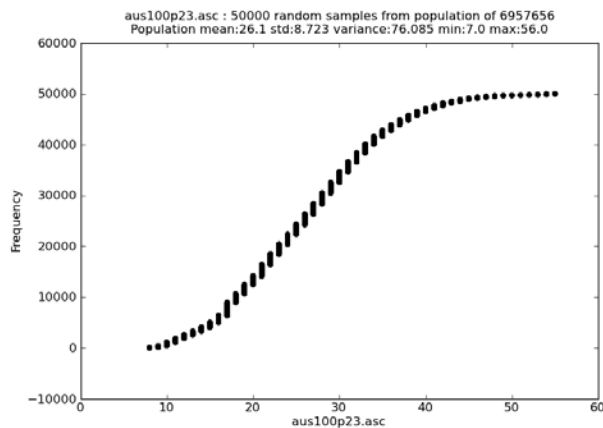
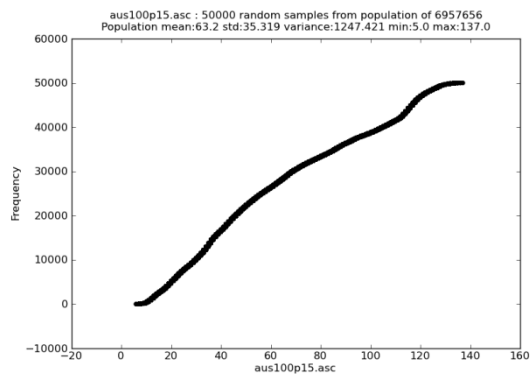
Black Grasswren: p12/20 – falls over

# Environmental Surface Issues

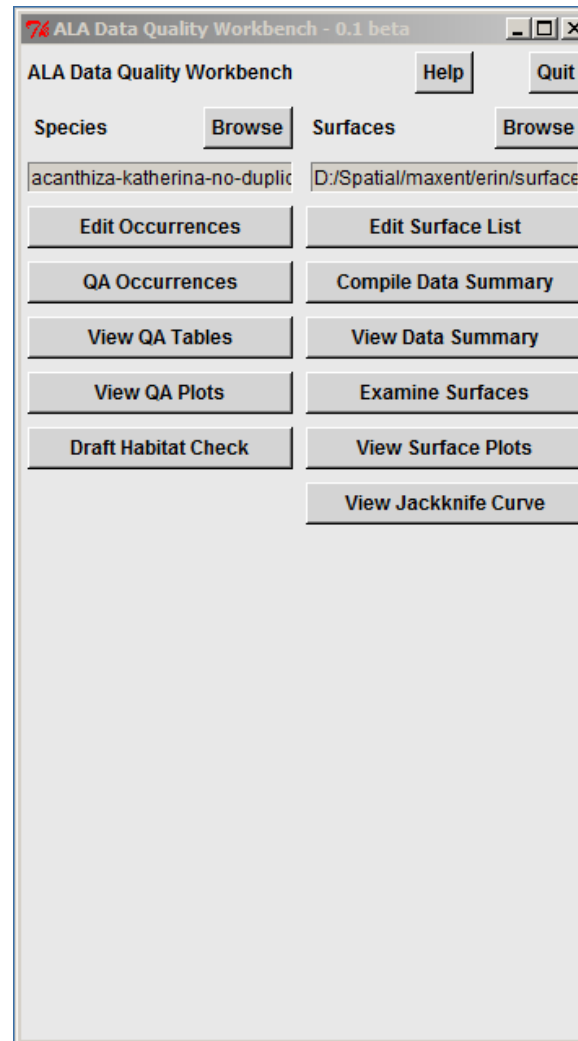
- Many 100's of environmental surface
- The various outlier checks exposing issue surfaces
- Will need an optimal set for each check



# Lee's minimal set



# Python Tool



# Outlier – Maxent Pair wise Analysis



With Simon Barry and Warren Muller, CSIRO Mathematics, Informatics and Statistics

Jack-knife occurrence records through Maxent:

- a) for all data
- b) removing each record in turn,
- c) then removing each pair of records in turns

Obtain Maxent prediction for each occurrence for every combination

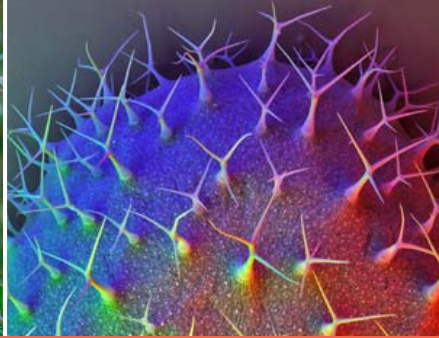
Mountain Thornbill – 189 records became 15000 Maxent runs!!

Output presently being analysed by CSIRO.



# Next Steps

- Complete documentation of precision/uncertainty issues
- Propose extent of location be added as Darwin Core field
- Write up progress with outlier detection work and refine methodology e.g. remove outliers to minimise hull area
- Determine optimal environmental surfaces to use with each test
- Continue work with Data Quality Wiki
-



# The Atlas of Living Australia Participants

[www.ala.org.au](http://www.ala.org.au)



Tasmanian Museum & Art Gallery



Australian Government  
Department of the Environment,  
Water, Heritage and the Arts



Australian Government  
Department of Agriculture,  
Fisheries and Forestry



THE UNIVERSITY  
OF ADELAIDE  
AUSTRALIA



The Council of Heads of Australian  
Faunal Collections (CHAFC)  
The Council of Heads of Australian  
Entomological Collections (CHAEC)

The Australian Microbial Resources  
Research Network (AMRRN)  
The Council of Heads of Australasian  
Museum Directors (CAMD)



**An Australian Government Initiative**  
**National Collaborative Research**  
**Infrastructure Strategy**



ATLAS OF **LIVING**  
**AUSTRALIA**  
sharing biodiversity knowledge

The Atlas is funded by the  
Australian Government under the  
National Collaborative Research  
Infrastructure Strategy  
and further supported by the  
Super Science Initiative of the  
Education Investment Fund