

TG3: Use Case Library

Data Quality

Miles Nicholls

Data Manager, Atlas of Living Australia

miles.nicholls@csiro.au

Context and purpose

- The third task group of the GBIF TDWG biodiversity data quality interest group
- Tasked to assemble a library of use cases describing specific examples of data selection

Anticipated Process

1. Use the framework from TG₁ to understand the information that needs to be captured in a data quality use case
2. Create a model to store use case descriptions
3. Capture use cases
4. Analyse the use case elements
5. Cross reference the use cases against the tools from TG₂
6. Identify areas for additional quality tools

Actual process

- The first few steps were relatively straight forward:
- Running through a few examples helped in understanding the practical application of the framework and led to the development of a use case description store – a series of related worksheets to capture a use case:
 - Sheets that collect the elements from the framework: Use Case Title, Information Elements, Data Quality Dimensions, Data Quality Criteria and Data Quality Enhancements
- Allan Koch Veiga added a DQ profile tool which produces a report based on the information in the worksheets

Data Quality Use Cases - Copy

File Edit View Insert Format Data Tools Add-ons Help DQ Profile Last edit was made on February 4 by anonymous

ID	Use Case	Group	Usage examples	Owner
1	ecological niche modelling	distribution modelling	Modelling using MaxEnt, GARP, ANN...	example
2	ecological gap analysis	distribution modelling		
3	evaluation of the diversity of species of a country	species list		
4	update status of the name of taxa	taxonomy		
5	update status of the name of taxa	taxonomy		
6	list of species in an area e.g. to validate a sighting or provide a list	species list		
7	list a species for conservation sensitivity	species status		
8	list a species as sensitive	species status		
9	look at where a pest could spread to	distribution modelling		
10	where could something be replanted	distribution modelling		
11	what to replant here	distribution modelling		
12	where to release a captive raised species	distribution modelling		
13	model of species in an area	distribution modelling		
14	whether a pest is present in a country	species list		
15	associated species	species characters		
16	migration patterns	distribution modelling		
17	provide a list of likely species in an area	distribution modelling		
18	identify a species	species identification		
19	identify population trends (rising or falling)	population modelling		
20	assess whether a development should take place at a location	species list		
21	identify the impact of clearing a region of vegetation			Australian Department of the Environment
22	Australian Government Threatened species distribution modelling			Australian Department of the Environment
23	Australian Government Threatened ecological community mapping			Australian Department of the Environment
24	Australian Government Exotic			Australian Department of the Environment
25	Australian Government Name not matched			Australian Department of the Environment

Profile: ecological niche modelling

Valuable Information Elements:

- Decimal coordinates
- Location
- Event date
- Event time
- Taxonomic information
- Occurrence
- Country
- Basis of record

DQ Validation Policy:

- Numerical coordinates precision must be higher than 0.001
- Stated coordinates precision must be higher than 0.001
- Coordinates must be completed
- Basis of Record must be well formed
- Coordinates must be 100% complete
- Event date precision must be at level month or higher.

Data Quality Use Cases - Copy

File Edit View Insert Format Data Tools Add-ons Help DQ Profile Last edit was made on February 4 by anonymous

ID	Use Case	Name	Composed of	Type	Description
1	ie:coordinates	Decimal coordinates	dwc:decimalLatitude + dwc:decimalLongitude	Composed	
2	ie:georeferenceProtocol	Georeference protocol	dwc:georeferenceProtocol	Single	http://rs.tdwg.org/dwc/terms/index.htm#georeferenceProto
3	ie:location	Location	ie:coordinates + dwc:locality + dwc:municipality...	Composed	
4	ie:eventDate	Event date	dwc:eventDate	Single	http://rs.tdwg.org/dwc/terms/index.htm#eventDate
5	ie:evenTime	Event time	dwc:eventTime	Single	http://rs.tdwg.org/dwc/terms/index.htm#eventTime
6	ie:eventMoment	Event moment	dwc:eventDate + dwc:eventTime	Composed	
7	ie:taxonomicInfo	Taxonomic information	dwc:scientificName + dwc:taxonRank + taxonHigher...	Composed	
8	ie:taxonName	Taxon name	dwc:scientificName + dwc:taxonRank	Composed	
9	ie:occurrence	Occurrence	ie:taxonInfo + ie:eventMoment + ie:location	Composed	
10	ie:institution	Institution	dwc:institutionCode + dwc:institutionID	Composed	
11	ie:country	Country	dwc:country + dwc:countryCode	Composed	
12	ie:vernacularName	Vernacular name	dwc:vernacularName	Single	http://rs.tdwg.org/dwc/terms/index.htm#vernacularName
13	ie:basisOfRecord	Basis of record	dwc:basisOfRecord	Single	http://rs.tdwg.org/dwc/terms/index.htm#basisOfRecord
14	ie:geodeticDatum	Geodetic datum	dwc:geodeticDatum	Single	http://rs.tdwg.org/dwc/terms/index.htm#geodeticDatum
15	ie:scientificNameAuthorship	Scientific name authorship	dwc:scientificNameAuthorship	Single	http://rs.tdwg.org/dwc/terms/index.htm#scientificNameAut
16	ie:namePublishedInYear	Name published in year	dwc:namePublishedInYear	Single	http://rs.tdwg.org/dwc/terms/index.htm#namePublishedInY

Data Quality Use Cases - Copy

File Edit View Insert Format Data Tools Add-ons Help DQ Profile Last edit was made on February 4 by anonymous

ID	Measure the DQ in the Use Case	Contextualized Dimension Name	Fundamental Dimension	Information Element target Resource Type target	Measure description	
1	d.coordinatesPrecision_0	Coordinates precision of single records	Precision	ie:coordinates	Single Record	A decimal representation of the p...
2	gap.niche	Stated coordinates precision of single records	Precision	ie:coordinates	Single Record	Assume the value of "dwc: coordin...
3	d.coordinatesPrecision_1	Coordinates completeness of datasets	Completeness	ie:coordinates	Dataset	Proportion of records with supplie...
4	d.coordinatesCompleteness_ds	Coordinates completeness of single records	Completeness	ie:coordinates	Single Record	Is "complete" if both latitude and l...
5	d.coordinatesCompleteness_sr	Coordinate certainty of single records	Precision/Accuracy/Certain	ie:coordinates	Single Record	Assume the value of "dwc: coordin...
6	dimension4	Georeference protocol completeness of single reco	Completeness	ie:georeferenceProtocol	Single Record	The quality is better as nearer of...
7	dimension5	Event date precision of single records	Precision	ie:eventDate	Single Record	Is "complete" if georeference prot...
8	dimension6	Event date completeness of single records	Completeness	ie:eventDate	Single Record	Precision is 0 if any value is supp...
9	dimension7	Event time precision of single records	Precision	ie:eventTime	Single Record	2 if year and month are supplied, ...
10	dimension8	Event time completeness of single records	Completeness	ie:eventTime	Single Record	
11	dimension9	Moment precision of single records	Precision	ie:precision	Single Record	
12	dimension10	Taxonomic information precision of single records	Precision	ie:taxonomicInfo	Single Record	
13	dimension11	Occurrence completeness of single records	Completeness	ie:occurrence	Single Record	
14	dimension12	Occurrence completeness of datasets	Completeness	ie:occurrence	Dataset	Measure how comprehensive is th...
15	dimension13	Occurrence uniqueness of datasets	Uniqueness	ie:occurrence	Dataset	Count any duplicates
16	dimension14	Force of evidence of single records	Credibility	ie:occurrence	Single Record	Measure how much the occuranc...
17	dimension15	Completeness of institution code of single records	Completeness	ie:institution	Single Record	
18	dimension16	Completeness of country of single records	Completeness	ie:country	Single Record	
19	dimension17	Completeness of common name of single records	Completeness	ie:vernacularName	Single Record	
20	dimension18	Completeness of geodetic datum of single records	Completeness	ie:geodeticDatum	Single Record	
21	dimension19	Completeness of scientific name of single records	Completeness	ie:scientificName	Single Record	
22	dimension20	Completeness of scientific name authorship of sing	Completeness	ie:scientificNameAuthorship	Single Record	
23	dimension21	Completeness of name published in year of single	Completeness	ie:namePublishedInYear	Single Record	
24	dimension22					

Data Quality Use Cases - Copy

File Edit View Insert Format Data Tools Add-ons Help DQ Profile Last edit was made on February 4 by anonymous

ID	Validate the DQ in the Use Case	Contextualized Criterion Statement	Fundamental Criterion	Information Element target Resource Type target
1	niche	Numerical coordinates precision must be higher than 0.001	DQ measure must be in the range	ie:coordinates
2	c.thresholdCoordinatesPrecision_1	Stated coordinates precision must be higher than 0.001	DQ measure must be in the range	ie:coordinates
3	c.thresholdCoordinatesPrecision_2	Coordinates must be completed	DQ measure must be in the range	ie:coordinates
4	c.thresholdCoordinatesMustBeComp	Coordinates must be completed	DQ measure must be in the range	ie:coordinates
5	c.basisOfRecordWellFormed	Basis of Record must be well formed	Data must be according to some standard	ie:basisOfRecord
6	gap	Numerical coordinates precision must be higher than 0.01	DQ measure must be in the range	ie:coordinates
7	c.thresholdCoordinatesPrecision_3	Coordinates must be 100% complete	DQ measure must be in the range	ie:coordinates
8	c.threshold100Complete	Coordinates must be 100% complete	DQ measure must be in the range	ie:coordinates
9	c.thresholdDatePrecision	Event date precision must be at level month or higher.	DQ measure must be in the range	ie:eventDate

Data Quality Use Cases - Copy

File Edit View Insert Format Data Tools Add-ons Help DQ Profile Last edit was made on February 4 by anonymous

ID	Improve the DQ in the Use Case	Contextualized Enhancement Statement	Fundamental Enhancement	Information Element target Resource Type target	Improvement description	
1	ie:recommendName	Recommend closest valid scientific name	Recommend name based on similarity	ie:scientificName	Single Record	Recommend the most similar
2	niche.gap					
3						
4						
5						
6						
7						
8						
9						
10						
11						
12						
13						
14						
15						
16						
17						
18						
19						
20						
21						
22						
23						
24						
25						
26						
27						
28						
29						
30						
31						
32						
33						
34						
35						

Actual process

- Review of the initial worksheets indicated that although it collected the required information the level of technical understanding required to fill in the series of related sheets was likely too high for it to be widely used.
- That the worksheets were online and didn't export well was also a barrier to data entry
- So an additional offline sheet was developed focussed on capturing a single use case as simply as possible

	A	B	C	D	E	F	G
1	Use Case name						
2	Owner name						
3	Contact email						
4							
5	Information/Data element/Field name	Description/Link to standard	Criteria	Relates to single record or data set	Criteria description/remarks	Enhancement	Example implementation
6	The name of the field or fields that need to meet a particular quality metric for the data to be fit for this use case	A description of the field(s) or link to the standard where one exists	The criteria the field needs to meet	Is the criteria relevant to a single record or an entire data set	Additional information about the criteria	Possible processing to improve the quality	Example of a tool or system that provides this check, if known
7							
8							
9							
10							
11							
12							
13							
14							
15							
16							
17							
18							

- Once the template was filled in it would be converted to the more complex data structure manually

Actual Process

- The template was then posted up on the GBIF community site with some help material, a worked example and a request for contributions
- ...
- A further post was added to indicate that an email with a text description would be enough to start the process
- ...
- Unfortunately there has been very little participation

Key challenge so far

- How to address the limited number of contributions:
 - Is the community aware of the working groups and what we're doing?
 - More and varied communication channels
 - More frequent updates on progress
 - More communication on the purpose and outcomes of the work to encourage participation
 - Are the proposed tools too difficult to use?
 - Needs feedback on the usability of the framework and worksheets from the community using them

What's next

- Continue to collect use cases
 - Improve communication
 - Publicise contributions
- Transfer use cases to the use case store
- Improve use case collection tools based on feedback

The future

- Encourage use of the use case library as a resource to research how data is selected for particular purposes
- Begin to build profiles of the most used information elements, dimensions and criteria
 - Feed these back to tools and initiatives that collect data
- Feedback use cases into the development of data quality assessment tools
- Begin to develop automated data selection for use cases

Thank You

Do you have a data quality use case?

[Task Group 3 - Use Case Library for Data Quality](#)