



Datamart  
by Information Sciences Institute

The ISI Datamart project is building technology to create the largest publicly available knowledge graph to power data-driven models in a wide variety of domains. At the core of Datamart is Wikidata, a publicly available knowledge graph that already contains over 93 million entities. Datamart will enable communities of interest to build satellite knowledge graphs that contain detailed knowledge in domains of interest. The enabler technologies includes an architecture for combining public data with private data kept within an organization's firewall, and tools to automate the ETL process required to syntactically and semantically align the data in millions of spreadsheets and CSV files to the Wikidata semantic representation.

The ISI Datamart project webpage is here <https://usc-isi-i2.github.io/datamart/>.

## Semi-Automated Data Curation

Datamart provides a semi-automated tool T2WML (Table to Wikidata Mapping Language) to transform structured data to canonical tabular format.

- T2WML automatically detects syntactic and semantic regions of spreadsheets and CSV files. The regions include main subject, property, value, unit and time. See the color coded cells in the middle pane below. The user can adjust the regions as needed.
- T2WML transforms these annotated regions into canonical format where each row is a single record, see the Output Preview pane on the right.
- The ingested data is aligned with Wikidata to enable joins across multiple data sources. Notice in the Output pane in the lower right should that the string ARGENTINA is wikified and mapped on the Wikidata entity Q414.

The screenshot shows the T2WML version 2.9.4 interface. On the left is a File Tree with folders for 'AMIS.xlsx' and 'fsi-2012.csv'. The main area displays a table with columns A-F and rows 1-41. Row 10 is highlighted, showing 'TOTAL SUPPLY' for 'Argentina' in 'Million tonnes' with a value of 18.10 in 2006. On the right, the 'Output Preview' shows a table with columns A-E and rows 1-16. Row 10 is highlighted, showing 'Argentina' with 'TOTAL SUPPLY' of 18.10 Million tonnes in 2006. Below the preview is a 'YAML Editor' and an 'Output' section showing the structured data for 'Argentina (Q14)'.

A	B	C	D	E	F
1	Data Last ...	03 Apr 20...			
2	COUNTRY:	ARGENTINA			
3	COMMOD...	WHEAT			
4	SOURCE	FAO-AMIS...			
5					
6					
7	Supply an...				
8					
9	National ...	December...	2006	2007	2008
10	TOTAL SU...	Million to...	18.10	19.94	12.76
11	Opening S...	Million to...	3.55	3.45	4.25
12	Production	Million to...	14.55	16.49	8.51
13	Imports (L...	Million to...	0	0	0
14	TOTAL UTL...	Million to...	18.10	19.94	12.76
15	Domestic ...	Million to...	4.92	5.07	5.04
16	Food Use	Million to...	4.59	4.6	4.69
17	Feed Use	Million to...	0.03	0.17	0.2
18	Other Uses	Million to...	0.3	0.3	0.15
19	Exports (N...	Million to...	9.73	10.62	6.32
20	Closing St...	Million to...	3.45	4.25	1.4
21	UNBALAN...	Million to...	0	0	0
22	UNBALAN...				
23	For exam...				
24					
25	Internatio...	July/June			
26					
27	Imports (L...	Million to...	0	0	0
28	Exports (IT...	Million to...	11.3	9.37	8.1
29	For exam...				
30					
31	Other				
32					
33			2006	2007	2008
34	Population	1000s	39559	39970	40382
35	Per Capita...	Kg/Yr	116	115.1	116.1
36	Area Plant...	Million Ha			
37	Area Harv...	Million Ha	5.54	5.83	4.34
38	Yield	Tonnes/Ha	2.63	2.83	1.96
39					
40					
41					

A	B	C	D	E
1	subject	property	value	unit
2	Argentina	TOTAL SUPPLY	18.10	Million tonnes
3	Argentina	TOTAL SUPPLY	19.94	Million tonnes
4	Argentina	TOTAL SUPPLY	12.76	Million tonnes
5	Argentina	TOTAL SUPPLY	10.42	Million tonnes
6	Argentina	TOTAL SUPPLY	17.14	Million tonnes
7	Argentina	TOTAL SUPPLY	17.99	Million tonnes
8	Argentina	TOTAL SUPPLY	9.19	Million tonnes
9	Argentina	TOTAL SUPPLY	9.95	Million tonnes
10	Argentina	TOTAL SUPPLY	16.48	Million tonnes
11	Argentina	TOTAL SUPPLY	17.10	Million tonnes
12	Argentina	TOTAL SUPPLY	20.51	Million tonnes
13	Argentina	TOTAL SUPPLY	19.80	Million tonnes
14	Argentina	TOTAL SUPPLY	21.01	Million tonnes
15	Argentina	Opening Stocks	3.55	Million tonnes
16	Argentina	Opening Stocks	3.45	Million tonnes

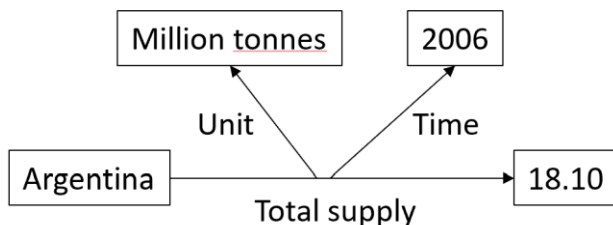
**Output**

Argentina (Q14)

TOTAL SUPPLY (Pcausx3-total\_supply\_) 18.10 ( Million tonnes)  
- point in time: 2006-01-01T00:00:00

## Datamart Core

Internally Datamart stores datasets as part of the Wikidata knowledge graph. Statements extracted from tabular data are stored as edges, e.g. the Total supply for Argentina is 18.10. Additional context about statements are stored as qualifiers on the edges, e.g. the Unit is Million tonnes and the Time is 2006.



The backend storage engine for Datamart is Postgres. The Datamart database schema is designed to be compatible with SPARQL [Blazegraph](#) and the [Knowledge Graph Toolkit \(KGTK\)](#).

Datamart always outputs structured data in canonical format with one record per row, see below. The canonical data format is defined [here](#).

	A	B	C	D	E
1	subject	property	value	unit	point in time
2	ARGENTINA	TOTAL SUPPLY	18.1	Million tonnes	2006
3	ARGENTINA	TOTAL SUPPLY	19.94	Million tonnes	2007
4	ARGENTINA	TOTAL SUPPLY	12.76	Million tonnes	2008
5	ARGENTINA	TOTAL SUPPLY	10.42	Million tonnes	2009
6	ARGENTINA	TOTAL SUPPLY	17.14	Million tonnes	2010
7	ARGENTINA	TOTAL SUPPLY	17.99	Million tonnes	2011
8	ARGENTINA	TOTAL SUPPLY	9.19	Million tonnes	2012

## Datamart Fuzzy Search

ISI Datamart provides a graphical user interface to search for datasets. The functionality include:

- Approximate string matches of query terms to dataset metadata description
- Geographical search of countries and administrative regions.
- Search results returned as scrollable tiles
- Time series overview of datasets

The screenshot shows the ISI Datamart search interface. At the top, there is a search bar with the query 'food' and a 'Choose Regions' button. Below the search bar, a list of regions is displayed: Ethiopia, Abazinia, Afghanistan, Albania, Algeria... A 'Debug queries' button is also visible.

The search results are displayed in a grid of tiles. Each tile contains a dataset title, ID, and a relevance score. The tiles shown are:

- Food, beverages and t...** (P110763) 0.046: Value added in manufacturing is the sum of gross output less the value of intermediate inputs used in production for industries classified in ISIC major division D. Food, beverages, and tobacco correspond to ISIC divisions 15 and 16. Time range: 1990 - 2015, Time precision: year, Count: 26 records.
- Travel services (% of c...** (P111435) 0.035: Travel services (% of commercial service imports) covers goods and services acquired from an economy by travelers in that economy for their own use during visits of less than one year for business or personal purposes. Time range: 1977 - 2017, Time precision: year, Count: 41 records.
- Travel services (% of c...** (P111464) 0.035: Travel services (% of commercial service exports) covers goods and services acquired from an economy by travelers in that economy for their own use during visits of less than one year for business or personal purposes. Time range: 1977 - 2017, Time precision: year, Count: 41 records.
- CPI (Alcoholic Bevera...** (P1200000) 0.023: CPI (Alcoholic Beverages and Tobacco). Time range: 2000-12 - 2017-6, Time precision: month, Count: 199 records.
- maximum food additi...** (P4851) 0.020: maximum allowed level of food additive permitted in a quantity of food.
- Women who believe a...** (P111036) 0.018: Percentage of women ages 15-49 who believe a husband/partner is justified in hitting or beating his wife/partner when she burns the food. Time range: 2000 - 2016, Time precision: year, Count: 4 records.

On the right side of the interface, a time series chart is displayed for the dataset 'P110026 Food production Index (2004-2006 = 100)'. The chart shows the index value over time from 1960 to 2010. The y-axis ranges from 0 to 1000. The x-axis shows years from 1960 to 2010. The chart features a prominent green line for Afghanistan, which peaks around 1980 at approximately 1000. Other countries are represented by various colored lines, mostly clustered below 200. A legend at the bottom identifies the countries: Afghanistan (green), Albania (purple), Algeria (blue), Angola (light blue), Antigua and Barbuda (dark blue), Argentina (orange), and Armenia (red).

## System

Datamart is dockerized so that it can be easily deployed. A sample Datamart deployment is here: <https://dsbox02.isi.edu:8888/datamart-api-dev/>

The ISI Datamart dataset metadata schema and data schema are defined here: <https://datamart-upload.readthedocs.io/en/latest/>

The ISI Datamart API is defined here: <https://datamart-upload.readthedocs.io/en/latest/api/>. This Jupyter notebook demonstrates how to use the the API (<https://github.com/usc-isi-i2/datamart-api/blob/master/Datamart%20Data%20API%20Demo.ipynb>).

Datamart repository: <https://github.com/usc-isi-i2/datamart-api>

T2WML repository: <https://github.com/usc-isi-i2/t2wml>

Fuzzy search repository: <https://github.com/usc-isi-i2/wikidata-fuzzy-search>

