

CSCI-548: Data Cleaning Spring 2015

Bo Wu

University of Southern California/ISI

Outline

- Introduction
- OpenRefine
- Data Wrangler
- FlashFill

Problem

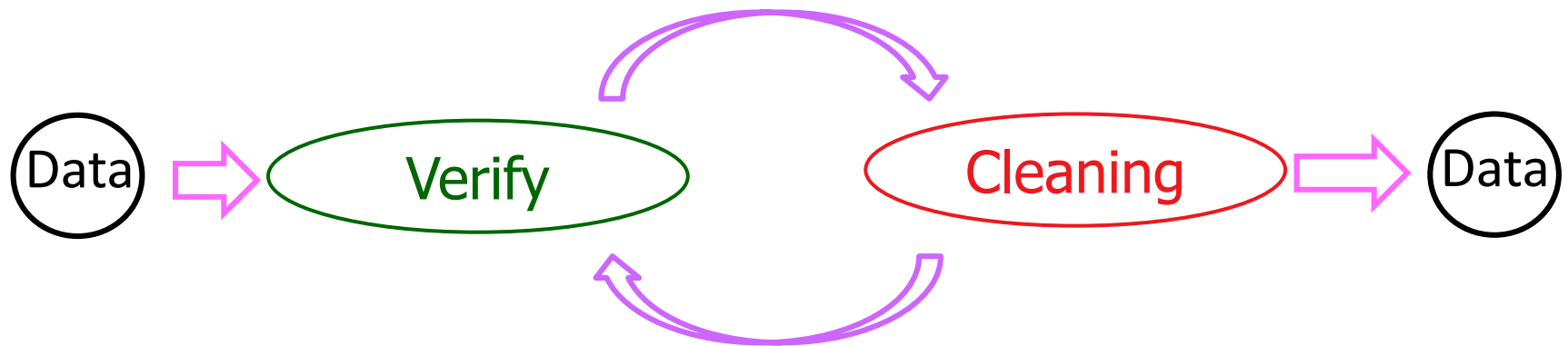
- Data is not in the expected format

artist	artistyear	artyear	credittext	dimensions	photo
FRANK STELLA	born 1936	1970	Gift of Alice Pratt Brown	120 x 600 inches	Oil on canvas
JASPER JOHNS	born 1930		Museum purchase with funds provided by; the Agnes Cullen Arnold Endowment Fund	75 x 50 inches	http://www.mfah.org/site_media/cache/e5/37/e537510de3215396c6c0336b17eab82b.jpg
GEORGE BELLOWS	1882 - 1925	1914	Gift of Mr. and Mrs. Meredith Long in memory of Agnes Cullen Arnold	38 x 30 inches	http://www.mfah.org/site_media/cache/72/7c/727c9b61b85b16d432a79031137e3d6b.jpg

Solutions

- Prevent dirty data from getting into the system
 - Enforce source integrity constraints
 - Not allow “null” for a field
 - Only allow numbers
 - ...
- Clean data
 - Manually clean the data
 - Create clean scripts

Data Cleaning Workflow



Data Transformation Operations

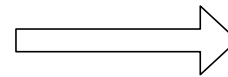
- One to one mapping
- One to many/ many to one mapping
- Look up and join
- Positional
- Filter
- Common functions: sum, min, max, avg ...

One Row to One Row Mappings

- Drop Column
- Copy Column
- Add Column
 - Constant value, random number, serial number
- Merge Columns with Glue
- Value update

Name
George
Paul
Ann

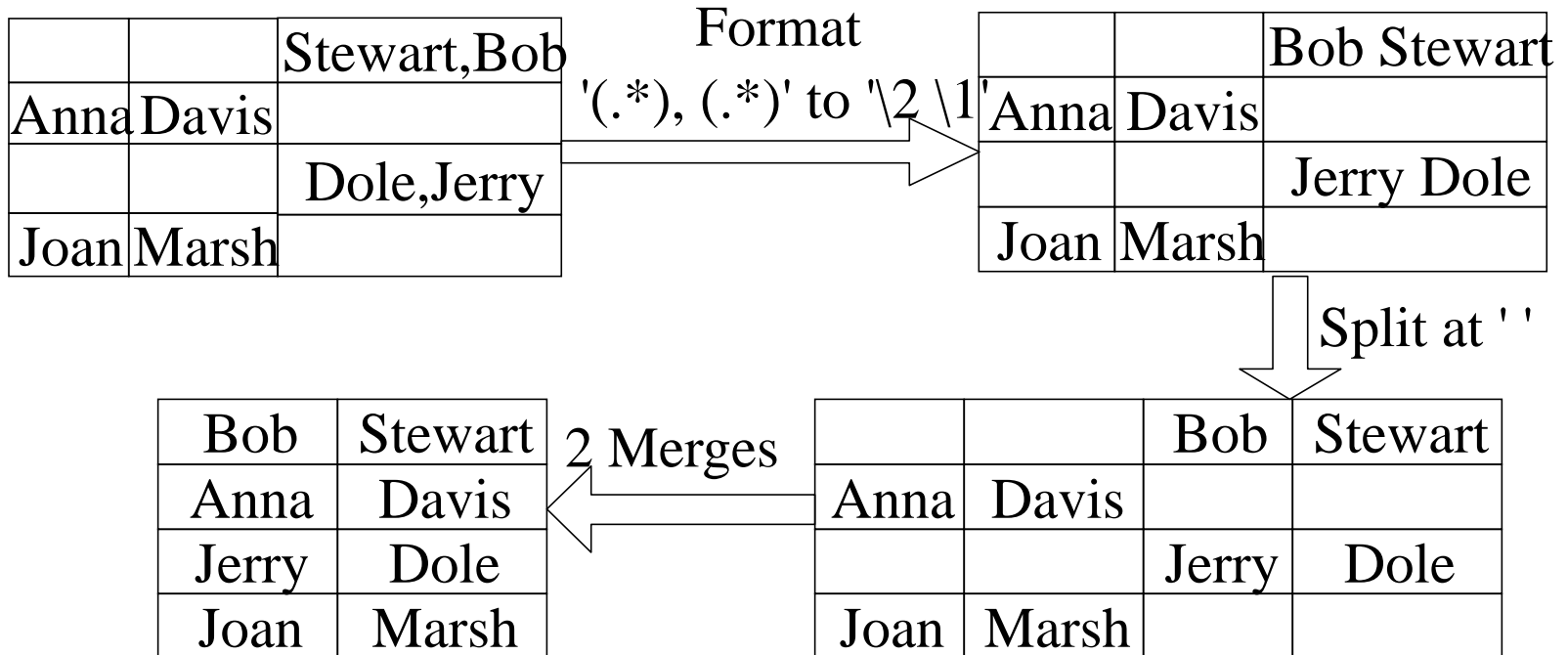
(.*) to `<\C>1</\C>`



Name
<code><Name>George</Name></code>
<code><Name>Paul</Name></code>
<code><Name>Ann</Name></code>

- Split Column
 - by position
 - by regular expression (first match)

Example



One Row To Many Rows: Fold

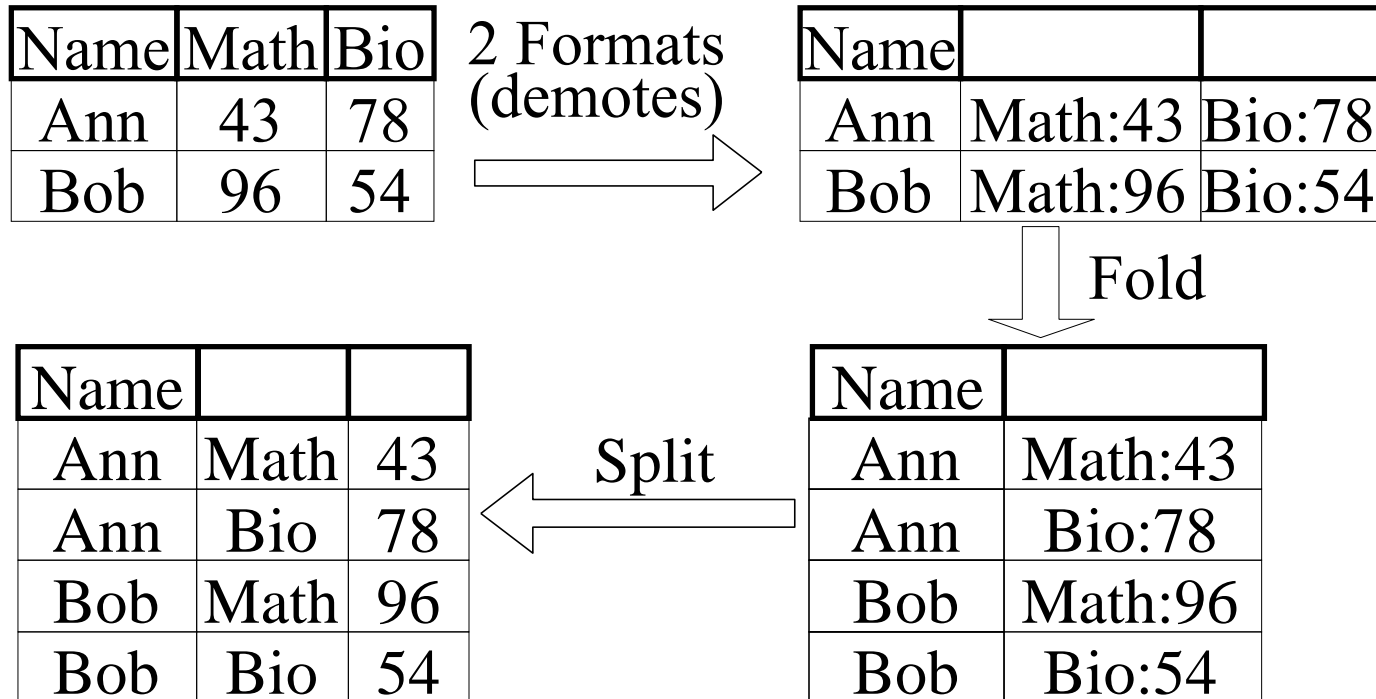
A	B	C	D	E	F
a1	b1	c1	d1	e1	f1



fold(D, E, F)

A	B	C	
a1	b1	c1	d1
a1	b1	c1	e1
a1	b1	c1	f1

One Row To Many Rows: Fold



Many Rows to One Row: Unfold

unfold(col_1, col_2)

Name		
George	Math	65
George	French	42
Anna	Math	43
Anna	French	78
Bob	English	96
Bob	French	54
Joan	English	79

unfold(2,3)

Name	Math	French	English
George	65	42	
Anna	43	78	
Bob		54	96
Joan			79

OpenRefine

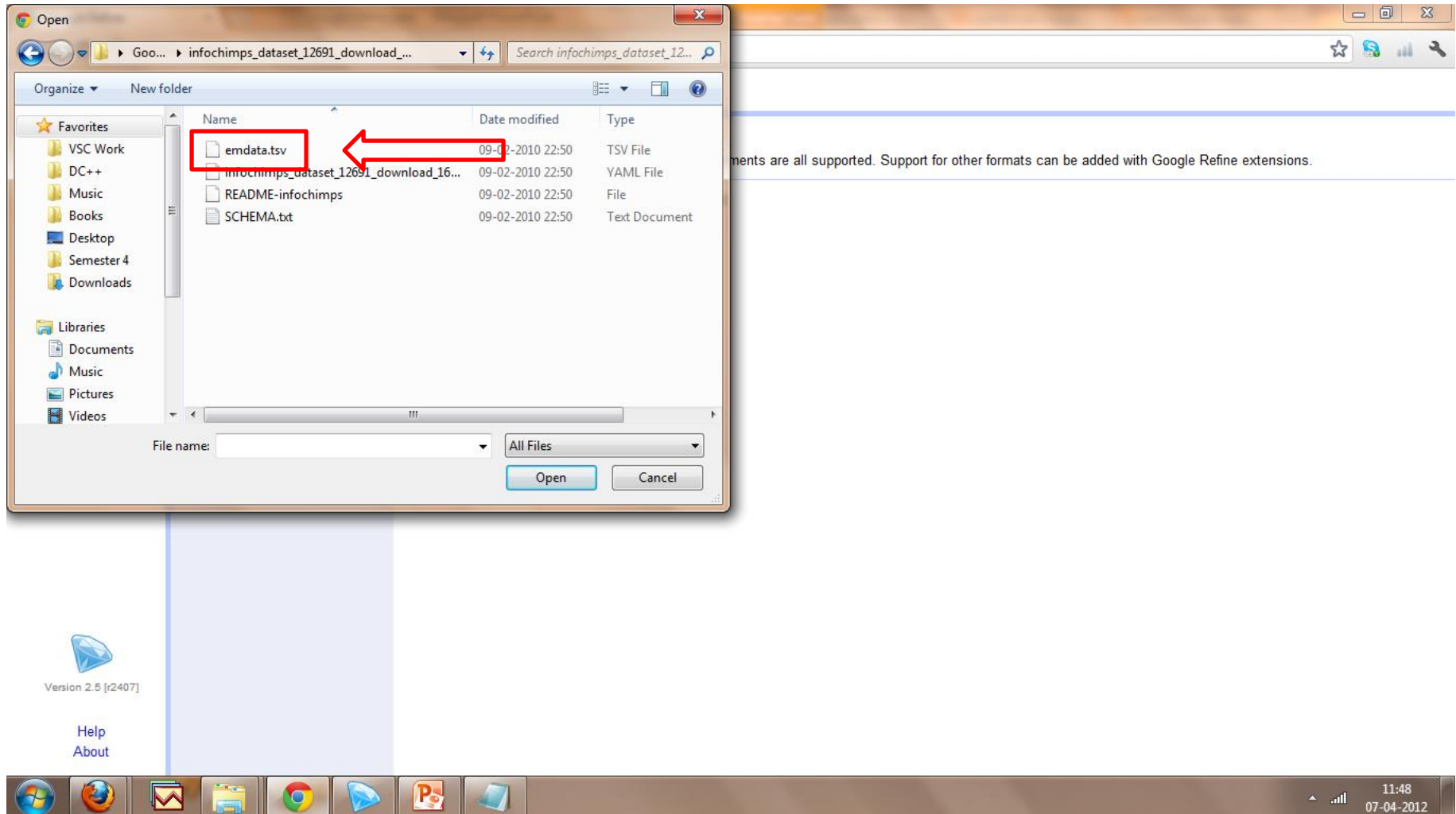


- A powerful tool that can be effectively used for data cleansing.
- It helps in working with raw data, cleaning it up, transforming from one format to other, encompassing it with web services and linking it to databases.
- It is very easy to use and has a web interface.
- It is freely available and works well with any browser.
- a desktop application and it runs a small web server on your system and we need to point our browser to the server to use refine.

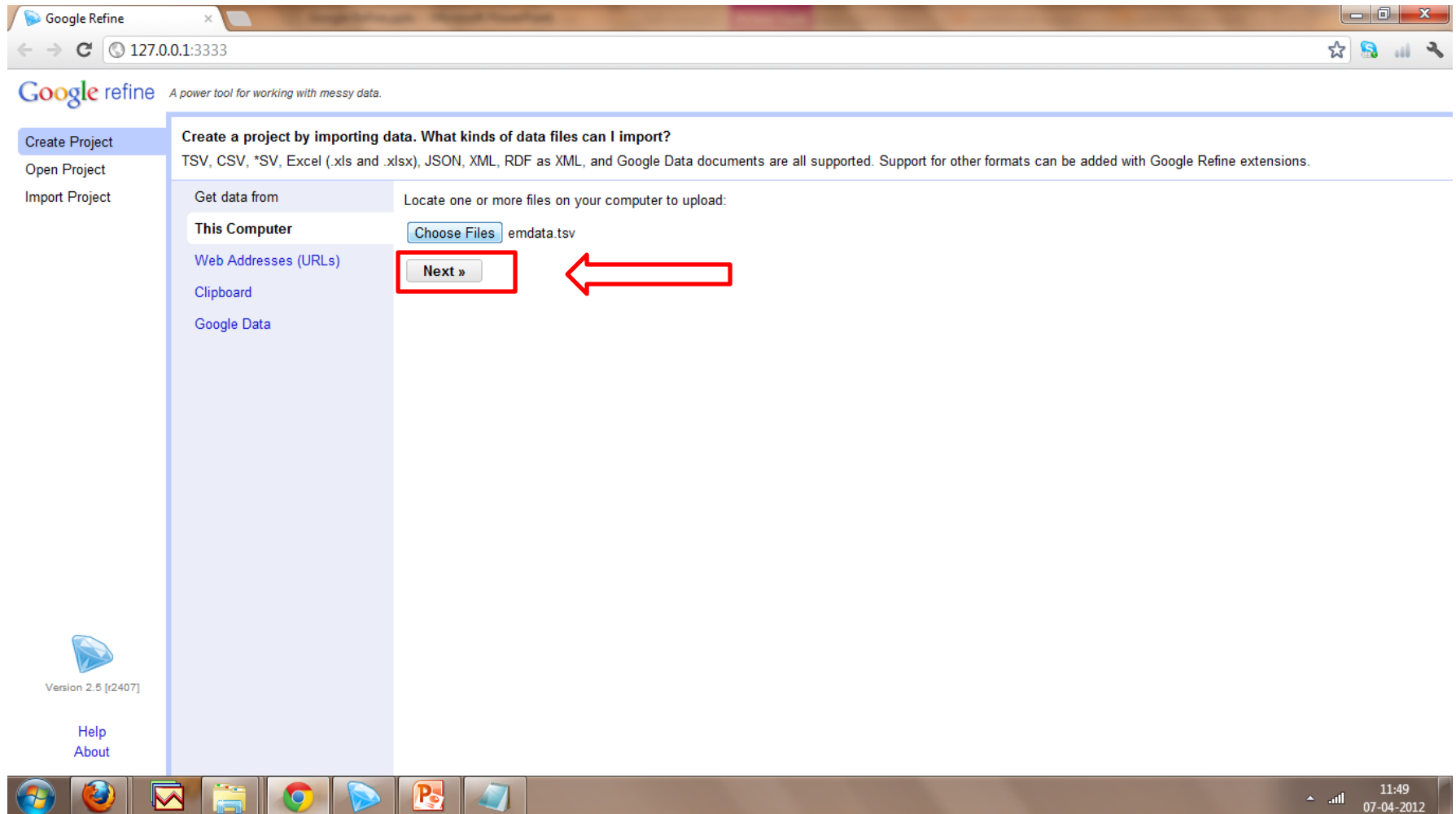
Importing Data

- Google Refine supports TSV, CSV, Excel (.xls and .xlsx), JSON, XML, and Google data document formats.
- Once imported the data is in Google Refine's own data format.
- We have used TSV data on Disasters worldwide from 1900-2008 available from <http://www.infochimps.com/datasets/disasters-worldwide-from-1900-2008> for the tutorial.

Importing Data



Importing Data



Data
Uploaded

Creating Project

Google Refine 127.0.0.1:3333

Google refine A power tool for working with messy data.

Create Project Open Project Import Project

« Start Over Configure Parsing Options Project name emdata tsv Create Project »

	Start	End	Country	Location	Type	Sub_Type	Name	Killed	Cost	Affected	Id	Column 12
1.	102008	102008	Afghanistan	Kunduz, Balkh, Faryab, Ba	Drought	Drought			280000		2008-9475	
2.	72006	2006	Afghanistan		Drought	Drought			1900000		2006-9570	
3.	52000	2002	Afghanistan	Kandahar, Helmand, Nimroz ...	Drought	Drought		37	2580000	0.05	2000-9186	
4.	81971	1973	Afghanistan	Central, North-West, Nort ...	Drought	Drought					1971-9085	
5.	11969	1969	Afghanistan	Paktia province	Drought	Drought			48000	0.2	1969-9007	
6.	29072006	29072006	Afghanistan	Emam Sahib (Kunduz distri	Earthquake (seismic activity)	Earthquake (ground shaking)		1	935		2006-0405	
7.	13122005	13122005	Afghanistan	Hindu Kush	Earthquake (seismic activity)	Earthquake (ground shaking)		5	501		2005-0686	
8.	8102005	8102005	Afghanistan	Nangarhar, Jalalabad prov	Earthquake (seismic activity)	Earthquake (ground shaking)		1		0.05	2005-0575	

Parse data as

CSV / TSV / separator-based files

- Line-based text files
- Fixed-width field text files
- PC-Axis text files
- JSON files
- RDF/N3 files
- XML files
- Open Document Format spreadsheets (.ods)
- RDF/XML files

Character encoding

Columns are separated by

- ☐ commas (CSV)
- ☒ tabs (TSV)
- ☐ custom \t

Escape special characters with \

☐ Ignore first 0 line(s) at beginning of file

☒ Parse next 1 line(s) as column headers

☐ Discard initial 0 row(s) of data

☐ Load at most 0 row(s) of data

☒ Parse cell text into numbers, dates, ...

☒ Quotation marks are used to enclose cells containing column separators

☒ Store blank rows

☒ Store blank cells as nulls

☐ Store file source (file names, URLs) in each row

Update Preview

Version 2.5 [r2407]

Help About

Creating Project

Project
Created

17828 rows

Extensions: Freebase

Show as: rows records

Show: 5 10 25 50 rows

« first « previous 1 - 10 next » las

▼ All	▼ Start	▼ End	▼ Country	▼ Location	▼ Type	▼ Sub_Type	▼ Name	▼ Killed	▼ Cost	▼ Affected	▼ Id
☆	1.	102008	102008	Afghanistan	Kunduz, Balkh, Faryab, Ba ...	Drought	Drought		280000		2008-9475
☆	2.	72006	2006	Afghanistan		Drought	Drought		1900000		2006-9570
☆	3.	52000	2002	Afghanistan	Kandahar, Helmand, Nimroz ...	Drought	Drought	37	2580000	0.05	2000-9186
☆	4.	81971	1973	Afghanistan	Central, North-West, Nort ...	Drought	Drought				1971-9085
☆	5.	11969	1969	Afghanistan	Paktia province	Drought	Drought		48000	0.2	1969-9007
☆	6.	29072006	29072006	Afghanistan	Emam Sahib (Kunduz distri ...	Earthquake (seismic activity)	Earthquake (ground shaking)	1	935		2006-0405
☆	7.	13122005	13122005	Afghanistan	Hindu Kush	Earthquake (seismic activity)	Earthquake (ground shaking)	5	501		2005-0686
☆	8.	8102005	8102005	Afghanistan	Nangarhar, Jalalabad prov ...	Earthquake (seismic activity)	Earthquake (ground shaking)	1		0.05	2005-0575
☆	9.	18072004	18072004	Afghanistan	Paktia province	Earthquake (seismic activity)	Earthquake (ground shaking)	2	1040		2004-0436
☆	10.	10042003	10042003	Afghanistan	Yakabagh (Takhar province ...	Earthquake (seismic activity)	Earthquake (ground shaking)	1	1001		2003-0236

Faceting

- Faceting is about seeing the big picture and filtering based on rows to work on data you want to change in bulk.
- We can create a facet for a column to get the details about that column and then we can filter to a subset of rows with a constraint.
- We can perform text facet, Numeric facet, timeline facet and scatterplot facet. Also various customized facets can be designed.

Faceting

17828 rows										
Show as: rows records		Show: 5 10 25 50 rows		« first						
▼ All	▼ Start	▼ End	▼ Country	▼ Location	▼ Type	▼ Sub_Type	▼ Name	▼ Killed	▼	
☆	1.	102008	102008	Afghanistan	Kunduz, Balkh, Faryab, Ba ...	Facet ▶	Text facet			2
☆	2.	72006	2006	Afghanistan		Text filter	Numeric facet			19
☆	3.	52000	2002	Afghanistan	Kandahar, Helmand, Nimroz ...	Edit cells ▶	Timeline facet		37	25
☆	4.	81971	1973	Afghanistan	Central, North-West, Nort ...	Edit column ▶	Scatterplot facet			
☆	5.	11969	1969	Afghanistan	Paktia province	Transpose ▶	Custom text facet...			
☆	6.	29072006	29072006	Afghanistan	Emam Sahib (Kunduz distri ...	Sort...	Custom numeric facet...			
☆	7.	13122005	13122005	Afghanistan	Hindu Kush	View ▶	Customized facets ▶		1	
☆					Reconcile ▶	ity) Earthquake (ground shaking)			5	

Faceting

emdata tsv - Google Refine x

127.0.0.1:3333/project?project=1912038646727

Google refine emdata tsv Permalink

Facet / Filter Undo / Redo 0

Refresh Reset All Remove All

17828 rows

Show as: rows records Show: 5 10 25 50 rows

Type 18 choices Sort by: name count

The Column Type has 18 unique options

Type	All	Start	End	Country	Location
Complex Disasters 12	102008	102008	Afghanistan	Kunduz, Balkh,	
Drought 561	72006	2006	Afghanistan		
Earthquake (seismic activity) 1120	22000	2002	Afghanistan	Kandahar, Helm	
Epidemic 1179	4.	81971	1973	Afghanistan	Central, North-V
Extreme temperature 361	5.	11969	1969	Afghanistan	Paktia province
Flood 3512	6.	29072006	29072006	Afghanistan	Emam Sahib (Ki
Industrial Accident 1190	7.	13122005	13122005	Afghanistan	Hindu Kush
Insect infestation 83	8.	8102005	8102005	Afghanistan	Nangarhar, Jala
Mass movement dry 48	9.	18072004	18072004	Afghanistan	Paktia province
Mass Movement Dry 4					
Mass Movement Wet 12					

Removing Redundancy

Google refine emdata tsv Permalink

Facet / Filter Undo / Redo

Refresh Reset All Remove All

Type change

18 choices Sort by: name count Cluster

Industrial Accident 1190

Insect infestation 83

Mass movement dry 48

Mass Movement Dry 4

Mass Movement Wet 12

Mass movement wet 105

Miscellaneous accident 1133

Storm 3206

Transport accident 196

Transport Accident 4155

Volcano 211

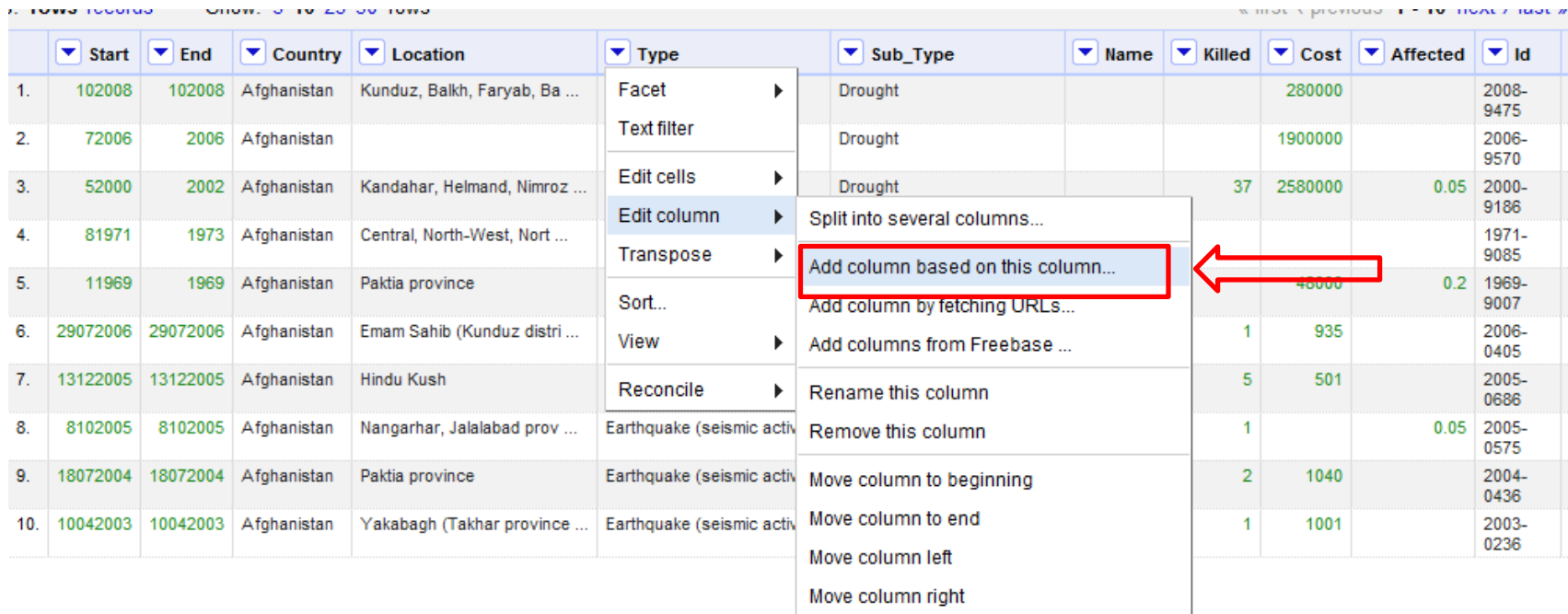
17828 rows

Show as: rows records S

			All	Start	End
☆	🗨	1.	102008	10200	
☆	🗨	2.	72006	200	
☆	🗨	3.	52000	200	
☆	🗨	4.	81971	197	
☆	🗨	5.	11969	196	
☆	🗨	6.	29072006	2907200	
☆	🗨	7.	13122005	1312200	
☆	🗨	8.	8102005	810200	
☆	🗨	9.	18072004	1807200	

Even though they are of same type, shows as different options due to case

Removing Redundancy



	Start	End	Country	Location	Type	Sub_Type	Name	Killed	Cost	Affected	Id
1.	102008	102008	Afghanistan	Kunduz, Balkh, Faryab, Ba ...	Facet	Drought			280000		2008-9475
2.	72006	2006	Afghanistan		Text filter	Drought			1900000		2006-9570
3.	52000	2002	Afghanistan	Kandahar, Helmand, Nimroz ...	Edit cells	Drought		37	2580000	0.05	2000-9186
4.	81971	1973	Afghanistan	Central, North-West, Nort ...	Edit column						1971-9085
5.	11969	1969	Afghanistan	Paktia province	Transpose				48000	0.2	1969-9007
6.	29072006	29072006	Afghanistan	Emam Sahib (Kunduz distri ...	Sort...			1	935		2006-0405
7.	13122005	13122005	Afghanistan	Hindu Kush	View			5	501		2005-0686
8.	8102005	8102005	Afghanistan	Nangarhar, Jalalabad prov ...	Reconcile			1		0.05	2005-0575
9.	18072004	18072004	Afghanistan	Paktia province	Earthquake (seismic activ			2	1040		2004-0436
10.	10042003	10042003	Afghanistan	Yakabagh (Takhar province ...	Earthquake (seismic activ			1	1001		2003-0236

Removing Redundancy

data tsv [Permalink](#)

Redo 0

Reset All Remove All

change

count Cluster

33

17

▼ Name ▼ K

Add column based on column Type

New column name

On error ☒ set to blank ☐ store error ☐ copy value from original column

Expression Language Google Refine Expression Language (GREL) ▼

No syntax error.

Preview History Starred Help

row	value	value
1.	Drought	Drought
2.	Drought	Drought
3.	Drought	Drought
4.	Drought	Drought
5.	Drought	Drought
6.	Earthquake (seismic activity)	Earthquake (seismic activity)

OK Cancel

Removing Redundancy

The screenshot shows a dialog box titled "Add column based on column Type". It has a light blue header. Below the header, there is a text input field for "New column name" containing "Type_UC". Under the "On error" section, three radio buttons are present: "set to blank" (selected), "store error", and "copy value from original column". The "Expression" field contains the text `toUppercase(value)`, which is underlined in red. To the right of the expression field is a "Language" dropdown menu set to "Google Refine Expression Language (GREL)". To the right of the expression field, the text "No syntax error." is displayed. At the bottom of the dialog, there are four tabs: "Preview", "History", "Starred", and "Help". The "Preview" tab is active, showing a table with two columns: "row value" and "toUppercase(value)".

Add column based on column Type

New column name

On error ☒ set to blank ☐ store error ☐ copy value from original column

Expression Language Google Refine Expression Language (GREL) ▼

No syntax error.

Preview History Starred Help

row value	toUppercase(value)
-----------	--------------------

Removing Redundancy

The screenshot shows the Google Refine interface with a dataset of 17,828 rows. The 'Facet / Filter' sidebar on the left shows a facet for 'Type_UC' with 15 unique choices. A red box highlights the '15 choices' text, and a red arrow points from this box to a callout bubble that says 'Reduced to 15 unique options'. The main table displays 10 rows of data, including columns for Start, End, Country, Location, Type, Type_UC, Sub_Type, Name, Killed, and Cost. The 'Type_UC' column shows values like 'DROUGHT' and 'EARTHQUAKE (SEISMIC ACTIVITY)'.

emdata tsv - Google Refine

127.0.0.1:3333/project?project=1912038646727

Google refine emdata tsv Permalink

Open... Export Help

Facet / Filter Undo / Redo

Refresh Reset All Remove All

17828 rows

Show as: rows records Show: 5 10 25 50 rows

Extensions: Freebase

« first < previous 1 - 10 next > last »

Type_UC change

15 choices Sort by: name count Cluster

COMPLEX DISASTERS 12

DROUGHT 561

EARTHQUAKE (SEISMIC ACTIVITY) 1120

EPIDEMIC 1179

EXTREME TEMPERATURE 361

FLOOD 351

INDUSTRIAL ACCIDENT 1190

INSECT INFESTATION 83

MASS MOVEMENT DRY 52

MASS MOVEMENT WET 517

1. 102008 102008 Afghanistan Kunduz, Balkh, Faryab, Ba Drought DROUGHT Drought 280000

2. 72006 2006 Afghanistan Drought DROUGHT 1900000

3. 52000 2002 Afghanistan Kandahar, Helmand, Nimroz Drought DROUGHT Drought 37 2580000

4. 81971 1973 Afghanistan Central, North-West, Nort... Drought DROUGHT Drought

5. 11969 1969 Afghanistan Paktia province Drought DROUGHT Drought 48000

6. 29072006 29072006 Afghanistan Emam Sahib (Kunduz distri Earthquake (seismic activity) EARTHQUAKE (SEISMIC ACTIVITY) Earthquake (ground shaking) 1 935

7. 13122005 13122005 Afghanistan Hindu Kush Earthquake (seismic activity) EARTHQUAKE (SEISMIC ACTIVITY) Earthquake (ground shaking) 5 501

8. 8102005 8102005 Afghanistan Nangarhar, Jalalabad prov Earthquake (seismic activity) EARTHQUAKE (SEISMIC ACTIVITY) Earthquake (ground shaking) 1

9. 18072004 18072004 Afghanistan Paktia province Earthquake (seismic activity) EARTHQUAKE (SEISMIC ACTIVITY) Earthquake (ground shaking) 2 1040

10. 10042003 10042003 Afghanistan Yakabagh (Takhar province Earthquake (seismic activity) EARTHQUAKE (SEISMIC ACTIVITY) Earthquake (ground shaking) 1 1001

javascript:{} 12:37 07-04-2012

Numeric Faceting

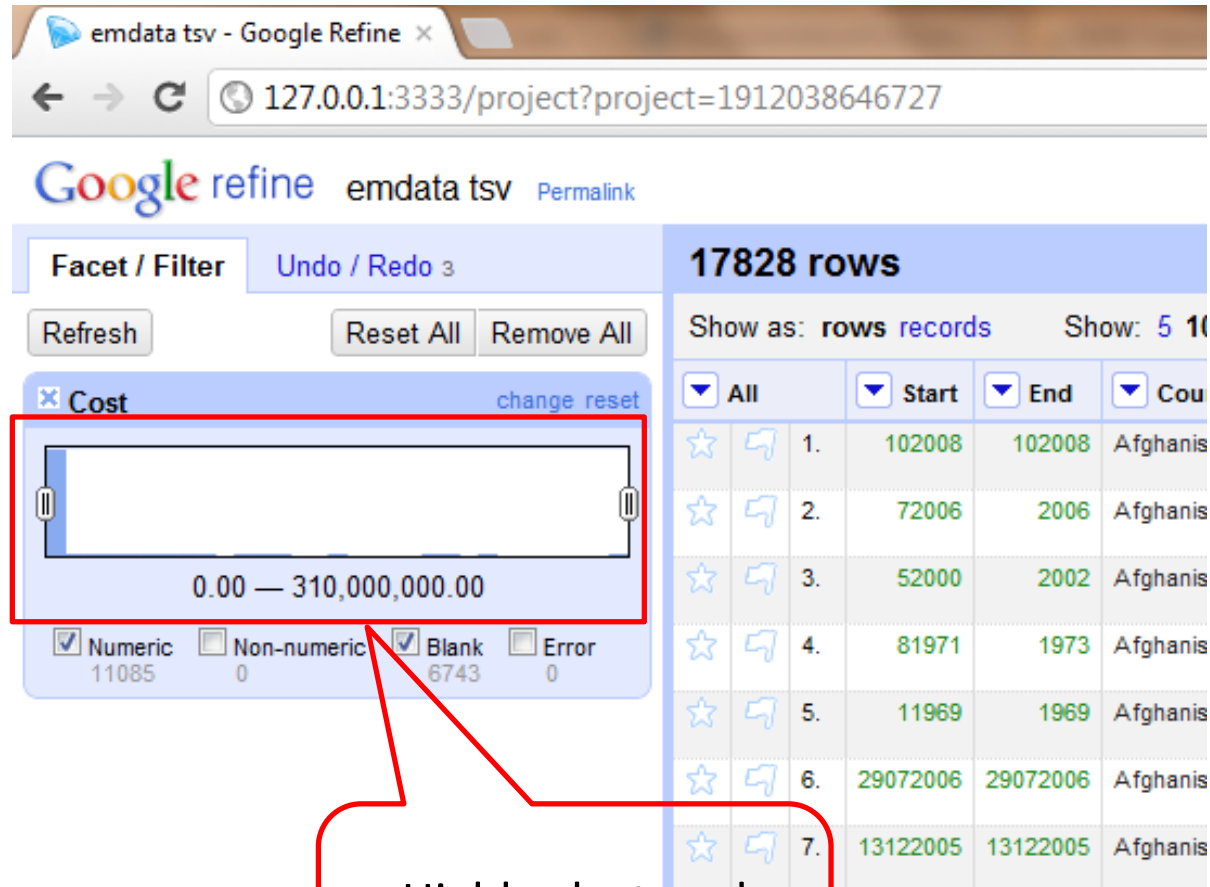
Open... Export ▾ Help

Extensions: Freebase ▾

50 rows « first ◀ previous 1 - 10 next ▶ last »

Location	Type	Type_UC	Sub_Type	Name	Killed	Cost
Kunduz, Balkh, Faryab, Ba ...	Drought	DROUGHT	Text facet	Facet		00000
	Drought	DROUGHT	Numeric facet	Text filter		00000
Kandahar, Helmand, Nimroz ...	Drought	DROUGHT	Timeline facet	Edit cells		80000
Central, North-West, Nort ...	Drought	DROUGHT	Scatterplot facet	Edit column		
Paktia province	Drought	DROUGHT	Custom text facet...	Transpose		48000
Emam Sahib (Kunduz distri ...	Earthquake (seismic activity)	EARTHQUAKE (SEISMIC ACTIVITY)	Custom numeric facet...	Sort...		935
Hindu Kush	Earthquake (seismic activity)	EARTHQUAKE (SEISMIC ACTIVITY)	Customized facets ▶	View		501
Nangarhar, Jalalabad prov ...	Earthquake (seismic activity)	EARTHQUAKE (SEISMIC ACTIVITY)	Earthquake (ground shaking)		1	
Paktia province	Earthquake (seismic activity)	EARTHQUAKE (SEISMIC ACTIVITY)	Earthquake (ground shaking)		2	1040
Yakabagh (Takhar province ...	Earthquake (seismic activity)	EARTHQUAKE (SEISMIC ACTIVITY)	Earthquake (ground shaking)		1	1001

Numeric Faceting



Highly clustered
towards low values

Numeric Faceting

The screenshot shows the Google Refine web interface. A dialog box titled "Edit Facet's Expression based on Column Cost" is open. The "Expression" field contains `value.log()`. The "Language" dropdown is set to "Google Refine Expression Language (GREL)". Below the expression field is a "Preview" section showing a table with 6 rows. The first three rows show the log of the 'Cost' column values (280000, 1900000, 2580000). The fourth row shows an error: "Error: log expects a number" for a null value. The fifth and sixth rows show the log of 48000 and 935 respectively. At the bottom of the dialog, the "OK" button is highlighted with a red box, and a red arrow points to it from the right. The background shows a facet on the "Cost" column with a range of 0.00 to 310,000,000.00. The main data table on the right has columns: Sub_Type, Name, Killed, and Cost. It lists several rows, including Drought and Earthquake (ground shaking) events.

row	value	value.log()
1.	280000	5.447158031342219
2.	1900000	6.278753600952829
3.	2580000	6.41161970596323
4.	null	Error: log expects a number
5.	48000	4.681241237375588
6.	935	2.9708116108725178

Numeric Faceting

Google refine emdata tsv [Permalink](#)

Facet / Filter

Undo / Redo 3

Refresh

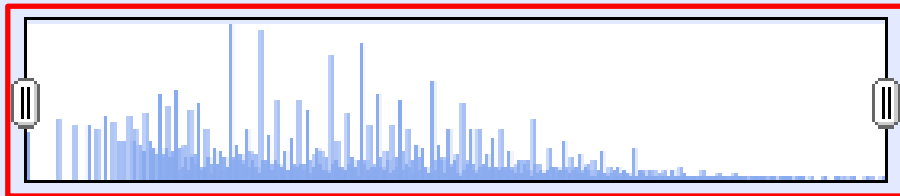
Reset All

Remove All

☒ Cost

[change](#) [reset](#)

grel:value.log()



0.00 — 8.48

☒ Numeric 11085 ☐ Non-numeric 0 ☐ Blank 0 ☒ Error 6743

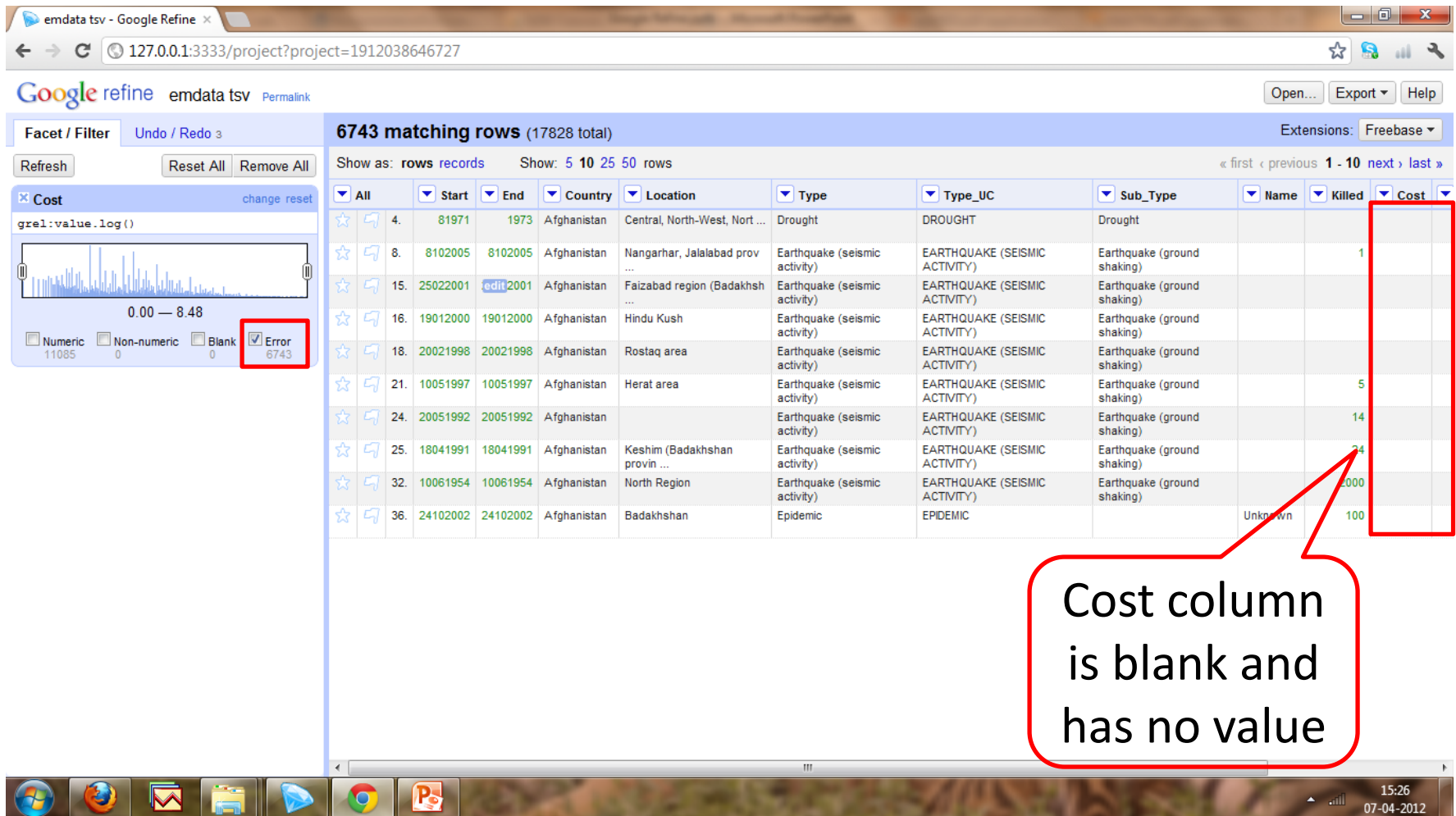
17828 rows

Show as: [rows](#) [records](#)

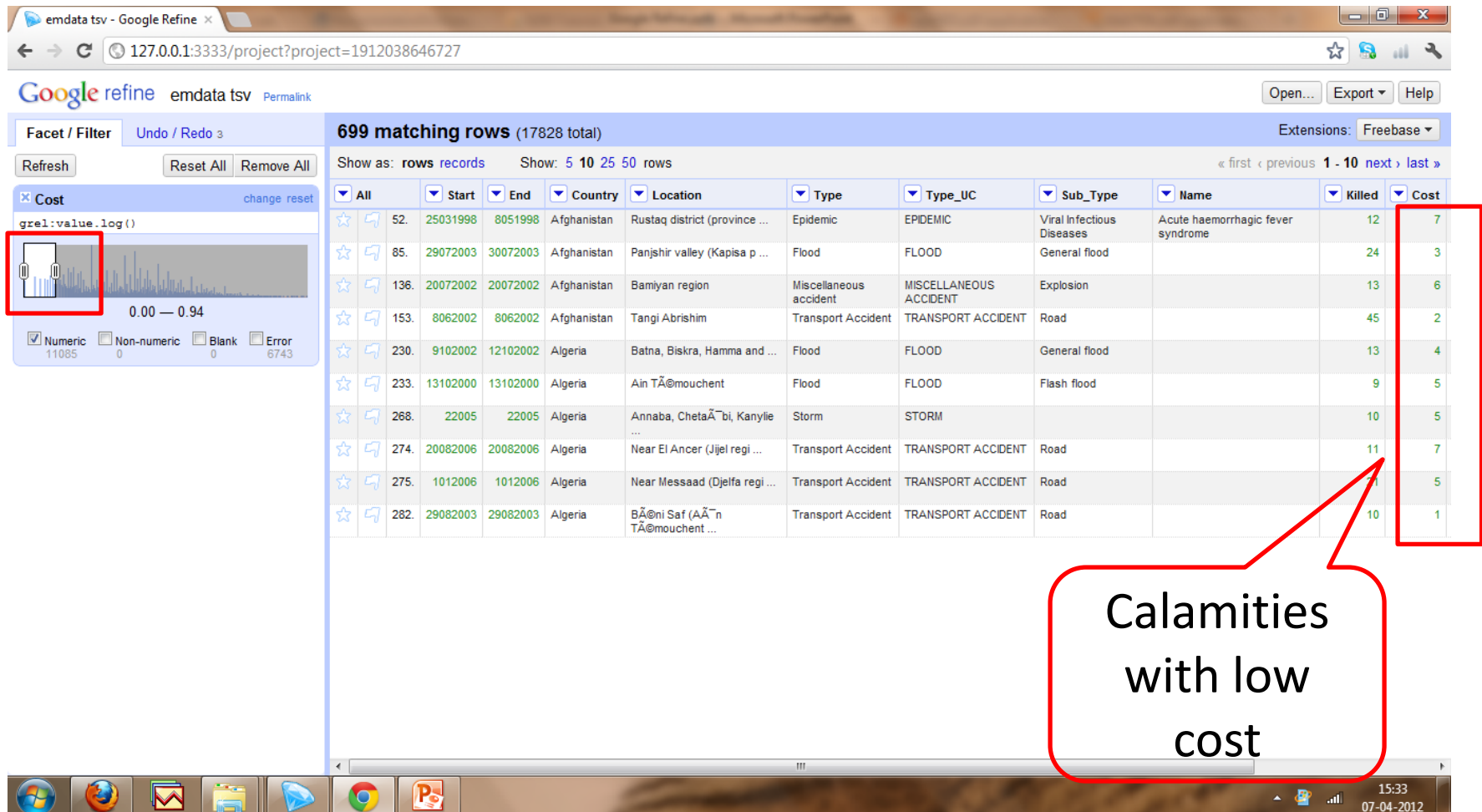
Show: 5 10

▼ All			▼ Start	▼ End	▼ Cou
☆	🗨	1.	102008	102008	Afghanis
☆	🗨	2.	72006	2006	Afghanis
☆	🗨	3.	52000	2002	Afghanis
☆	🗨	4.	81971	1973	Afghanis
☆	🗨	5.	11969	1969	Afghanis
☆	🗨	6.	29072006	29072006	Afghanis

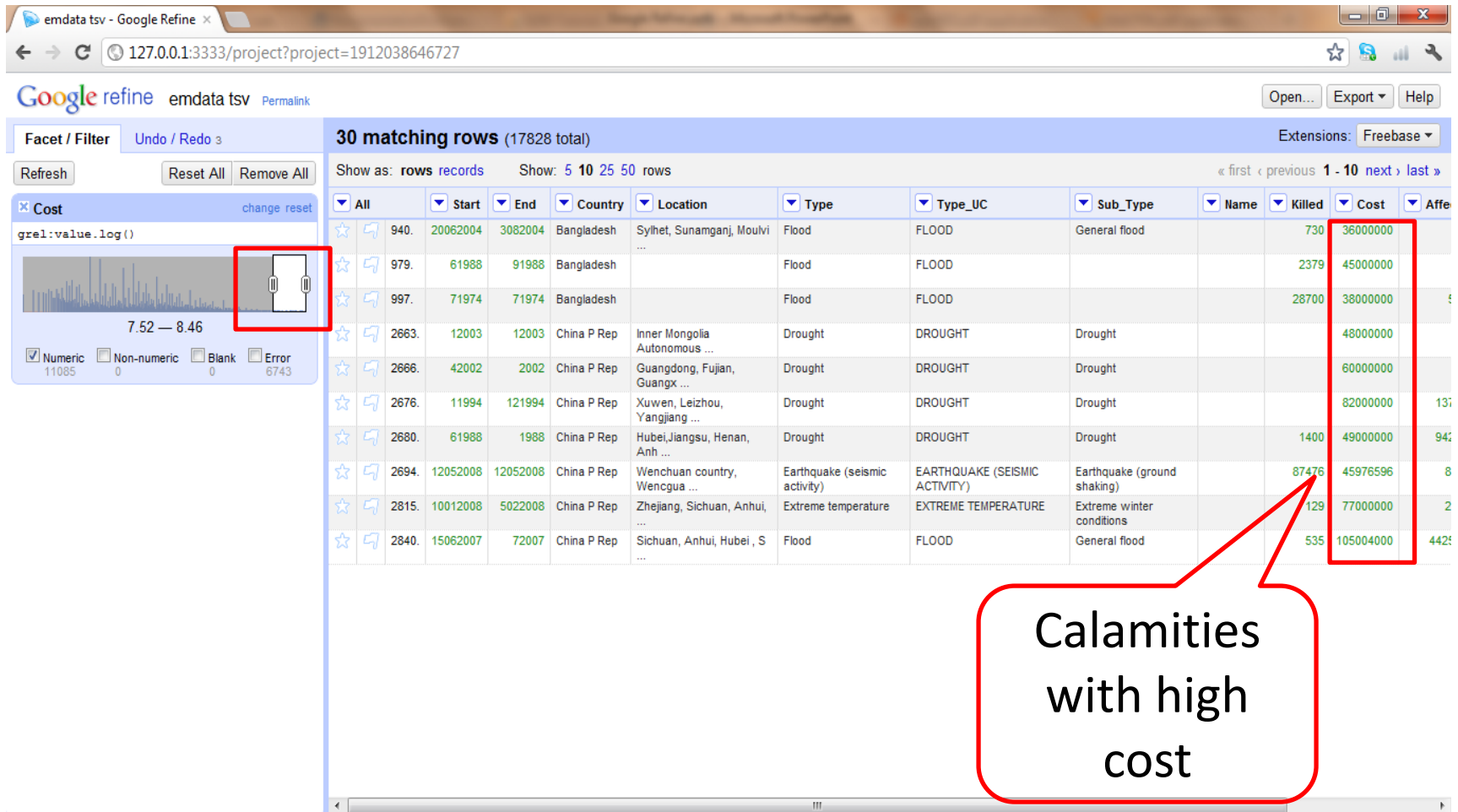
Numeric Faceting



Numeric Faceting



Numeric Faceting



Clustering

- Clustering is used to merge choices which look similar.

The screenshot shows the Google Refine interface with a data table containing 17828 rows. The interface includes a 'Facet / Filter' sidebar on the left, a 'Show as: rows records' dropdown, and a 'Show: 5 10 25 50 rows' dropdown. The table has columns for Country, Location, Type, Type_UC, Sub_Type, Name, Killed, and Cost. A red box highlights the 'Cluster' button in the Country facet, and a red arrow points to the first row of the table.

	All	Start	End	Country	Location	Type	Type_UC	Sub_Type	Name	Killed	Cost
1.	102000	102008		Afghanistan	Kunduz, Balkh, Faryab, Ba ...	Drought	DROUGHT	Drought			280000
2.	72006	2006		Afghanistan		Drought	DROUGHT	Drought			1900000
3.	52000	2002		Afghanistan	Kandahar, Helmand, Nimroz ...	Drought	DROUGHT	Drought	37		2580000
4.	81971	1973		Afghanistan	Central, North-West, Nort ...	Drought	DROUGHT	Drought			
5.	11969	1969		Afghanistan	Paktia province	Drought	DROUGHT	Drought			48000
6.	29072006	29072006		Afghanistan	Emam Sahib (Kunduz distri ...	Earthquake (seismic activity)	EARTHQUAKE (SEISMIC ACTIVITY)	Earthquake (ground shaking)	1		935
7.	13122005	13122005		Afghanistan	Hindu Kush	Earthquake (seismic activity)	EARTHQUAKE (SEISMIC ACTIVITY)	Earthquake (ground shaking)	5		501
8.	8102005	8102005		Afghanistan	Nangarhar, Jalalabad prov ...	Earthquake (seismic activity)	EARTHQUAKE (SEISMIC ACTIVITY)	Earthquake (ground shaking)	1		
9.	18072004	18072004		Afghanistan	Paktia province	Earthquake (seismic activity)	EARTHQUAKE (SEISMIC ACTIVITY)	Earthquake (ground shaking)	2		1040
10.	10042003	edit 2003		Afghanistan	Yakabagh (Takhar province ...	Earthquake (seismic activity)	EARTHQUAKE (SEISMIC ACTIVITY)	Earthquake (ground shaking)	1		1001

Clustering

Cluster & Edit column "Country"

This feature helps you find groups of different cell values that might be alternative representations of the same thing. For example, the two strings "New York" and "new york" are very likely to refer to the same concept and just have capitalization differences, and "Gödel" and "Godel" probably refer to the same person. [Find out more ...](#)

Method: **key collision** Keying Function: **fingerprint** 1 cluster found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
2	15	<ul style="list-style-type: none">St Vincent and The Grenadines (14 rows)St Vincent and the Grenadines (1 rows)	<input checked="" type="checkbox"/>	St Vincent and The Grenadines

Select All Deselect All Merge Selected & Re-Cluster Merge Selected & Close Close

Using Expressions

Custom text transform on column Country

Expression Language Google Refine Expression Language (GREL)

`value.trim()` No syntax error.

Preview History Starred Help

row	value	value.trim()
1.	Afghanistan	Afghanistan
2.	Afghanistan	Afghanistan
3.	Afghanistan	Afghanistan

Using Expressions

emdata tsv - Google Refine

127.0.0.1:3333/project?project=1912038646727

Google refine emdata tsv Permalink

Text transform on 1 cells in column Country:
grel:value.trim() Undo

Open... Export Help

Facet / Filter Undo / Redo

17828 rows Extensions: Freebase

Show as: rows records Show: 5 10 25 50 rows « first < previous 1 - 10 next > last »

Using facets and filters

Use facets and filters to select subsets of your data to act on. Choose facet and filter methods from the menus at the top of each data column.

Not sure how to get started?
[Watch these screencasts](#)

All	Start	End	Country	Location	Type	Type_UC	Sub_Type	Name	Killed	Cost
1.	102008	102008	Afghanistan	Kunduz, Balkh, Faryab, Ba	Drought	DROUGHT	Drought			280000
2.	72006	2006	Afghanistan		Drought	DROUGHT	Drought			1900000
3.	52000	2002	Afghanistan	Kandahar, Helmand, Nimroz	Drought	DROUGHT	Drought		37	2580000
4.	81971	1973	Afghanistan	Central, North-West, Nort ...	Drought	DROUGHT	Drought			
5.	11969	1969	Afghanistan	Paktia province	Drought	DROUGHT	Drought			48000
6.	29072006	29072006	Afghanistan	Emam Sahib (Kunduz distri	Earthquake (seismic activity)	EARTHQUAKE (SEISMIC ACTIVITY)	Earthquake (ground shaking)		1	935
7.	13122005	13122005	Afghanistan	Hindu Kush	Earthquake (seismic activity)	EARTHQUAKE (SEISMIC ACTIVITY)	Earthquake (ground shaking)		5	501
8.	8102005	8102005	Afghanistan	Nangarhar, Jalalabad prov	Earthquake (seismic activity)	EARTHQUAKE (SEISMIC ACTIVITY)	Earthquake (ground shaking)		1	
9.	18072004	18072004	Afghanistan	Paktia province	Earthquake (seismic activity)	EARTHQUAKE (SEISMIC ACTIVITY)	Earthquake (ground shaking)		2	1040
10.	10042003	10042003	Afghanistan	Yakabagh (Takhar province	Earthquake (seismic activity)	EARTHQUAKE (SEISMIC ACTIVITY)	Earthquake (ground shaking)		1	1001

14:37 07-04-2012

OpenRefine

- Nice tutorial
 - http://www.intersect.org.au/docs/GoogleRefine_Slides.pdf
 - <http://www.intersect.org.au/docs/GoogleRefineExercises.pdf>

Data Wrangler

- An interactive system for creating data transformations
 - Suggest transforms
 - Iteratively explore the space of applicable operations

The screenshot displays the Data Wrangler interface. On the left, the 'Transform Script' panel lists three operations: 'Split data repeatedly on newline into rows', 'Split split repeatedly on \',', and 'Promote row 0 to header'. Below the script, there are buttons for 'Text', 'Columns', 'Rows', 'Table', and 'Clear'. Further down, there are options to 'Delete row 7', 'Delete empty rows', and 'Fill row 7 by copying values from above'. On the right, a table is shown with two columns: 'Year' and 'Property_crime_rate'. The table contains data for 'Reported crime in Alabama' and 'Reported crime in Alaska' for the years 2004 through 2006. Row 7 is highlighted in light blue.

	Year	Property_crime_rate
0	Reported crime in Alabama	
1		
2	2004	4029.3
3	2005	3900
4	2006	3937
5	2007	3974.9
6	2008	4081.9
7		
8	Reported crime in Alaska	
9		
10	2004	3370.9
11	2005	3615
12	2006	3582

(Video) <http://vimeo.com/19185801>

Inference Engine

- Inference parameter sets based on users' input
 - Row, type and text selection
- Identify all compatible operations
- Rank the operations

An Example

Transform Script		Import	Export
<ul style="list-style-type: none"> Split data repeatedly on newline into rows Split split repeatedly on ' ' Promote row 0 to header Delete empty rows Extract from Year after 'In ' Set extract's name to State Fill State by copying values from above 			
Text	Columns	Rows	Table
Clear			
Delete rows where Year starts with 'Reported'			
Delete rows where Year contains 'Reported'			
Extract from Year between positions 0, 8			

	Year	State	#	Property
0	Reported crime in Alabama	Alabama		
1	2004	Alabama	4029.3	
2	2005	Alabama	3900	
3	2006	Alabama	3937	
4	2007	Alabama	3974.9	
5	2008	Alabama	4081.9	
6	Reported crime in Alaska	Alaska		
7	2004	Alaska	3370.9	
8	2005	Alaska	3615	
9	2006	Alaska	3582	
10	2007	Alaska	3373.9	
11	2008	Alaska	2928.3	
12	Reported crime in Arizona	Arizona		
13	2004	Arizona	5073.3	
14	2005	Arizona	4827	
15	2006	Arizona	4741.6	
16	2007	Arizona	4502.6	
17	2008	Arizona	4087.3	
18	Reported crime in Arkansas	Arkansas		
19	2004	Arkansas	4033.1	
20	2005	Arkansas	4058	

FlashFill in Excel

- Transforming by entering input-output examples

(<https://www.youtube.com/watch?v=aMdnbMQImVg>)

A Data Table

Accession	Credit	Dimensions	Medium	Name
01.2	Gift of the artist	5.25 in HIGH x 9.375 in WIDE	Oil on canvas	John Mix Stanley
05.411	Gift of James L. Edison	20 in HIGH x 24 in WIDE	Oil on canvas	Mortimer L. Smith
06.1	Gift of the artist	Image: 20.5 in. HIGH x 17.5 in. WIDE	Oil on canvas	Theodore Scott Dabo
06.2	Gift of the artist	9.75 in 16 in HIGH x 13.75 in 19.5 in WIDE	Oil on canvas	Leon Dabo
...				
09.8	Gift of the artist	12 in 14 in HIGH x 16 in 18 in WIDE	Oil on canvas	Gari Melchers

Programming by Example

	Raw Value	Target Value
R1	5.25 in HIGH x 9.375 in WIDE	9.375
R2	20 in HIGH x 24 in WIDE	24
R3	Image: 20.5 in. HIGH x 17.5 in. WIDE	17.5
R4	9.75 in 16 in HIGH x 13.75 in 19.5 in WIDE	13.75 13.75
...		
R5	12 in 14 in HIGH x 16 in 18 in WIDE	18 16

Transformation Program

Transform(value)

Conditional
Statement

```
label = classify(value)
```

```
switch label:
```

```
case "partition1":
```

Partition
Transformation
Program

```
pos1 = value.indexOf('BNK', 'NUM', -1)
```

```
pos2 = value.indexOf('NUM','BNK',2)
```

```
output = value.substring(pos1,pos2)
```

```
case "partition2":
```

Partition
Transformation
Program

```
pos1 = value.indexOf('BNK', 'NUM', 1)
```

```
pos2 = value.indexOf('NUM','BNK',3)
```

```
output = value.substring(pos1,pos2)
```

```
return output
```

BNK: blankspace
NUM[0-9]+: 98
UWRD[A-Z]: I
LWRD[a-z]: mae
WORD[a-zA-Z]
START:
END:

Creating Hypothesis Spaces

- Segmenting the Outputs ($n!$ traces)

1st I: 5.25 in HIGH x 9.375 in WIDE
O: 9.375

Constant

2nd 9.375

3rd 9.375

...

$O(N^2)$ unique substrings

- Create Hypothesis Space

Substring

1. subString(startPosition, endPosition)
2. Constant String

StartPostion

1. 15
2. (BNK, NUM, 1)
3. (BNK, 9, -1)
4. (LWRD BNK, NUM ' ', 1)

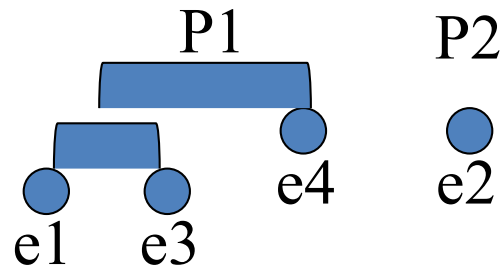
...

Generating Programs

- Generate programs from Hypothesis Space
 - Generate and then verify
 - Select values for non-terminals based on a partial order
 - Number of the Segments
 - Length of the context
 - Token class or actual text

Learn Conditional Statement

- Clustering:
 1. Agglomerative Clustering



2. Compatibility Score

- Thank you