



OpenRefine

'a free, open source, powerful tool for working with messy data'



Jeff Moon
Data Librarian &
Academic Director, Queen's Research Data Centre
Queen's University Library





What is OpenRefine?

- Formerly known as "Google Refine"
- Free & open-source
- Tool for cleaning large datasets
- Can create links between metadata sets
- Desktop application run locally

http://openrefine.org

openrefine.org





Home

Download

Documentation

Community

Post archive

Welcome!

OpenRefine (formerly Google Refine) is a powerful tool for working with messy data: cleaning it; transforming it from one format into another; extending it with web services; and linking it to databases like Freebase.

Please note that since October 2nd, 2012, Google is not actively supporting this project, which has now been rebranded to OpenRefine. Project development, documentation and promotion is now fully supported by volunteers. Find out more about the history of OpenRefine and how you can help the community.

Using OpenRefine - The Book

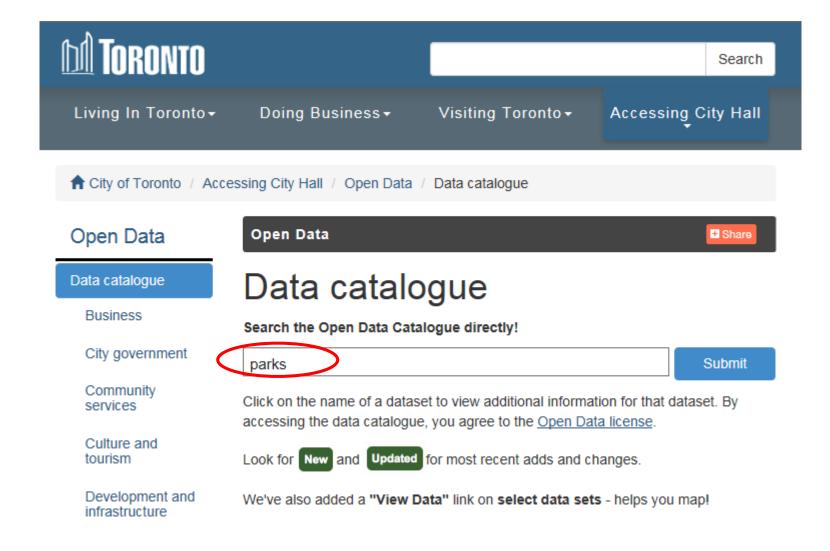
OpenRefine Core

Google Refine 2.5 - Stable version

- **Windows kit**, Download, unzip, and double-click on *google-refine.exe*. If you're having issues with the above, try double-clicking on *refine.bat* instead.
- Mac kit, Download, open, drag icon into the Applications folder and double click on it. NOTE: If you have issues installing Refine on Mac, please refer to issue 590
- Linux kit, Download, extract, then type ./refine to start.



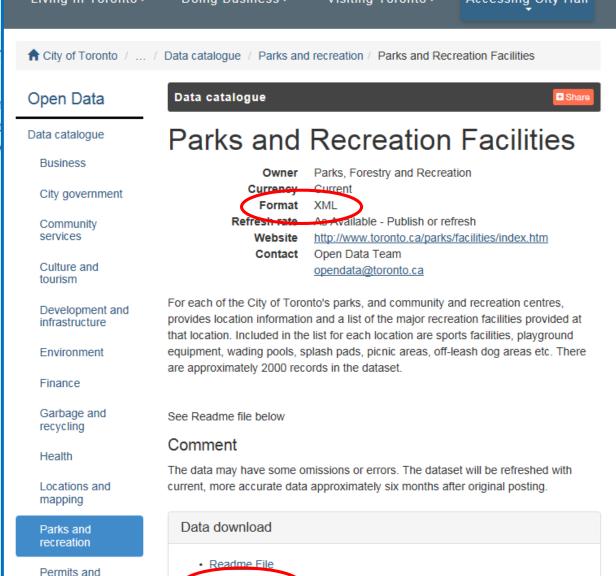
Real-world question...





Parks and Recreation F

... Live With Culture; Parks and Toronto Transit Commission; Tr www1.toronto.ca/wps/portal/content vgnextoid=d28b12f464151310VgnV



Facilities Data

licenses

```
http://www1.toronto.ca/City_Of_Toronto/Information_Technol 🔎 🔻 🖒
                                                               m www1.toronto.
      View Favorites Tools Help
<?xml version="1.0" encoding="UTF-8" standalone="true"?>
<Locations xmlns="http://www.example.org/PFRMapData">
   <Location>
       <LocationID > 1 < /LocationID >
       <LocationName>ASHBRIDGE'S BAY PARK//LocationName>
       <Address>Lake Shore Blvd Est</Address>
       <PostalCode>M4L 3W6</PostalCode>
       <Facilities>

    <Facility>

              <FacilityID>42376</FacilityID>
              <FacilityType>Skateboard Park</FacilityType>
              <FacilityName>Skatebord Park/FacilityName>
             <FacilityDisplayName>Skateboard Park</FacilityDisplayName>
           </Facility>
         <Facility>
              <FacilityID>529P1-Playground-0</FacilityID>
              <FacilityType>Playground</FacilityType>
```

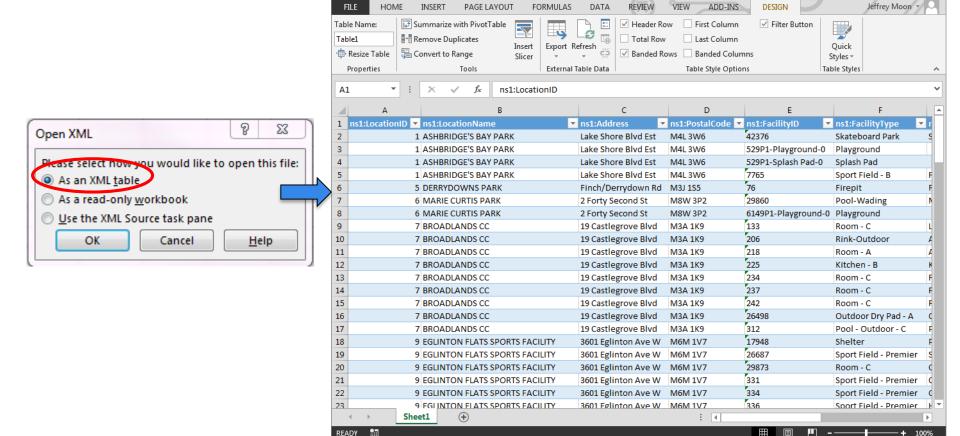
So, how do we get these data into a spreadsheet and manipulate them?



XI .



TABLE TOOLS



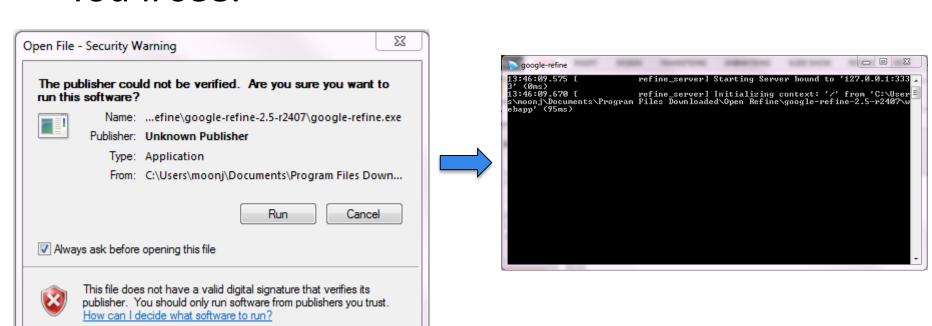
Book2 - Excel

But let's try OpenRefine



- Start the program

 Still called 'google-refine'
- You'll see:





Warning to use a 'modern browser'

close

Google Chrome Frame

Thank you for stopping by

We've wound down Chrome Frame and ceased support. Please read our previous announcement for more information. If you haven't already, consider installing and using a modern browser to get the best of the web.

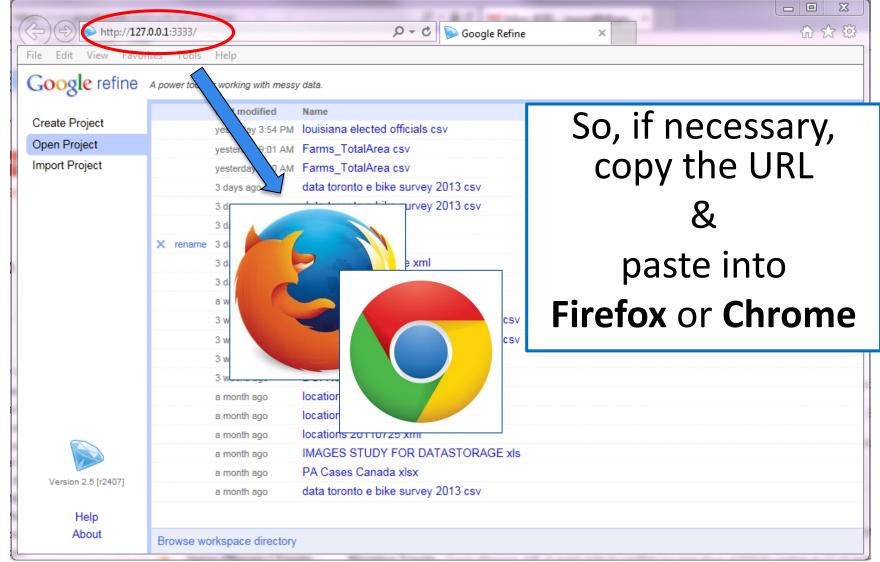
For more developer information please see the FAQ.

Sincerely,

The Google Chrome Frame team

On my PC, Internet Explorer is the default, BUT, IE is not the best choice for OpenRefine

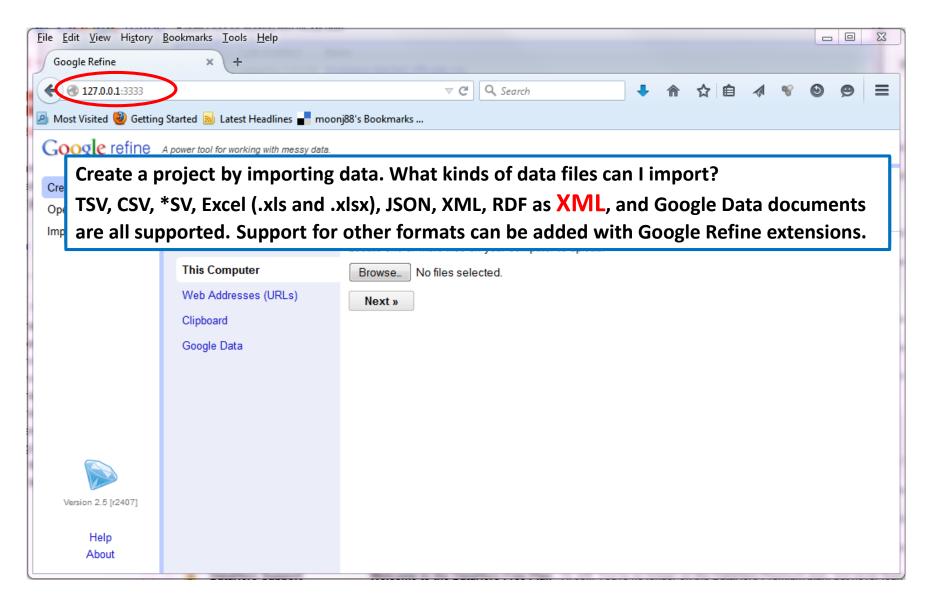






OpenRefine in Firefox

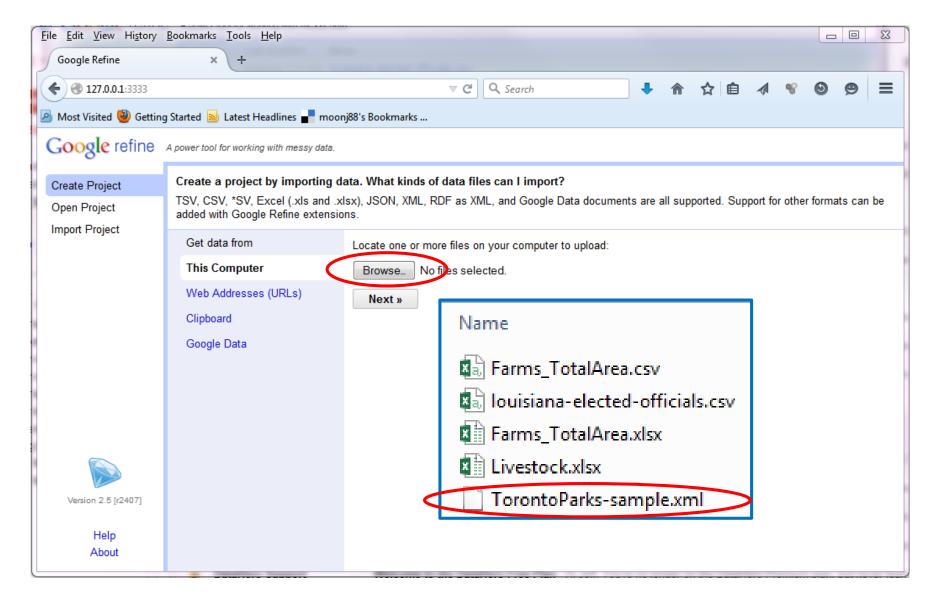






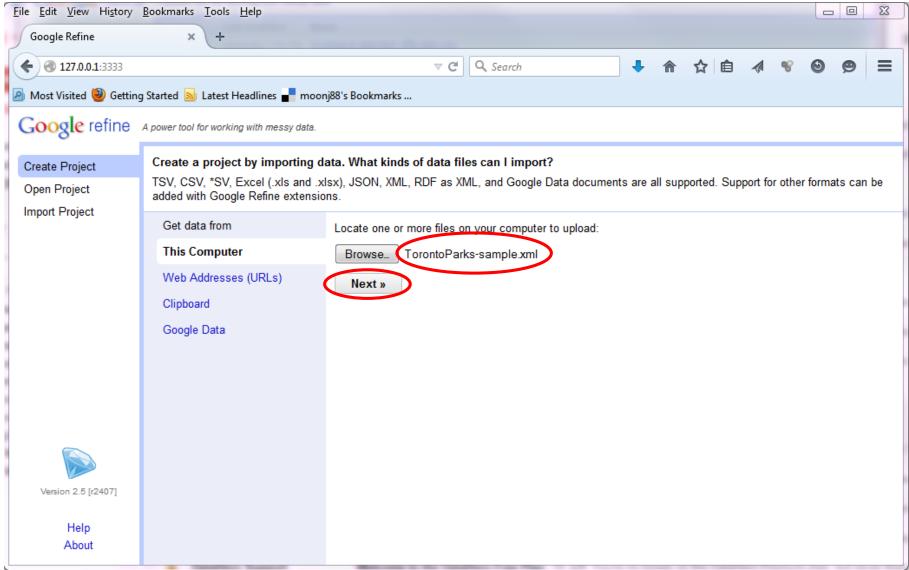
Browse to the XML file





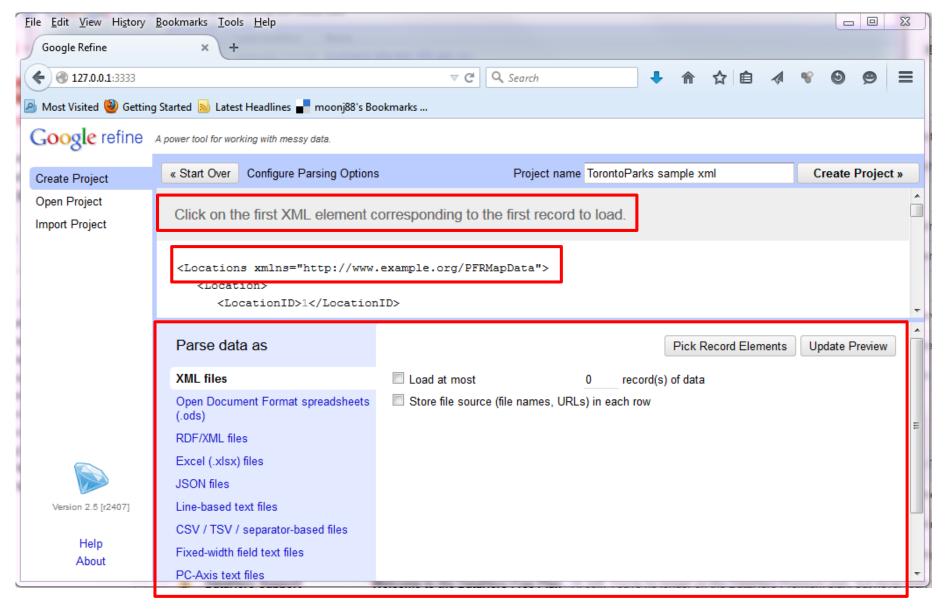
Click 'Next' when ready...





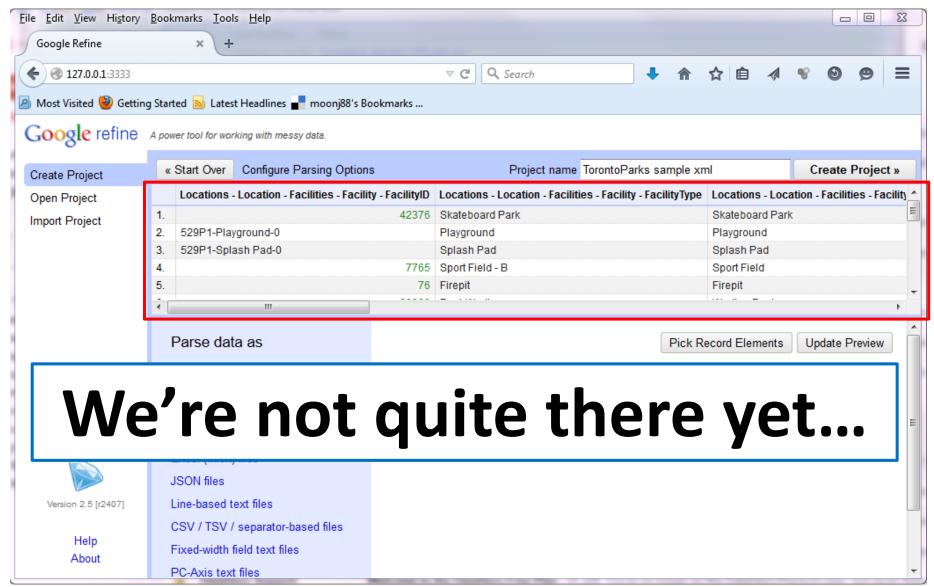
OpenRefine recognizes your file





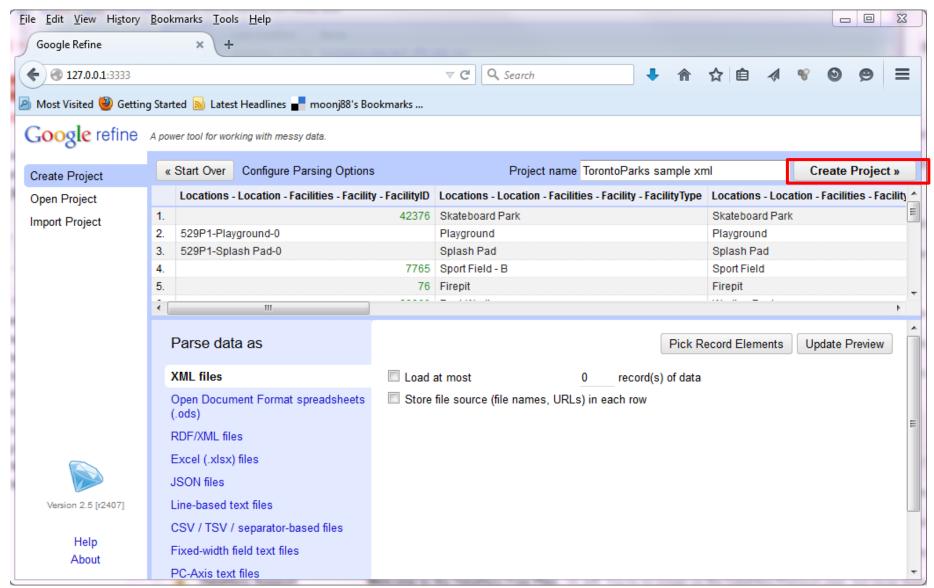
We get a preview of the data...





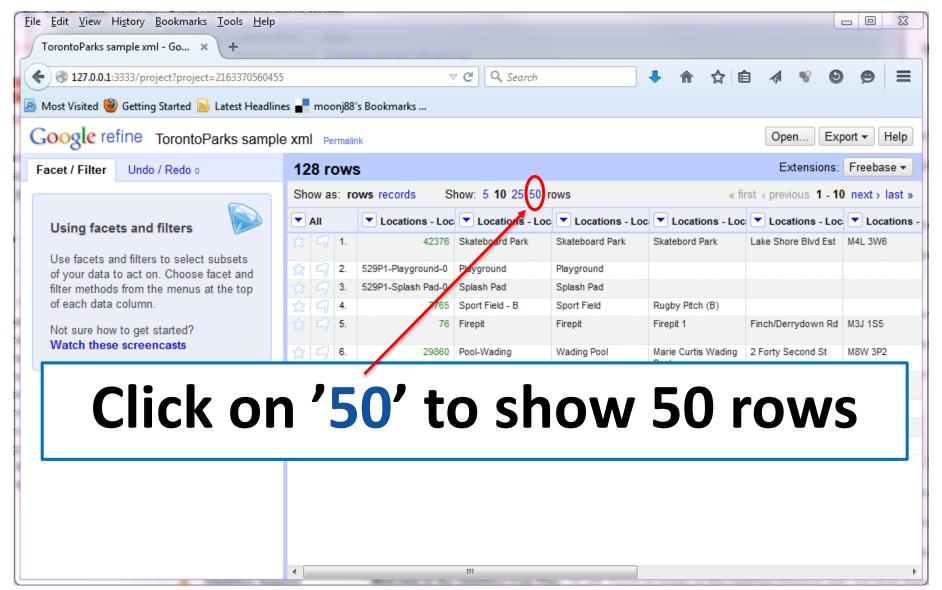
We're ready to 'create' our project...





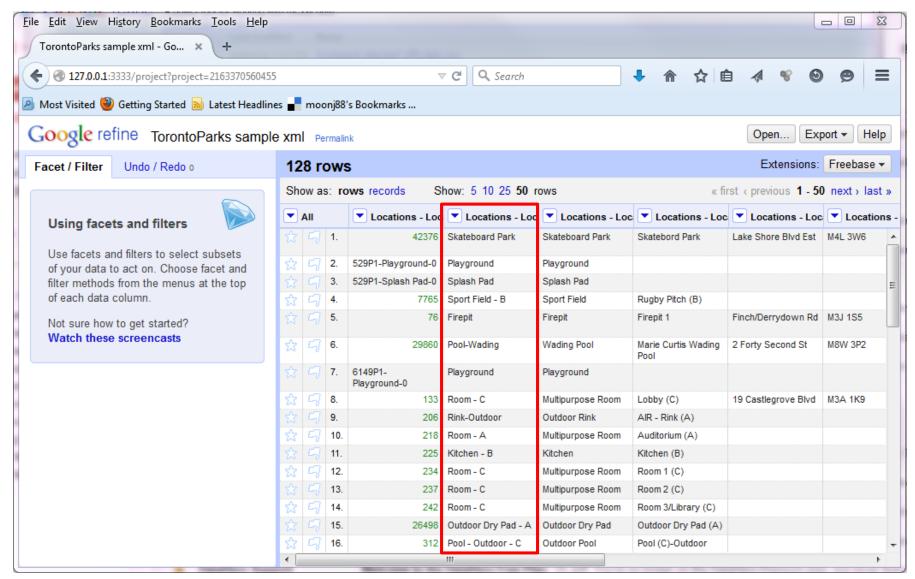
Data is now in OpenRefine





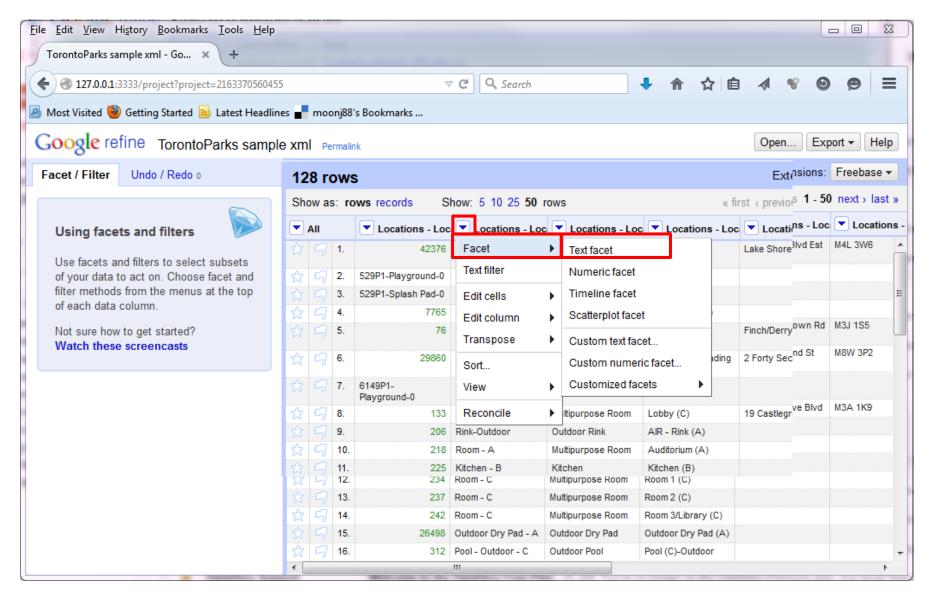
Text facets – let's look at the 2nd 'Location' column...





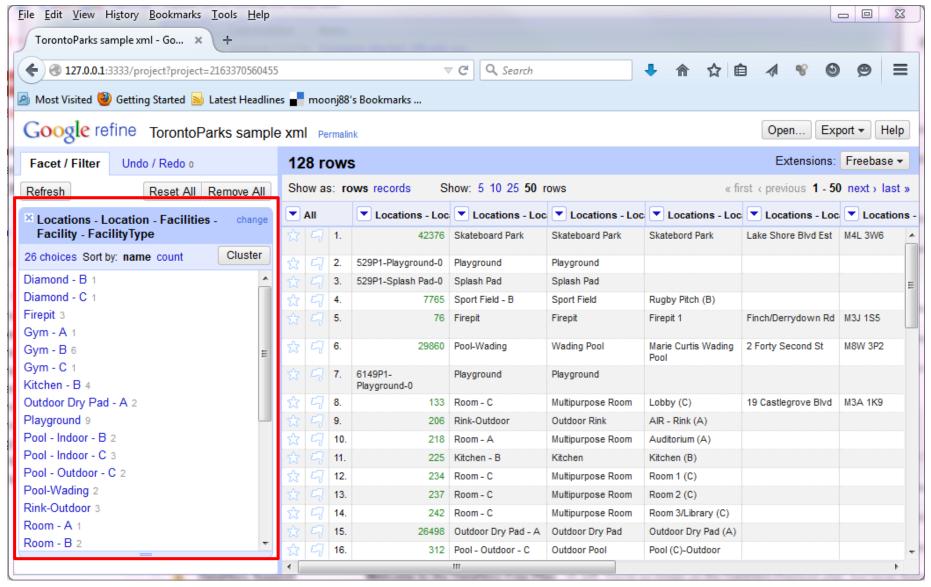
Dropdown: Facet \rightarrow **Text facet**



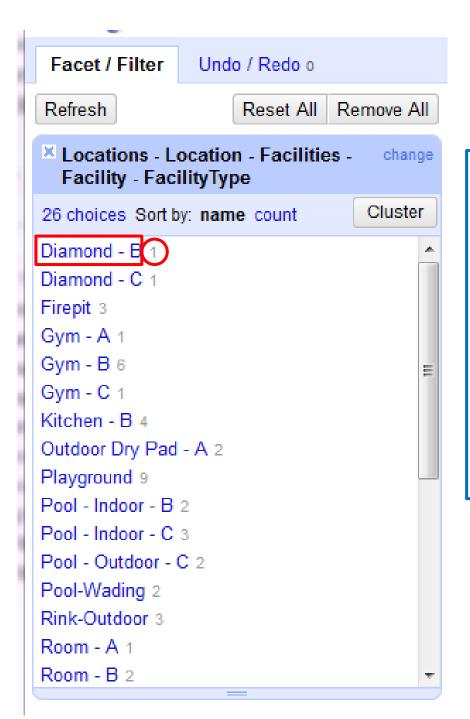


Facets show all unique entries, with frequencies











Unique values

Frequencies



Let's look at another example

Canadian Century Research Infrastructure
Infrastructure de recherche sur le Canada au 20e siècle



1911 Census Microdata

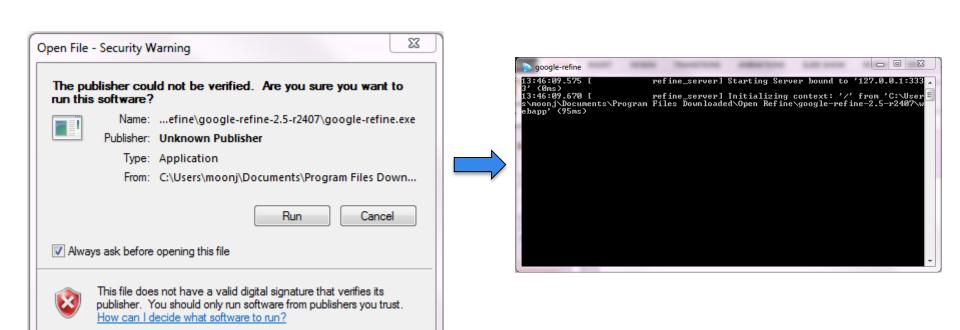
For this workshop I created a subset of ~100 records from each Province/Territory





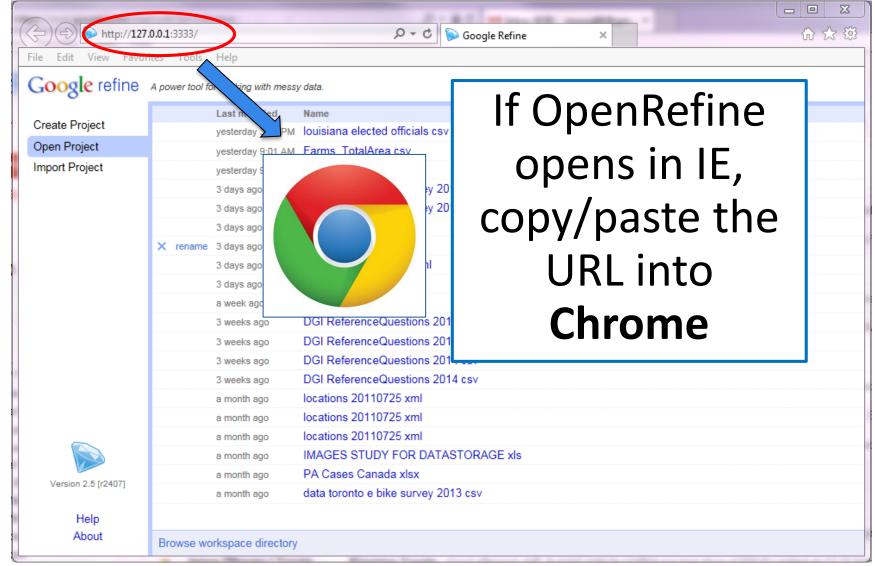
- Start the program

 Still called 'google-refine'
- You'll see:



If needed, copy the URL to Chrome...

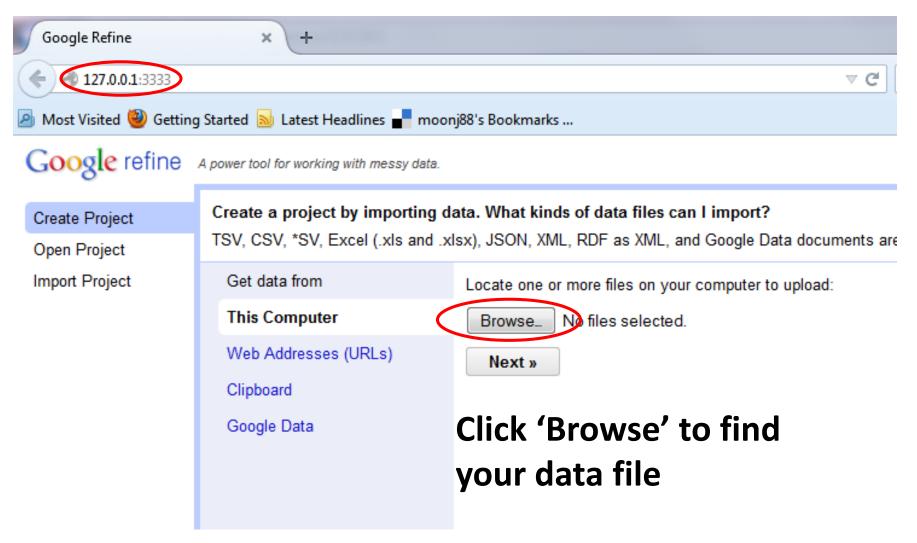




OpenRefine start screen



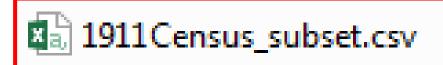
→ Click Browse



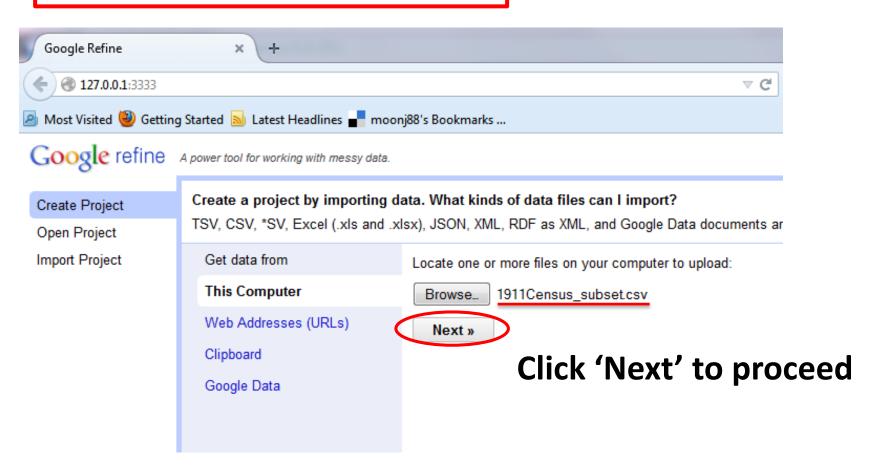
Navigate to the '.csv' file

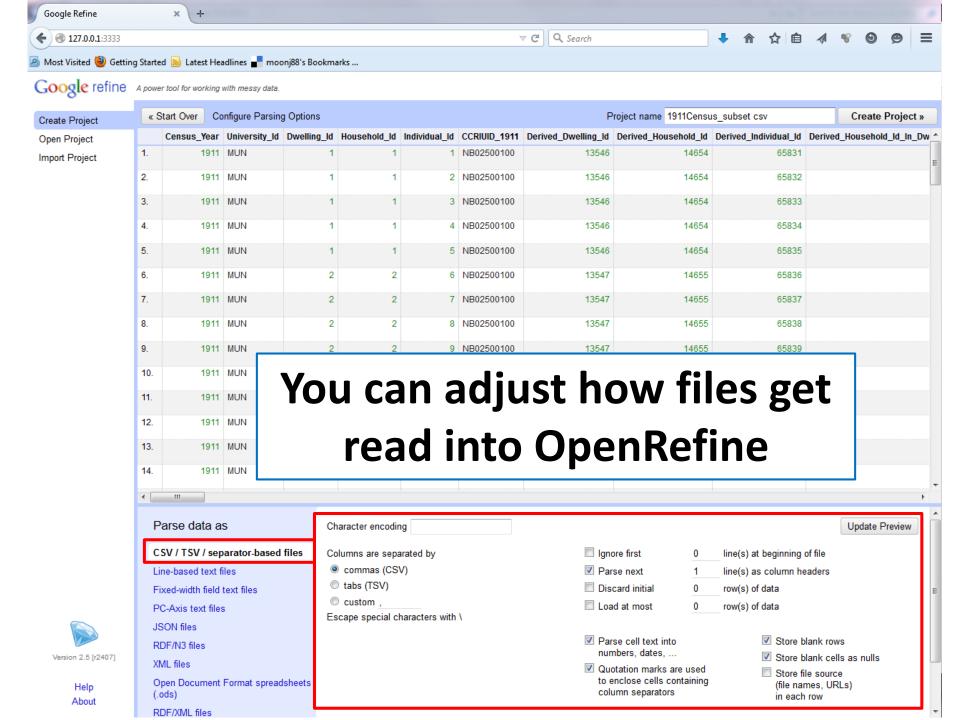


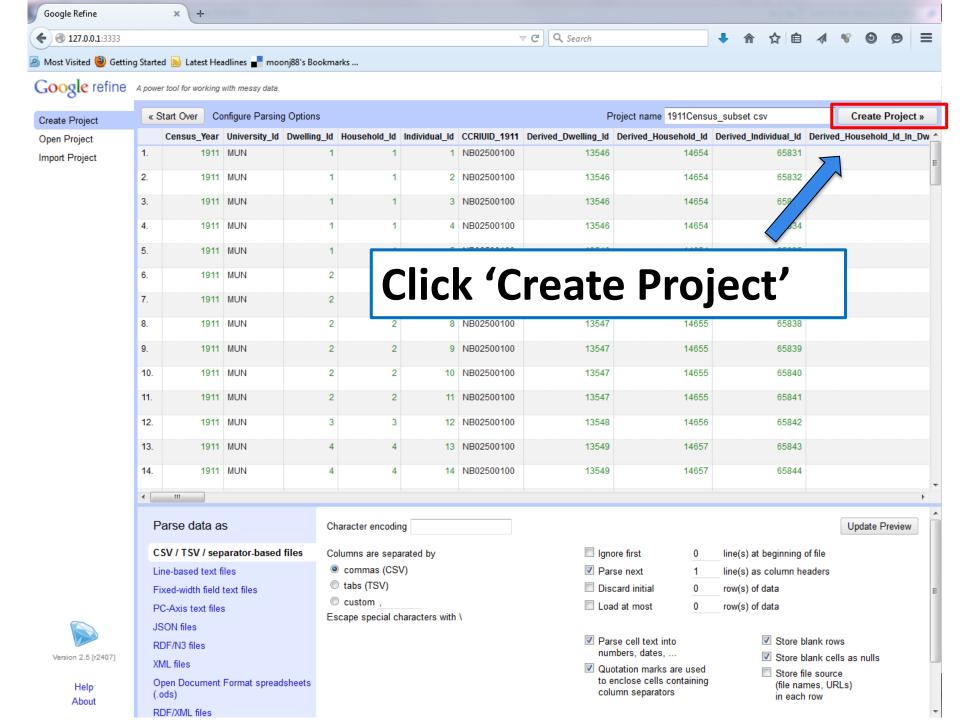
D:\pccf\1911Census_subset.csv

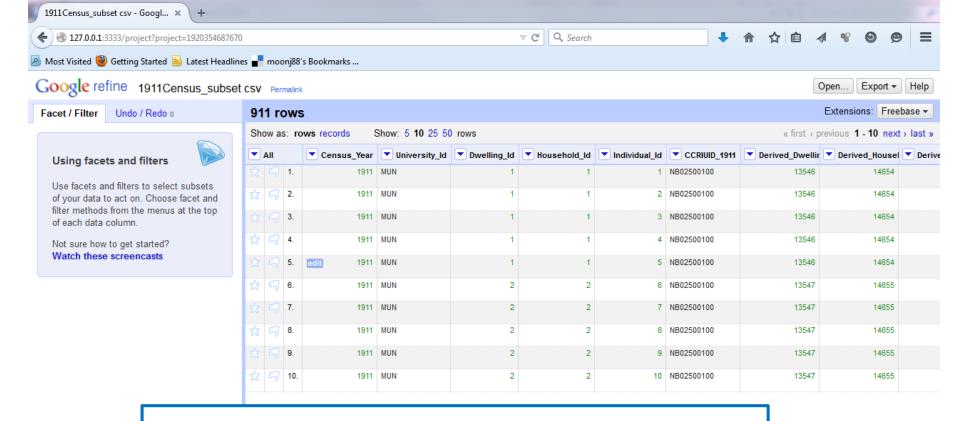


And click 'Open'

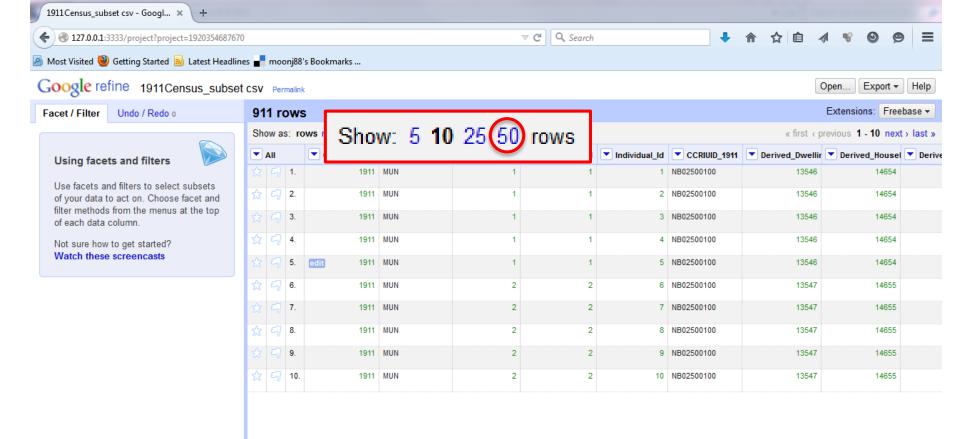




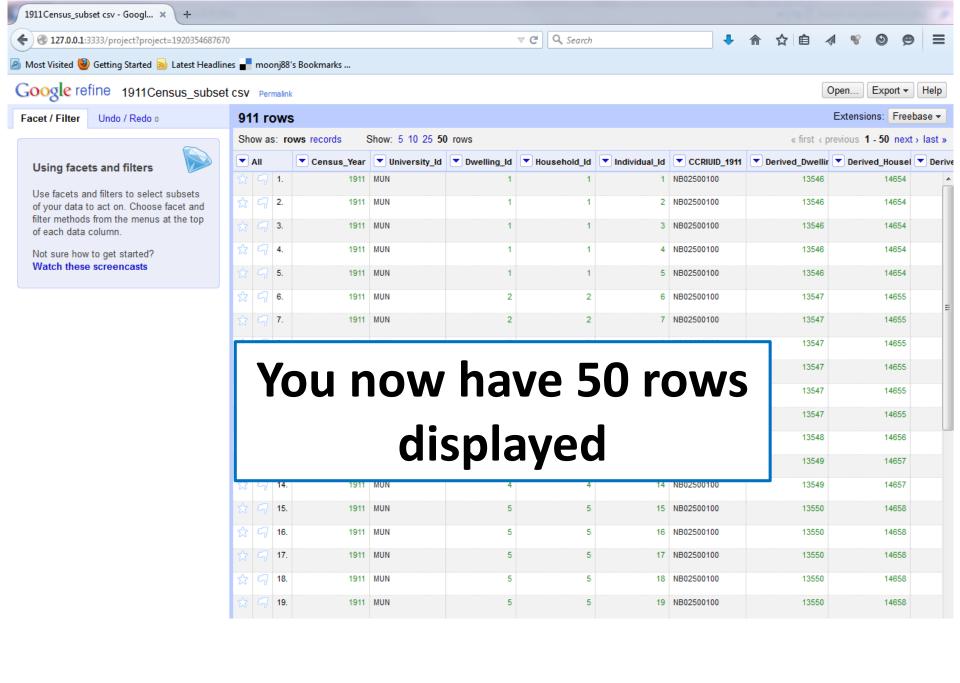




Your data is now open as a 'Project' in OpenRefine

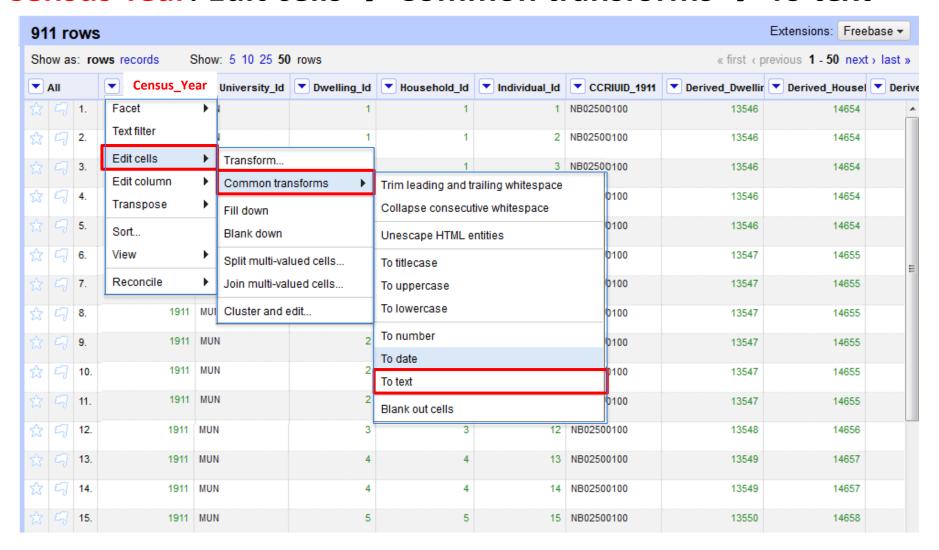


You can specify how many rows to show



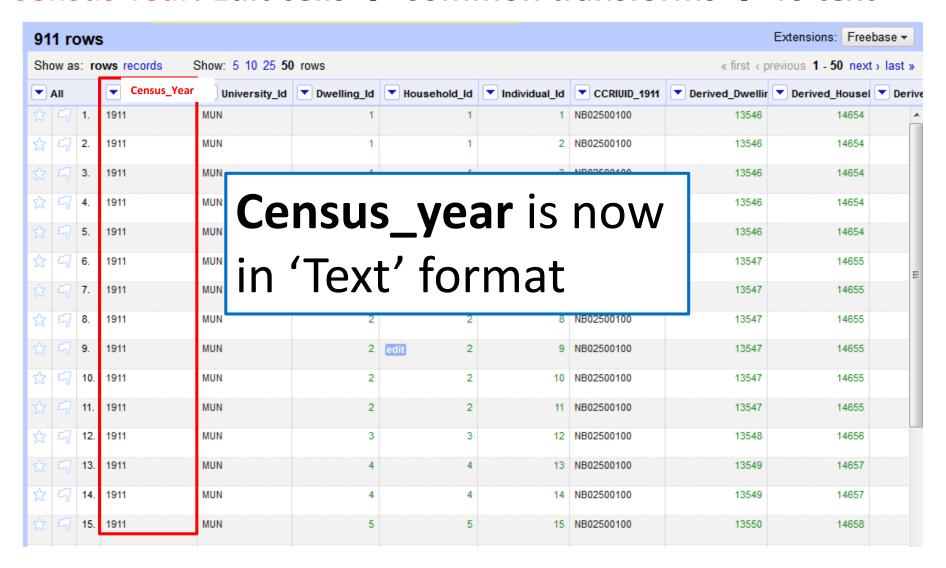
Exercise 1. Transform data

Census Year: Edit cells → Common transforms → To text



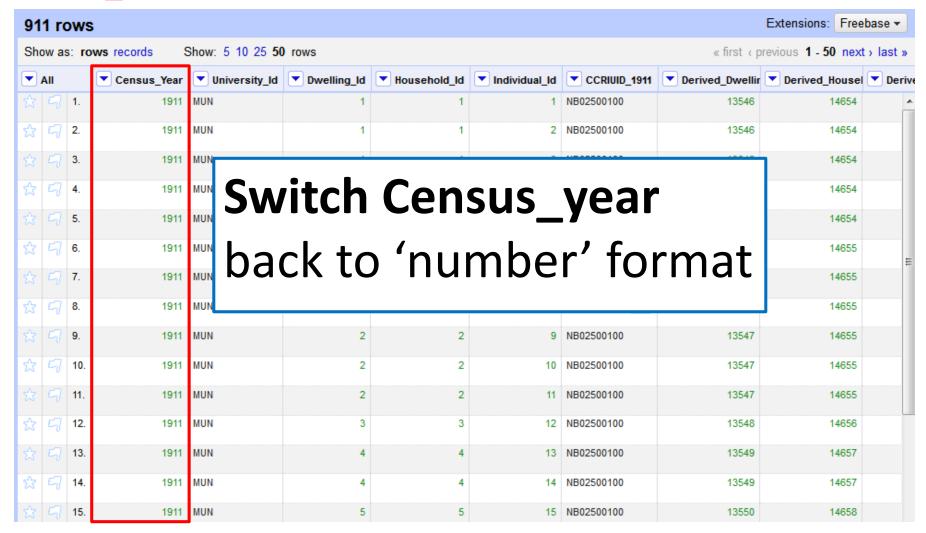
Exercise 1. Transform data

Census Year: Edit cells > Common transforms > To text



Exercise 2. Transform back to number

Census_Year: Edit cells → Common transforms → To number



Exercise 3. Text facet First, SCROLL to Province column



how as	: rows records		Show: 5 10	25 50 rows				« first	c previous 1 - 50	next > last »
n_Size	Dwelling_Unit_1	▼	Province	District_Numbe	Sub_District_Nu	Census_Form_	▼ District_Name	Sub_District_Na	Enumerator_Fir	Enumer
U	U	NB		25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
UI	U	NB	1	25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
U	U	NB		25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
UI	U	NB		25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
UI	U	NB		25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
UI	U	NB		25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
UI	U	NB		25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
UI	U	NB		25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
UI	U	NB	edit	25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
UI	U	NB		25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
UI	U	NB		25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
UI	U	NB		25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
UI	U	NB		25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
U	U	NB		25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
UI	U	NB		25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON

Exercise 3. Text facet

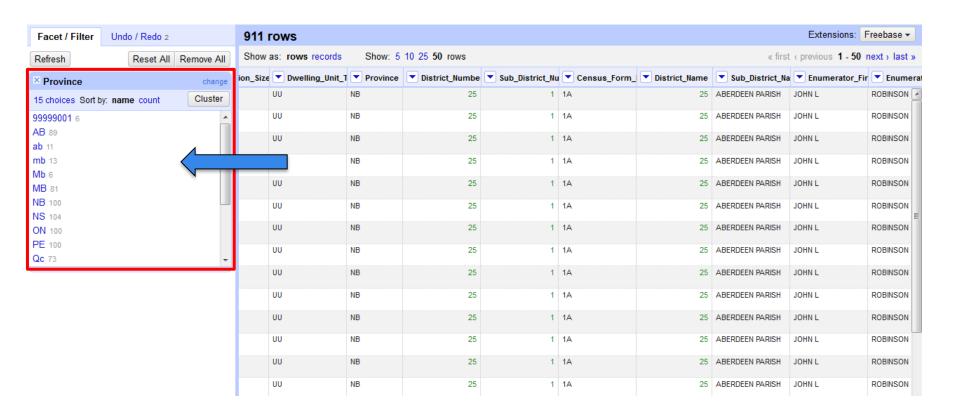
Province: Facet → Text facet



911 rows						Extensions:	reebase 🕶						
Show as: rows records Show: 5 10 25 50 rows « first < previous 1 - 50 next > last »													
ion_Size 💌 Dwelling_Un	t_1 Province	District_Numbe Sub_District_Nu	Census_Form_	▼ District_Name	Sub_District_Na	Enumerator_Fir	Enumera						
UU	Facet	Text facet	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON A						
UU	Text filter	Numeric facet	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON						
UU	Edit cells	Timeline facet	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON						
UU	Edit column	Scatterplot facet	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON						
UU	Transpose Sort	Custom text facet Custom numeric facet	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON						
UU	View	Customized facets	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON						
UU	Reconcile) 25 1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON						
UU	NB	25 1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON						
UU	NB edit	25 1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON						
UU	NB	25 1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON						
UU	NB	25 1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON						
UU	NB	25 1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON						
UU	NB	25 1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON						
UU	NB	25 1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON						
UU	NB	25 1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON						

Exercise 3. Facets Facets appear at the left of the screen

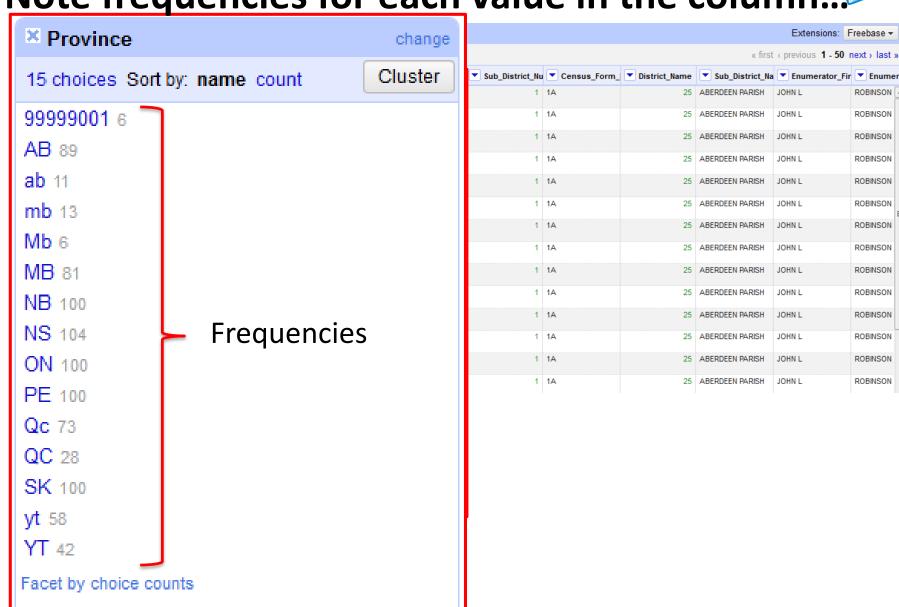




Exercise 3. Facets

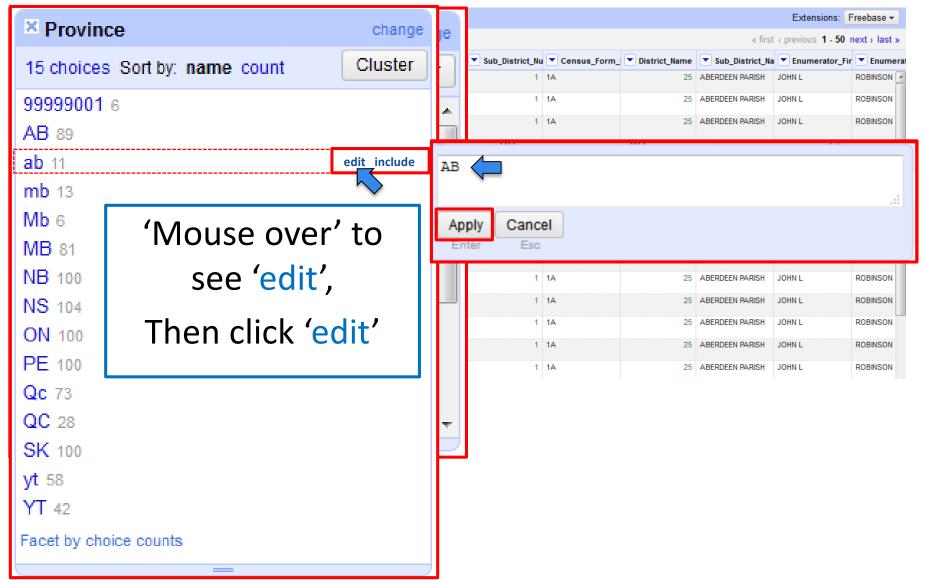
Note frequencies for each value in the column.

ROBINSON



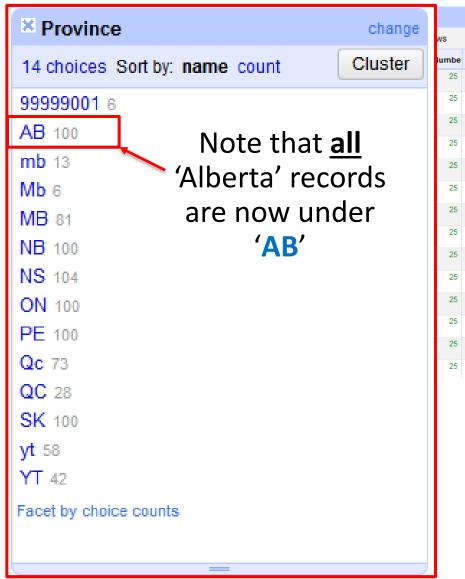
Exercise 4. EDIT entry for Alberta [ab]

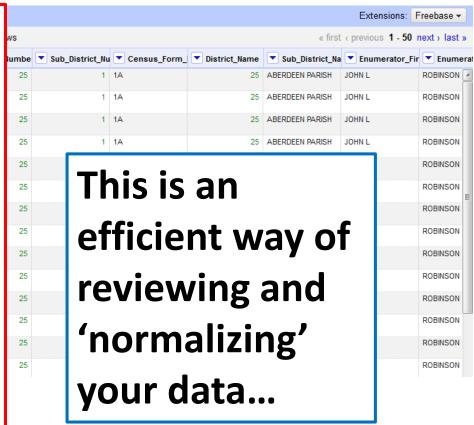




Exercise 4. EDIT entry for Alberta [ab]







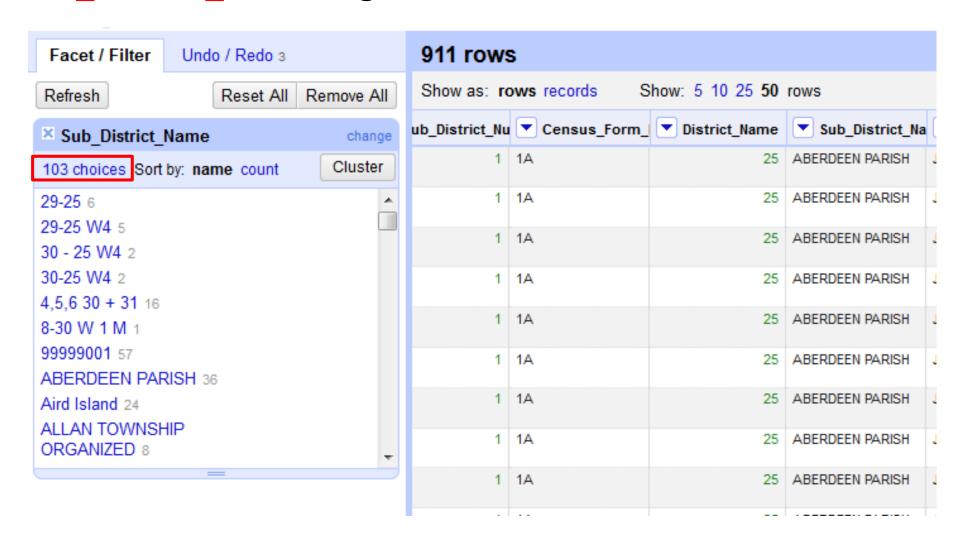
Exercise 5. Text facets & Clustering Sub_District_Name: Facet → Text facet



now a	s: rows records	Show: 5	10 25 50 rows				« first	c previous 1 - 50	next > last
n_Size	Dwelling_Unit_1	Province	▼ District_Numbe	Sub_District_Nu	Census_Form_	▼ District_Name	Sub_District_Na	▼ Enumerator_Fi	r 💌 Enume
ı	UU	NB	25	1	1A	25	Facet)	Text facet	
	UU	NB	25	1	1A	25	Text filter	Numeric facet	
	UU	NB	25	1	1A	25	Edit cells	Timeline facet	
							Edit column	Scatterplot facet	
	UU	NB	25	1	1A	25	Transpose	Custom text facet	Ĺ
1	UU	NB	25	1	1A	25	Sort	Custom numeric	facet
	υυ	NB	25	1	1A	25	View	Customized face	ts 🕨
	υu	NB	25	1	1A	25	Reconcile	HN L I	ROBINSON
	υυ	NB	25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
	UU	NB edit	25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
	υu	NB	25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
I	UU	NB	25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
1	UU	NB	25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
1	UU	NB	25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSON
I	UU	NB	25	1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSOI
	UU	NB	25	-1	1A	25	ABERDEEN PARISH	JOHN L	ROBINSO

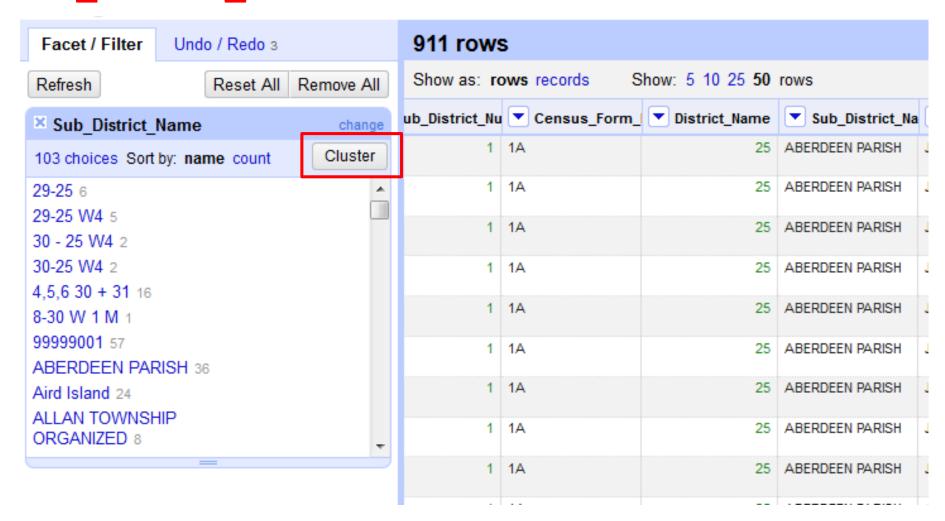


Sub_District_Name: Again, facets shown on left side



Exercise 5. Text facets & Clustering **Sub_District_Name:** Click on Cluster



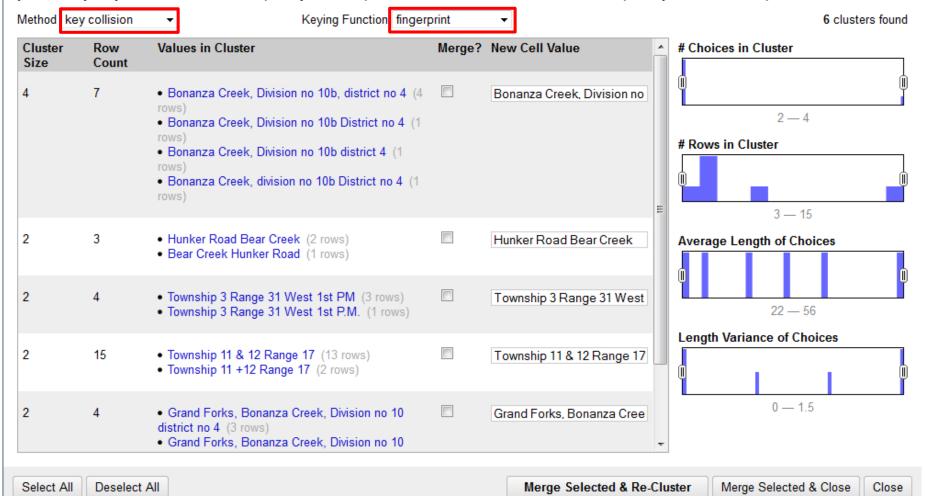


Exercise 5. Text facets & Clustering Note different 'methods' and 'keying



Cluster & Edit column "Sub_District_Name"

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 ...



Exercise 5. Text facets & Clustering Note suggested 'clusters'...



Cluster & Edit column "Sub_District_Name"

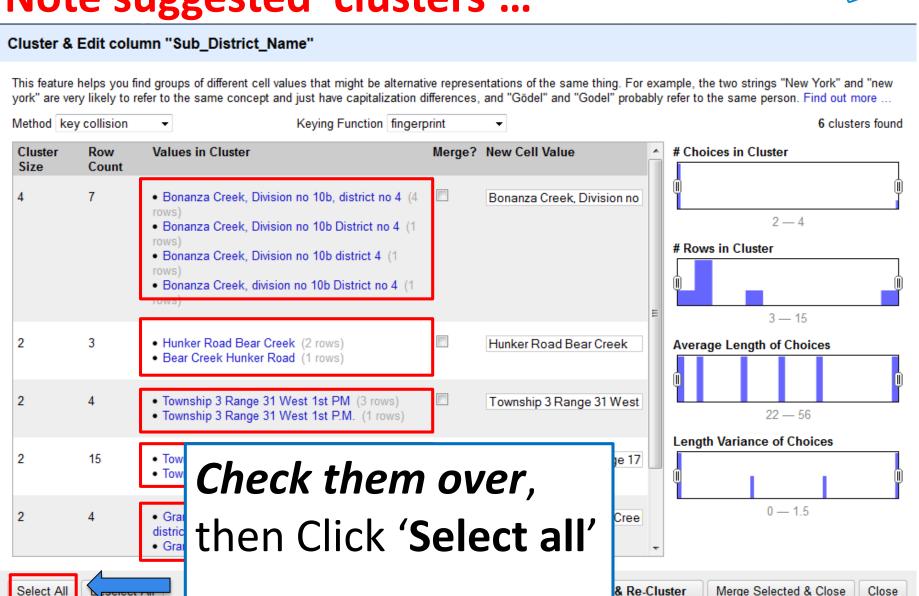
This feature helps you find groups of different cell values that might be alternative representations of the same thing. For example, the two strings "New York" and "new



Exercise 5. Text facets & Clustering Note suggested 'clusters'...



Close



Exercise 5. Text facets & Clustering Now we're ready to 'Re-Cluster'



Cluster & Edit column "Sub_District_Name" 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 6 clusters found Merge? New Cell Value Cluster Row Values in Cluster # Choices in Cluster Size Count . Bonanza Creek, Division no 10b, district no 4 (4) Bonanza Creek, Division no 2 - 4. Bonanza Creek, Division no 10b District no 4 (1 # Rows in Cluster . Bonanza Creek, Division no 10b district 4 (1 . Bonanza Creek, division no 10b District no 4 (1 3 - 15Click Average Length of Choices eek 'Merge Selected & Re-Cluster West 22 - 56Length Variance of Choices 1 2 15 Township 11 & 12 Range 17 (13 rows) Township 11 & 12 Range 17 Township 11 +12 Range 17 (2 rows) 0 - 1.5. Grand Forks, Bonanza Creek, Division no 10 Grand Forks, Bonanza Cree

Select All Deselect All

district no 4 (3 rows)

. Grand Forks, Bonanza Creek, Division no 10

Merge Selected & Re-Cluster

Merge Selected & Close

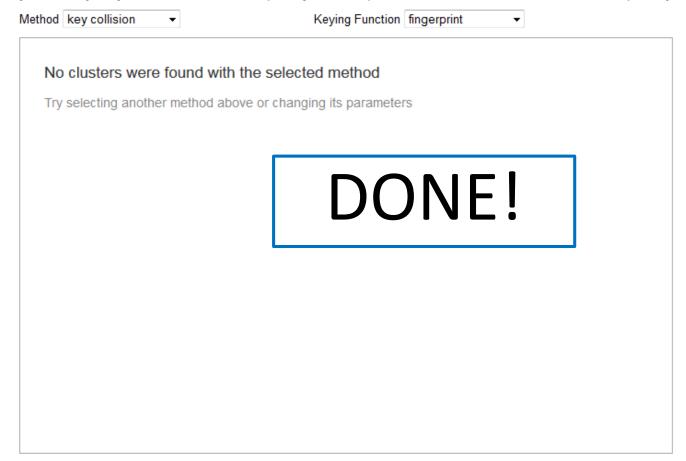
Close



'fingerprint' clustering done!

Cluster & Edit column "Sub_District_Name"

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 ...

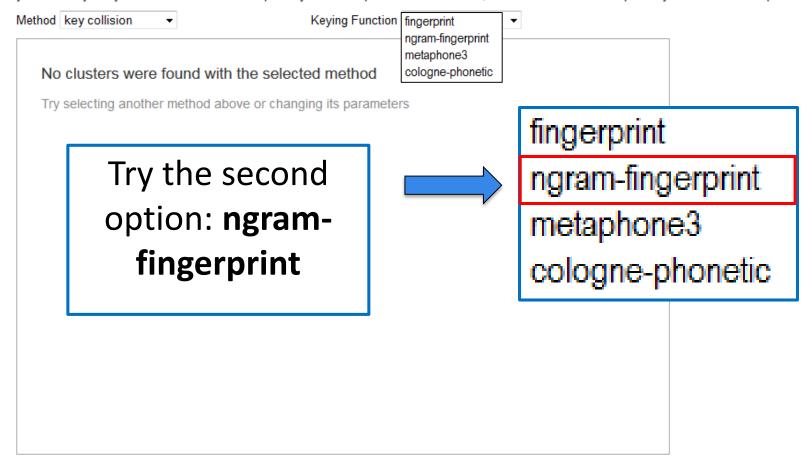




Now try the next function, 'ngram-fingerprint'

Cluster & Edit column "Sub_District_Name"

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 ...







Cluster & Edit column "Sub_District_Name"

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 ...

Ngram Size 2 Method key collision Keying Function | ngram-fingerprint | ▼ 2 clusters found Cluster Size Row Count Values in Cluster **New Cell Value** # Rows in Cluster Merge? 2 • Township 11+12 Range 17 (21 rows) 36 Township 11+12 Range 17 Township 11 & 12 Range 17 (15 rows) 4 - 36 30 - 25 W4 (2 rows) 30 - 25 W4 Average Length of Choices 30-25 W4 (2 rows) 9 - 24Review, Two more Select All, 'clusters' then Click Merge Selected & Re-Cluster

Exercise 5. Text facets & Clustering 'ngram-fingerprint' clustering done!



Cluster & Edit column "Sub_District_Name"

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 ...

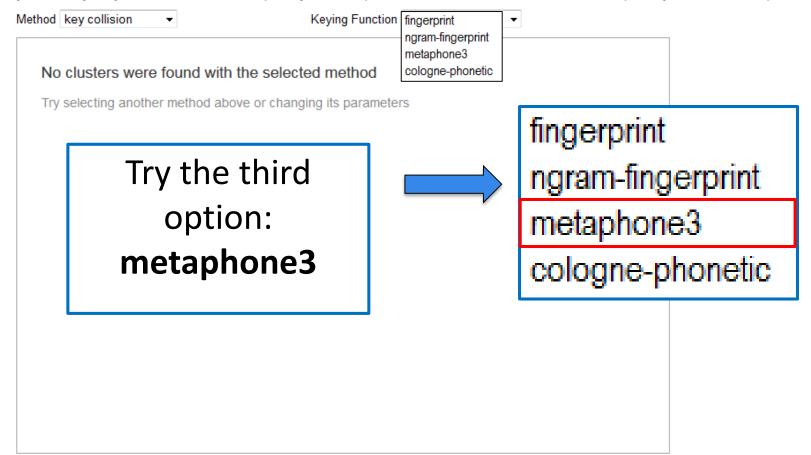
Keying Function | ngram-fingerprint ▼ Method key collision No clusters were found with the selected method Try selecting another method above or changing its parameters DONE!



Now try the next keying function, 'metaphone3'

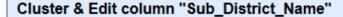
Cluster & Edit column "Sub_District_Name"

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 ...

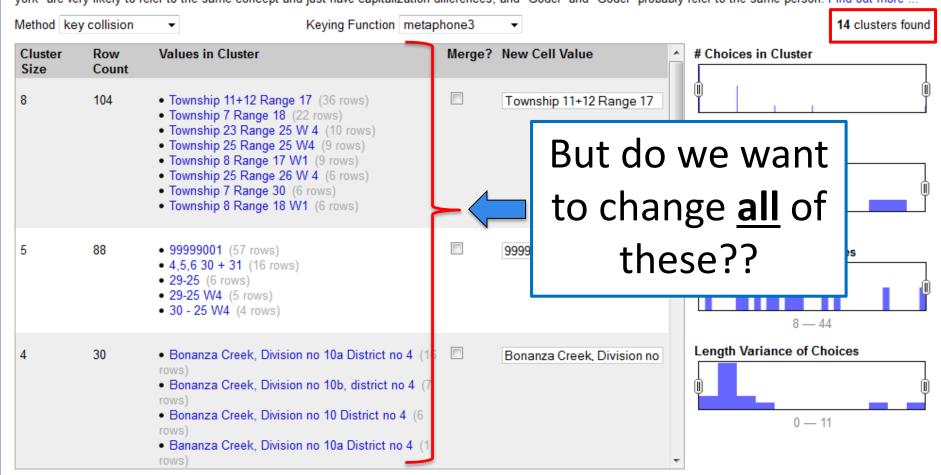




We get a lot of new clusters...



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 ...



Select All Deselect All

Merge Selected & Re-Cluster

Merge Selected & Close

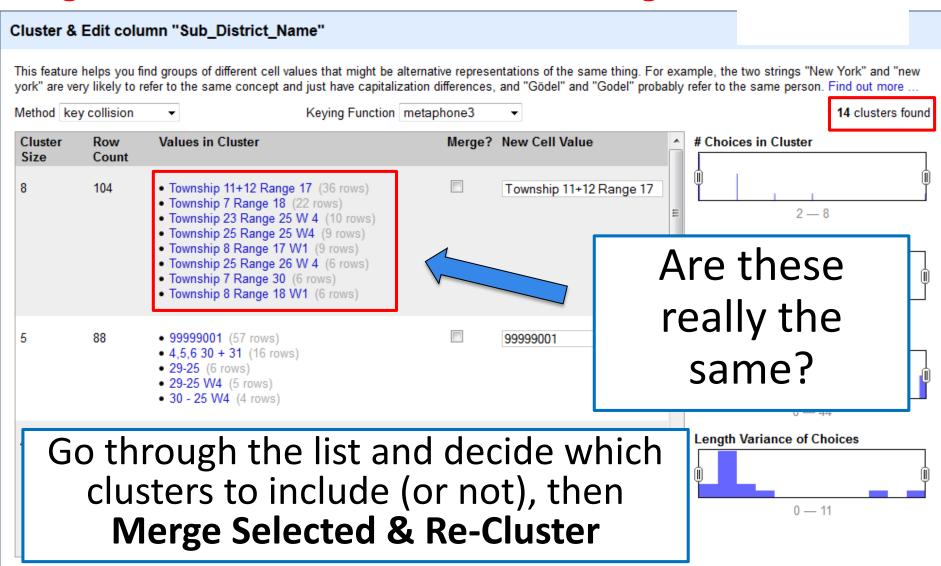
Close

Deselect All

Select All



We get a lot of new clusters... scroll through the list



Merge Selected & Re-Cluster

Merge Selected & Close

Close

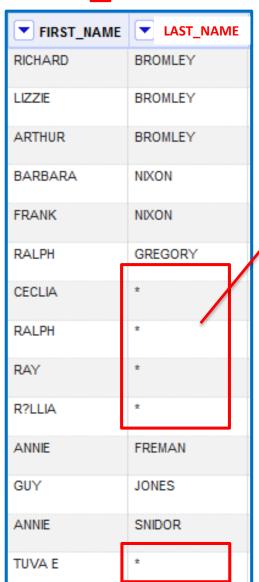
Exercise 5. Text facets & Clustering Clustering... a powerful clean-up tool

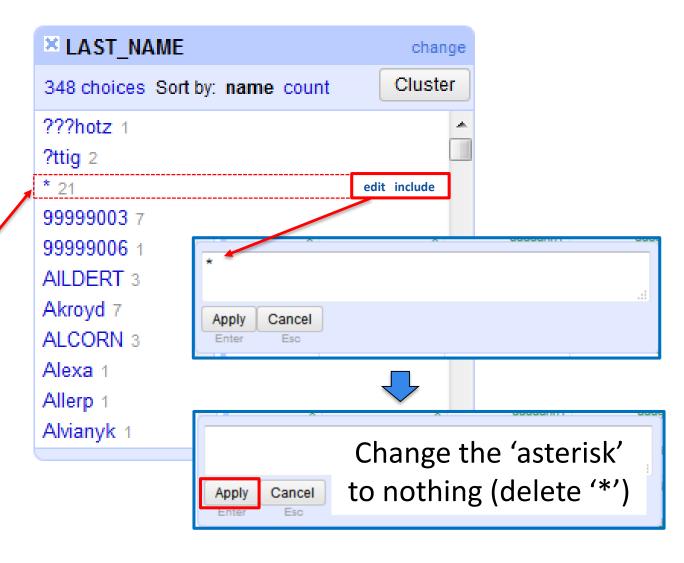


You can use, and re-use clustering techniques until your data is 'clean'

Exercise 6. Editing data in columns LAST_NAME: Facet \rightarrow Text facet

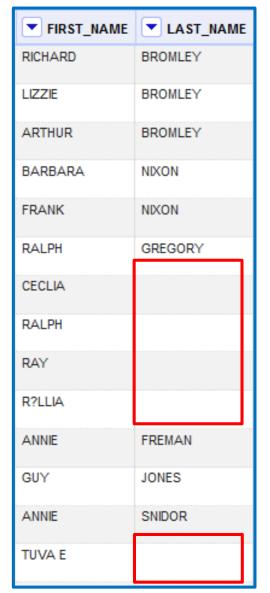


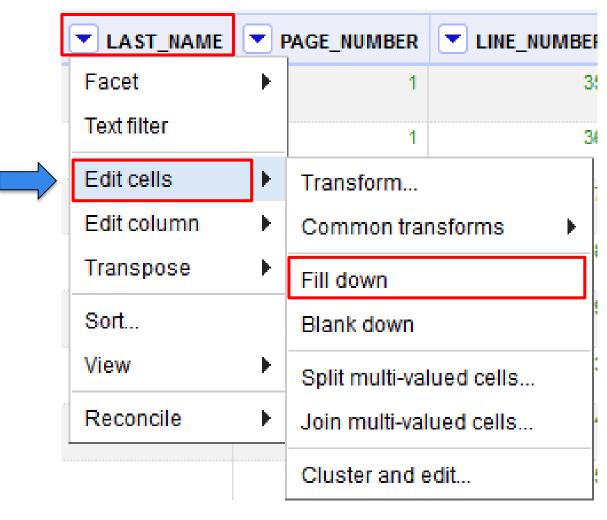




Exercise 6. Editing data in columns LAST_NAME: Edit cells -> Fill down



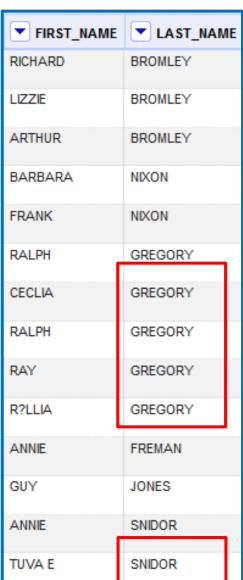




Exercise 6. Editing data in columns LAST_NAME: Edit cells → Fill down



FIRST_NAME	LAST_NAME
RICHARD	BROMLEY
LIZZIE	BROMLEY
ARTHUR	BROMLEY
BARBARA	NIXON
FRANK	NIXON
RALPH	GREGORY
CECLIA	
RALPH	
RAY	
R?LLIA	
ANNIE	FREMAN
GUY	JONES
ANNIE	SNIDOR
TUVA E	

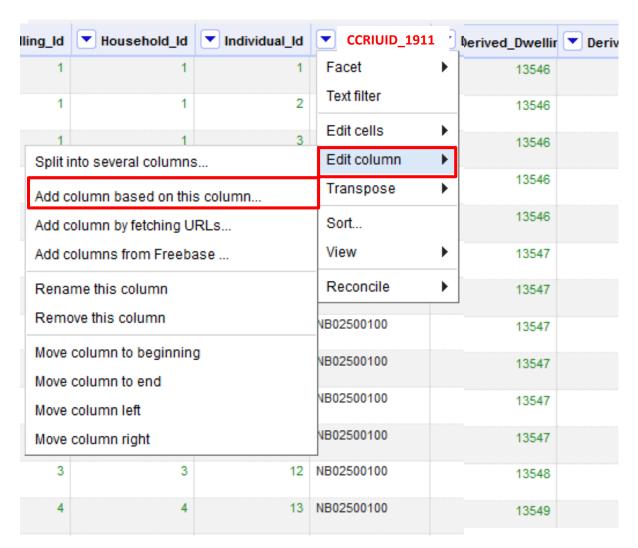


Not so easy to do in Excel...

Exercise 7. Creating new columns CCRIUID_1911: Edit column

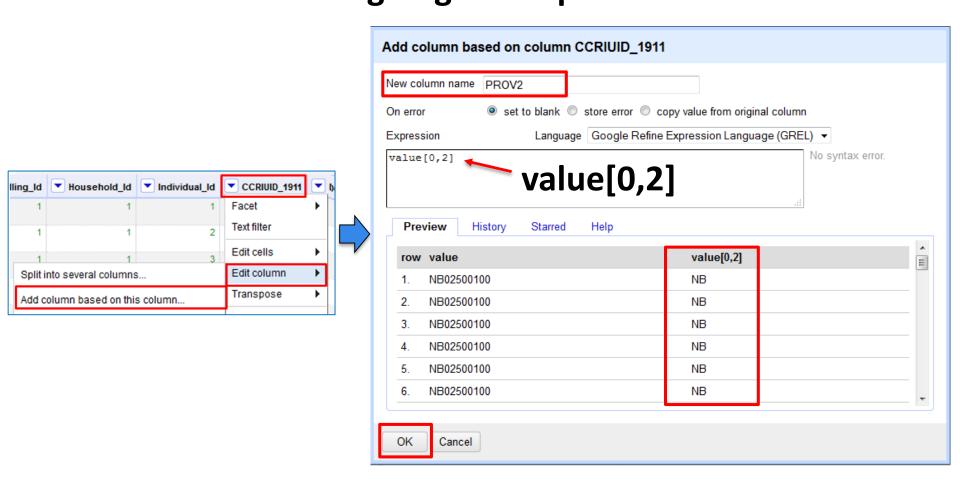


→ Add column based on this column



Exercise 7. Creating new columns CCRIUID_1911: Define values for new column using 'regular expression'





A very powerful tool!



Create new column PROV2 based on column CCRIUID_1911 by filling 911 rows with grel:value[0,2] Undo



'grel'

OpenRefine [formerly Google]
Regular Expression Language

'A regular expression is a string that describes a text pattern occurring in other strings.'

CCRIUID_1911	▼ PROV2
NB02500100	NB

WHENEVER I LEARN A
NEW SKILL I CONCOCT
ELABORATE FANTASY
SCENARIOS WHERE IT
LETS ME SAVE THE DAY.

OH NO! THE KILLER MUST HAVE ROLLOWED HER ON VACATION!



BUT TO FIND THEM WE'D HAVE TO SEARCH THROUGH 200 MB OF EMAILS LOOKING FOR SOMETHING FORMATTED LIKE AN ADDRESS!



IT'S HOPELESS!

















In Transformations:

Edit cells > Transform:

- value.toNumber()
- value.toString()
- value.replace("+", "")
- value.replace("~", "").replace(",","")
- value.unescape('url')
- value.toUppercase()

← change to 'number'

← change to 'string'

← replace one thing

← replace more than 1 thing

← remove 'odd' characters

← change to Uppercase

In Facets

Find entries based on what they contain:

Facet → Custom text facet:

value.contains("million")

← Facet based on presence/absence of a "string": returns 'True' or 'False'

value.replace("F", "f").replace("a","A")



Custom text transf	form on column OCCUPATION_CHIE	F_OCC_IND_CL
Expression	Language Google Refine Expression	Language (GREL) ▼
value.replace("F	No syntax en	
Preview Histo	ry Starred Help	
row value	value.replace("F", "f").replace("a","A")	Â
1. Farmer	fArmer	
2. None	None	=
3. null	Error: replace expects 3 strings, or 1 string,	regex, and 1 string
4. None	None	
5. None	None	
6. Farmer	fArmer	Ţ
	blank	mes until no change
3.1		

value.contains("Farmer")



Cı	ustor	n Facet on c	olumn OCCUF	PATION_CHIEF_OCC_I	ND_CL		
E	Language (GREL) ▼						
V	alue	.contains("F	No syntax error.				
ļ	Prev	view Histor	y Starred	Help			
	row	value	value.co	ntains("Farmer")	ŕ		
	1.	Farmer	true				
	2.	None	false				
	3.	null	null				
	4.	None	false				
	5.	None	false				
	6.	Farmer	true				
OK Cancel							

Source: http://enipedia.tudelft.nl/wiki/OpenRefine_Tutorial

There is so much more... Edit cells → Common transforms



Trim leading and trailing whitespace	Just what it says
Collapse consecutive whitespace	Collapses embedded >1 whitespace
Unescape HTML entities	Removes HTML or XML character references and entities from a text string
To titlecase]]
To uppercase]
To lowercase	lust what they say
To number	Just what they say
To date	
To text	

There is so much more...

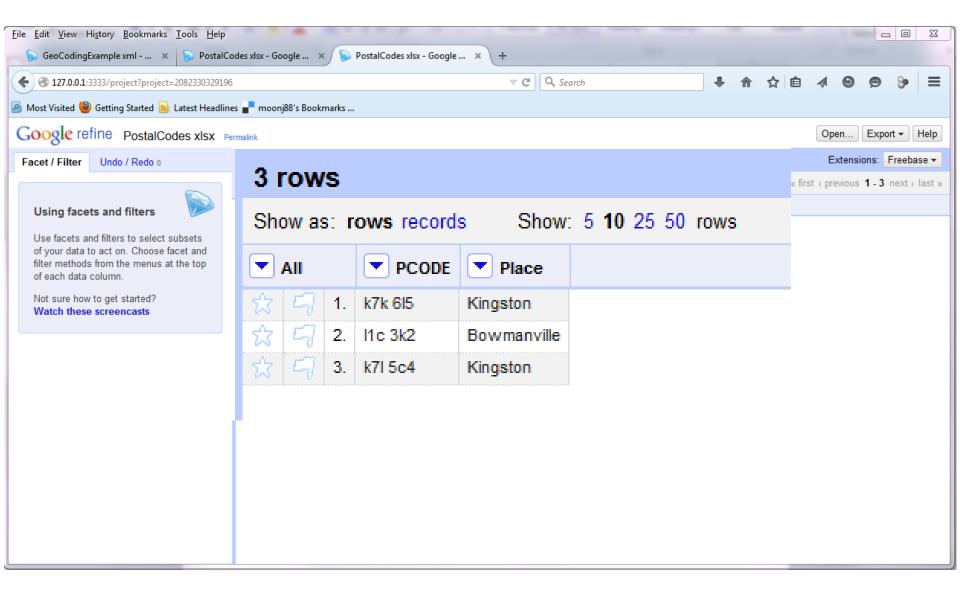


Use Google Regular Expression
Language (GREL) to query external
data sources

For example, using Postal Codes to extract information from Google Maps

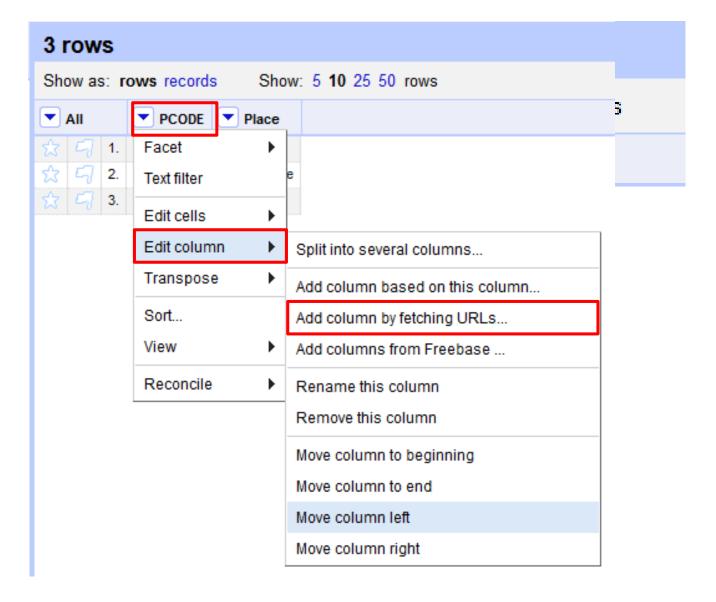
Here's a very simple example...





Use PCODE to retrieve external data...



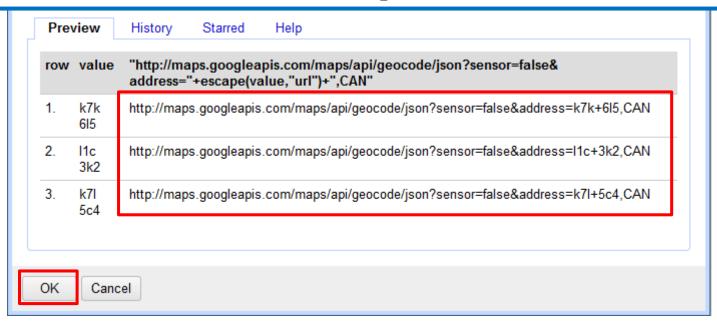


Add GREL: Google Regular Expression Language



Add column by fetching URLs based on column PCODE						
New column name	GoogleInfo Throttle delay 5000 milliseconds					
On error	set to blank store error					
Formulate the URLs to fetch:						
Expression Language Google Refine Expression Language (GREL) ▼						

"http://maps.googleapis.com/maps/api/geocode/json?sensor=false&address="+escape(value, "url")+", CAN"



Google Map API is queried and sends back results...





Embedded in this data, you'll find 'latitude' & 'longitude'

JSON output can be formatted...

```
{ "results" : [ { "address_components" : [ {
"long name": "K7L 5C4", "short name": "K7L 5C4",
"types" : [ "postal code" ] }, { "long name" :
"Kingston", "short name": "Kingston", "types": [
"locality", "political" ] }, { "long_name" : "Frontenac
County", "short_name": "Frontenac County",
"types": [ "administrative area level 2", "political"
] }, { "long_name" : "Ontario", "short_name" : "ON",
"types": [ "administrative area level 1", "political"
] }, { "long name" : "Canada", "short name" : "CA",
"types" : [ "country", "political" ] } ],
"formatted address": "Kingston, ON K7L 5C4,
Canada", "geometry" : { "bounds" : { "northeast" : {
"lat": 44.2275953, "lng": -76.4946668 },
"southwest" : { "lat" : 44.2269466. "lng" : -
76.4954425 } }, "location" : { "lat" : 44.2272811,
"Ing": -76.49506989999999 }, "location type":
"APPROXIMATE", "viewport" : { "northeast" : { "lat"
: 44.22861993029149, "lng" : -76.49370566970849 },
"southwest" : { "lat" : 44.22592196970849, "lng" : -
76.4964036302915 } } }, "partial_match" : true,
"place id": "ChlJo1lXewSr0kwRqBwN8Nxp4HM",
"types" : [ "postal_code" ] } ], "status" : "OK" }
```

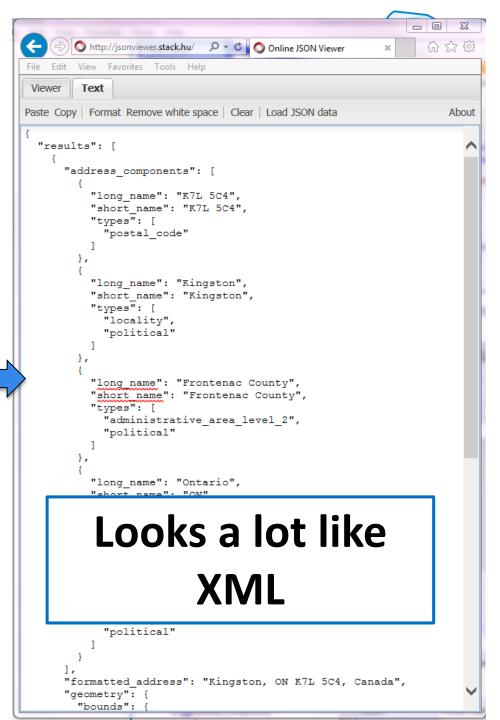


Online JSON Viewer

http://jsonviewer.stack.hu

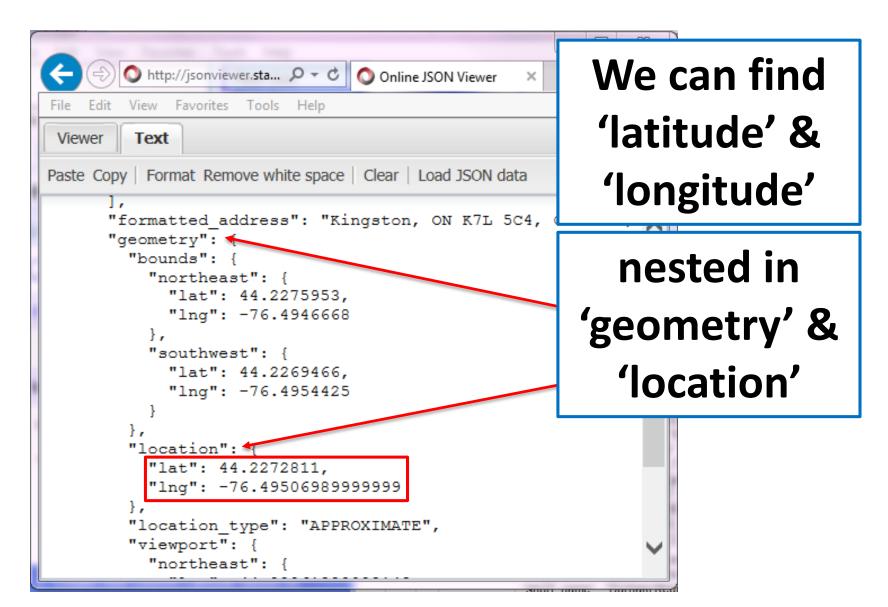
JSON output can be formatted...

```
{ "results" : [ { "address_components" : [ {
"long name": "K7L 5C4", "short name": "K7L 5C4",
"types" : [ "postal code" ] }, { "long name" :
"Kingston", "short name": "Kingston", "types": [
"locality", "political" ] }, { "long name" : "Frontenac
County", "short_name": "Frontenac County",
"types": [ "administrative area level 2", "political"
] }, { "long_name" : "Ontario", "short_name" : "ON",
"types": [ "administrative area level 1", "political"
] }, { "long name" : "Canada", "short name" : "CA",
"types" : [ "country", "political" ] } ],
"formatted address": "Kingston, ON K7L 5C4,
Canada", "geometry" : { "bounds" : { "northeast" : {
"lat": 44.2275953, "lng": -76.4946668 },
"southwest" : { "lat" : 44.2269466, "lng" : -
76.4954425 } }, "location" : { "lat" : 44.2272811,
"Ing": -76.49506989999999 }, "location type":
"APPROXIMATE", "viewport" : { "northeast" : { "lat"
: 44.22861993029149, "lng" : -76.49370566970849 },
"southwest" : { "lat" : 44.22592196970849, "lng" : -
76.4964036302915 } } }, "partial_match" : true,
"place id": "ChIJo1lXewSr0kwRqBwN8Nxp4HM",
"types" : [ "postal_code" ] } ], "status" : "OK" }
```



Now it is easier to find elements we might want to use





Edit column > Add column based on this column

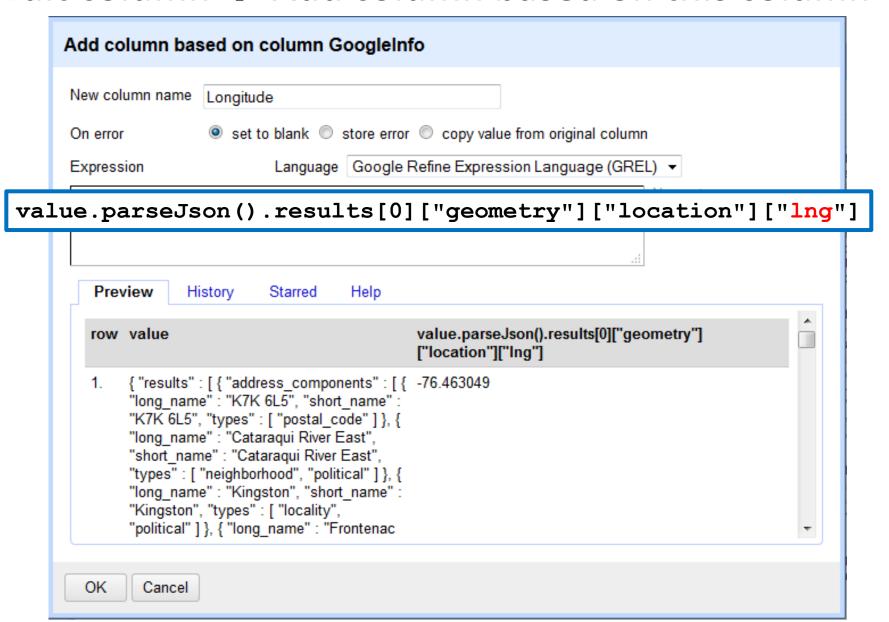
	Add column based on column GoogleIn	fo
	New column name Latitude	
	On error	copy value from original column
	Expression Language Google I	Refine Expression Language (GREL) ▼
valu	e.parseJson().results[0]	["geometry"]["location"]["lat"]
		.41
	Preview History Starred Help	
	row value	value.parseJson().results[0]["geometry"] ["location"]["lat"]
	1. { "results" : [{ "address_components" : [{ "long_name" : "K7K 6L5", "short_name" : "K7K 6L5", "types" : ["postal_code"] }, { "long_name" : "Cataraqui River East", "short_name" : "Cataraqui River East", "types" : ["neighborhood", "political"] }, { "long_name" : "Kingston", "short_name" : "Kingston", "types" : ["locality", "political"] }, { "long_name" : "Frontenac	44.2540515
	OK Cancel	

New 'Latitude' column



3 rc	ows	;		Extensions: Free	base ▼
Shov	w as:	rows record	ds Show: 5 10 25 50 rows « first «	previous 1 - 3 next) last »
▼ Al	▼ AII ▼ PCODE ▼		▼ GoogleInfo	▼ Latitude	Place
	57	l. k7k 6l5	{"results" : [{"address_components" : [{"long_name" : "K7K 6L5", "short_name" : "K7K 6L5", "types" : ["postal_code"]}, { "long_name" : "Cataraqui River East", "short_name" : "Cataraqui River East", "types" : ["neighborhood", "political"]}, { "long_name" : "Kingston", "short_name" : "Kingston", "types" : ["locality", "political"]}, { "long_name" : "Frontenac County", "short_name" : "Frontenac County", "types" : ["administrative_area_level_2", "political"]}, { "long_name" : "Ontario", "short_name" : "ON", "types" : ["administrative_area_level_1", "political"]}, { "long_name" : "Canada", "short_name" : "CA", "types" : ["country", "political"]}], "formatted_address" : "Kingston, ON K7K 6L5, Canada", "geometry" : { "bounds" : { "northeast" : { "lat" : 44.2550538, "lng" : -76.4626056}, "southwest" : { "lat" : 44.2535827, "lng" : -76.4638978}}, "location" : { "lat" : 44.2540515, "lng" : -76.463049}, "location_type" : "APPROXIMATE", "viewport" : { "northeast" : { "lat" : 44.25566723029149, "lng" : -76.4619027197085}, "southwest" : { "lat" : 44.2529692697085, "lng" : -76.46460068029151}}}, "partial_match" : true, "place_id" : "ChlJf723pzCq0kwRvM3KrXzj0", "types" : ["postal_code"]}], "status" : "OK"}	44.2540515	Kingston
公	9:	2. I1c 3k2edit	{"results" : [{ "address_components" : [{ "long_name" : "L1C 3K2", "short_name" : "L1C 3K2", "types" : ["postal_code"] }, { "long_name" : "Haydon", "short_name" : "Haydon", "types" : ["locality", "political"] }, { "long_name" : "Durham Regional Municipality", "short_name" : "Durham Regional Municipality", "types" : ["administrative_area_level_2", "political"] }, { "long_name" : "Ontario", "short_name" : "ON", "types" : ["administrative_area_level_1", "political"] }, { "long_name" : "Canada", "short_name" : "CA", "types" : "country", "political"] }, "formatted_address" : "Haydon, ON L1C 3K2, Canada", "geometry" : { "bounds" : { "northeast" : { "lat" : 44.066854, "lng" : -78.69502659999999 }, "southwest" : { "lat" : 43.9239747, "lng" : -78.779425 } }, "location" : { "lat" : 43.97851499999999, "lng" : -78.7306158 }, "location_type" : "APPROXIMATE", "viewport" : { "northeast" : { "lat" : 44.066854, "lng" : -78.69502659999999 }, "southwest" : { "lat" : 43.9239747, "lng" : -78.779425 } }, "partial_match" : true, "place_id" : "Chllp9BwSfcL1YkRr0ANR05zWyE", "postcode_localities" : ["Bowmanville", "Haydon"], "types" : ["postal_code"] }], "status" : "OK" }	43.97851499999999	Bowmanv
垃	5	3. k7l 5c4	{"results" : [{ "address_components" : [{ "long_name" : "K7L 5C4", "short_name" : "K7L 5C4", "types" : ["postal_code"] }, { "long_name" : "Kingston", "short_name" : "Kingston", "types" : ["locality", "political"] }, { "long_name" : "Frontenac County", "short_name" : "Frontenac County", "types" : ["administrative_area_level_2", "political"] }, { "long_name" : "Ontario", "short_name" : "ON", "types" : ["administrative_area_level_1", "political"] }, { "long_name" : "Canada", "short_name" : "CA", "types" : ["country", "political"] }], "formatted_address" : "Kingston, ON K7L 5C4, Canada", "geometry" : { "bounds" : { "northeast" : { "lat" : 44.2275953, "lng" : -76.4946668 }, "southwest" : { "lat" : 44.2269466, "lng" : -76.4954425 } }, "location" : { "lat" : 44.2272811, "lng" : -76.49506989999999 }, "location_type" : "APPROXIMATE", "viewport" : { "northeast" : { "lat" : 44.22861993029149, "lng" : -76.49370566970849 }, "southwest" : { "lat" : 44.22592196970849, "lng" : -76.4964036302915 } } }, "partial_match" : true, "place_id" "ChIJo1IXewSr0kwRqBwN8Nxp4HM", "types" : ["postal_code"] }], "status" : "OK" }	44.2272811	Kingston

Edit column > Add column based on this column



New 'Longitude' column

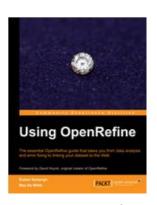


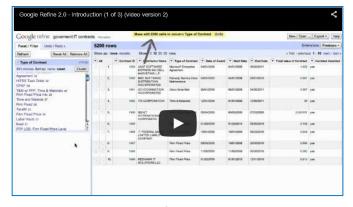
Sh	ow a	s: r	ows records	s Show: 5 10 25 50 rows	« first	c previous 1 - 3	
•	▼ AII ▼ PCODE		▼ PCODE	CODE GoogleInfo		▼ Latitude	
☆	5	1.	k7k 6l5	<pre>{ "results" : [{ "address_components" : [{ "long_name" : "K7K 6L5", "short_name" : "K7K 6L5", "types" : ["postal_code"] }, { "long_name" : "Cataraqui River East", "short_name" : "Cataraqui River East", "types" : ["neighborhood", "political"] }, { "long_name" : "Kingston", "short_name" : "Kingston", "types" : ["locality", "political"] }, { "long_name" : "Frontenac County", "short_name" : "Frontenac County", "types" : ["administrative_area_level_2", "political"] }, { "long_name" : "Ontario", "short_name" : "ON", "types" : ["administrative_area_level_1", "political"] }, { "long_name" : "Canada", "short_name" : "CA", "types" : ["country", "political"] } }, "formatted_address" : "Kingston, ON K7K 6L5, Canada", "geometry" : { "bounds" : { "northeast" : { "lat" : 44.2550538, "lng" : -76.4626056 }, "southwest" : { "lat" : 44.2535827, "lng" : -76.4638978 } }, "location" : { "lat" : 44.2540515, "lng" : -76.463049 }, "location_type" : "APPROXIMATE", "viewport" : { "northeast" : { "lat" : 44.25566723029149, "lng" : -76.4619027197085 }, "southwest" : { "lat" : 44.2529692697085, "lng" : -76.46460068029151 } } }, "partial_match" : true, "place_id" : "ChIJf723pzCq0kwRvM3KrXzj_0", "types" : ["postal_code"] }], "status" : "OK" }</pre>	-76.463049	44.2540515	
☆	9	2.	I1c 3k2	<pre>{"results":[{"address_components":[{"long_name":"L1C 3K2", "short_name":"L1C 3K2", "types":["postal_code"]}, {"long_name":"Haydon", "short_name":"Haydon", "types":["locality", "political"]}, { "long_name":"Durham Regional Municipality", "short_name":"Durham Regional Municipality", "types":["administrative_area_level_2", "political"]}, {"long_name":"Ontario", "short_name":"ON", "types":["administrative_area_level_1", "political"]}, {"long_name":"Canada", "short_name":"CA", "types":["country", "political"]}], "formatted_address":"Haydon, ON L1C 3K2, Canada", "geometry": {"bounds": {"northeast":{"lat":44.066854, "lng":-78.69502659999999}, "southwest":{"lat":43.9239747, "lng": -78.779425}}, "location":{"lat":43.978514999999999, "lng":-78.7306158}, "location_type": "APPROXIMATE", "viewport":{"northeast":{"lat":44.066854, "lng":-78.69502659999999}, "southwest":{"lat":43.9239747, "lng":-78.779425}}}, "partial_match": true, "place_id": "ChUp9BwSfcL1YkRr0ANR05zWyE", "postcode_localities":["Bowmanville", "Haydon"], "types":["postal_code"]}], "status":"OK"}</pre>	-78.7306158	43.97851499999999	
	9	3.	k7l 5c4	{"results" : [{ "address_components" : [{ "long_name" : "K7L 5C4", "short_name" : "K7L 5C4", "types" : ["postal_code"] }, { "long_name" : "Kingston", "short_name" : "Kingston", "types" : ["locality", "political"] }, { "long_name" : "Frontenac County", "short_name" : "Frontenac County", "types" : ["administrative_area_level_2", "political"] }, { "long_name" : "Ontario", "short_name" : "ON", "types" : ["administrative_area_level_1", "political"] }, { "long_name" : "Canada", "short_name" : "CA", "types" : ["country", "political"] } }, "formatted_address" : "Kingston, ON K7L 5C4, Canada", "geometry" : { "bounds" : { "lat" : 44.2269466, "lng" :	-76.49506989999999	44.2272811	

Closing comments



- Just scratched the surface of what OpenRefine can do
- More resources can be found at OpenRefine.org





Manual

Videos

Regular Expressions Cheat Sheet:

http://arcadiafalcone.net/GoogleRefineCheatSheets.pdf

Using OpenRefine to Update, Clean up, and Link your Metadata to the Wider World, Sarah Weeks, St. Olaf College.

Youtube: https://www.youtube.com/watch?v=E-NbMR3 MRw

Data: https://raw.githubusercontent.com/PaulMakepeace/refine-client-py/master/tests/data/louisiana-elected-officials.csv

