

Open Refine Tutorial 1 of 3: Cleaning Messy Data

In this tutorial on cleaning messy data you'll resolve variations of titles into a single identity using Open Refine. You'll also learn how to undo functions and export your refined spreadsheet to your desktop.

Getting Started

First you need to download a dataset to work with. This tutorial uses open data on public housing from the City of Los Angeles' Housing and Community Development Department, found here [<https://data.lacity.org/A-Livable-and-Sustainable-City/HCIDLA-Affordable-Housing-Projects-Catalog-And-Lis/u4mj-cwbz>].

On the webpage, click the blue 'export' box and export the file as a 'CSV for Excel' to your computer.

Now download Open Refine's installation package [<http://openrefine.org/download.html>]. Once installed, click on its application icon, and it will pop open in your default web browser. Create a new project and click 'browse' to upload your file. Click 'Next'. On the following screen, check that the box, "Parse cell text into numbers, dates, etc", is ticked, then click 'create project' on the following screen.

You'll see a familiar looking spreadsheet format. You can also see that you're able to edit each individual cell, and that each of the column headings has an arrow – clicking on the arrow brings up a dropdown menu of operations.

Resolving Variations in a Name

Open Refine helps find and clean variations or misspellings of titles, resolving them into a single term. This function is useful because spreadsheets often contain human error or suffer from poorly designed records keeping, with multiple name variations for the same entity. For instance, the author William Burroughs might be entered as "William Burroughs", "William S. Burroughs", "William Seward Burroughs", or "William Seward Burroughs II."

Furthermore, computers can read two titles as distinct if one is capitalized and the other isn't, or if one entry has a blank space at the end and another doesn't, a difference that isn't even apparent to the human eye. Consequently, when you sort or search for a term, you may miss variants in the retrieved results.

To locate and correct these inconsistencies, you'll learn two features: faceting and clustering.

Faceting

You'll be faceting the 'Developer' column to see if there are any inconsistencies in the titles.

Refine OPEN HCIDLA_Affordable_Housing_Projects_Catalog_And_Listing_2003_To_Present_csv Permalink Open... Export Help

Facet / Filter Undo / Redo 0 **341 rows** Extensions:

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

	JOB#	DEVELOPER	FINISH DATE(YE	MANAGEMENT	CONTACT PHON	DATE_STAMP	LONGITUDE	LATITU
to/CD_1_-_Barbizon_Hotel.png	47				(323) 937-6468	Nov 30 2015 5:25AM	-118.273227219	34.058957
to	146				(714) 282-2520	Nov 30 2015 5:25AM	-118.417627842	34.201029
to 124.jpg	169			LC	(818) 887-6920	Nov 30 2015 5:25AM	-118.598286345	34.208140
to	29				(714) 533-3450	Nov 30 2015 5:25AM	-118.257500077	34.091310
to/CD_9_-	53	BEYOND SHELTER HOUSING,	2006	WNC & ASSOCIATES, INC.	(213) 251-2111	Nov 30 2015 5:25AM	-118.278625635	33.996336
to 7%20acre%20-%20010139.JPG	7	NATIONAL FOUNDATION FOR ... (AKA PENNY LANE),	2003	PENNY LANE TRANSITIONAL HOUSING PROGRAM	(818) 892-3423	Nov 30 2015 5:25AM	-118.465337052	34.229852
to%20Site2	7	NATIONAL FOUNDATION FOR ... (AKA PENNY LANE),	2003	PENNY LANE TRANSITIONAL HOUSING PROGRAM	(818) 892-3423	Nov 30 2015 5:25AM	-118.465325782	34.231033
to	75	SRO HOUSING CORPORATION,	2005	SRO HOUSING CORPORATION	(213) 229-9641	Nov 30 2015 5:25AM	-118.2471476	34.041155
to/CD_13_-	65	ST ANNE'S TRANSITIONAL HOUSING,	2005	ST ANNE'S TRANSITIONAL HOUSING	(213) 381-2931	Nov 30 2015 5:25AM	-118.27894062	34.071697

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

Facet menu options: Text facet, Numeric facet, Timeline facet, Scatterplot facet, Custom text facet..., Custom Numeric Facet..., Customized facets

In the dropdown menu of the ‘Developer’ column, you’ll see the ‘Facet’ feature. In the submenu, click on the ‘facet text’ option. Now, in the panel to the left of the spreadsheet, you’ll now see a list of all the unique values within that column – 136 in total. Each unique value is assigned a number, meaning the number of times it appears in the column. By default the faceted list is sorted in alphabetical order. You can now easily compare the values to see if there are any variations of the same name.

Clearly this list has problems with consistent data entry. 1755 EFM is most likely the same entity as 1755 EFM, LLC, and the name Community of Friends has been entered both with a comma at the end and without one.

To resolve these differences, we will next use the *clustering* feature.

Clustering

Now click the ‘Cluster’ button. A pop-up window shows all groups of different cell values that might be alternative representations of the same thing:

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 9 clusters found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
3	8	<ul style="list-style-type: none"> SKID ROW HOUSING TRUST, (5 rows) SKID ROW HOUSING TRUST (2 rows) SKID ROW HOUSING TRUST (1 rows) 	<input checked="" type="checkbox"/>	SKID ROW HOUSING TRU...
2	6	<ul style="list-style-type: none"> THOMAS SAFRAN & ASSOCIATES (5 rows) THOMAS SAFRAN & ASSOCIATES, (1 rows) 	<input checked="" type="checkbox"/>	THOMAS SAFRAN & ASSO...
2	4	<ul style="list-style-type: none"> EAST L.A. COMMUNITY CORPORATION (3 rows) EAST LA COMMUNITY CORPORATION (1 rows) 	<input checked="" type="checkbox"/>	EAST L.A. COMMUNITY CC...
2	4	<ul style="list-style-type: none"> RETIREMENT HOUSING FOUNDATION, (3 rows) RETIREMENT HOUSING FOUNDATION (1 rows) 	<input checked="" type="checkbox"/>	RETIREMENT HOUSING FC...
2	2	<ul style="list-style-type: none"> MCCORMACK BARON SALAZAR (1 rows) MCCORMACK BARON SALAZAR, (1 rows) 	<input checked="" type="checkbox"/>	MCCORMACK BARON SAL...
2	18	<ul style="list-style-type: none"> META HOUSING CORPORATION, (11 rows) META HOUSING CORPORATION (7 rows) 	<input checked="" type="checkbox"/>	META HOUSING CORPORA...

Choices in Cluster: 2 - 3

Rows in Cluster: 2 - 18

Average Length of Choices: 22.66 - 31.5

Length Variance of Choices: 0.4710000000000003 - 1

The cluster function relies on a few different algorithms to determine how aggressively it will apply criteria to determine sameness. You can switch between these algorithms by selecting within the 'Keying Functions' drop down menu. Fingerprinting will be the strictest heuristic – this formula detects values that differ only in punctuation and capitalization. For instance, 'Atwood, Margaret' and 'MARGARET ATWOOD' would both be read as equivalent with the fingerprinting formula, but 'Margaret E. Atwood' would not.

The metaphone function is looser, so it will return more results; it groups entities together by how they sound. Try selecting among each keying function and go through the clusters, using your best judgement to determine if the algorithm has found similarities or delivered false positives. Click the box under 'Merge?' next to any clusters that need resolving, and adjust the title if necessary.

Some cases are less clear. When you use the Metaphone keying function, for example, you get the result of MCCORMACK BARON SALAZAR, MCCORMACK BARON & ASSOCIATES, and MCCORMACK BARON SALAZAR, INC. Should we assume that these are all the same company?

For clues, you can point your mouse over the row and click on 'browse this cluster'. Clicking this function allows you to isolate these rows and scrutinize them further. You can see that all three of these companies have the same Management company in the 'Management' column, so it's likely that they can be resolved as one entity.

Manual Cleaning

Now examine your facet again. If you have resolved all clusters retrieved by the cluster function, that will have narrowed down your list of unique entities. But the cluster function doesn't catch all similar matches. There still are some entities that should be merged that will need to be edited by hand. For example, "Women Organizing Resources Knowledge & Services" is most likely

the same as “W.O.R.K.S.”. To resolve these, click on the ‘edit’ function next to one of these names, and rename it with the same title as its match.

Additional functions: Creating a custom facet

You can use facets also to isolate rows of specific items. For instance, perhaps we want to find out all of the developers of public housing in Hollywood. To do this, go to the “Community” column and click ‘Facet’ → ‘Custom facet’. In the text box write:

```
value.contains("HOLLYWOOD")
```

The command is case sensitive, so you’ll need to keep the text in all-caps.

The results generates a boolean value of either ‘true’ or ‘false’ depending on whether the current value in the cell contains the text ‘HOLLYWOOD’ anywhere.

NOTE: A Boolean is a data type with only two possible values: true or false.

Next, you’ll note that the results also return ‘E. HOLLYWOOD’ and ‘N. HOLLYWOOD.’ To refine the results further, go to the drop down menu for the “Community” column and click ‘Facet’ → ‘Text Facet’. Now, in your facet box, you can select ‘HOLLYWOOD’ by clicking on the text ‘include’. The results are now isolated to 14 rows.

Note: This command is actually an example of a programming language called Google Refine Expression Language (GREL). Several advanced features like this are available. If you want to explore more commands visit the Google Refine Expression Language (GREL) reference [<https://github.com/OpenRefine/OpenRefine/wiki/Documentation-For-Users#reference>] or this tutorial [http://www.meanboyfriend.com/overdue_ideas/wp-content/uploads/2014/11/Introduction-to-OpenRefine-handout-CC-BY.pdf].

Note also: You can get to the same results by using the Text Filter. In the “Community” column click the drop down menu, then ‘Text filter’. Type in ‘hollywood’ – in this instance the text does not need to be case sensitive.

Other Types of Facets

You’ll have seen in the drop down menu that Open Refine allows other types of facets, including Numeric facets, Timeline facets for dates, Custom facets, and Scatterplot facets.

- Numeric and Timeline facets display graphs instead of lists of values.
- Scatterplot facets, which allow you to visualize the relationships between numeric values in columns. For more information, you can visit this tutorial:

http://enipedia.tudelft.nl/wiki/OpenRefine_Tutorial#Exploring_the_data_with_scatter_plots

- Custom facets allow a range of commands, including isolating cells in a column that contain a particular word (as shown above), finding duplicates, finding blanks, and faceting by length of characters in a cell.

Note: See the 2nd tutorial in this series for more information on Numeric facets.

Undo/Redo

With the undo and Redo function you can go back to any past action you've made on your dataset and restore the data to its state at that point. In this way you can experiment on your data but still access all prior versions.

'Undo' and 'Redo' are accessed via the lefthand panel. Through the 'extract' and 'apply' functions you can also save the operations you've carried out on one data set and apply them to another data set by copying and pasting them.

Exporting Your Cleaned Data

Click on the 'Export' box on the upper right hand corner. You can export as HTML, Excel, csv (comma separated values), tsv (tab separated values), or a custom export.

Additional Resources

Video tutorials on Open Refine:
<http://openrefine.org/>

Text tutorials:

"OpenRefine Tutorial" by David François Huynh
http://enipedia.tudelft.nl/wiki/OpenRefine_Tutorial

"Cleaning Spending Data with OpenRefine," by School of Data
<http://schoolofdata.org/handbook/recipes/cleaning-spending-data-open-refine/>

"Using Google Refine to Clean Messy Data," by Dan Nguyen for ProPublica
<https://www.propublica.org/nerds/item/using-google-refine-for-data-cleaning>

"Introduction to OpenRefine," developed by Owen Stephens on behalf of the British Library.
http://www.meanboyfriend.com/overdue_ideas/wp-content/uploads/2014/11/Introduction-to-OpenRefine-handout-CC-BY.pdf

"Cleaning Data with OpenRefine," by Seth van Hooland, Ruben Verborgh and Max De Wildo
<http://programminghistorian.org/lessons/cleaning-data-with-openrefine>

Credits:
Written by Morgan Currie, September 2016.