

Open Refine Tutorial 3 of 3: Sorting, Column Splits, Dates, and Geocoding

In this tutorial on Open Refine we'll learn how to reorder and sort columns, create new columns from information in existing cells, manipulate dates, and create geo-coordinates from addresses.

Getting Started

First you need to download a dataset to work with. This tutorial uses open data on festivals from the City of Los Angeles' Department of Cultural Affairs, found here [<https://data.lacity.org/A-Livable-and-Sustainable-City/Events-from-LA-Festival-Guide-2014/acy8-72w9>].

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

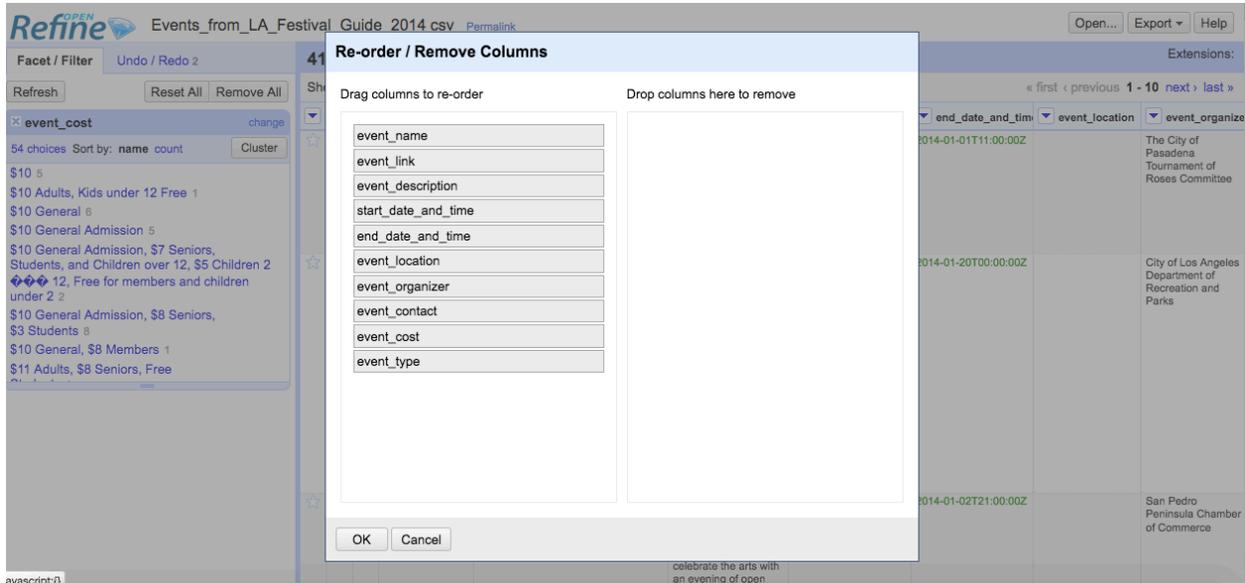
In Open Refine create a new project and click 'browse' to upload your file. Click 'Next'. On the following screen, check the box "commas (CSV)". *Make sure that the box, "Parse cell text into numbers, dates, etc", is unchecked.* Then click 'create project' on the following screen.

Note: For instructions on downloading Open Refine, see Tutorial 1.

Reordering Columns

You may want to reorder or even remove columns in Open Refine. To perform these actions, click the drop down menu at the top of the 'All' column, and choose 'Edit columns' → 'Re-order / remove columns ...'

You'll find a list of all columns in the spreadsheet, and you can drag and drop these in the order you desire, or remove columns completely. Go ahead and resort the columns so that 'Event Cost' is next to the 'Event Location' column, or any combination you choose.



Sorting Data

You can also sort data in OpenRefine. Currently this spreadsheet is sorted by date of festival, but say you want to sort by price. Click on the drop down menu for the ‘Price’ column and choose ‘Sort’.

Notice that a new ‘Sort’ drop down menu is now displayed above the column headings. The ‘Sort’ drop down menu allows you to reverse the sort order, remove existing sorts, and make sorts permanent. (Unlike Excel ‘Sorts’, sorts in OpenRefine are temporary). Click on the ‘Sort’ drop down menu and click ‘Reorder rows permanently’ to keep the sort you just made.

You can also sort on multiple columns at the same time, so try a few sorting options if you like.

Splitting Columns

Open Refine allows you to make new columns out of data from a single column. To try this function, first make sure your dates are read as text and display like this:

Refine OPEN Events_from_LA_Festival_Guide_2014 csv Permalink

Open... Export Help

Facet / Filter Undo / Redo 0 412 rows Extensions:

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

Filter:	event_description	start_date_and_time	end_date_and_time	event_location	event_organizer	event_contact	event_cost	event_type
0. Create project 1. Split 411 cell(s) in column start_date_and_time into several columns by separator	amentofroses.com This historic and traditional California celebration features majestic floral floats, high-stepping equestrian units, and spirited marching bands from across the nation.	01/01/2014 08:00:00 AM	01/01/2014 11:00:00 AM		The City of Pasadena Tournament of Roses Committee	626-449-ROSE, 626-449-4100	Free	Family Friendly, Holiday Celebration, Fairs & Festivals
	ingsquareicerink.com Situating amid the Downtown Los Angeles skyscrapers, the rink will be open for ice skating daily, including holidays. There will also be a variety of free activities, including live concerts, youth programs, and special events. Event times: Mondays through Thursdays, 11:30 a.m. - 9:30 p.m., Fridays through Sundays, 10:00 a.m. - 11:00 p.m.	01/01/2014 12:00:00 AM	01/20/2014 12:00:00 AM		City of Los Angeles Department of Recreation and Parks	213-624-4289	\$8 per hour	Family Friendly, Great Outdoors
	ursday.com On the first Thursday of each month, the artists, entertainers, and business people of San Pedro celebrate the arts with an evening of open	01/02/2014 05:00:00 PM	01/02/2014 09:00:00 PM		San Pedro Peninsula Chamber of Commerce		free	Family Friendly, Culture & Community, Arts & Antiques

Note: If the dates are read as 'dates', you may need to reupload your data, making sure that the box, "Parse cell text into numbers, dates, etc", is unchecked.

Now go to the drop down menu of the 'start_date_and_time' column to 'Edit column' → 'Split into several columns'. In this case, we want to make a separate column for the festivals start date times.

In the pop up window that appears, make sure 'by separator' is checked, and use a blank space in the text box next to 'Separator' (you won't see anything in the box when you do this). Then you'll want to split into 2 columns. These action tells Open Refine that you want to take all text after the first space in the cells and put this data into a new column. Now click 'OK.'

Now rename the columns now to 'Festival start date' and 'Festival start time.' Go to 'Edit column' → 'Rename this column' on both the old and new columns to perform this action. Go through the same steps now on the 'end_date_and_time' column.

Cleaning and Displaying Dates

Open Refine allows you to display dates in different formats. In this case, you'll use a Transformation command. In the first set of parentheses, you have the current date format. In the second set of parentheses, you have the format you want to convert to:

```
value.toDate('MM/dd/yy').toString('MMMM dd yyyy')
```

There are many different ways to display dates. After the transformation, try playing around with combinations of these commands to see how dates are differently displayed:

```
value.toDate('MMMM dd yy').toString('MM dd yy')
```

```
value.toDate('MMMM dd yy').toString('MMM dd yy')
value.toDate('MMMM dd yy').toString('MM-dd-yy')
```

Timeline Faceting

We can also tell Open Refine to read this column as a date. Go to 'Edit cells' → 'Common Transformations' → 'To date'. This conversion undoes the formatting from the previous exercise. However, with the date format we can now perform a timeline facet on the column. This action allows some quick analysis, such as seeing which months attract the most festivals.

To create the timeline facet, go to the 'Festival start date' column, click 'Facet' → 'Timeline facet.' Now you can toggle the histogram sliders so that they bracket the months with the most festivals:

The screenshot shows the Open Refine interface with a dataset named 'Events_from_LA_Festival_Guide_2014.csv'. A facet is applied to the 'start_date_and_time' column, showing a histogram of event counts by month. The histogram has a prominent peak in April. The main table below the facet shows 73 matching rows. The first row is for 'Pomona Art Walk' on 2014-04-26. The second row is for 'City of Lights, City of Angels (COLCOA) - A Week of French Film Premieres in Hollywood' on 2014-04-21. The third row is for 'City of Lights, City of Angels' on 2014-04-23.

	event_name	event_link	event_description	start_date_and_time	start_date_and_time	end_date_and_time	event_location
94.	Pomona Art Walk	http://www.metropomona.com	The Downtown Pomona Arts Colony has an art walk every second and last Saturday of the month that brings big crowds to stroll the streets. Over a dozen galleries host receptions and open houses to showcase their latest exhibits and artists. Music, food, wine, and art throughout the evening. This event takes place on April 12 & 26	2014-04-26T00:00:00Z	03:00:00 PM	04/26/2014 11:00:00 PM	
101.	City of Lights, City of Angels (COLCOA) - A Week of French Film Premieres in Hollywood	http://www.colcoa.org	COLCOA, a week of French film premieres in Hollywood, is one of the biggest events dedicated to French cinema in the world with an exclusive program of 50 films, including world and North American premieres. All films presented with English subtitles. This event takes place April 21 - 28	2014-04-21T00:00:00Z	09:00:00 AM	04/21/2014 11:59:00 PM	
102.	City of Lights, City of Angels	http://www.colcoa.org	COLCOA, a week of French film premieres	2014-04-23T00:00:00Z	09:00:00 AM	04/23/2014 11:59:00 PM	

You can see that April, May, June, and July are the months that attract the most festivals in Los Angeles.

Geocoding

Finally, here are the steps to take to create geocode values of text addresses. Geocodes are very powerful to have in your dataset, since they will allow you to map locations into GIS software. These functions use Google Maps behind the scenes to figure out your most likely geographic location.

To begin, we need to isolate all the cells that have addresses (you'll see most are blank). Go to the column 'Event location'. In the drop down menu click 'Facet' → 'Facet by text.' Scroll down to the bottom of your list. You'll see that you have 383 blanks. To facet the blank cells out, click 'Facet by choice counts' at the very bottom of the facet list. This shows a histogram; you can now toggle the sliders so that they exclude the tallest column, which represents the blanks.

Refine **Events_from_LA_Festival_Guide_2014.csv** [Permalink](#) Open... Export Help

Facet / Filter Undo / Redo **29 matching rows (412 total)** Extensions:

Refresh Reset All Remove All Show as: rows records Show: 5 10 25 50 rows « first < previous 1 - 29 next > last »

event_location	event_descriptive	start_date_and_	start_date_and_	end_date_and_t	event_location	event_organizer	event_contact	event_cost	event_type
20 choices Sort by: name count Cluster	Ring in the New Year and the Year of the Horse with fun arts and crafts, food, and exciting performances. In Japan the most important and elaborate holiday is Oshogatsu - the celebration of the New Year.	Jan 5, 2014	11:00:00 AM	01/05/2014 05:00:00 PM	100 N Central Ave	Japanese American National Museum	213-625-0414	Free	Holiday Celebration, Family Friendly, Culture & Community
501 S Diamond Bar Blvd 1	The International Los Angeles photographic art exposition features the finest photographic art from the earliest 19th century photographic experiments to the most contemporary photography and photo-based art. More than 80 premiere galleries and private dealers present international and US artists. Event times: Friday and Saturday 11:00 a.m. - 7:00 p.m., Sunday 11:00 a.m. - 6:00 p.m.	Jan 16, 2014	12:00:00 AM	01/19/2014 12:00:00 AM	1933 Broadway Blvd	ArtLA	323-937-5523	General \$20, Students and seniors \$15	Culture & Community, Art & Antiques

Facet by choice counts

event_location change reset

0.00 — 380.00

You should now have 29 rows left, all with location information. This amount is a good size for the next function, which can take some time when working with larger datasets.

Now go to 'Edit column' → 'Add column by fetching URLs'. This action will create a new column with cell values that come from the internet, in this case Google maps. Enter this command into the text box:

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

You'll need to name the new column and set the throttle delay to around 500 milliseconds. The new column should look like this:

Refine Events_from_LA_Festival_Guide_2014 csv Permalink

Open... Export Help

Facet / Filter Undo / Redo 6

Refresh Reset All Remove All

29 matching rows (412 total)

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

Extensions: « first < previous 1 - 29 next > last »

event_location	event_descriptic	start_date_and_	start_date_and_	end_date_and_t	event_location	Geocode_Set
20 choices Sort by: name count Cluster	Ring in the New Year and the Year of the Horse with fun arts and crafts, food, and exciting performances. In Japan the most important and elaborate holiday is Oshogatsu - the celebration of the New Year.	Jan 5, 2014	11:00:00 AM	01/05/2014 05:00:00 PM	100 N Central Ave	{ "results": [{ "address_components": [{ "long_name": "100", "short_name": "100", "types": ["street_number"] }, { "long_name": "North Central Avenue", "short_name": "N Central Ave", "types": ["route"] }, { "long_name": "Umatilla", "short_name": "Umatilla", "types": ["locality", "political"] }, { "long_name": "Lake County", "short_name": "Lake County", "types": ["administrative_area_level_2", "political"] }, { "long_name": "Florida", "short_name": "FL", "types": ["administrative_area_level_1", "political"] }, { "long_name": "United States", "short_name": "US", "types": ["country", "political"] }, { "long_name": "32784", "short_name": "32784", "types": ["postal_code"] }, { "long_name": "7514", "short_name": "7514", "types": ["postal_code_suffix"] }], "formatted_address": "100 N Central Ave, Umatilla, FL 32784, USA", "geometry": { "bounds": { "lat": 28.9276567, "lng": -81.6680648 }, "location": { "lat": 28.927529799999999, "lng": -81.668243399999999 }, "location_type": "ROOFTOP", "viewport": { "northeast": { "lat": 28.92894232915, "lng": -81.66680511970848 }, "southwest": { "lat": 28.9262442697085, "lng": -81.66950308029151 } } }, "place_id": "ChJ94IW4R-I54gRsnwKkD6MkUQ", "types": ["premise"] }, { "address_components": [{ "long_name": "100", "short_name": "100", "types": ["street_number"] }, { "long_name": "North Central Avenue", "short_name": "N Central Ave", "types": ["route"] }, { "long_name": "Hartsdale", "short_name": "Hartsdale", "types": ["locality", "political"] }, { "long_name": "Greenburgh", "short_name": "Greenburgh", "types": ["administrative_area_level_3", "political"] }, { "long_name": "Westchester County", "short_name": "Westchester County", "types": ["administrative_area_level_2", "political"] }, { "long_name": "New York", "short_name": "NY", "types": ["administrative_area_level_1", "political"] }, { "long_name": "United States", "short_name": "US", "types": ["country", "political"] }, { "long_name": "10530", "short_name": "10530", "types": ["postal_code"] }, { "long_name": "1910", "short_name": "1910", "types": ["postal_code_suffix"] }], "formatted_address": "100 N Central Ave, Hartsdale, NY 10530, USA", "geometry": { "bounds": { "lat": 41.0217295, "lng": -73.7949527 }, "location": { "lat": 41.02153, "lng": -73.795197199999999 }, "location": { "lat": 41.0216297, "lng":

To convert this into a more readable format, you need to click on the new column, and then select 'Edit column' → 'Add column based on this column.' Create a new name for the column, such as 'Geocoordinates', and enter in the expression below:

```
with(value.parseJson().results[0].geometry.location, pair, pair.lat + ", " + pair.lng)
```

You now have a column containing the geocoordinates.

Credits:

Written by Morgan Currie, September 2016.