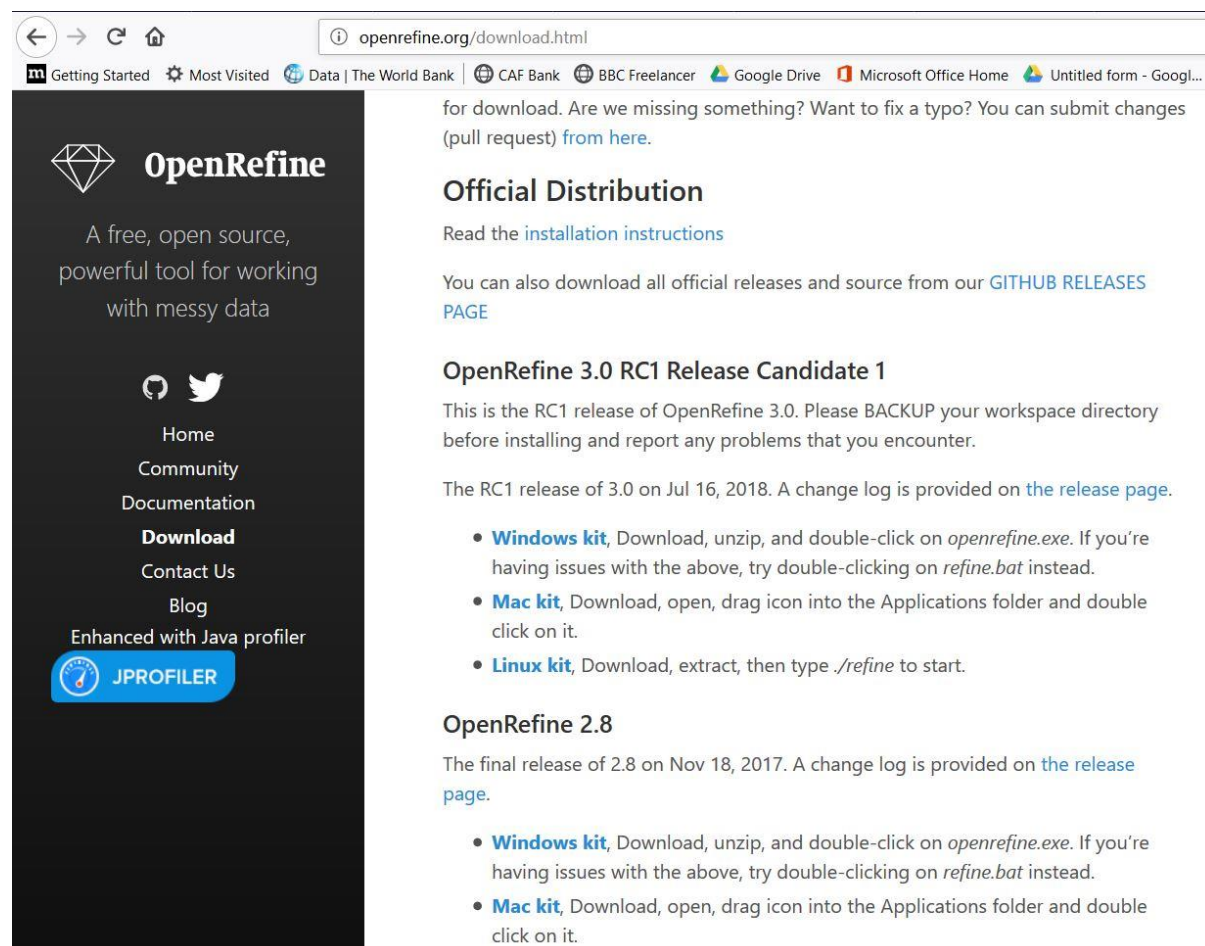OPEN REFINE – CLEAN DATA PAINLESSLY

Open Refine is a powerful tool for cleaning large amounts of data quickly and efficiently. And it's free.

You download a copy from openrefine.org/download.html – at the time of writing the most recent stable version is 2.8



Follow the instructions for installing onto your system – Mac/PC/Linux. As a program, OpenRefine runs within a browser window, and so will force your default browser to open (if it's already open, it will force a new tab to open). It works well with Firefox or Chrome, but not Microsoft Edge or Explorer.

When you run the program for the first time you may get error messages – your virus software may not like being forced to open a browser, or you may need to update your version of Java. Just do what it asks.

When you go to the new tab it will look like this



 The main feature is a fairly standard file choosing menu >browse>select file > next. But take a closer look when you have a moment sometime – there are different ways of importing data, and on the far left, there are options to create projects, open existing projects, or import projects.

Every file you work on is seen as a "project". OpenRefine does all its work on the file as a virtual copy, which you never actually open, but make edit decisions, which are recorded by the program, viewed within your browser. This means you never change the original data. It also means that when you have finished cleaning the dataset, and want to work on it back in Excel or R or whatever, then you will need to export it (in whatever format you need) – it then gets a new title, and you save it just as you would any other file you downloaded or imported to your hard drive.

By the way – the fact that you are using a web browser does not mean your data is on the internet: you don't need a connection to the web, and your data remains on your hard drive, not in the cloud. It's just Refine's way of presenting the data to you.

Projects can be shared too. The 3rd option – import project – is very useful in a newsroom or team context. Projects can be exported from one person's hard drive and imported by another – the 2nd user acquires the original file, and the editing record created by the 1st user, so when they import a project, they can see what their colleague has already done, and take on the file and all the editing decisions as if they had done it all along. Very powerful.

Back to creating a project (ie opening a dataset in OpenRefine for the first time):

The first thing you see after selecting your file, and clicking next is this preview



You have various options – OpenRefine usually (99.9% of the time) guesses what the file format is and opens it correctly. You can choose the formatting (in this case UTF-8 is the only one which reads the £ correctly – just click on the encoding to see the other options:
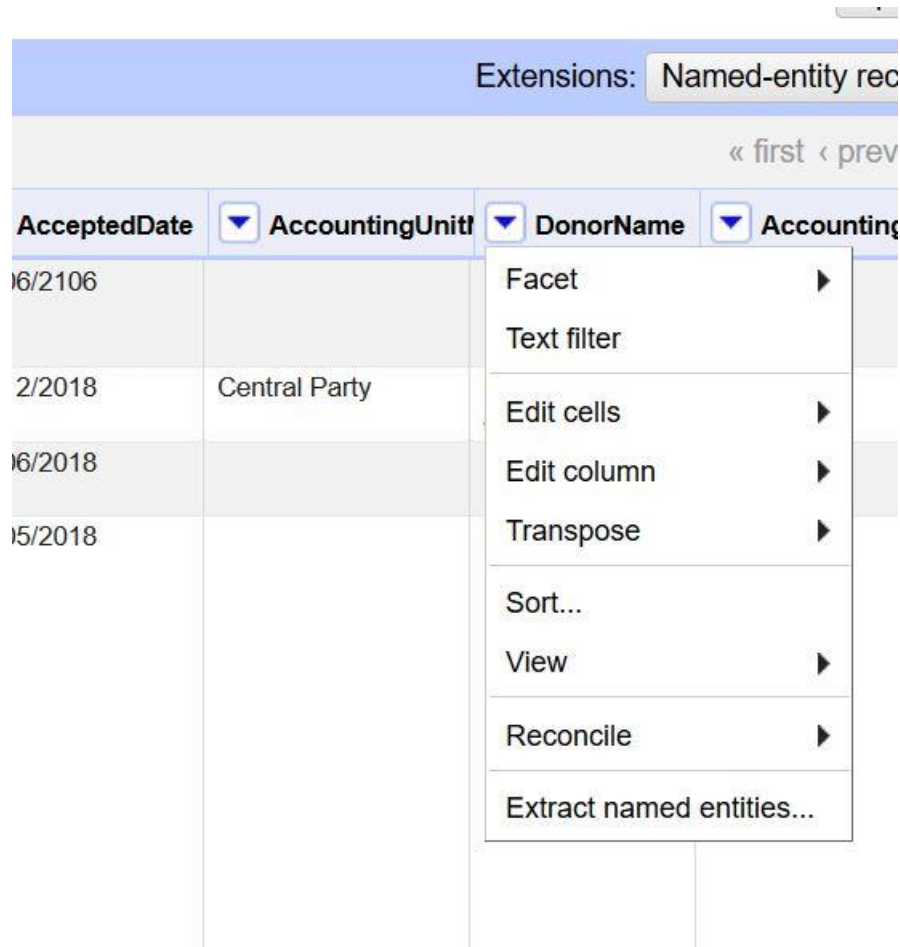


You can also choose options to help deal with blank cells, or column names, or ignore lines with logos or other redundant data in them. Remember, you're not editing the original file, and you can go back to this step at any point and restart a project – either from scratch, or by using the undo function, which we will come to shortly.
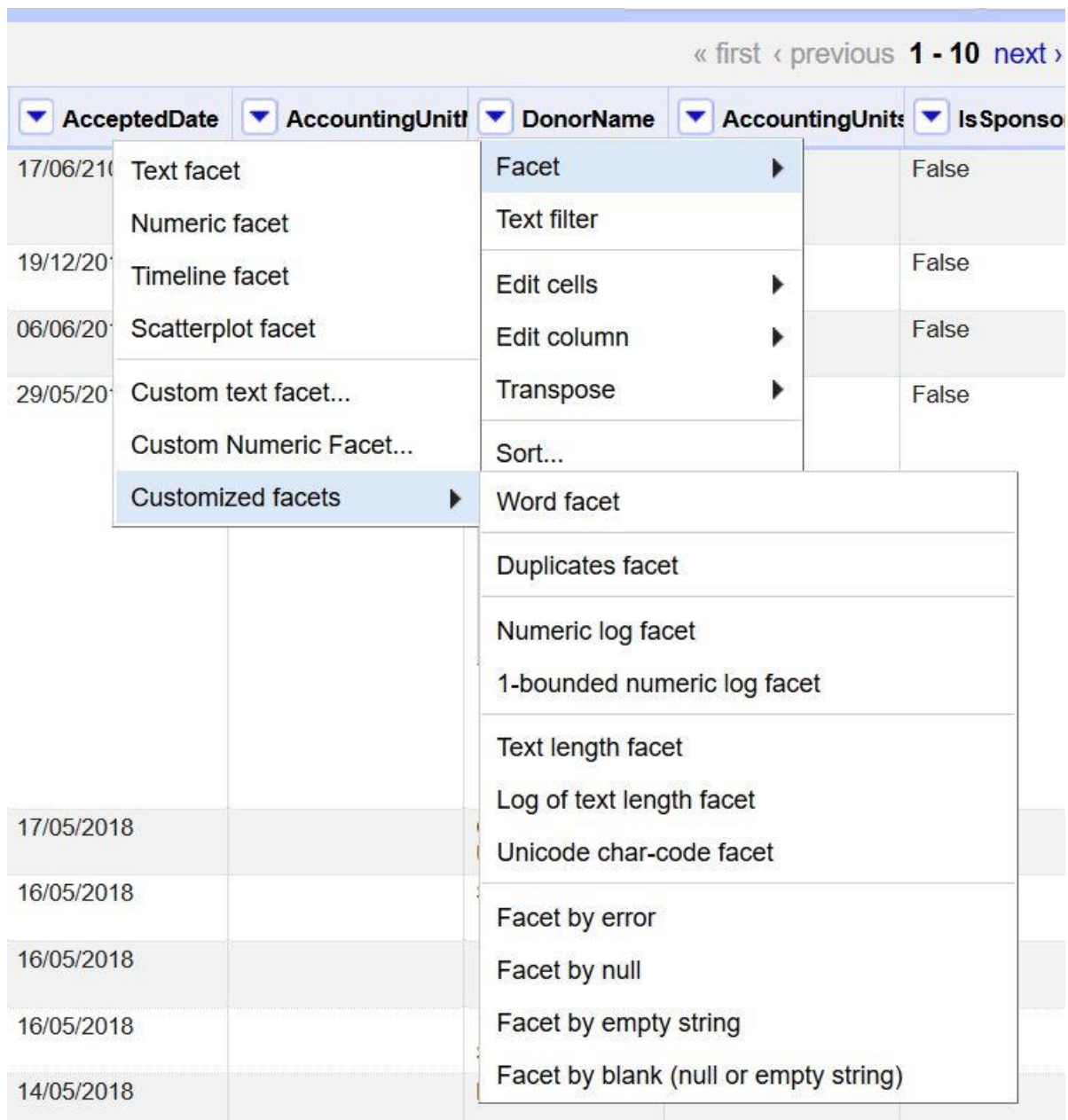
OpenRefine is very powerful and has many built-in functions which you may never use. When you have time, just play with it to see what happens. In this case we are going to use the most powerful function for datajournalists – cleaning one column of data by finding all the anomalies and getting names to be recorded in a uniform way.

As you will have seen, OpenRefine's logo and imagery is centred on a diamond, and the polishing of rough diamonds into perfect jewels. So we are going to use a function called "faceting" (polishing one side or facet of a diamond) and we are going to polish the text.

Every column has a dropdown arrow selection tool. We are working on the "Donor Name" column:



We choose "Facet" which brings up the next set of options:

| ▼ AcceptedDate | ▼ AccountingUnitH | ▼ DonorName | ▼ AccountingUnits | ▼ IsSponso |
|---|---|---|---|---|
| 17/06/21( | | | | False |
| 19/12/20 | | | | False |
| 06/06/20 | | | | False |
| 29/05/20 | | | | False |
| 17/05/2018 | | | | |
| 16/05/2018 | | | | |
| 16/05/2018 | | | | |
| 16/05/2018 | | | | |
| 14/05/2018 | | | | |

Menu items shown in the image:

Text facet
Numeric facet
Timeline facet
Scatterplot facet
Custom text facet...
Custom Numeric Facet...
Customized facets ▶

Facet ▶
Text filter
Edit cells ▶
Edit column ▶
Transpose ▶
Sort...

Word facet
Duplicates facet
Numeric log facet
1-bounded numeric log facet
Text length facet
Log of text length facet
Unicode char-code facet
Facet by error
Facet by null
Facet by empty string
Facet by blank (null or empty string)

We choose "Text Facet" – but do come back and take a look at some of the others. There are many options which may come in handy one day. Documentation is pretty good, but OpenRefine is a unique search string and you will be able to find a lot of tutorials and solutions by googling these options (including looking for them on youtube).

Choosing "text facet" makes OpenRefine do some work, and a few seconds later your screen will look like this – with a new box on the left, taking the name of the column you chose "Donor Name". It shows the names once each, in alphabetical order. You can also click on "count" and see them in order of frequency (slightly faster than pivoting them in a spreadsheet if you just want to see whose name pops up the most).

To see the real power of faceting, click on the "cluster" button on the right of the facet box. Now you see all the ways in which OpenRefine identifies entries which could be the same – but are entered differently ("dirty data") – the first one in this example shows the power of this tool – not only does it spot that "Ministry of Foreign Affairs, Qatar" and "Ministry of Foreign Affairs – Qatar" are probably the same, but it even finds another less obvious variant – "Qatar Ministry of Foreign Affairs" which is yet another way of saying the same thing.



Just above the table showing the clustered entries you will see two dropdown selectors – "Method" and "Keying function". These are the different ways – 6 in all – which OpenRefine assesses and groups possible errors in data entry which you can then use your own knowledge of the subject area to decide which entries to merge, and under which version to merge them (if you don't like the options you can just type your version into the box under the heading "New Cell Value".

Under the method "key collision" there are 4 different keying functions. Work through them in the order they are listed. The next one is "ngram-fingerprint" – which brings up a whole different set of options. As before you work through them, selecting or ignoring them depending on the accuracy of OpenRefine's assessment. When you're happy with a set of selections, click "Merge selected and re-cluster" and wait for OpenRefine to do its magic. Then carry on, either with a new batch, or a new keying function. (I tend to do about 10-15 mergers in one go, as the actual edit can take several

minutes. I did once crash OpenRefine by doing too many "mass edits" – as they're called – in one go!)
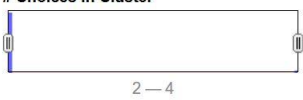
Take a moment to look at how these groups of results differ slightly from each other. Sometimes you can't see the difference, and it's usually down to an extra space somewhere in the string – white space is "seen" differently by the computer, and would become a different entry in a pivot table, for example.
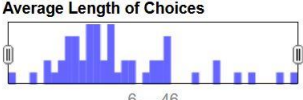


As you work through the different options, merging as you go, the number of options offered for merging in each cluster goes down – it's worth doing the clustering in the order suggested by the drop down options. I often find that "metaphone3" is the least useful analysis – OpenRefine sees all entries beginning with the same words: "Mr Christopher", for example, as being potentially the same as each other. Obviously this is nonsense. But it's worth looking through all the results just in case there is something you wouldn't otherwise have seen.

When you have done all the merging and clustering you can (or that you can stand – it can make your eyeballs burn after a while), close the clustering tool.

If you, or a colleague, wants to revise what you have done at any point, you can click on "undo/redo" at the top left of your screen
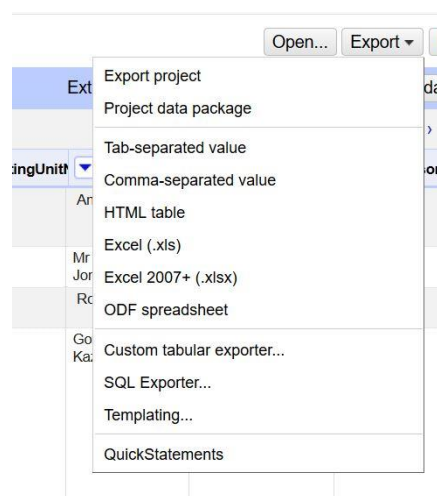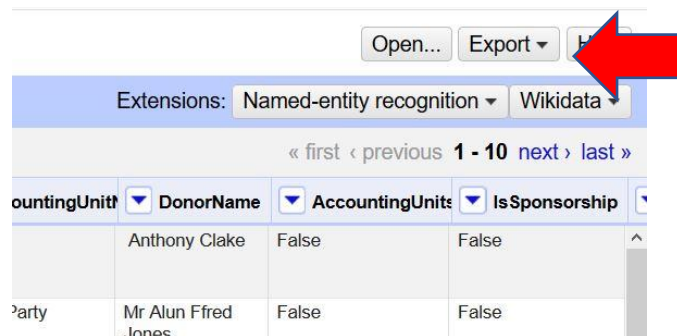


This brings up a list in chronological order of every edit you have made – in this case, I have done two things – opened the project, and then a mass edit of 44 cells in one column. If I decide I did the wrong thing, or perhaps new information comes in, instead of starting all over again, I could just go back one step – in this case, if I click on "0. Create Project" I will go back to the way the dataset looked at that point, and either export it, or re-cluster the columns and edit it differently.

(At that point, I can't unfortunately create two versions of history – my new set of edits will overwrite my previous set. If I want to preserve what I have done before I start re-editing it, then I can export the dataset – see all the various options.





Whatever I do at this point, one of the great strengths of OpenRefine is that what I have in my browser tab now will be there next time I reopen that project. It saves everything as I go along, and it records the edits I make forever.