# Exercise on the use of Open Refine

**Contents**
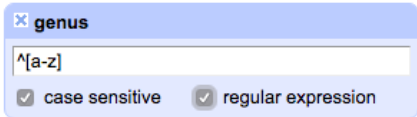
# CONVENTIONS

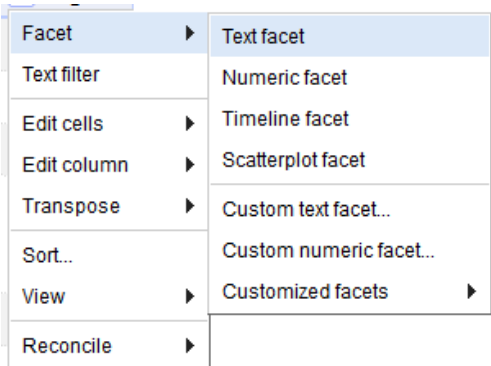| | |
|---|---|
| *Formulas (copy-paste)* | Text in blue |

**Example:** ...then paste the expression `^[a-z]`



*Commands in Refine*  Text in red

Example: ...and follow the route to Text facet



*Column names*  Text in green

Example: ...go to column Cat. Numb



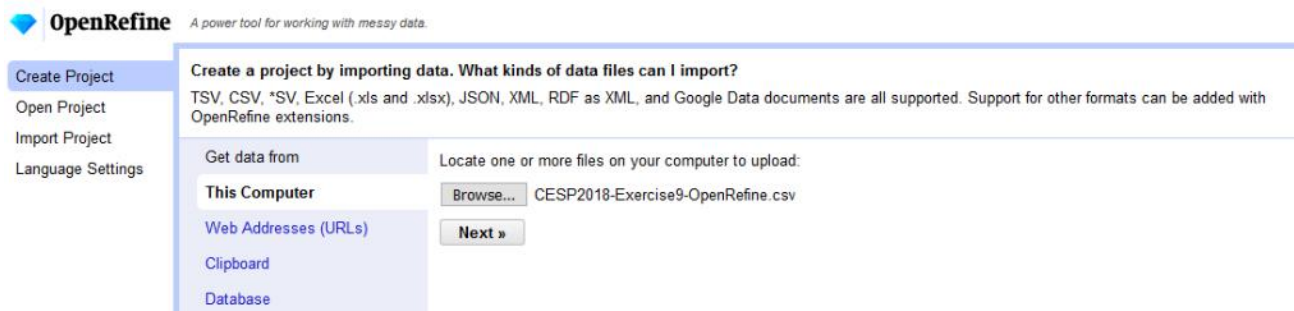*Hyperlinks*  www.gbif.org

*Column menu*

# 2. BASIC USE

## 2.1. FILE LOADING AND PROJECTS
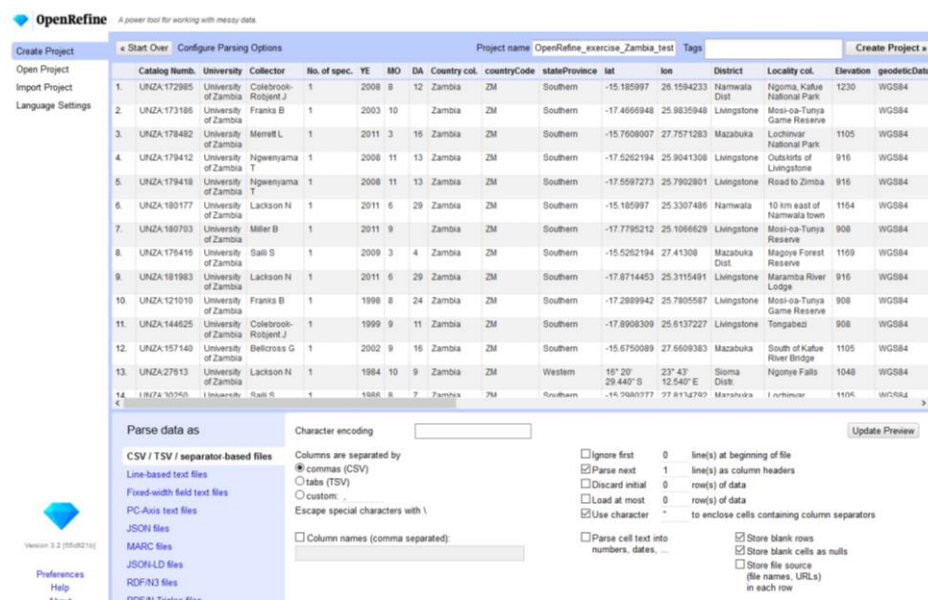
### 2.1.1. Before starting

Data can be loaded from various data sources: TSV, CSV, SV, Excel (.xls and .xlsx), JSON, XML, RDF and XML data as well as Google Docs. Loading data involves two steps: the first is selecting the file, and the second is the creation of the project.

### 2.1.2. Exercise 1. Create a project

1. Find the file Exercise-OpenRefine.csv in your course folder.
2. Open *OpenRefine* (using openrefine.exe), click on Create Project, and follow the route Get data from > This Computer, then click on Browse. Select the file. Click on Next.



3. A parsing options menu will appear. Fill in the options as shown in the picture below. Note that columns may be separated by tabs, commas or semicolons (;), so if you are not sure which delimiter your file uses, experiment with different options, until the data is split into columns correctly, as shown in the image below :

5. On the top centre you can provide a Project name, and click Create Project and you will be ready to work!
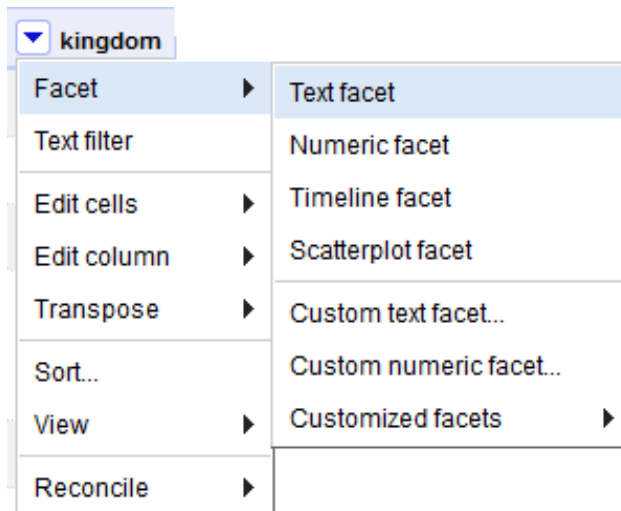
## 2.2. FACETING

### 2.2.1. Before starting

Faceting is a feature that will allow us to get a big picture overview of the data, and to filter down to just the subset of rows that we want to change or view in bulk. It facilitates the use and analysis of data and can be done with cells containing any kind of text, numbers and dates.

### 2.2.2. Exercise 2. Faceting and mass editing

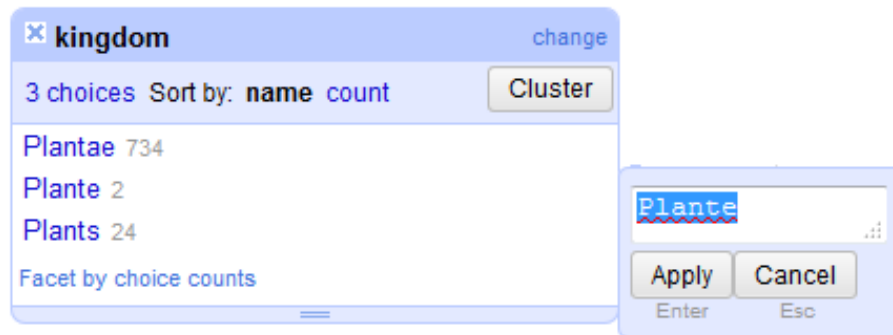1. Go to column kingdom, and then click on the column menu ▼ and follow the route to Text facet as shown below:



2. On the left a window with the name of the column will appear, which is the facet:



Click on count to sort by count, then click on name to sort alphabetically.

3. To fix the spelling mistakes (all values should be identical), place the cursor over the text in the window and click on Edit, then fix the error in the text box, and to save click on Apply.
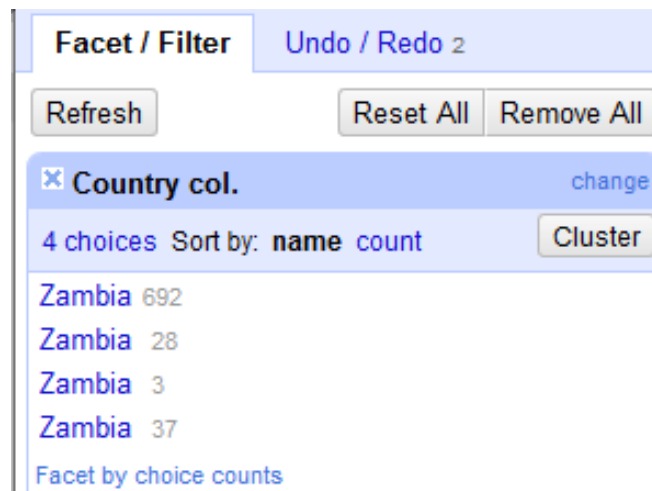
All the values will be fixed automatically in all records.

4. To close a facet, click on the **x** next to the facet's name.

## 2.2.3. Exercise 3. Faceting and white spaces I

1. Go to Country col. and click on column menu and perform a Text Facet.

It appears that the country name is spelled correctly, but the facet shows four different values. This is due to the extra spaces at the end of the text.
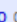
2. To fix the error in this column, go to Country col. and click on column menu ▼ > Edit Cells > Common transforms > Trim leading and trailing whitespace. You will see a notification message:

Text transform on 68 cells in column Country col.: value.trim()   Undo

3. Now check the facet window: only one value will remain.

## 2.2.4. Exercise 4. Faceting and white spaces II

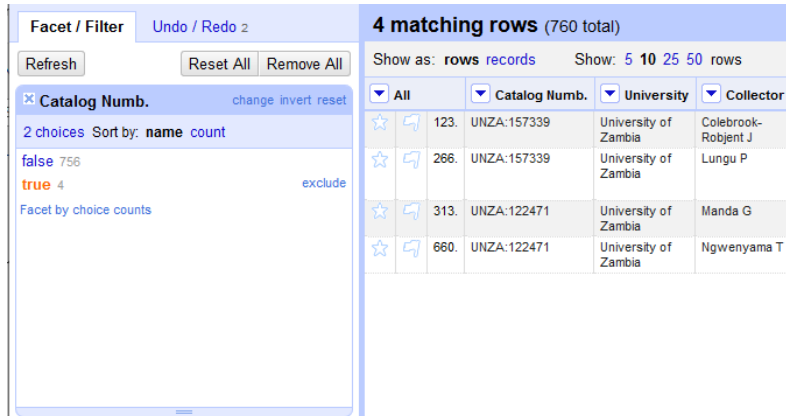1. Go to column Full name and click on ▼ then go to Text facet. Then click on count. The facet will show:

| Facet / Filter | Undo / Redo 0 |
|---|---|

Refresh        Reset All  Remove All

☒ Full name                              change

207 choices  Sort by: name **count**        Cluster

Brachystegia spiciformis 81
Brachystegia stipulata 41
Isoberlinia tomentosa 30
Brachystegia spiciformis 28
Brachystegia boehmii 22
Brachystegia longifolia 22
Brachystegia microphylla 21
Isoberlinia angolensis  20
Brachystegia woodiana 18
Caesalpinioideae 18
Brachystegia glaucescens 17
Colophospermum mopane 15
Brachystegia bakeriana 13
Tessmannia  10

As seen above, *Brachystegia spiciformis* is the first item in the list with 81 specimens, but it is also present in the 4th place with 28 specimens.  This is because there are additional white spaces between *Brachystegia* and *spiciformis.*

2. Fix the error from the Full name column menu, Edit Cells > Common transforms > Collapse consecutive whitespaces.
3. Once the whitespaces are removed, *Brachystegia spiciformis* should only appear once in the list with 109 records.

## 2.2.5. Exercise 5. Faceting and duplicates

1. Go to column catalog in Cat. Numb, and follow the route Facet > Customized facets > Duplicates facet. The facet will show 4 duplicates
2. Click on true, and you'll see the values in the main window:



3. Fix the catalog numbers using the correct values shown in the table below, by clicking on edit directly in the relevant cell. A window will open that allows editing of the value. Once it has been corrected, click Apply.

| | |
|---|---|
| UNZA:122470 | Manda G |
| UNZA:122471 | Ngwenyama T |
| UNZA:157351 | Colebrook-Robjent J |
| UNZA:157339 | Lungu P |

GBIF

## 2.3. FILTERING

### 2.3.1. Exercise 6. Basic filter

1. Go again to Full name column menu and perform a Text facet to visualize the values, then go again to ▼ and click on Text filter, perform the following filters and fix them as shown below:

| Filter | How to fix | Correct value |
|---|---|---|
| **⊠ Full name** <br> sp1 <br> ☐ case sensitive ☐ regular expression | Edit directly in the cell | Eragrostis |
| **⊠ Full name** <br> SP2 <br> ☑ case sensitive ☐ regular expression | Edit directly in the cell, check case sensitive | Julbernardia |
| **⊠ Full name** <br> spp <br> ☐ case sensitive ☐ regular expression | 1. Go to ▼ on Full name, then click Edit cells > Transform… <br> 2. In the text box paste the formula `value.replace(" spp","")` <br> 3. Click OK | Digitaria Diheteropogon Sporobolus Tessmannia |

2. Remember to always close filters before continuing to the next exercise.

### 2.3.2. Exercise 7. Advanced filter I

1. Go to column Full name and perform a Text filter.
2. Check regular expression and case sensitive, then paste the expression ^[a-z]



This regular expression filters the strings in which the first letter is lowercase.

**3. WRITE DOWN THE NUMBER OF RECORDS, AS WELL AS THE NAME IN THE FULL NAME COLUMN, THAT HAVE THE FIRST LETTER AS A LOWERCASE LETTER. YOU WILL REQUIRE THIS INFORMATION TO COMPLETE THE QUIZ.**

4. Perform a correction since the genus should be capitalized.

Note: If you want to know more about regular expressions click here
https://github.com/OpenRefine/OpenRefine/wiki/Understanding-Regular-Expressions

### 2.3.3. Exercise 8. Advanced filter II

1. Go to column Full name and perform a Text filter.
2. Check regular expression and case sensitive, then paste the expression ^[A-Z].*\s[A-Z]



This regular expression filters the strings that start with a capital letter followed by any character, then a space, then a capital letter.

**3. WRITE DOWN THE NUMBER OF RECORDS, AS WELL AS THE NAME IN THE FULL NAME COLUMN, THAT HAVE THE FIRST LETTER OF THE SECOND WORD AS AN UPPERCASE LETTER. YOU WILL REQUIRE THIS INFORMATION TO COMPLETE THE QUIZ.**

4. Perform a correction since the second word of the name should be lowercase.

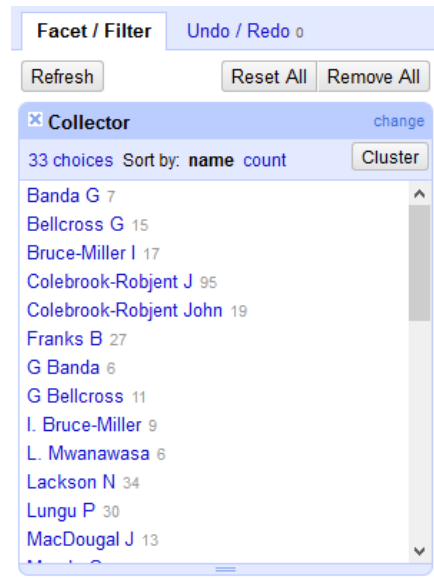Note: If you want to know more about regular expressions click here - https://github.com/OpenRefine/OpenRefine/wiki/Understanding-Regular-Expressions

### 2.3.4. Exercise 9. Advanced filter – additional exercise on advanced filter 1

1. Go to column genus and perform a Text filter.
2. Check regular expression and case sensitive, then paste the correct expression to be able to find the strings in which the first letter is lowercase.

3. **WRITE DOWN THE NUMBER OF RECORDS AS WELL AS THE NAME IN THE GENUS COLUMN THAT HAVE THE FIRST LETTER AS A LOWERCASE LETTER. YOU WILL REQUIRE THIS INFORMATION TO COMPLETE THE QUIZ.**

4. Perform a correction since the genus should be capitalized.

## 2.4. CLUSTERING

### 2.4.1. Exercise 9. Basic clustering

1. Go to Collector, then in the menu column click Text facet.



2. On the top right of the facet window click on Cluster, a new window will appear:



3. Now you can see information about the clusters:
   - Cluster size: the number of different versions that the clustering algorithm believes to be the same.
   - Row count: the number of records with any of the cluster values.
   - Values in cluster: the actual values that the algorithm believes to be the same. There is also the number of records with each particular value, and the possibility to browse the contents of the cluster in a different tab.
   - Merge?: check if values are to be merged into a single standard value.

- **New cell value:** the value to be applied to every record in the cluster. By default, it is the value with most records. You can also click on any value to apply that to the New cell value.

On the right will be a number of graphs providing additional information about the clusters. The top graph will be #Choices in Cluster. In our example, this will range from 1 to 2. Drag the slider to be greater than 1, so that only those 6 clusters that have more than one choice will show.

Note: If you want to know more about clustering click [here](here).

4. Click on Select All and then on Merge Selected & close, you will see a notification message:


Mass edit 167 cells in column Collector   Undo

5. To fix the remaining collectors go again to Cluster in the facet window of Collector.
6. In the Cluster and edit window, go to Keying Function, then select metaphone3. Again, the top graph on the right will be #Choices in Cluster. Drag the slider to be greater than 1, so that only those clusters that have more than one choice will show. The 3 clusters shown in the image below will be identified.



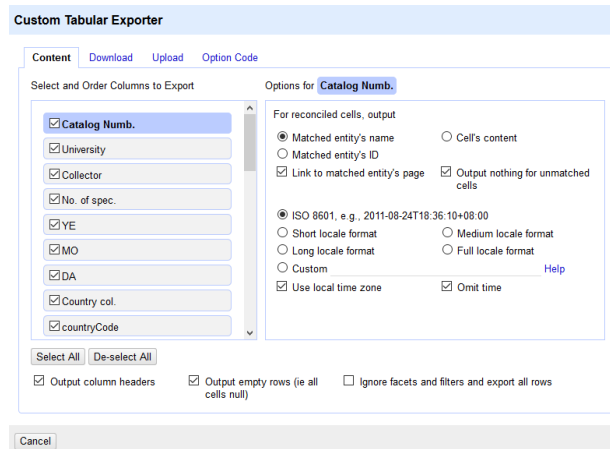7. Click on Select All and then on Merge Selected & close, you will see a notification message:


Mass edit 225 cells in column Collector   Undo

8. Your collectors are now fixed.

## 2.5. EXPORTING

You will have several options for exporting your cleaned data, but the following option is useful in most cases.

1. On the upper right corner click on Export and select Custom tabular exporter.
2. You will see the exportation window:



3. In the content tab you can choose the columns that you want to export, if you select Ignore facets and filters and export all rows, all facets and filterings will be ignored - this is useful if you forget to clear them before exporting.
4. Go to the Download tab and select the separator that you prefer, or whether you would like to download it as an Excel file. Don't modify the other options unless you need to.  You can also Upload the data to a Google spreadsheet.

You can also export the whole project to open it in OpenRefine on another computer by following the route Export > Export project. In this case you are not downloading a data file to open in a spreadsheet or text processor, but rather a GZIP file that will only be accessible through OpenRefine.

# 3. USEFUL LINKS AND REFERENCES
● Documentation
https://github.com/OpenRefine/OpenRefine/wiki/Documentation-For-Users
● Resources list for OpenRefine:
https://github.com/OpenRefine/OpenRefine/wiki/External-Resource

Exercise concept and content developed by Néstor Beltrán, and adapted for OpenRefine 3.2 and an African example by Lizanne Roxburgh.