

The strategy used to clean the data was first to get an idea of what it looked like overall. After looking at a few YouTube videos on how to use Open Refine it was clear that the best way to accomplish this quickly would be to use a text facet and quickly browse through the data of each column. My strategy for cleaning the data is to go column by column. Things I want to take care of are duplicates, error values, and consolidating data that is unnecessarily separated. I chose to keep the N/A values because I did not want to lose that data just yet. My thought process is that I don't want to get rid of any data that I may want later and if I really don't need it later I can discard it then.

The location column would be the most straight forward to clean up since the values contained in the column would be the most familiar and easiest for spotting errors and inefficiencies. This data was for the most part well formatted but a name of a single place could be typed out different ways causing more rows than necessary. Using the text facet, I quickly sifted through the data to edit certain columns to conform to others.

**Refine** OPEN data\_mining.txt Perm

**Facet / Filter** Undo / Redo 1

Refresh Reset All Remove All

**Barcelona, Spain** change invert reset

245 choices Sort by: name count Cluster

Tokyo, Japan	1
Ã...ÃÃÃÃÃÃÃÃÃÃÃÃÃÃÃÃ Ã¡Ã¢Ã£Ã¤Ã¥Ã¦Ã§Ã¨Ã©ÃªÃ«Ã¬Ã­Ã®Ã¯Ã°Ã±Ã²Ã³Ã´ÃµÃ¶Ã·Ã¸Ã¹ÃºÃ»Ã¼Ã½Ã¾Ã¿ÃÃÃÃÃÃÃÃÃÃÃÃÃÃÃÃÃ Ã¡Ã¢Ã£Ã¤Ã¥Ã¦Ã§Ã¨Ã©ÃªÃ«Ã¬Ã­Ã®Ã¯Ã°Ã±Ã²Ã³Ã´ÃµÃ¶Ã·Ã¸Ã¹ÃºÃ»Ã¼Ã½Ã¾Ã¿	1
<b>Aberdeen, Scotland</b>	1
Aberdeen, Scotland (UK)	1
Adelaide, Australia	1
Alicante, Spain	1
Anchorage, Alaska, USA	1
Arras, France	1
Astrophysical Journal - Supplement Serie	1
Athens, Greece	2
Atlanta, GA	1

Since this was my first time using the program I was curious to see how quickly I could go through this data and group up the necessary rows. The time it took was 10 minutes which was a little tedious but a rather quick way to group up the rows with the location. I had dropped down to 184 rows.



The second column to be cleaned was the acronym column because immediately I noticed extra information that could be easily removed. Many rows had acronym names but also a prefix including ACM. To clear these up, the transform function was used ( `value.replace(arguments)` )



This transformation is essentially the same as doing a find and replace all, like in Excel. There were a few instances that I used this approach to scrub the data in this column as well as the name column. For example, many of the values had the unnecessary information of the year which needed to be erased.

Custom text transform on column ICFIP--EI, Scopus 2018

Expression

Language General Refine Expression Language (GREL) ▼

value.replace("ACM--", "")

No syntax error.

Preview

History

Starred

Help

row	value	value.replace("ACM--", "")
1.	BDIOT - 2017	BDIOT - 2017
2.	EMSA 2018	EMSA 2018
3.	AI 2017	AI 2017
4.	ADCOM 2017	ADCOM 2017
5.	CCSEA 2018	CCSEA 2018
6.	OPTLJ 2017	OPTLJ 2017
7.	LINKP 2017	LINKP 2017

On error

☒ keep original
 ☐ set to blank
 ☐ store error

☐ Re-transform up to 10 times until no change

OK

Cancel

It was clear that some of the formatting in the values could be improved. Some of the cells had numbers and values beside one another with no space in between. Again, the transform function was used using the General Refine Expression Language. This time using regular expressions to match to letters and then create a space to the numbers.

Custom text transform on column ICFIP--EI, Scopus 2018

Expression

Language General Refine Expression Language (GREL) ▼

value.replace(/(\p{IsAlphabetic})(?=\d)/, '\$1 ')

No syntax error.

Preview

History

Starred

Help

Another instance of needing to scrub acronyms from the data.



A screenshot of a data table with a light blue header and a white body. The table has two columns: the first column contains text labels, and the second column contains numerical counts. The labels are in blue text, and the counts are in black text. The table is partially visible, showing the bottom portion of the data.

IEEE - ICBDA	1
IEEE - ICKEA	1
IEEE Big Data	1
IEEE DSAA'	1
IEEE FRUCT	1
IEEE HPSC	1
IEEE ISI	3
IEEE ITNEC	1
IEEE ITOEC	1
IEEE TDSC Journal SI	1
IEEE TETC SI	1
IEEE TETCI	1

Getting a little bit more complex and using regular expressions it was clear that values between parentheses were not needed. This was taken care of in two steps. The first step was to clean the data of the values in between parentheses followed by simply removing the parentheses that were left over. I struggled with writing the regular expression myself for a while until finally turning to the internet to help find a function that would accomplish what I needed to do as seen below.

## Custom text transform on column 2018 2nd International Conference on Frontiers of Image Processing (ICFIP 2018)--EI Compendex, Scopus, and ISI CPCS

Expression

Language General Refine Expression Language (GREL) ▾

`value.split('(')[0]+'('+value.split(')')[1]`

No syntax error.

**Preview**

History

Starred

Help

row	value	<code>value.split('(')[0]+'('+value.split(')')[1]</code>
1.	International Conference on Big Data and Internet of Things (BDIOT)--Ei Compendex, Scopus and ISI	International Conference on Big Data and Internet of Things ()--Ei Compendex, Scopus and ISI
2.	International Conference on Embedded Systems and Applications	International Conference on Embedded Systems and Applications() International Conference on Embedded Systems and Applications
3.	International Conference on Artificial Intelligence and Applications	International Conference on Artificial Intelligence and Applications() International Conference on Artificial Intelligence and Applications

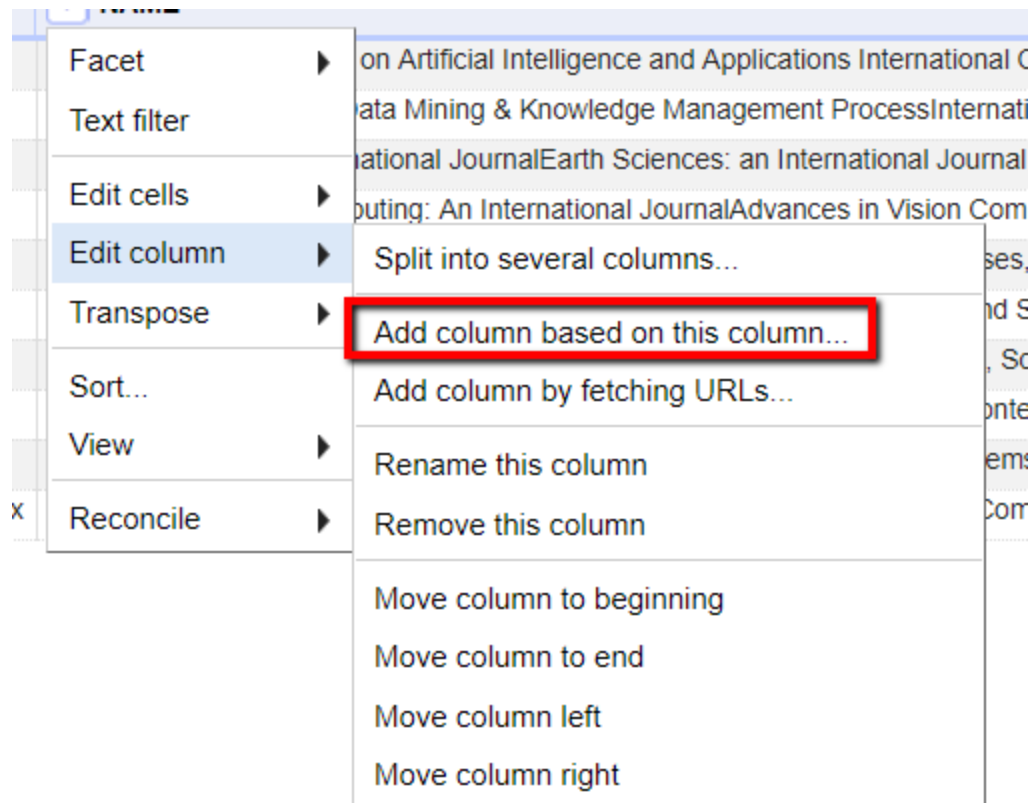
On error ☒ keep original ☐ Re-transform up to  times until no change  
☐ set to blank  
☐ store error

OK

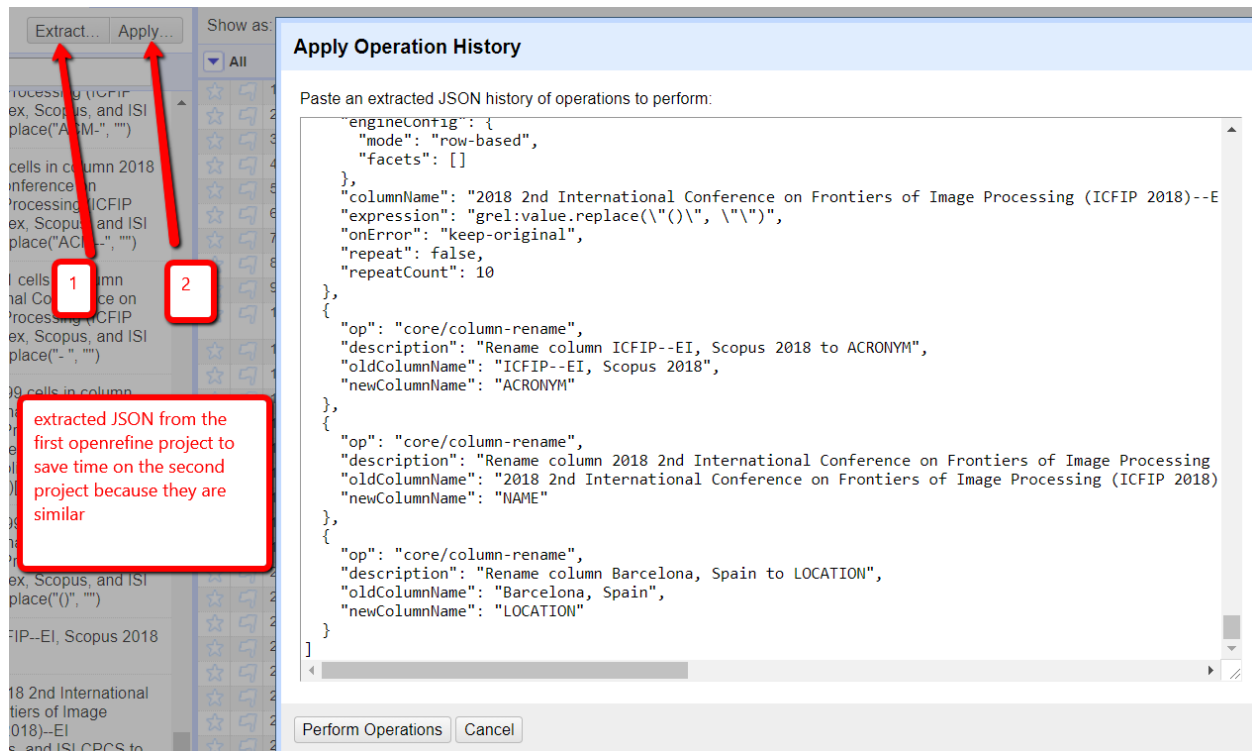
Cancel

**Text transform on 399 cells in column 2018 2nd International Conference on Frontiers of Image Processing (ICFIP 2018)--EI Compendex, Scopus, and ISI CPCS: `grel:value.replace("()", "")`** Undo

It was also to get rid of the .Journals since we are only interested in conferences. To accomplish this the functionality to create a new column bases on a column was used. The function was a .contains to check if the column contains the value I was looking for. From there it is possible to filter for only the value you needed and then delete the rest.



My crawler was run 4 times to crawl all the pages needed. This means that I have 4 very similar data sets. I used the functionality from Open Refine to extract the JSON code of instructions to clean the data set and applied it to my other outputs of the crawler. This helped but a large start on the data cleansing process. However to be thorough I followed a similar process manually on all subsequent datasets to make sure nothing was missed.



Ultimately, the hardest part of the phase 1 of this assignment was the crawling of the web page. I had a very difficult time extracting the data I needed from the html table because the table itself was not built in a very intuitive way. The rows were not specified as headers when necessary and the values were not stored by class name so it took a lot of iterations through each row to extract each piece of data needed. I assume this was the point of the exercise because data in the wild will never be perfect. Also, I am quite confident there is an easier way to clean the data using Open Refine and there is probably some advanced functionality that I may have missed out on.