

Data cleaning

Name: Xinyu Hu

Student ID: 29350379

Email: xh2n17@soton.ac.uk

Overview of the dataset

This dataset is about the US government spending on the project. However, it contains so many errors. In order to make this dataset available for further use, we have to detect all the errors and find out solutions to fix them. Here are some errors I find:

1. Multiple representation problem in 'Department Name'

When we go through the first few rows in the 'department name' column, we can discover there are some abbreviations like 'DoA', the 'Agriculture Department' or the 'Department of Agriculture', which all represent the same meaning and should be combined as one.

Solution:

We first apply the 'trim leading and trailing whitespace' method this column, and then access the 'cluster' bottom in the facet:

The screenshot shows the 'Cluster & Edit column "Column3"' interface. It includes a description of the clustering feature, settings for Method (nearest neighbor), Distance Function (levenshtein), Radius (1.0), and Block Chars (6). It indicates that 2 clusters were found. The main table lists the clusters:

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
2	162	<ul style="list-style-type: none">Department of Agriculture (153 rows)Department of Agraculture (9 rows)	<input checked="" type="checkbox"/>	Department of Agriculture
2	259	<ul style="list-style-type: none">Department of Commerce (247 rows)Department of Comerce (12 rows)	<input checked="" type="checkbox"/>	Department of Commerce

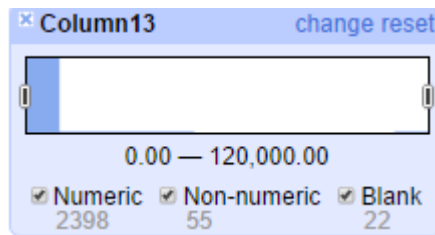
On the right side, there are three sliders for '# Rows in Cluster', 'Average Length of Choices', and 'Length Variance of Choices'. At the bottom, there are buttons for 'Select All', 'Unselect All', 'Export Clusters', 'Merge Selected & Re-Cluster', 'Merge Selected & Close', and 'Close'.

Finally, edit the abbreviations to the full name in the facet.

2. Redundant information and mixed use of numerical scales in 'Lifecycle cost'

As we can see, there are some redundant information like the '(\$m)' appearing so many times. Since we default all the numbers in the column are in the units of '\$million', we have to accurately point out each of them, and delete the '(\$m)' by using the tool.

Solution: as we apply the numeric facet to this column, we find out that there is no numeric cell, so we first apply ().toNumber action.



Then in order to get rid of the data with '(\$m)', we use the `value.replace()` function to deal with.

3. Mixed formats in the 'Completion date'

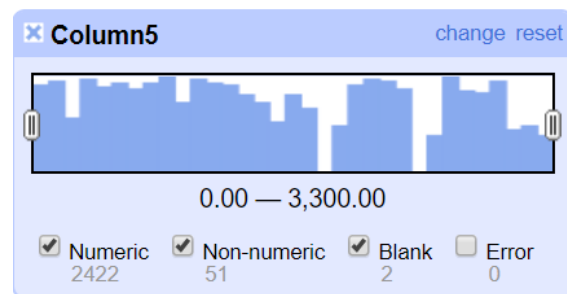
The mixed format like '31/03/2012' and '2012-30-09' are quite obvious in the 'Completion date' column. To deal with this problem, we just transfer the one with '-' to the one with '/', as the slash one is more common to see.

Solution: use the method `toString(toDate(value),"dd/MM/yyyy")` to make the transformation.

4. Mixed formats in the project ID

When we check the inconsistency in the 'project ID' column, it's easy to see there are different kinds of formats like '1,923' and '1025'. To uniform the data in this column, we have to remove the comma.

Solution: we first check the numeric facet of this part, realizing that there are 50 non-numeric rows. Therefore we use the `value.replace()` function to cope with it.



5. Duplicated formation in the project ID

Noticing that each ID may be unique, we have to exclude all duplicated ones. To complete this step, we first check how many rows are identical. And then use the 'Blank down' method to cope with it.

Column5	change
2 choices	Sort by: name count
false	2465
true	10
Facet by choice counts	

6. Summation records

For the summation records in the dataset, we just delete them directly.

Solution: to remove them, we first choose the text facet in the second column, and then 'star' them all and choose the 'remove all matching rows'.