

Adam Reilly
Lab 8

1. 5,377 rows do not have states (This is the number of "blank" states listed in the text facet).
2. 4,362 rows do not have zip codes
3. There are 380,136 valid zip codes and 4,362 blank zip codes (which I would call invalid)
4. This is what occurs when radius is changed to 3.

Cluster & Edit column "Local"

This feature helps you find groups of different cell values that might be alternative representations of the same thing. For example, the two strings "New York" and "new york" are very likely to refer to the same concept and just have capitalization differences, and "Gödel" and "Godel" probably refer to the same person. [Find out more ...](#)

Method: Distance Function: Radius: Block Chars: 4 clusters found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
2	85	<ul style="list-style-type: none">California (84 rows)Caifornia (1 rows)	<input type="checkbox"/>	<input type="text" value="California"/>
2	795	<ul style="list-style-type: none">Alaska (791 rows)alaska (4 rows)	<input type="checkbox"/>	<input type="text" value="Alaska"/>
2	61	<ul style="list-style-type: none">Tajikistan (36 rows)Pakistan (25 rows)	<input type="checkbox"/>	<input type="text" value="Tajikistan"/>
2	805	<ul style="list-style-type: none">Indonesia (797 rows)Micronesia (8 rows)	<input type="checkbox"/>	<input type="text" value="Indonesia"/>

Rows in Cluster: —

Average Length of Choices: —

Length Variance of Choices: —

Select All Unselect All Export Clusters Merge Selected & Re-Cluster Merge Selected & Close Close

I would merge the first and second clusters, but not the third or fourth.

5. Changing the block characters to 2 greatly expands the number of clusters. Two clusters that I would choose to change are:

4	87	<ul style="list-style-type: none">California (84 rows)Caliofrnia (1 rows)Calfiornia (1 rows)Caifornia (1 rows)	<input type="checkbox"/>	<input type="text" value="California"/>
3	796	<ul style="list-style-type: none">Alaska (791 rows)alaska (4 rows)Alska (1 rows)	<input type="checkbox"/>	<input type="text" value="Alaska"/>

6. When you cluster on a place column, OpenFile groups together words based on letter similarity and proximity. Block will determine the number of letters that need to match in words for them to be included. I believe this is because we opted for Levenhstein for “distance”, and the radius determines the maximum allowable distance. An additional functionality that could be provided is to pick and choose words to match when instead of making you pick a whole cluster. For example of the flaw, see below. You can only merge the entire cluster, not just the words that match.

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
10	809	<ul style="list-style-type: none"> Alaska (791 rows) Alabama (6 rows) alaska (4 rows) Candaa (2 rows) Alaksa (1 rows) Albania (1 rows) Alaa (1 rows) Alaka (1 rows) Malawi (1 rows) Alska (1 rows) 	<input type="checkbox"/>	Alaska

7. Here is a schema of the Levenshtein's distance for number 7.

		1	2	3	4	5	6	7	8	9	10
			G	U	M	B	A	R	R	E	L
1	+	0	1	2	3	4	5	6	7	8	9
2	G	1	0	1	2	3	4	5	6	7	8
3	U	2	1	0	1	2	3	4	5	6	7
4	N	3	2	1	1	2	3	4	5	6	7
5	B	4	3	2	2	1	2	3	4	5	6
6	A	5	4	3	3	2	1	2	3	4	5
7	R	6	5	4	4	3	2	1	2	3	4
8	E	7	6	5	5	4	3	2	2	2	3
9	L	8	7	6	6	5	4	3	3	3	2
10	L	9	8	7	7	6	5	4	4	4	3