# Open Refine

for efficient data cleaning

# The Goal: Get More Reliable Answers

Standardization is how we do that. Let's take a look at our PPP data (use your own copy) and find some problems of consistency. Use what you've learned from the Evaluating Data presentation.

- Where/how do similar values look different?
- Are there abbreviations or truncated data?
- Weird/missing punctuation or characters?

Enter some examples (and describe them) in this spreadsheet.

# Get Started with OpenRefine

Go to the applications folder and click Open Refine to open. This should automatically launch a browser window. If it doesn't you can go here: http://127.0.0.1:3333/

Then download a CSV version of the PPP data from your sheet.

# Import the PPP data

**OpenRefine**

Create Project

Open Project

**Create a project by importing data. What kinds of data files can I import?**

TSV, CSV, *SV, Excel (.xls and .xlsx), JSON, XML, RDF as XML, and Google Data documents are all supported.

Get data from

**This Computer**

Web Addresses (URLs)

Clipboard

Database

Google Data

Locate one or more files on your computer to upload:

Choose Files   No file chosen

Next »

# Check the Results

Once uploaded, you should see a preview. Check it out: Does it look right?

Note the options in the lower right-hand corner of the screen. Has OpenRefine selected the right ones?

# How It Should Look

# If It Looks Right

Click on "Create Project" and get started. It'll look like this:

# Facets

Open Refine relies on things called 'facets' to help us clean up data.

Click on the triangle next to the column header 'CITY' and select 'Text Facet' from the dropdown menu.

# Facet Box

A 'facet box' will appear on the left side of your workspace. Click the Cluster button to use one of Refine's most powerful features.

**Facet / Filter**    Undo / Redo 0 / 0

Refresh      Reset All   Remove All

☒ ▬ **city**      change

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

100 East Main St  1
1706 HAMPSHIRE GREEN LANE APT 23  1
aberdeen  1
Aberdeen  34
ABERDEEN  12
Abingdon  36
ABINGDON  19
Abington  1
Accident  9
ACCIDENT  1

# What is Clustering?

Clustering involves using different characteristics of words to group likely identical ones together. Some clustering techniques rely on having letters in common. Others group together words that sound alike even if they are spelled differently. Each method has strengths and weaknesses, so it's useful to try more than one.

# Putting Like Things Together

## Cluster & Edit column "city"

This feature helps you find groups of different cell values that might be alternative representations of the same thing. For ex: york" are very likely to refer to the same concept and just have capitalization differences, and "Gödel" and "Godel" probably

Method [ key collision        ⌄ ]          Keying Function [ fingerprint        ⌄ ]

| Cluster Size | Row Count | Values in Cluster | Merge? | New Cell Value |
|---|---|---|---|---|
| 5 | 277 | • Upper Marlboro (220 rows)<br>• UPPER MARLBORO (53 rows)<br>• upper marlboro (2 rows)<br>• Upper marlboro (1 rows)<br>• upper Marlboro (1 rows) | ☐ | Upper Marlboro |
| 5 | 115 | • Fort Washington (96 rows)<br>• FORT WASHINGTON (16 rows)<br>• Fort WASHINGTON (1 rows)<br>• fort Washington (1 rows)<br>• fort washington (1 rows) | ☐ | Fort Washington |
| 5 | 725 | • Silver Spring (546 rows)<br>• SILVER SPRING (170 rows)<br>• silver spring (6 rows)<br>• Silver spring (2 rows)<br>• silver Spring (1 rows) | ☐ | Silver Spring |

# Don't Trust, Verify

Is this a match we want? Probably, and we can check. If you hover, you get the option to 'Browse this Cluster'. Click that. A new window will pop open showing just the rows that would be included in that potential cluster.

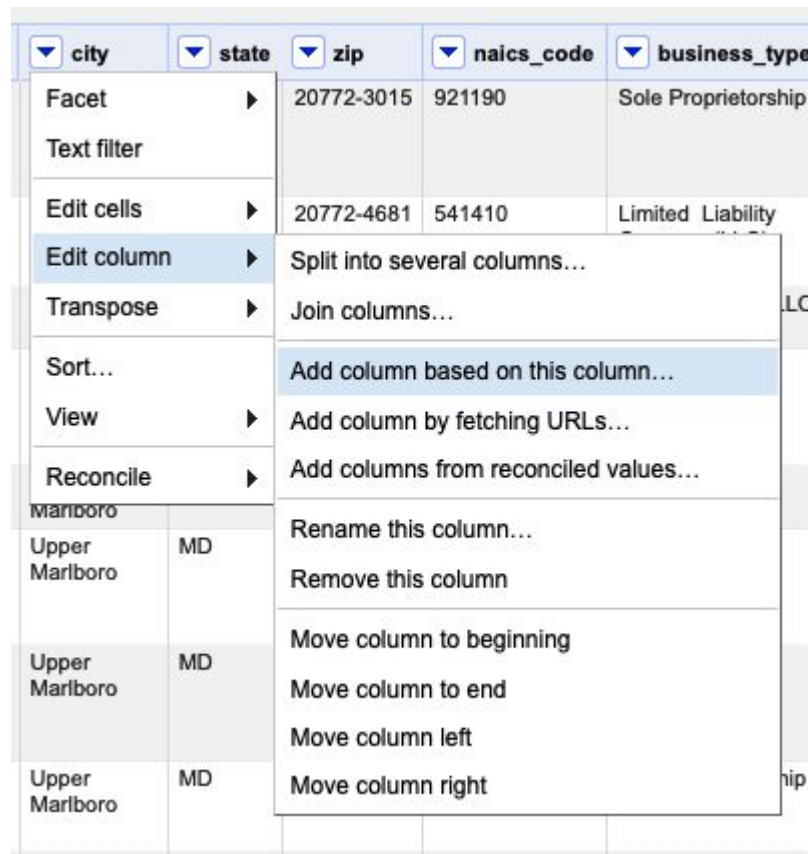In this case we can see that the cities have very similar zip-codes. They're probably a match.

Close the additional window to return to the complete data.

# Make a Copy First!

Remember how Sophie said "make a copy of the original data"? We can do that with columns, too, so we can see what edits we make.

Go to the "city" column and choose Edit column -> Add column based on this column...

Then give the new column a name like "city_clean" and cluster on that column.

# Everyday I'm Clustering

Browse the suggested clusters. Some will look good. Some won't. If you see one that makes sense, check the 'merge' checkbox. Then make sure the 'New cell value' is appropriate. If not, you can edit it. When you're finished with a cluster, or clusters, you can click 'Merge selected and re-cluster' or 'Merge selected and close.' It's a good idea to work your way through the different clustering options.

# But Wait, There's More!

OpenRefine also has some other useful functions you can explore.

- Splitting cells into multiple columns
- Trimming unneeded white space
- Changing column data types
- Changing cases