

CS 513 Theory and Practice of Data Cleaning

Final Project Report

Team Name: Dream of Data

Mondal, Valentina vmondal2@illinois.edu

Panchal, Harsh Mahendrakumar hmpanch2@illinois.edu

Balasubramanian, Sembian sembian2@illinois.edu

Table of Contents

1. <i>Overview and initial assessment of the dataset [25%]</i>	3
a. A clear description of the structure and content of the dataset [3%]	3
b. A comprehensive list of data quality issues [7%].....	3
c. Identifying a feasible use case and the essential data cleaning goals [10%]	4
d. Fitness of Use.....	8
e. Identifying use cases for which the dataset is already clean and use cases for which it will never be clean enough or usable [5%]	9
2. <i>Data cleaning with OpenRefine (and other tools) [40%]</i>	10
a. Identifying the appropriate data cleaning steps for the use case	10
b. A clear description of the data cleaning steps with supplemental information.....	10
c. Quantifying the results of cleaning (e.g., provide a table of changes along with appropriate quantification) [5%]	38
3. <i>Developing a relational schema [15%]</i>	42
a. Identifying the appropriate integrity constraints [5%]	42
b. Loading data into a database with proper schema [3%]	43
c. Writing queries to check the integrity constraints [7%]	46
4. <i>Creating a workflow model [10%]</i>	51
a. Identifying the key inputs, outputs of your workflow along with the dependencies [3%]	51
b. A visual representation of your overall workflow, e.g., using YesWorkflow [4%]	53
c. A visual representation of your OpenRefine workflow, e.g., using OR2YWTool [3%]	54
5. <i>Other factors [10%]</i>	54
a. Further analysis/takeaways/challenges [5%] (In Progress)	54
REFERENCES:	54
APPENDIX	54

1. Overview and initial assessment of the dataset [25%]

a. A clear description of the structure and content of the dataset [3%]

Organization:

U.S. Agricultural Marketing Service creates domestic and international marketing opportunities for U.S. producers of food, fiber, and specialty crops. AMS also provides the agriculture industry with valuable services to ensure the quality and availability of wholesome food for consumers across the country.

The Farmers Market Directory lists (<https://www.ams.usda.gov/local-food-directories/farmersmarkets>) markets that feature two or more farm vendors selling agricultural products directly to customers at a common, recurrent physical location. The U.S. Agricultural Marketing Service maintains the directory listing designed to provide customers with convenient access to farmers' market listings to find operations offering locally grown markers:

The listings include some key details like the market locations, directions, operating times, product offerings, accepted forms of payment, and more. The Farmers Market Directory listing website (Farmers Market, 2020) provides an option for export to CSV format. Additionally, it also has API access, and the listing does not provide a dictionary of all fields or details about how the data is entered and validated.

Architecture:

The Farmers market listing contains 59 columns. The first column FMID is an integer column and can be used as a unique identifier for each farmers' market. The second column MarketName is a text column which allows special characters as well. The next set of columns – Website, Facebook, Twitter, Youtube and OtherMedia allows URL to their website and/or social media pages/handles and contains NULL values. The data in these columns looks legible but contains blank or NULL values as well. The next five columns are related to location of farmers' market – street, city, country, state and zip. Next set of columns are related to seasons – date and time period, this provides market hours and days when it is open. The next three columns – x, y and location, represents latitude, longitude and additional details about market's location. The rest columns except last one (updateTime) are Boolean values (Y or N) representing availability of different products in a market. The last column updateTime is a timestamp column, represents when a record was inserted/updated.

b. A comprehensive list of data quality issues [7%]

This dataset contains lot of quality issues like inconsistent formatting and data values, repetitive information with different column names and redundant values etc. This might be because of lack of validation rules for each column and not following any ISO formatting. We could classify data quality issues using dimensions of data quality:

1. Uniqueness: FMID seems to be unique id associated with each row in the dataset, but it needs to be validated
2. Validity: The columns related to seasons have inconsistent granularity of time period – certain values are date1 to date2 and other month1 to month2.
3. Accuracy:
 - i. MarketName is a free text column and allows special characters and hence, this column contains consecutive spaces, leading and trailing spaces and names which are technically same but have few different characters making them look different.

- ii. Latitude and Longitude are stored with column names x and y, column names aren't meaningful.
- 4. Completeness
The columns related to URL don't have validation; making them allow blank, NULL and text instead of URL.
- 5. Consistency
The set of columns related to location have leading and/or trailing white spaces, values in column zip contains more than 940 blank values and two letter state names or "-".



Location column supposed to provide specific geographic location of a market, but it contains irrational values like other, educational institution, healthcare institution and blanks.

The columns related to availability of a product contains not only Y or N values, but it contains blank or “-”.

- 6. Timeliness
updateTime columns should represent timestamp but it has inconsistent formatting and values like year, MON DD YYYY HH: MM: SS AM/PM, MM/DD/YY HH: MM: SS AM/PM. Moreover, we don't have enough information on how this dataset can be consumed (i.e. API) and can be used in data pipeline as real-time feed.

c. Identifying a feasible use case and the essential data cleaning goals [10%]

This dataset can be used to answer questions like availability of Farmers' Market and markets with availability of group of grocery items like fruits, vegetables, meat etc.

Data Cleaning Goals-

ID	Column Name	Steps	Desired outcome
1	FMID	1.1 Convert to Numbers	FMID needs to be 7-digit unique number to use it as Primary Key
2	Market Name	2.1 Remove trailing and leading spaces including consecutive spaces	Market Names are accurate

		2.2 Remove duplicates using clustering 2.3 Set data type to Text. 2.4 Convert special character ‘&’ to ‘and’	
3	Website	3.1 Remove trailing and leading spaces including consecutive spaces 3.2 Remove duplicates using clustering 3.3 Convert all URL to lowercase 3.4 Validate URL	Consistent web URL structure “http(s)://(www.)” Spaces in URL are ignored Leave blank URL’s
4	Facebook	4.1 Remove trailing and leading spaces including consecutive spaces 4.3 Remove all duplicates using clustering	Leave blank URL’s Spaces in URL are ignored
5	Twitter	5.1 Remove trailing and leading space including consecutive spaces 5.3 Remove all duplicates using clustering	Leave blank URL’s Spaces in URL are ignored
6	Youtube	6.1 Remove trailing and leading spaces including consecutive spaces 6.2 Remove duplicates using clustering	Leave blank URL’s Spaces in URL are ignored
7	OtherMedia	7.1 Remove trailing and leading spaces including consecutive spaces	Leave blank URL’s Spaces in URL are ignored
8	street	8.1 Remove trailing and leading spaces including consecutive spaces 8.2 Identifying redundant 8.3 Change street names case sensitive to title case 8.3 Use clusters to standardize street names with common names (e.g. St > Street) 8.4 Rename column “street” to “Street” to be consistent with other column names	Street names are consistent Ignore or replace blank street names with NULL
9	city	9.1 Remove trailing and leading spaces including consecutive spaces 9.2 Remove duplicates and case sensitive using clustering	City names are consistent
10	County	10.1 Remove trailing and leading spaces including consecutive spaces 10.2 Remove duplicates and case sensitive using clustering	Country names are consistent
11	State	11.1 Remove trailing and leading spaces including consecutive spaces 11.2 Remove duplicates and case sensitive using clustering	State names are consistent
12	zip	12.1 Convert to Numbers 12.2 Remove values other than numbers (10 nonnumeric value like FL, IL, OR and 947 blank value) 12.3. Standardize Zip code to 5 digit and split “99999-99999” > (99999)	The values in Zip codes must be numbers (Assumption - The dataset is of the United States)
13	Season1Date	13.1 Remove trailing and leading spaces including consecutive spaces 13.2 Create new columns From and To based on	Consistent From and To Dates

		separator value "to" 13.3 Convert to ISO Standard Date (or only Month) and remove default timestamp	New columns "Season1FromDate" & "Season1ToDate"
14	Season1Time	14.1 Remove trailing and leading spaces 14.2 Convert to lowercase 14.3 Split into several columns based on the separator ";"	Consistent hours of store
15	Season2Date	15.1 Remove trailing and leading spaces including consecutive spaces 15.2 Create new columns From and To 15.3 Convert to ISO Standard Date (or only Month) and remove default timestamp	Consistent From and To Dates New columns "Season2FromDate" & "Season2ToDate"
16	Season2Time	16.1 Remove trailing and leading spaces 16.2 Convert to lowercase 16.3 Split into several columns based on the separator ";"	Consistent hours of store
17	Season3Date	17.1 Remove trailing and leading spaces including consecutive spaces 17.2 Create new columns From and To 17.3 Convert to ISO Standard Date (or only Month) and remove default timestamp	Consistent From and To Dates New columns "Season3FromDate" & "Season3ToDate"
18	Season3Time	18.1 Remove trailing and leading spaces 18.2 Convert to lowercase 18.3 Split into several columns based on the separator ";"	Consistent hours of store
19	Season4Date	19.1 Remove trailing and leading spaces including consecutive spaces 19.2 Create new columns From and To 19.3 Convert to ISO Standard Date and remove default timestamp	Consistent From and To Dates New columns "Season4FromDate" & "Season4ToDate"
20	Season4Time	20.1 Remove trailing and leading spaces 20.2 Convert to lowercase 20.3 Split into several columns based on the separator ":"	Consistent hours of store
21	x	21.1 Convert to Numbers 21.2 Rename column to Latitude	Accurate Latitude
22	y	22.1 Convert to Numbers 22.2 Rename column to Longitude	Accurate Longitude
23	Location	23.1 Remove trailing and leading spaces including consecutive spaces	Consistent location description values
24	Credit	24.1 Remove trailing and leading spaces 24.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
25	WIC	25.1 Remove trailing and leading spaces 25.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
26	WICCash	26.1 Remove trailing and leading spaces 26.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
27	SFMNP	27.1 Remove trailing and leading spaces 27.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
28	SNAP	28.1 Remove trailing and leading spaces	Consistent Boolean values

		28.2 Replace values to NULL if other than "Y" or "N"	
29	Organic	29.1 Remove trailing and leading spaces 29.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
30	Bakedgoods	30.1 Remove trailing and leading spaces 30.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
31	Cheese	31.1 Remove trailing and leading spaces 31.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
32	Crafts	32.1 Remove trailing and leading spaces 24.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
33	Flowers	33.1 Remove trailing and leading spaces 33.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
34	Eggs	34.1 Remove trailing and leading spaces 34.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
35	Seafood	35.1 Remove trailing and leading spaces 35.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
36	Herbs	36.1 Remove trailing and leading spaces 36.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
37	Vegetables	37.1 Remove trailing and leading spaces 37.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
38	Honey	38.1 Remove trailing and leading spaces 38.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
39	Jams	39.1 Remove trailing and leading spaces 39.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
40	Maple	40.1 Remove trailing and leading spaces 40.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
41	Meat	41.1 Remove trailing and leading spaces 41.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
42	Nursey	42.1 Remove trailing and leading spaces 42.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
43	Nuts	43.1 Remove trailing and leading spaces 43.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
44	Plants	44.1 Remove trailing and leading spaces 44.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
45	Poultry	45.1 Remove trailing and leading spaces 45.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
46	Prepared	46.1 Remove trailing and leading spaces 46.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
47	Soap	47.1 Remove trailing and leading spaces 47.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
48	Trees	48.1 Remove trailing and leading spaces 48.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
49	Wine	49.1 Remove trailing and leading spaces 49.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
50	Coffee	50.1 Remove trailing and leading spaces 50.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
51	Beans	51.1 Remove trailing and leading spaces 51.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
52	Fruits	52.1 Remove trailing and leading spaces 52.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values

53	Grains	53.1 Remove trailing and leading spaces 53.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
54	Juices	54.1 Remove trailing and leading spaces 54.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
55	Mushrooms	55.1 Remove trailing and leading spaces 56.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
56	PetFood	56.1 Remove trailing and leading spaces 56.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
57	Tofu	57.1 Remove trailing and leading spaces 57.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
58	WildHarvested	58.1 Remove trailing and leading spaces 58.2 Replace values to NULL if other than "Y" or "N"	Consistent Boolean values
59	updateTime	59.1 Remove trailing and leading spaces 59.2 Convert to ISO Standard Date	Consistent timestamp value

d. Fitness of Use

Columns	Use/Application
FMID Market Name Season1Time Season2Time Season3Time Season4Time	Market Hour of Operation Per Season?
FMID MarketName Website Facebook Twitter Youtube OtherMedia	Market affiliation in online?
FMID Market Name Season1Date Season2Date Season3Date Season4Date	Markets by Seasonality (Summer/Winter/Fall)?
FMID MarketName Latitude Longitude Zip	Availability of Market in Proximity?

FMID MarketName Credit WIC WICash SFMNP SNAP	Market having specific payment type?
FMID MarketName Organic Bakedgoods Cheese Crafts Flowers Eggs Seafood Herbs Vegetables Honey Jams Maple Meat Nursery Nuts Plants Poultry Prepared Soap Trees Wine Coffee Beans Fruits Grains Juices Mushrooms PetFood Tofu WildHarvested	Products and their type of Products in a specific Market? (Products can be classified into the following categories - a) Organic b) Beverages c) Dairy d) Meat e) Produce f) Dry Goods g) Cleaners h) Canned/Jarred i) Garden j) Other
FMID MarketName updateTime	Market Updated Latest or Oldest?

- e. Identifying use cases for which the dataset is already clean and use cases for which it will never be clean enough or usable [5%]

Overall, column zip, x and y are quite clean, and it can be used to map out density/ heat map of availability of Farmers' Market across the United States.

The columns related to Market's website and social media pages are not reliable and can't be used for any kind of analysis. If we want to derive some additional information like how many people have liked a page on Facebook and Twitter, promotional offers on social media, pickup availability, etc. would be difficult.

2. Data cleaning with OpenRefine (and other tools) [40%]

- Identifying the appropriate data cleaning steps for the use case
- A clear description of the data cleaning steps with supplemental information

1. FMID

FMID seems to be unique for each row. To validate it, we converted this column into Number and analyzed the results using text facet. The count of each number is "1". The total number of FMID were 8806 and all of them were converted to number, yielding in 8806 unique FMID.

The screenshot shows the OpenRefine interface with a table of data. The first column contains FMIDs (e.g., 1011100, 1009845, 1005586). The second column contains market names and addresses (e.g., '12 South Farmers Market', '125th Street Fresh Connect Farmers' Market', '12th & Brandywine Urban Farm'). The third column contains URLs (e.g., 'https://sites.google.com/site/aledoniafarmersmarket/', 'http://www.125thStreetFarm.com', 'https://www.facebook.com/StearnsHomesteadFarm'). A context menu is open over the URL column, specifically over the first URL cell. The menu path 'Edit cells' -> 'Transform...' -> 'Common transforms' is visible. A sub-menu is open under 'Common transforms' with the following options: 'Trim leading and trailing whitespace', 'Collapse consecutive whitespace', 'Unescape HTML entities', 'Replace Smart quotes with ascii', 'To titlecase', 'To uppercase', 'To lowercase', 'To number' (which is highlighted in blue), 'To date', 'To text', 'To null', and 'To empty string'. The 'To number' option is currently selected.

Facet	▶	https://sites.google.com/site/aledoniafarmersmarket/	https://www.facebook.co
Text filter	▶		
Edit cells	▶	Transform...	ad.com StearnsHomesteadFarm
Edit column	▶	Common transforms	Trim leading and trailing whitespace Collapse consecutive whitespace Unescape HTML entities Replace Smart quotes with ascii
Transpose	▶		To titlecase To uppercase To lowercase
Sort...	▶		To number To date To text
View	▶		To null To empty string
Reconcile	▶		
Avenue			
1011100	12 South Farmers Market	http://www.125thStreetFarm.com	
1009845	125th Street Fresh Connect Farmers' Market		
1005586	12th & Brandywine Urban Farm		

2. Market Name

- Trim leading and trailing whitespace – The number of cells affected by this transformation were 410
- Collapse consecutive whitespace – The number of cells affected by this transformation were 51

Facet	►	://sites.google.com/site/caledoniafarmersmarket/	https://www.facebook.com/Danville.v
Text filter			
Edit cells	►	Transform...	StearnsHomesteadFarmersMarket
Edit column	►	Common transforms ►	Trim leading and trailing whitespace
Transpose	►	Fill down	Collapse consecutive whitespace
Sort...		Blank down	Unescape HTML entities
View	►	Split multi-valued cells...	Replace Smart quotes with ascii
Reconcile	►	Join multi-valued cells...	To titlecase
Avenue		Cluster and edit...	To uppercase
12 South Farmers Market	http:/	Replace	To lowercase
125th Street Fresh Connect Farmers' Market			To number
12th & Brandywine			To date
			To text
			To null
			To empty string

- iii) Special character “&” in MarketName values were replaced with “and”. This transformation affected 192 cells.

Custom text transform on column MarketName

Expression: value.replace("&","and")

Language: General Refine Expression Language (GREL) ▾ No syntax error.

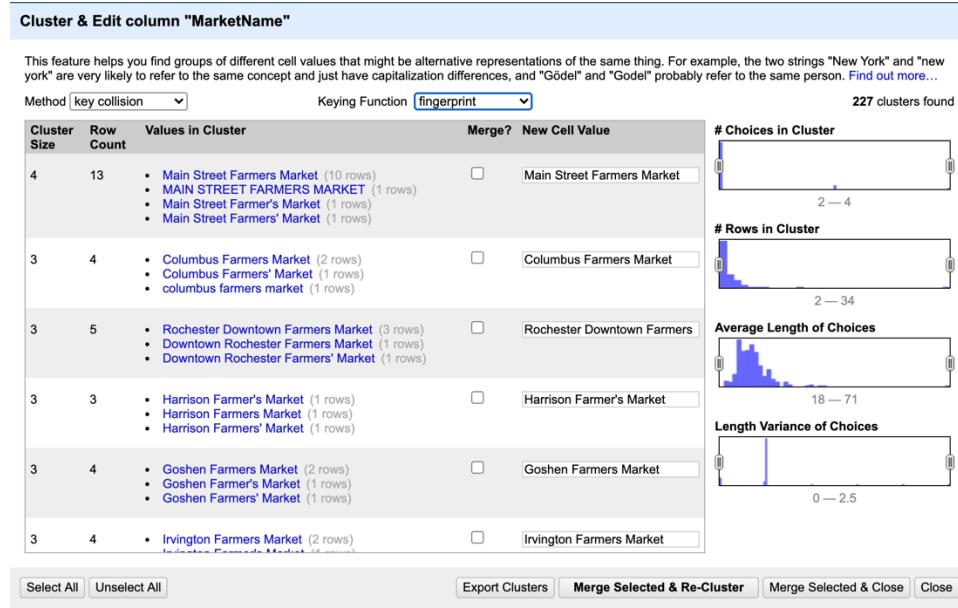
row	value	value.replace("&","and")
1.	Caledonia Farmers Market Association - Danville	Caledonia Farmers Market Association - Danville
2.	Stearns Homestead Farmers' Market	Stearns Homestead Farmers' Market
3.	106 S. Main Street Farmers Market	106 S. Main Street Farmers Market
4.	10th Street Community Farmers Market	10th Street Community Farmers Market
5.	112st Madison Avenue	112st Madison Avenue
6.	12 South Farmers Market	12 South Farmers Market

On error: keep original set to blank store error

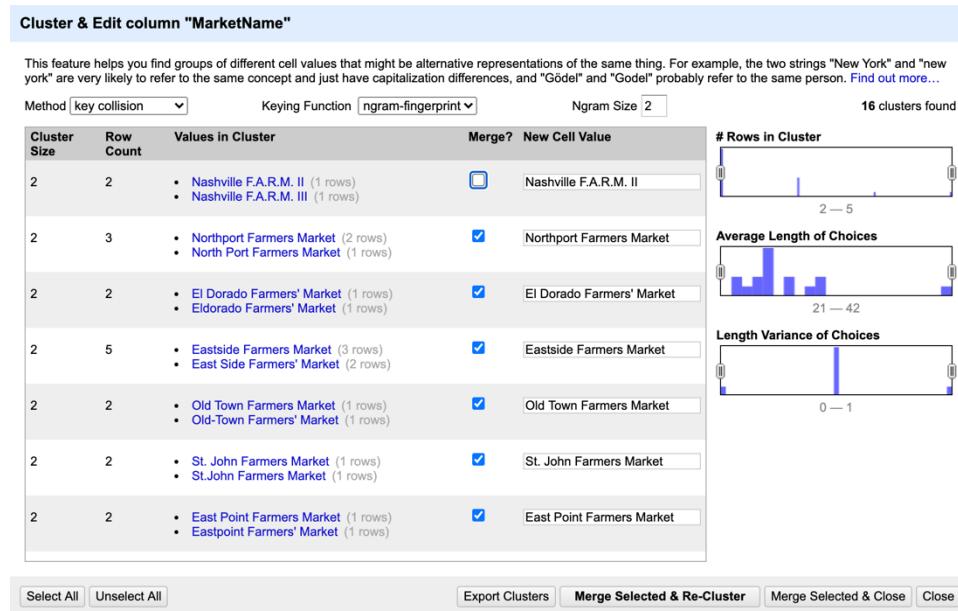
Re-transform up to 10 times until no change

OK Cancel

- iv) Since this is a text column, there's a possibility of having same names with minor difference. To remove these kinds of duplicates, we used clustering transformation using Key Collision method and Fingerprint keying function. This transformation affected 661 cells.



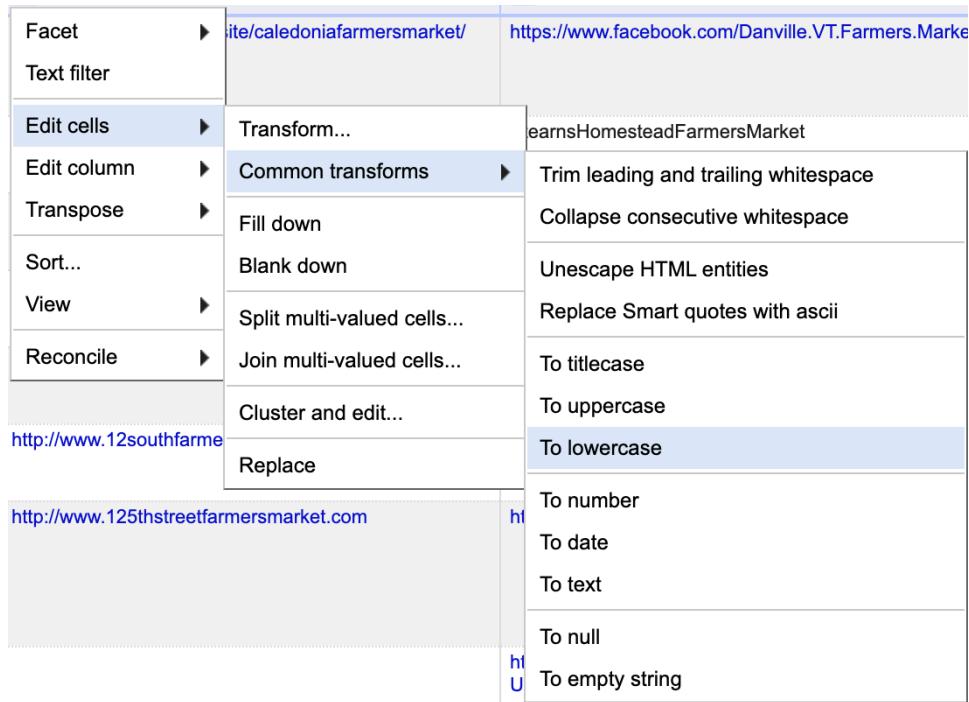
- v) We also transformed this column with Key Collision method but with a variant of fingerprint keying function – ngram-fingerprint with Ngram Size 2. This affected 30 cells.



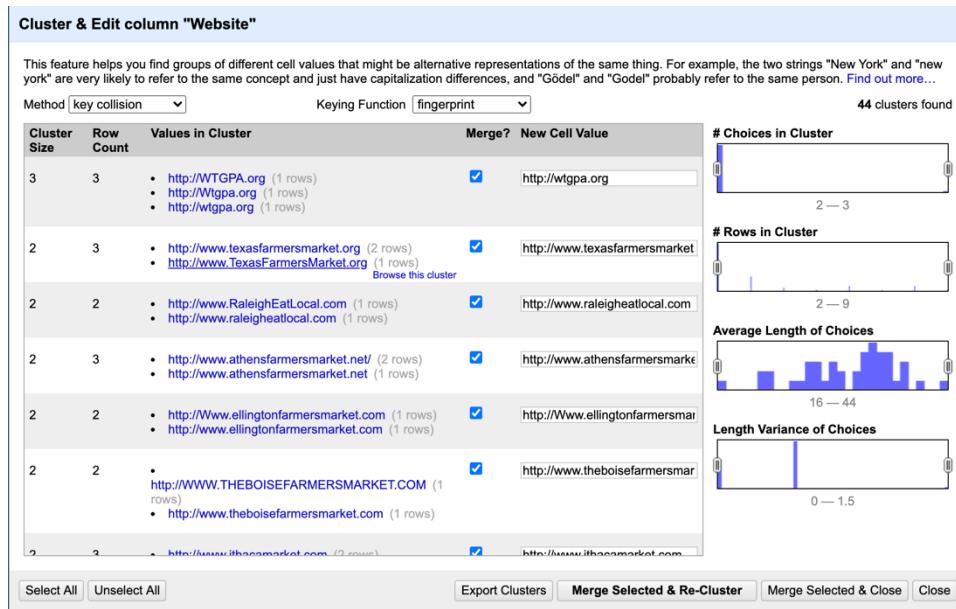
After applying all these steps, now we have clean and consistent MarketName without any duplicates for further use.

3. Website

- i) Trim leading and trailing whitespace – The number of cells affected by this transformation were 27.
- ii) Collapse consecutive whitespace – There were no cells affected by this transformation.
- iii) In order to do further cleaning, we converted all URLs to lowercase. This affected 606 cells.



- iv) We tried to validate if any URLs are duplicate. To remove these kinds of duplicates, we used clustering transformation using Key Collision method and Fingerprint keying function. This transformation affected 145 cells.



- v) We also transformed this column with Key Collision method but with different keying function called - ngram with Ngram Size 2. This affected 23 cells.

Cluster & Edit column "Website"

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 key collision Keying Function ngram-fingerprint Ngram Size 2 3 clusters found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
2	3	<ul style="list-style-type: none"> • http://www.41markets.com (2 rows) • http://www.41markets.com (1 rows) 	<input checked="" type="checkbox"/>	http://www.41markets.com
2	2	<ul style="list-style-type: none"> • http://www.richlandareafarmersmarket.org (1 rows) • http://www.richlandareafarmersmarket.org (1 rows) 	<input checked="" type="checkbox"/>	http://www.richlandareafarmersmarket.org
2	18	<ul style="list-style-type: none"> • http://www.recworcester.org/#markets/c1m3b (17 rows) • http://www.recworcester.org/#markets/c1m3b (1 rows) 	<input checked="" type="checkbox"/>	http://www.recworcester.org/#markets/c1m3b

Rows in Cluster
2 — 18

Average Length of Choices
24 — 44

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

- vi) At this point, we have unique URLs in Website column. To validate every URL contains https:// or http://, we used GREL function. We created a new column called "Check Website" and populated 0 and 1 values. 1 means valid URL and 0 means invalid URL or NULL value in Website. This step confirmed that all URLs are valid.

Add column based on column Website

New column name Check Website

On error set to blank store error copy value from original column

Expression Language General Refine Expression Language (GREL)

```
if(value==null,0,if(or(value.contains('https://'),value.contains('http://')),1,0))
```

No syntax error.

G

Preview History Starred Help

row	value	if(value==null,0,if(or(value.c ...
1.	https://sites.google.com/site/caledoniafarmersmarket/	1
2.	http://www.stearnsfarmstead.com	1
3.	http://thetownofsixmile.wordpress.com/	1
4.	null	0
5.	null	0
6.	http://www.12southfarmersmarket.com	1
7.	http://www.12southfarmersmarket.com	1

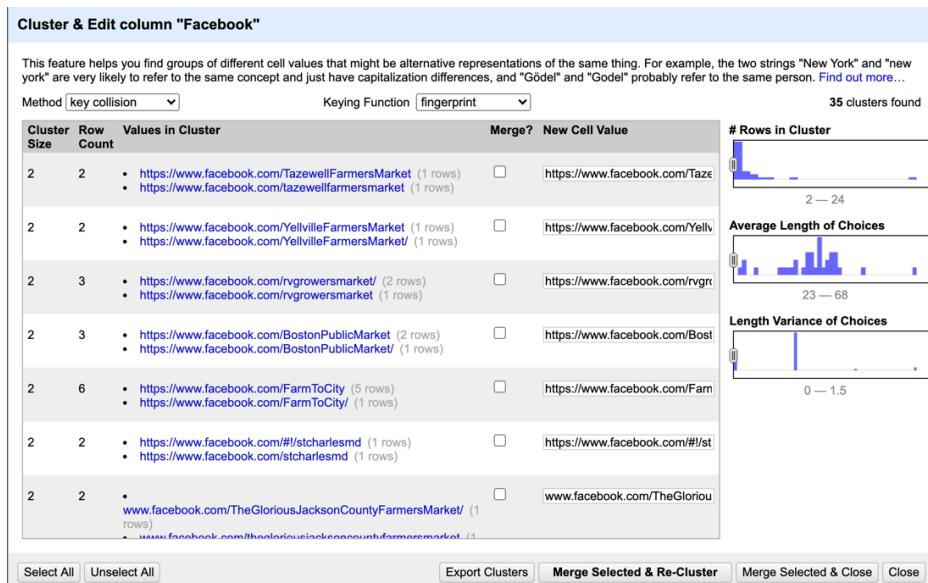
OK Cancel

The top screenshot shows a 'Check Website' interface with 2 choices. The first choice, '0 3514', is highlighted with a red box. The bottom screenshot shows a 'Website' interface with 4332 choices, sorted by count. The choice '(blank) 3514' is highlighted with a red box.

After applying all these steps, now we have valid URLs in Website column without any duplicates.

4. Facebook

- i) Trim leading and trailing whitespace – The number of cells affected by this transformation were 40.
- ii) Collapse consecutive whitespace – There were 3 cells affected by this transformation.
- iii) We tried to validate if any Facebook links are duplicate. To remove these kinds of duplicates, we used clustering transformation using Key Collision method and Fingerprint keying function. This transformation affected 117 cells.

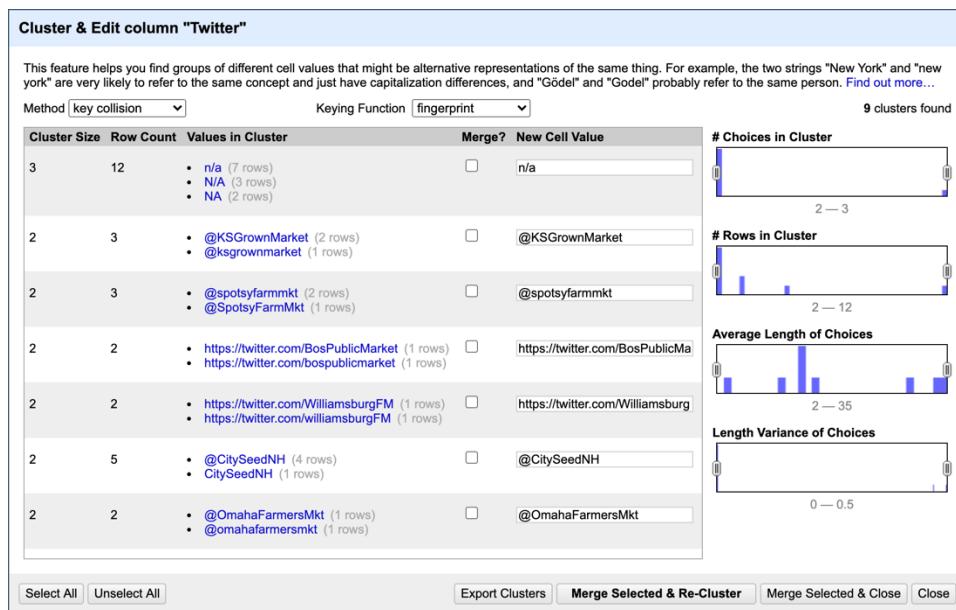


- iv) We tried transforming this column with Key Collision method but with different keying function called - ngram with Ngram Size 2 but no clusters were found.

After applying all these steps, now we have clean values without any duplicates in Facebook column.

5. Twitter

- i) Trim leading and trailing whitespace – The number of cells affected by this transformation were 9.
- ii) Collapse consecutive whitespace – There was 1 cell affected by this transformation.
- iii) We tried to validate if any Twitter links are duplicate. To remove these kinds of duplicates, we used clustering transformation using Key Collision method and Fingerprint keying function. This transformation affected 33 cells.



- iv) We tried transforming this column with Key Collision method but with different keying function called - ngram with Ngram Size 2 but no clusters were found.

After applying all these steps, now we have clean values without any duplicates in Twitter column.

6. YouTube

- i) Trim leading and trailing whitespace – The number of cells affected by this transformation were 4.
- ii) Collapse consecutive whitespace – There were no cells affected by this transformation.
- iii) We tried to validate if any YouTube links are duplicate. To remove these kinds of duplicates, we used clustering transformation using Key Collision method and Fingerprint keying function. This transformation affected 17 cells.

Cluster & Edit column "Youtube"

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	key collision	Keying Function	fingerprint	1 cluster found
Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
3	17	<ul style="list-style-type: none"> n/a (8 rows) N/A (7 rows) NA (2 rows) 	<input type="checkbox"/>	n/a

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

- iv) We tried transforming this column with Key Collision method but with different keying function called - ngram with Ngram Size 2 but no clusters were found.

After applying all these steps, now we have clean values without any duplicates in Youtube column.

7. OtherMedia

- i) Trim leading and trailing whitespace – The number of cells affected by this transformation were 18.
- ii) Collapse consecutive whitespace – There were 22 cells affected by this transformation.
- iii) We tried to validate if any OtherMedia links are duplicate. To remove these kinds of duplicates, we used clustering transformation using Key Collision method and Fingerprint keying function. This transformation affected 47 cells.

Cluster & Edit column "OtherMedia"

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	key collision	Keying Function	fingerprint	8 clusters found
Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
3	11	<ul style="list-style-type: none"> n/a (7 rows) N/A (3 rows) NA (1 rows) 	<input checked="" type="checkbox"/>	n/a
3	15	<ul style="list-style-type: none"> Instagram (11 rows) instagram (3 rows) Instagram: (1 rows) 	<input checked="" type="checkbox"/>	Instagram
2	2	<ul style="list-style-type: none"> Twitter , instagram (1 rows) Twitter, Instagram (1 rows) 	<input checked="" type="checkbox"/>	Twitter , instagram
2	4	<ul style="list-style-type: none"> instagram: visaliafarmersmarket (3 rows) visaliafarmersmarket (Instagram) (1 rows) 	<input checked="" type="checkbox"/>	instagram: visaliafarmersmarke
2	2	<ul style="list-style-type: none"> http://pinterest.com/morgantownmkt (1 rows) http://pinterest.com/morgantownmkt/ (1 rows) 	<input checked="" type="checkbox"/>	http://pinterest.com/morgantow
2	3	<ul style="list-style-type: none"> smfms (Instagram) (2 rows) Instagram: SMFMs (1 rows) 	<input checked="" type="checkbox"/>	smfms (Instagram)
2	5	<ul style="list-style-type: none"> instagram: @cityseedhaven (4 rows) Instagram: cityseedhaven (1 rows) 	<input checked="" type="checkbox"/>	instagram: @cityseedhaven

Choices in Cluster
2 — 3

Rows in Cluster
2 — 15

Average Length of Choices
2 — 35

Length Variance of Choices
0 — 0.5

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

- iv) We also transformed this column with Key Collision method but with different keying function called - ngram with Ngram Size 2. This affected 16 cells.

Cluster & Edit column "OtherMedia"

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	key collision	Keying Function	ngram-fingerprint	Ngram Size	2	1 cluster found
Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value		
2	16	<ul style="list-style-type: none"> Instagram (15 rows) insta gram (1 rows) 	<input type="checkbox"/>	Instagram		

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

After applying all these steps, now we have clean values without any duplicates in OtherMedia column.

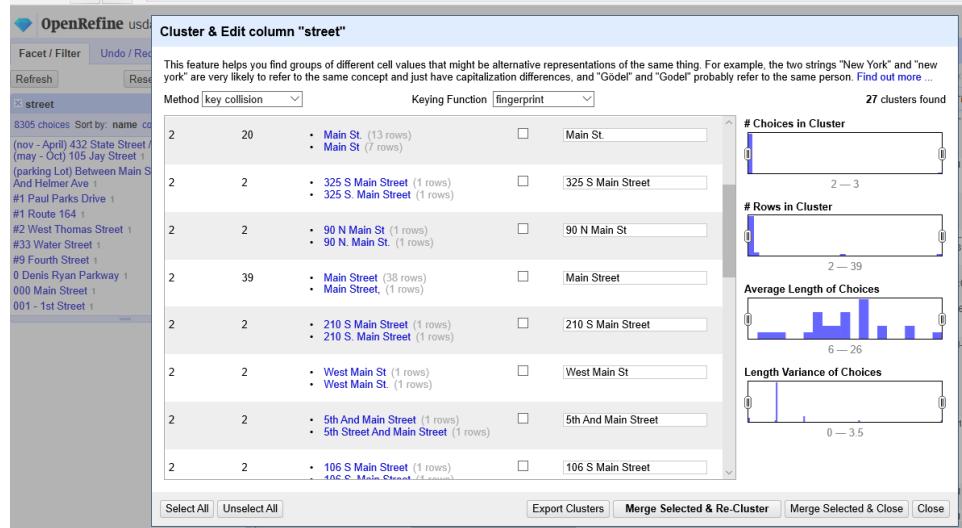
8. street

- i) Trim leading and trailing whitespace – The number of cells affected by this were 317
- ii) Collapse consecutive whitespace – The number of cells affected by this were 90
- iii) To further make the Street name consistent, we convert Street names to title case. This affected 2507 cells

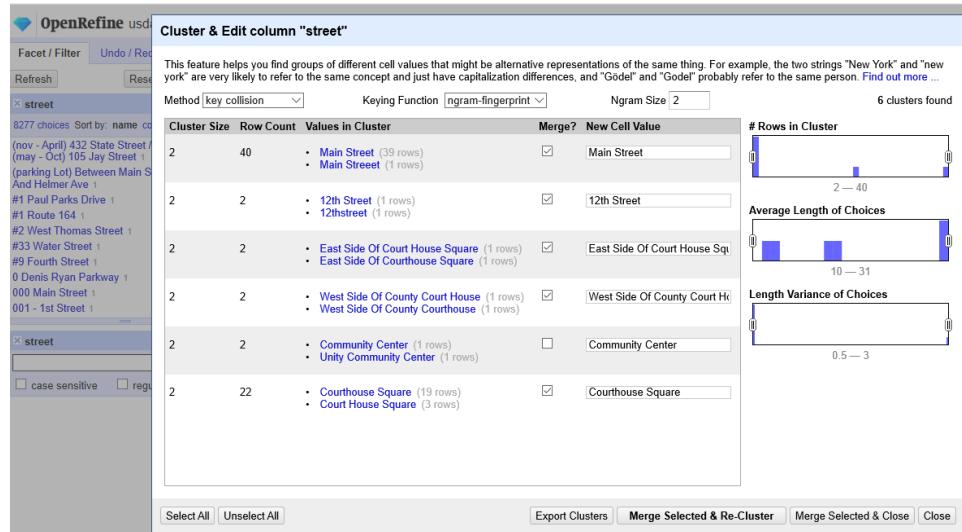
Show as: rows records Show: 5 10 25 50 rows first previous 1 - 10 next last

OtherMedia	street	city	County	State	zip	Season1Date	Season1Time	Season2Date	Season2Time	Search
http://agrimissouri.com/mo-grown/grodetail.php?type=mo-grown&ID=275	Facet	Caledonia	Vermont	05828	06/14/2017 to 08/30/2017	Wed: 9:00 AM-1:00 PM;	09/06/2017 to 10/18/2017	Wed: 2:00 PM-6:00 PM;		
@12southfrmsmkt	Text filter									
Instagram--> 125thStreetFarmersMarket	Edit cells	Transform...			06/24/2017 to Sat: 9:00 AM-1:00					
	Edit column	Common transforms				Trim leading and trailing whitespace				
	Transpose	Fill down				Collapse consecutive whitespace				
	Sort...	Blank down				Unescape HTML entities				
	View	SPLIT multi-valued cells...				To titlecase				
	Reconcile	Join multi-valued cells...				To uppercase				
		Cluster and edit...				To lowercase				
		Replace				To number				
						To date				
						To text				
						To null				
						To empty string				
https://www.facebook.com/delawareurbanfarmcoalition	163 West 125th Street and Adam Clayton Powell, Jr. Blvd.	New York	New York	New York	1002					
instagram:14kenfm	12th & Brandywine Streets	Wilmington	New Castle	Delaware	19801	05/16/2014 to 10/17/2014	Fri: 8:00 AM-11:00 AM;			
	1400 U Street NW	Washington	District of Columbia	District of Columbia	20009	05/03/2014 to 11/22/2014	Sat: 9:00 AM-1:00 PM;			
	5500	Washington	District of	District of	20011	04/09/2016 to	Sat: 9:00 AM-1:00			

- iv) Since this is a text column, there's a possibility of having same names with minor difference. To remove these kinds of duplicates, we used clustering transformation using Key Collision method and Fingerprint keying function. This transformation affected 111 cells.



- v) We also transformed this column with Key Collision method but with different keying function called - ngram with Ngram Size 2. This affected 68 cells. Note that have intentionally not merged the values Community Center and Unity Community Center as they are different values.



- vi) To keep the column street to be consistent with other column names, we have renamed the column "street" to titlecase "Street"

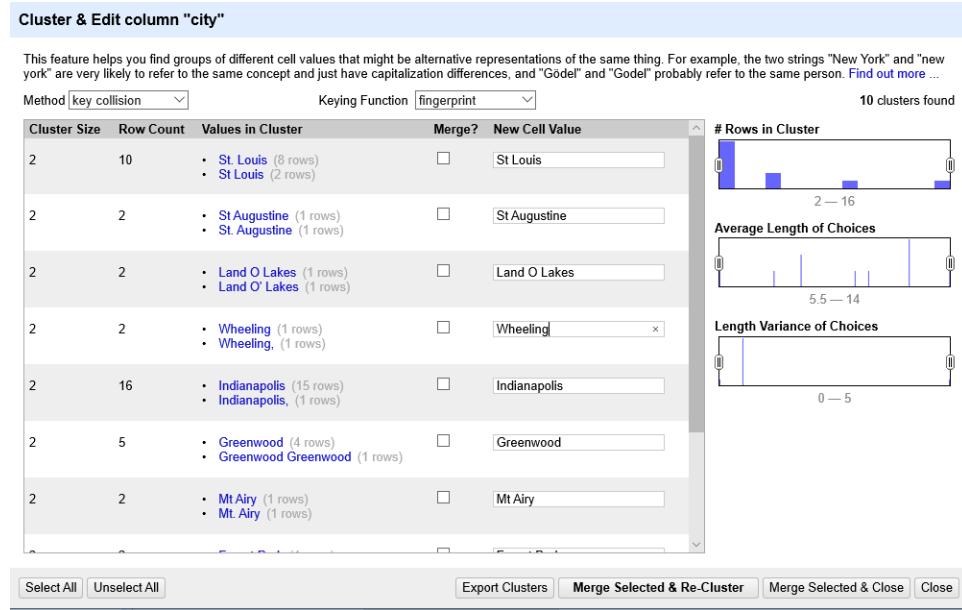
8806 rows											Extensions: Wikidata ▾	
Show as: rows records		Show: 5 10 25 50 rows									< first < previous 1 - 10 next > last >	
	OtherMedia	street	city	County	State	zip	Season1Date	Season1Time	Season2Date	Season2Time	Sea	
	http://agrimissouri.com/mo-grown/grodetall.php? type=mo-grown&ID=275	Facet	Caledonia	Vermont	05628	06/14/2017 to 08/30/2017	Wed: 9:00 AM-1:00 PM,	09/06/2017 to 10/18/2017	Wed: 2:00 PM-6:00 PM,			
		Text filter	Cuyahoga	Ohio		06/24/2017 to	Sat: 9:00 AM-1:00 PM,					
		Edit cells										
		Edit column										
		Transpose										
		Sort...										
		View										
		Reconcile										
		Madison Avenue	Nashville									
	@12southfrmsmkt	3000 Granny White Pike										
	Instagram--> 125thStreetFarmersMarket	163 West 125th Street And Adam Clayton Powell, Jr. Blvd.	New York									
	https://www.facebook.com/delawareurbanfarmcoalition	12th & Brandywine Streets	Wilmington	New Castle	Delaware	19801	05/16/2014 to 10/17/2014	Fri: 8:00 AM-11:00 AM;				
	instagram:14kentfm	1400 U Street Nw	Washington	District of Columbia	District of Columbia	20009	05/03/2014 to 11/22/2014	Sat: 9:00 AM-1:00 PM;				
		5500	Washington	District of	District of	20011	04/09/2016 to	Sat: 9:00 AM-1:00				

9. city

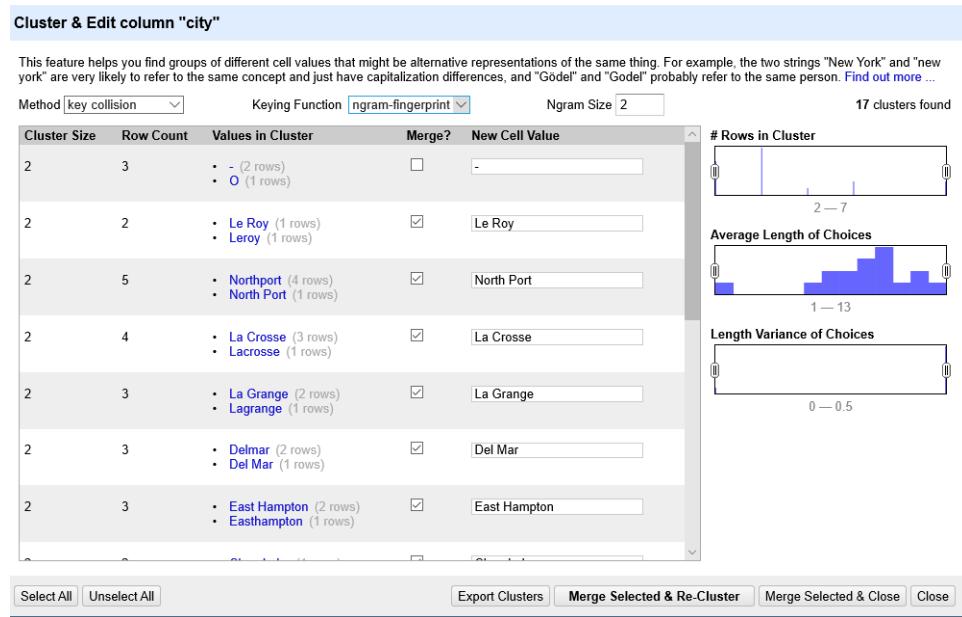
- i) Trim leading and trailing whitespace – The number of cells affected by this were 967
- ii) Collapse consecutive whitespace – The number of cells affected by this were 2
- iii) To further make the City name consistent, we convert city names to title case. This affected 342 cells

8806 rows											Extensions: Wikidata ▾	
Show as: rows records		Show: 5 10 25 50 rows									< first < previous 1 - 10 next > last >	
	OtherMedia	Street	city	County	State	zip	Season1Date	Season1Time	Season2Date	Season2Time	Se	
	http://agrimissouri.com/mo-grown/grodetall.php? type=mo-grown&ID=275	Facet	Vermont	05628	06/14/2017 to 08/30/2017		Wed: 9:00 AM-1:00 PM,	Sat: 9:00 AM-1:00	09/06/2017 to 10/18/2017	Wed: 2:00 PM-6:00 PM,		
		Text filter										
		Edit cells										
		Edit column										
		Transpose										
		Sort...										
		View										
		Reconcile										
		6975 Ridge Road										
		106 S Main Street	10th Street And Poplar									
	Instagram--> 125thStreetFarmersMarket	112th Madison Avenue	Nashville	Davidson								
	@12southfrmsmkt	3000 Granny White Pike										
	https://www.facebook.com/delawareurbanfarmcoalition	12th & Brandywine Streets	Wilmington	New Castle	Delaware	19801	05/16/2014 to 10/17/2014	Fri: 8:00 AM-11:00 AM;				
	instagram:14kentfm	1400 U Street Nw	Washington	District of Columbia	District of Columbia	20009	05/03/2014 to 11/22/2014	Sat: 9:00 AM-1:00 PM;				
		5500	Washington	District of	District of	20011	04/09/2016 to	Sat: 9:00 AM-1:00				

- iv) Since this is a text column, there's a possibility of having same names with minor difference. To remove these kinds of duplicates, we used clustering transformation using Key Collision method and Fingerprint keying function. This transformation affected 48 cells.



- v) We also transformed this column with Key Collision method but with different keying function called - ngram with Ngram Size 2. This affected 56 cells. Note that we have intentionally not merged the values blank with O as they are different values.



- vi) To keep the column city to be consistent with other column names, we have renamed the column "city" to titlecase "City"

OtherMedia	Street	city	County	State	zip	Season1Date	Season1Time	Season2Date	Season2Time	Se
		Vermont	05828	06/14/2017 to 08/30/2017		Wed: 9:00 AM-1:00 PM,	09/06/2017 to 10/18/2017	Wed: 2:00 PM-6:00 PM,		
	6975 Ridge Road	Ohio				Sat: 9:00 AM-1:00				
	106 S Main Street									
http://agrimissouri.com/mo-grown/grodetail.php?type=mo-grown&ID=275	10th Street And Poplar									
	112th Madison Avenue									
@12southfrmsmkt	3000 Granny White Pike	Nashville	Davidson							
Instagram--> 125thStreetFarmersMarket	163 West 125th Street And Adam Clayton Powell, Jr. Blvd.	New York	New York							
https://www.facebook.com/delawareurbanfarmcoalition	12th & Brandywine Streets	Wilmington	New Castle	Delaware	19801	05/16/2014 to 10/17/2014	Fri: 8:00 AM-11:00 AM;			
instagram: 14kenfm	1400 U Street Nw	Washington	District of Columbia	District of Columbia	20009	05/03/2014 to 11/22/2014	Sat: 9:00 AM-1:00 PM;			
	5500	Washington	District of	District of	20011	04/09/2016 to	Sat: 9:00 AM-1:00			

10. County

- i) Trim leading and trailing whitespace – The number of cells affected by this were 0
- ii) Collapse consecutive whitespace – The number of cells affected by this were 0
- iii) To further make the County name consistent, we convert County names to title case. This affected 232 cells

OtherMedia	Street	City	County	State	zip	Season1Date	Season1Time	Season2Date	Season2Time	Se
		Danville	Facet	05828	06/14/2017 to 08/30/2017	Wed: 9:00 AM-1:00 PM,	09/06/2017 to 10/18/2017	Wed: 2:00 PM-6:00 PM,		
	6975 Ridge Road	Parma	Text filter			Sat: 9:00 AM-1:00				
	106 S Main Street	Six Mile	Edit cells	Transform...						
	10th Street And Poplar	Lamar	Edit column	Common transforms		Trim leading and trailing whitespace				
	112th Madison Avenue	New York	Transpose	Fill down		Collapse consecutive whitespace				
@12southfrmsmkt	3000 Granny White Pike	Nashville	Davidson	Tennessee		Blank down				
Instagram--> 125thStreetFarmersMarket	163 West 125th Street And Adam Clayton Powell, Jr. Blvd.	New York	New York	10027	06/10/2014 to 11/25/2014	Unescape HTML entities	To titlecase			
https://www.facebook.com/delawareurbanfarmcoalition	12th & Brandywine Streets	Wilmington	New Castle	Delaware	19801	05/16/2014 to 10/17/2014	To uppercase			
instagram: 14kenfm	1400 U Street Nw	Washington	District Of Columbia	District of Columbia	20009	05/03/2014 to 11/22/2014	To lowercase			
	5500	Washington	District Of	District of	20011	04/09/2016 to	Cluster and edit..			
						Replace	To number			
							To date			
							To text			
							To null			
							To empty string			

- iv) Since this is a text column, there's a possibility of having same names with minor difference. To remove these kinds of duplicates, we used clustering transformation using Key Collision method and n-gram-fingerprint keying function. This transformation affected 13 cells.
- v) We also checked for clustering transformation using fingerprinting keying function and there were no duplicates found.

Cluster & Edit column "County"

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 key collision Keying Function ngram-fingerprint Ngram Size 2 2 clusters found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
2	11	• Stanislaus (10 rows) • St. Anislaus (1 rows)	<input type="checkbox"/>	St. Ahslaus
2	2	• De Witt (1 rows) • Dewitt (1 rows)	<input type="checkbox"/>	De Witt

Rows in Cluster
2 — 11

Average Length of Choices
6.5 — 11

Length Variance of Choices
0.5 — 1

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

11. State

- i) Trim leading and trailing whitespace – The number of cells affected by this were 0
- ii) Collapse consecutive whitespace – The number of cells affected by this were 0
- iii) To further make the State name consistent, we convert State names to title case. This affected 60 cells

Show as: rows records Show: 5 10 25 50 rows « first < previous 1 - 10 next > last »

OtherMedia	Street	City	County	State	zip	Season1Date	Season1Time	Season2Date	Season2Time	Se
	Danville	Caledonia		Facet Text filter		11/14/2017 to 30/2017	Wed: 9:00 AM-1:00 PM;	09/06/2017 to 10/18/2017	Wed: 2:00 PM-6:00 PM;	
	6975 Ridge Road	Parma	Cuyahoga	Edit cells		Transform...				
	106 S Main Street	Six Mile		Edit column		Common transforms		Trim leading and trailing whitespace		
	10th Street And Poplar	Lamar	Barton	Transpose		Fill down		Collapse consecutive whitespace		
http://agrimissouri.com/mo-grown/grodetail.php? type=mo-grown&ID=275				Sort...		Blank down		Unescape HTML entities		
@12southfrmsmkt	112th Madison Avenue	New York	New York	View		Split multi-valued cells...		To titlecase		
	3000 Granny White Pike	Nashville	Davidson	Reconcile		Join multi-valued cells...		To uppercase		
Instagram--> 125thStreetFarmersMarket	163 West 125th Street And Adam Clayton Powell, Jr. Blvd.	New York	New York			Cluster and edit...		To lowercase		
						Replace		To number		
https://www.facebook.com/delawareurbanfarmcoalition	12th & Brandywine Streets	Wilmington	New Castle	Delaware	19801	05/16/2014 to 10/17/2014	Tue: 10:00 AM-7:00 PM;	To date		
instagram:14kenfm	1400 U Street Nw 5500	Washington	District Of Columbia	District of Columbia	20009	05/03/2014 to 11/22/2014	Fri: 8:00 AM-11:00 AM;	To text		
								To null		
								To empty string		

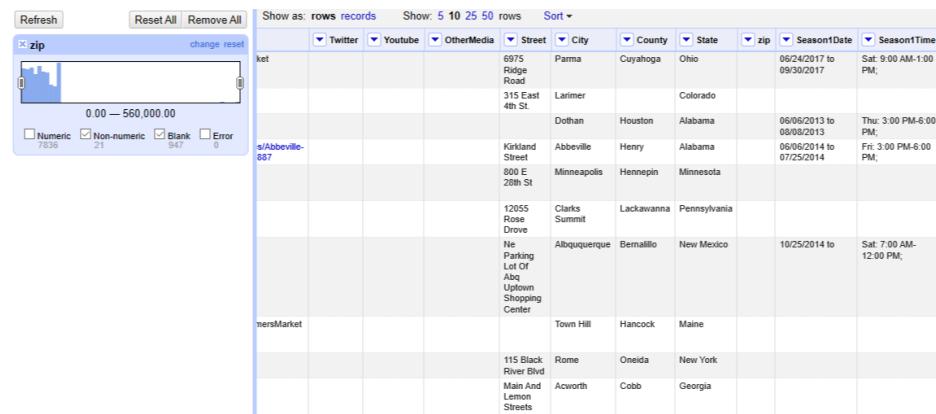
- iv) Since this is a text column, there's a possibility of having same names with minor difference. To remove these kinds of duplicates, we used clustering transformation using Key Collision method and functions – fingerprinting and n-gram-fingerprint keying function. This transformation affected 0 cells. This is an indication of consistent State names provided in the data. It can be inferred that the information in State column has overall high accuracy.

12.Zip

- i) Convert to Number – Since zip is supposed to be a numeric column, column is converted to number. 7838 cells were affected.

Show as: rows records Show: 5 10 25 50 rows										
OtherMedia	Street	City	County	State	zip	Season1Date	Season1Time	Season2Date	Season2Time	Season3Time
	Danville	Caledonia	Vermont			Facet Text filter		Wed: 9:00 AM-1:00 PM,	09/06/2017 to 10/18/2017	Wed: 2:00 PM-6:00 PM,
	6975 Ridge Road	Parma	Cuyahoga	Ohio		Edit cells	Transform...			
	106 S Main Street	Six Mile	Lamar			Trim leading and trailing whitespace	Common transforms			
http://agrimissouri.com/mo-grown/grodetail.php? type=mo-grown&ID=275	10th Street And Poplar					Collapse consecutive whitespace				
	112th Madison Avenue	New York				Unescape HTML entities				
@12southfrmsmkt	3000 Grayway White Pike	Nashville				To titlecase				
Instagram--> 125thStreetFarmersMarket	151 West 125th Street And Adam Clayton Powell, Jr. Blvd.	New York				To uppercase				
	12th & Brandywine Streets	Wilmington	New Castle	Delaware	19801	05/16/2014 to 05/17/2014				
	1400 U Street Nw	Washington	District Of Columbia	District Of Columbia	20009	05/03/2014 to 11/22/2014				
instagram:14kenfm	5500	Washington	District Of	District Of	20011	04/09/2016 to 04/09/2016				
						To number				
						To date				
						To text				
						To null				
						To empty string				

There is total 21 non-numeric values and 947 blank values as seen in below screenshot. This shows zip column is inconsistent and not completely reliable.



- ii) We used Google API to Recover blank Zip codes using location and compared it to existing city values and recovered 584 zip codes matching the exact city and latitude and longitude values as provided by Google API

We initially filtered thru applying a Text Facet and included the blank zip codes

Facet / Filter Undo / Redo 0 / 217

Refresh Reset All Remove All

zip change

6415 choices Sort by: name count Cluster

- FL 1
- IL 1
- KS 1
- Ks 1
- MA 1
- n/a 1
- NY 2
- OR 1
- TN 1
- (blank) 947

Facet by choice counts

Google API URL using a API Key [https://maps.googleapis.com/maps/api/geocode/json?latlng="](https://maps.googleapis.com/maps/api/geocode/json?latlng=) + value + "&key=AIzaSyCM0SHbIG6XqWqhDaGPV-RoPiDvvGe0cD8&result_type=postal_code" – the value accepts "lat,long"

Created a new column based on y value and appended the x value using Transform Join columns using a comma (,) separator so it can be passed to the API

Join columns

Select and order columns to join

- long
- x
- FMID
- MarketName
- Website
- Facebook
- Twitter
- Youtube
- OtherMedia

Select options

Separator between the content of each column: ,

Enter one or more characters, or keep blank to join the columns without separator.

Replace nulls with...

Enter one or more characters, or keep blank to replace nulls with blank strings.

Skip nulls.

In separator and nulls substitutes, use \n for new lines, \t for tabulation, \\n for \n, \\t for \t.

Write result in selected column.

Write result in new column named...

Delete joined columns.

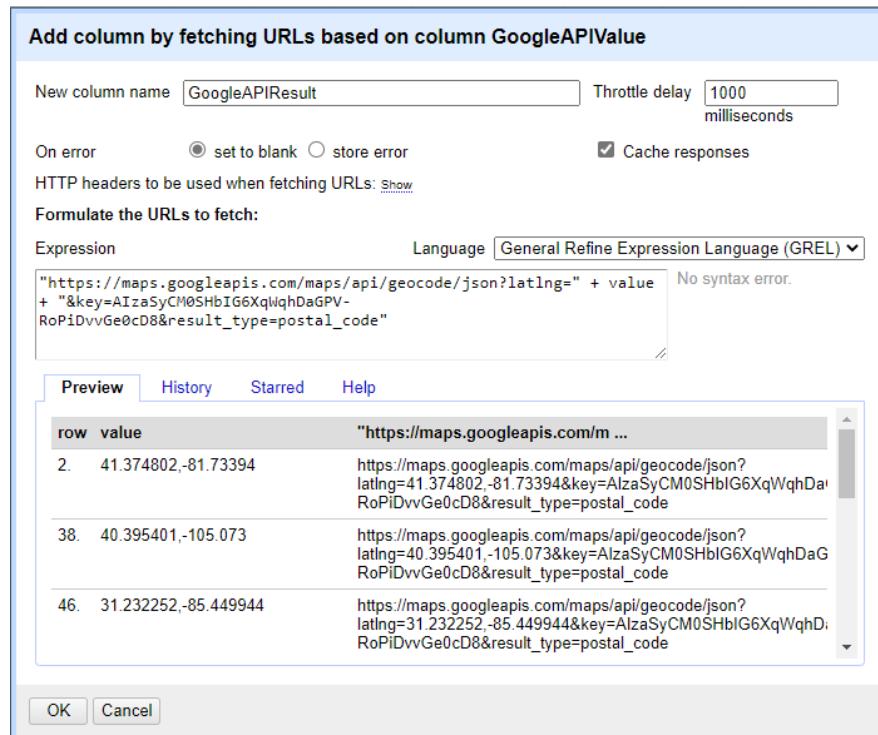
OK Cancel

The join operation combined the latitude and longitude in a format compliant with Google API requirements

x	y	long	GoogleAPIValue
-81.73394	41.374802	41.374802	41.374802,-81.73394
-105.073	40.395401	40.395401	40.395401,-105.073
-85.449944	31.232252	31.232252	31.232252,-85.449944
-85.250366	31.571272	31.571272	31.571272,-85.250366
-93.262901	44.955601	44.955601	44.955601,-93.262901
-75.774597	41.424999	41.424999	41.424999,-75.774597
-106.56584	35.103989	35.103989	35.103989,-106.56584

Using the GoogleAPIValue > Edit Column > Add Columns by fetching URLs passed in the value and concatenating the Google API URL we were able to generate the fetching URL including only the postal code details which returns a Json for the latitude and longitude and postal code along with city the final GREL expression with 10ms Throttle delay:

```
grel:"https://maps.googleapis.com/maps/api/geocode/json?latlng=" + value +
"&key=AIzaSyCM0SHbIG6XqWqhDaGPV-RoPiDvvGe0cD8&result_type=postal_code""
```



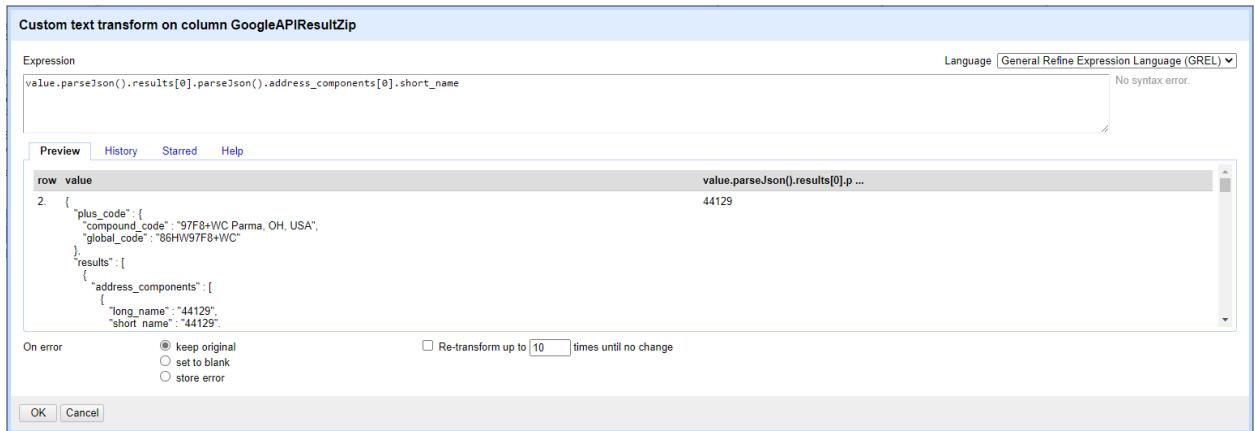
The API executed for around 30 minutes and returned zip files for all latitude longitude values resulting JSON file is now updated in the GoogleAPIResult column

GoogleAPISearch	GoogleAPIResult
41.374802,-81.73394	{ "plus_code": { "compound_code": "97F8+WC" }, "results": [{ "address_components": [{ "long_name": "44129", "short_name": "44129", "types": ["postal_code"] }, { "long_name": "Parma", "short_name": "Parma", "types": ["locality", "political"] }, { "long_name": "Cuyahoga County", "short_name": "Cuyahoga County", "types": ["administrative_area_level_2", "political"], "long_name": "Ohio", "short_name": "OH", "types": ["administrative_area_level_1", "political"] }] }] }

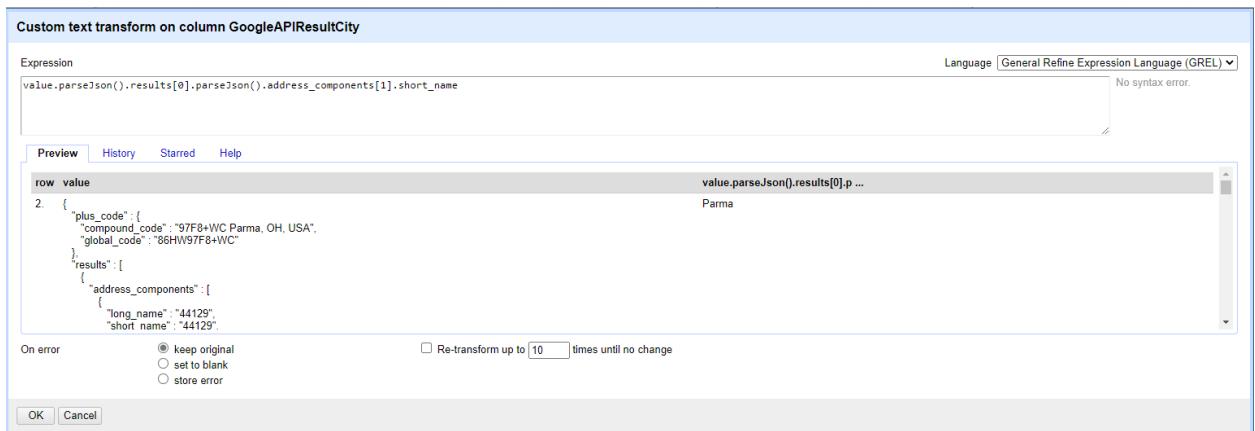
Duplicated the GoogleAPI Result to GoogleResultAPIZip and GoogleAPIResultCity using Add New Column based on the values to retain the JSON results for extracting Zip and City

GoogleAPIResult	GoogleAPIResultZip	GoogleAPIResultCity
{ "plus_code": { "compound_code": "97F8+WC" }, "results": [{ "address_components": [{ "long_name": "44129", "short_name": "44129", "types": ["postal_code"] }, { "long_name": "Parma", "short_name": "Parma", "types": ["locality", "political"] }] }, { "long_name": "Cuyahoga County", "short_name": "Cuyahoga County", "types": ["administrative_area_level_2", "political"], "long_name": "Ohio" }] }	{ "plus_code": { "compound_code": "97F8+WC" }, "results": [{ "address_components": [{ "long_name": "44129", "short_name": "44129", "types": ["postal_code"] }, { "long_name": "Parma", "short_name": "Parma", "types": ["locality", "political"] }] }, { "long_name": "Cuyahoga County", "short_name": "Cuyahoga County", "types": ["administrative_area_level_2", "political"], "long_name": "Ohio" }] }	{ "plus_code": { "compound_code": "97F8+WC" }, "results": [{ "address_components": [{ "long_name": "44129", "short_name": "44129", "types": ["postal_code"] }, { "long_name": "Parma", "short_name": "Parma", "types": ["locality", "political"] }] }, { "long_name": "Cuyahoga County", "short_name": "Cuyahoga County", "types": ["administrative_area_level_2", "political"], "long_name": "Ohio" }] }

Using the GoogleAPIResultZip column performed a transform to parse the JSON and extract the Zip by adding a custom text transform and using parseJson() GREL command
`value.parseJson().results[0].parseJson().address_components[0].short_name`



Using the GoogleAPIResultCity column performed another text transform to parse the JSON and extract the City by adding a custom transform GREL command:
`value.parseJson().results[0].parseJson().address_components[1].short_name`



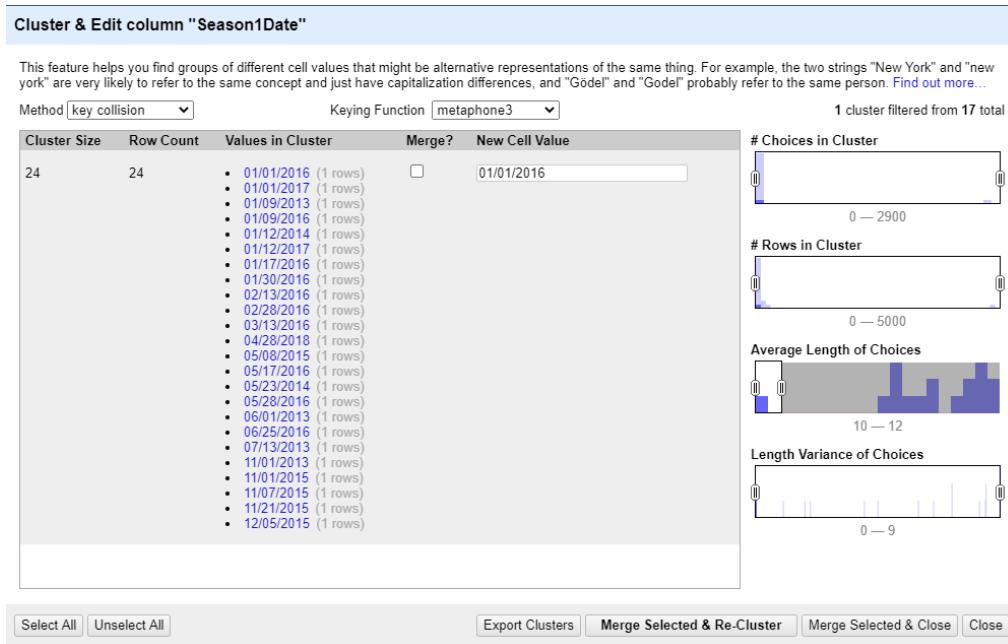
Removed the GoogleAPIResult column and used the Zip and City to extract FMID, city, zip GoogleAPIResultZip, GoogleAPIResultCity into a separate Excel and filtered using condition where the city column is equal to GoogleAPIResultCity.

The Final output we were able to recover 584 correct zip code matching the city names

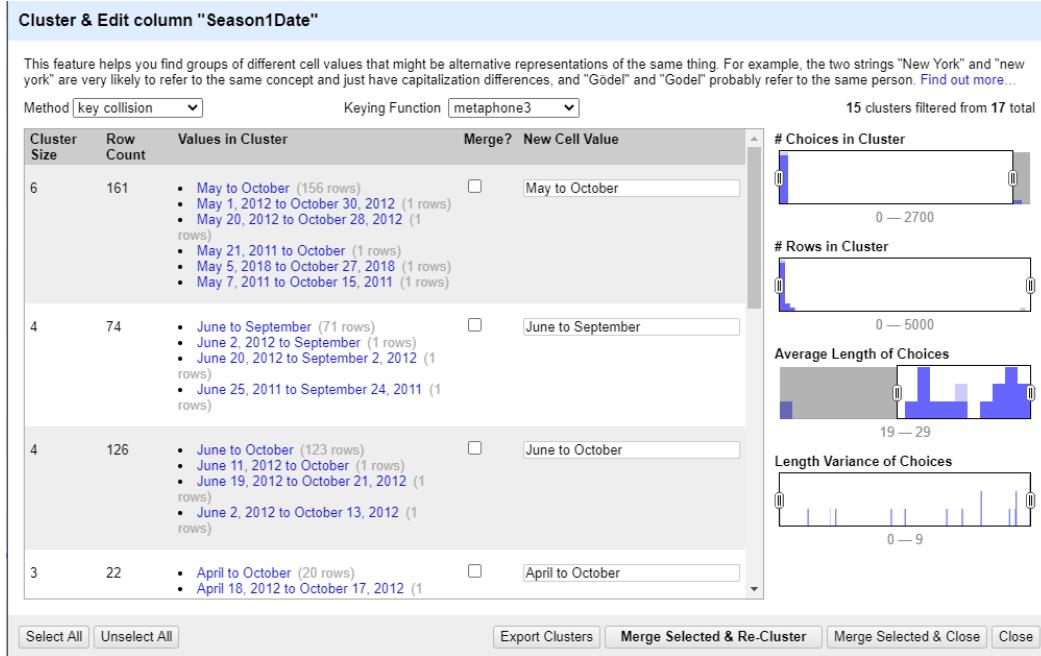
x	y	long	GoogleAPIValue	GoogleAPIResult	GoogleAPIResult
-81.73394	41.374802	41.374802	41.374802,-81.73394	44129	Parma
-105.073	40.395401	40.395401	40.395401,-105.073	80537	Loveland
-85.449944	31.232252	31.232252	31.232252,-85.449944	36305	Dothan
-85.250366	31.571272	31.571272	31.571272,-85.250366	36310	Abbeville
-93.262901	44.955601	44.955601	44.955601,-93.262901	55404	Minneapolis
-75.774597	41.424999	41.424999	41.424999,-75.774597	18411	S ABINGTN TWP
-106.56584	35.103989	35.103989	35.103989,-106.56584	87110	Albuquerque
-68.971603	44.590801	44.590801	44.590801,-68.971603	04438	Frankfort
-75.454102	43.209801	43.209801	43.209801,-75.454102	13440	Rome
-84.676697	34.065701	34.065701	34.065701,-84.676697	30101	Acworth

13. Season1Date, 15. Season2Date, 17. Season3Date, 19. Season4Date

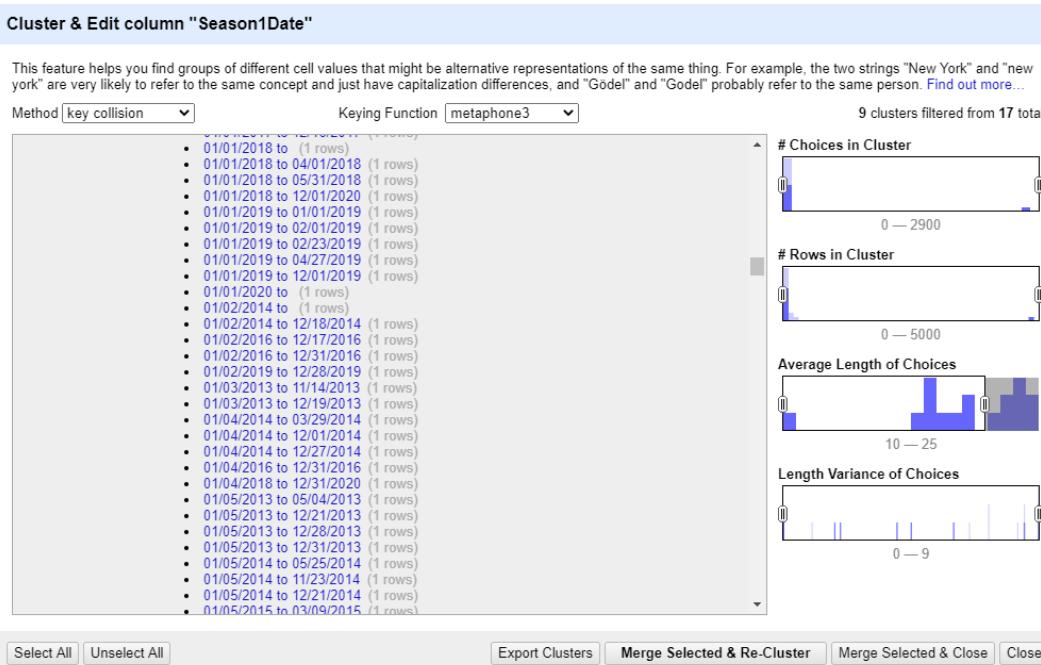
- i) The season1Date column has season start and end dates, using a text cluster using keying function metaphone3 we identified 24 single value dates



- ii) We then used the average length of choices showing dates that are long format



- iii) Majority of the dates around 4916 rows had common format “mm-dd-yyyy to mm-dd-yyyy” with a common separator “to”



- iv) Due the variations of patterns we decided to trim the Season1Date, Season2Date, Season3Date, Season4Date and used the split column transform to split values based on the separator “to”

The screenshot shows a data grid with columns for Season1Date, Season1Time, Season2Date, and Season2T. A context menu is open over the Season1Date column. The menu items include:

- Facet
- Text filter
- Edit cells
- Edit column
- Transpose
- Sort...
- View
- Reconcile

Below the menu, there are specific actions for the Season1Date column:

- 5:00 am: Rename this column, Remove this column
- 05/05/2015 to 10/27/2015: Move column to beginning, Move column to end, Move column left, Move column right
- 06/10/2014 to 11/25/2014: Move column to beginning, Move column to end, Move column left, Move column right

A secondary dialog box titled "Split column Season1Date into several columns" is overlaid. It contains the following fields:

- How to Split Column**
 - by separator
 - by field lengths
- After Splitting**
 - Guess cell type
 - Remove this column
- Separator: regular expression
- Split into columns at most (leave blank for no limit)
- List of integers separated by commas, e.g., 5, 7, 15

- v) The Season1Date, where separated into Season1Date 1, Season1Date 2 the two additional columns the Season1Date had 4940 rows with proper date formats (mm/dd/yyyy) the two new columns were then converted to Date using transform and changed the date to ISO format "YYYY-M-d" using the GREL operation value.toDate() and value.toString("YYYY-M-d") around 4967 Season1FromDate and 4983 Season1ToDate was converted to proper date formats.

Custom text transform on column Season1Date 1

Expression: `value.toDate()`

row	value	value.toDate()
1.	06/14/2017	[date 2017-06-14T00:00:00Z]
2.	06/24/2017	[date 2017-06-24T00:00:00Z]
3.	null	Error: Unable to parse as date
4.	04/02/2014	[date 2014-04-02T00:00:00Z]
5.	July	Error: Unable to convert to a date
6.	05/05/2015	[date 2015-05-05T00:00:00Z]
7.	06/10/2014	[date 2014-06-10T00:00:00Z]

On error: keep original, set to blank, store error

Re-transform up to times until no change

Custom text transform on column Season1Date 1

Expression: `value.toString("YYYY-M-d")`

row	value	value.toString("YYYY-M-d")
1.	2017-06-14T00:00:00Z	2017-6-14
2.	2017-06-24T00:00:00Z	2017-6-24
3.	null	Error: toString accepts an object and an optional second argument containing a date format string
4.	2014-04-02T00:00:00Z	2014-4-2
5.	July	Error: toString accepts an object and an optional second argument containing a date format string
6.	2015-05-05T00:00:00Z	2015-5-5
7.	2014-06-10T00:00:00Z	2014-6-10

On error: keep original, set to blank, store error

Re-transform up to times until no change

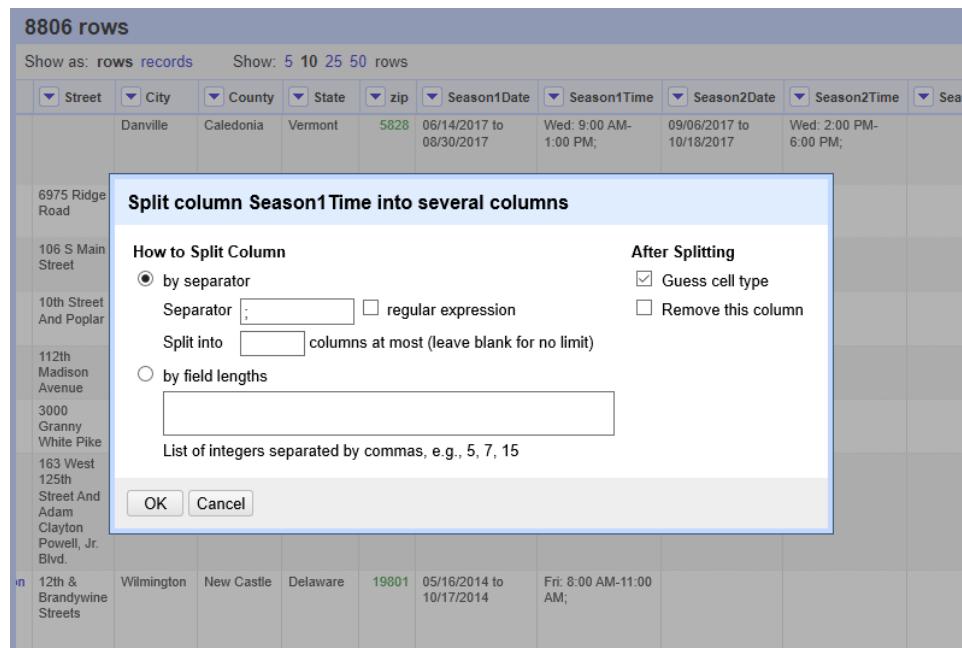
- vi) Steps i, ii, iii, iv, v where repeated for Season2Date, Season3Date, Season4Date resulting in additional columns Season2FromDate, Season2ToDate, Season3FromDate, Season3ToDate, Season4FromDate, Season4ToDate

14. Season1Time, 16. Season2Time, 18. Season3Time, 20. Season4Time

- i) Trim leading and trailing whitespace – There were no cells affected by this
- ii) To further make the Season1Time consistent, we convert Season1Time column to lower case. This affected 5868 cells
- iii) Now we split the Season1Time column into several columns based on the separator ";" but ensure to keep the original column by unchecking the "remove this column" option

This results in 8 new columns namely Season1Time1, Season1Time2 up till Season1Time8

There were 5868 cells that got split into several columns by separator ;. Here ; is the delimiter between the values such as "Wed: 3:00 Pm-6:00 Pm;sat: 8:00 Am-1:00 Pm;" of the Season1Time column.



- iv) Similar steps – 14 - i, ii and ii are repeated for columns – Season2Time, Season3Time and Season4Time. And the number of split columns vary based on the maximum number of multi values in the SeasonTime column.

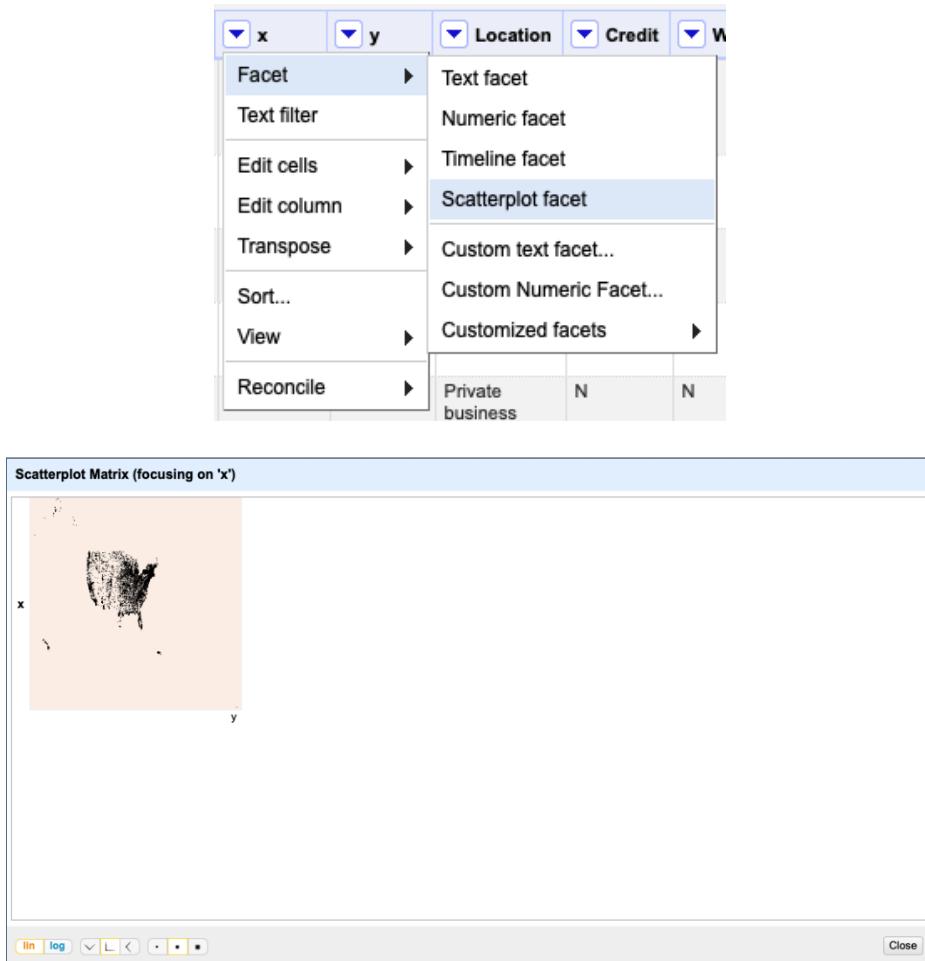
Season2Time – no cells affected by trim operation. 455 cells affected by lowercase () operation. And 455 cells were split into 8 columns.

Season3Time - no cells affected by trim operation. 75 cells affected by lowercase () operation. And 75 cells were split into 7 columns.

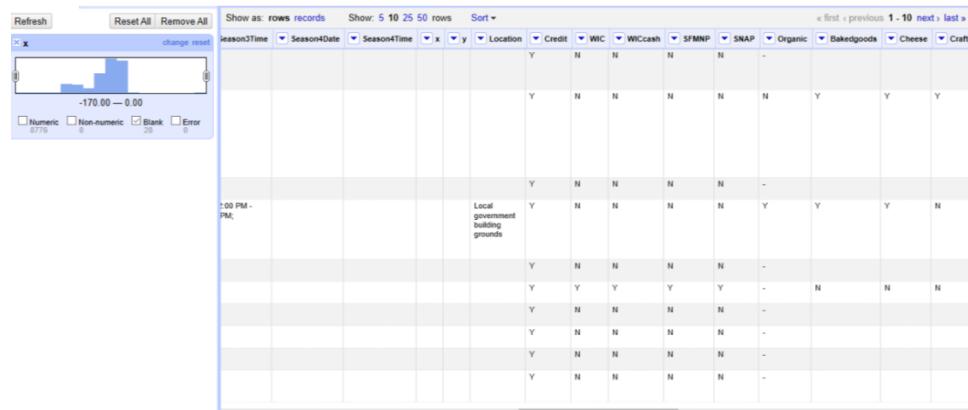
Season4Time - no cells affected by trim operation. 5 cells affected by lowercase () operation. And 5 cells were split into 2 columns.

21. X

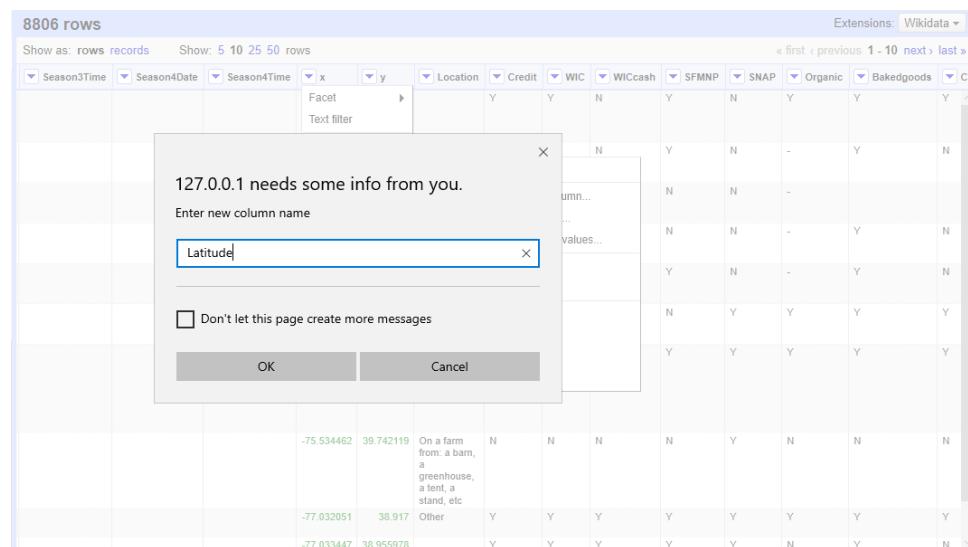
To verify the location accuracy, we did a Scatter Facet in OpenRefine to get a visual map of all the locations using the X and Y values to see an overview and observed that the data confirms to the shape of US with some outliers or errors



- i) Convert to number – This column seems to have numeric values. To validate it, Convert to Number transformation was applied and 8778 cells were affected.
- ii) Clustering using numeric facet shows there are 28 Blank values and rest are numeric.

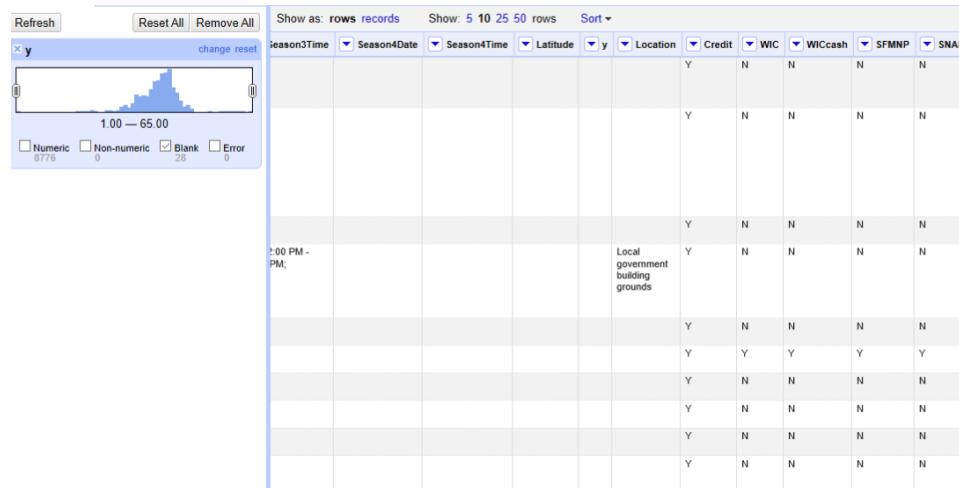


- iii) Seeing the values of column x and column y, it can be inferred that column "x" means latitude values for any given market. So, this column "x" is renamed to "Latitude" to be meaningful.

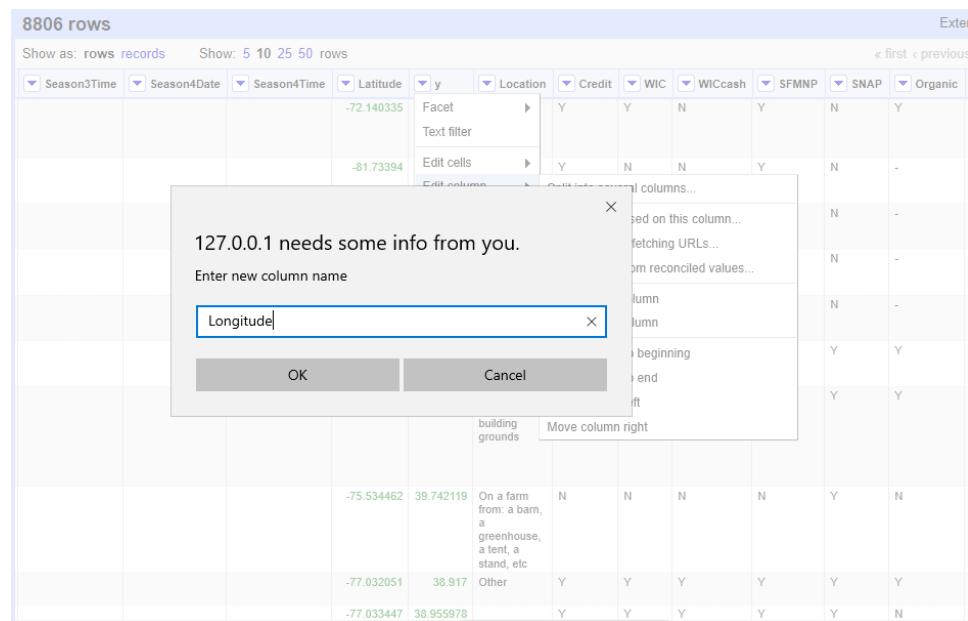


22. Y

- i) Convert to number - This column seems to have numeric values. To validate it, Convert to Number transformation was applied and 8778 cells were affected.
ii) Clustering using numeric facet shows there are 28 Blank values and rest are numeric.



- iii) Seeing the values of column x and column y, it can be inferred that column "y" means longitude values for any given market. So, this column "y" is renamed to "Longitude" to be meaningful.



23. Location

We performed a Transform > Trim leading and trailing spaces on the Location column resulting in cleaning any inconsistent spaces and also observed that there are 6592 blank values in the location field with 10 unique values was observed when performing a text facet.

Location		change
10 choices	Sort by: name count	Cluster
Closed-off public street	212	
Co-located with wholesale market facility	3	
Educational institution	80	
Faith-based institution (e.g., church, mosque, synagogue, temple)	101	
Federal/State government building grounds	52	
Healthcare Institution	60	
Local government building grounds	807	
On a farm from: a barn, a greenhouse, a tent, a stand, etc	73	
Other	487	
Private business parking lot	639	
(blank)	6292	
Facet by choice counts		

24-58 Columns

Credit, WIC, WICcash, SFMNP, SNAP, Organic, Bakedgoods, Cheese, Crafts, Flowers, Eggs, Seafood, Herbs, Vegetables, Honey, Jams, Maple, Meat, Nursery, Nuts, Plants, Poultry, Prepared, Soap, Trees, Wine, Coffee, Beans, Fruits, Grains, Juices, Mushrooms, PetFood, Tofu, WildHarvested – All columns had common theme of values having “Y”, “N” “-” and blank values – we applied a transform to trim Leading Spaces and replaced all blank Values and “-“ using a Text Facet and edit command the below steps indicated are for Organic and Vegetables the same steps was repeated for all the 24-58 columns.

Organic		change	APIResult	Location	Credit	WIC	WICcash	SFMNP
3 choices	Sort by: name count	Cluster			Y	Y	N	Y
-	5011							
N	1355							
Y	2440							
Facet by choice counts			UNKNOWN					
			Apply	Cancel	Enter	Esc		
					Y	N	N	N
				Private business parking lot	N	N	Y	Y
					Y	N	N	N

The screenshot shows a data editing interface. On the left, there's a facet titled 'Vegetables' with a count of 2 choices (N 159, Y 5826, (blank) 2821) and a 'Cluster' button. On the right, a 'Cluster' dialog box is open, showing a table with two columns: 'Federal/State government building grounds' and 'UNKNOWN'. The 'UNKNOWN' column contains four rows, each with a value 'Y'. At the bottom of the dialog are 'Apply' and 'Cancel' buttons.

The breakdown of all columns that had “-” or “blank” Values:

Credit, WIC, WICcash, SFMNP, SNAP, Organic[5011 “-” values] , Bakedgoods[2821 “blank” values], Cheese[2821 “blank” values], Crafts[2821 “blank” values], Flowers[2821 “blank” values], Eggs[2821 “blank” values], Seafood[2821 “blank” values], Herbs[2821 “blank” values], Vegetables[2821 “blank” values], Honey[2821 “blank” values], Jams[2821 “blank” values], Maple[2821 “blank” values], Meat[2821 “blank” values], Nursery[2821 “blank” values], Nuts[2821 “blank” values], Plants[2821 “blank” values], Poultry[2821 “blank” values], Prepared[2821 “blank” values], Soap[2821 “blank” values], Trees[2821 “blank” values], Wine[2821 “blank” values], Coffee[2821 “blank” values], Beans[2821 “blank” values], Fruits, Grains[2821 “blank” values], Juices[2821 “blank” values], Mushrooms[2821 “blank” values], PetFood[2821 “blank” values], Tofu[2821 “blank” values], WildHarvested[2821 “blank” values]

59. updateTime

- i) Trim leading and trailing whitespace – There were no leading and trailing space in any of the value of this column.
- ii) We converted this column into date column such that we have data in ISO format. This transformation affected 8251 cells.

The screenshot shows a context menu for a specific cell in a table. The menu includes options like 'Facet', 'Text filter', 'Edit cells' (which is selected), 'Edit column', 'Transpose', 'Sort...', 'View', 'Reconcile', 'Transform...', 'Common transforms', 'Fill down', 'Blank down', 'Split multi-valued cells...', 'Join multi-valued cells...', 'Cluster and edit...', and 'Replace'. The 'To date' option under 'Common transforms' is highlighted. In the background, there's a table with columns containing values like Y, N, N, Y, Y, N.

- c. Quantifying the results of cleaning (e.g., provide a table of changes along with appropriate quantification) [5%]

ID	Column Name	Steps	Cells Affected
1	FMID	1.1 Convert to Numbers	8806
2	Market Name	2.1 Remove trailing and leading spaces including consecutive spaces	461
		2.2 Convert special character '&' to 'and'	192
		2.3 Remove duplicates using clustering	691
3	Website	3.1 Remove trailing and leading spaces including consecutive spaces	27
		3.2 Convert to all URL to lowercase	606
		3.3 Remove duplicates using clustering	168
4	Facebook	4.1 Remove trailing and leading spaces including consecutive spaces	43
		4.2 Remove all duplicates using clustering	117
5	Twitter	5.1 Remove trailing and leading space including consecutive spaces	10
		5.2 Remove all duplicates using clustering	33
6	Youtube	6.1 Remove trailing and leading spaces including consecutive spaces	4
		6.2 Remove duplicates using clustering	17
7	OtherMedia	7.1 Remove trailing and leading spaces including consecutive spaces	40
		7.2 Remove duplicates using clustering	63
8	street	8.1 Remove trailing and leading spaces including consecutive spaces	407
		8.2 Change street names case sensitive to title case	2507
		8.3 Remove all duplicates using clustering	179
9	city	9.1 Remove trailing and leading spaces including consecutive spaces	969
		9.2 Change city names case sensitive to title case	342
		9.3 Remove duplicates using clustering	104
10	County	10.1 Remove trailing and leading spaces including consecutive spaces	0
		10.2 Change county names case sensitive to title case	232
		10.3 Remove duplicates using clustering	13
11	State	11.1 Remove trailing and leading spaces including consecutive spaces	0
		11.2 Change state names case sensitive to title case	60
		11.3 Remove duplicates using clustering	0
12	zip	12.1 Convert to Numbers	7838
		12.2 Remove values other than numbers (10 nonnumeric value like FL, IL, OR and 947 blank value)	21
		12.3. Retrieve zip using Latitude and Longitude values	584
13	Season1Date	13.1 Remove trailing and leading spaces on Season1Date	110

		13.2 Create new columns From and To based on separator value "to"	5703
		13.3 Remove trailing and leading spaces Season1FromDate	5673
		13.4 Remove trailing and leading spaces Season1ToDate	5571
		13.5 Convert Season1FromDate to Date format and convert to ISO standard "YYYY-M-d" format	4966
		13.6 Convert Season1ToDate to date format and convert to ISO standard "YYYY-M-d" format	4823
		13.7 Convert to ISO Standard Date (or only Month) and remove default timestamp	4940
14	Season1Time	14.1 Remove trailing and leading spaces	0
		14.2 Convert to lowercase	5868
		14.3 Split into several columns based on the separator ";"	5868
15	Season2Date	15.1 Remove trailing and leading spaces Season2Date	14
		15.2 Create new columns From and To based on separator value "to"	464
		15.3 Remove trailing and leading spaces Season2FromDate	459
		15.4 Remove trailing and leading spaces Season2ToDate	446
		15.5 Convert Season2FromDate to Date format and convert to ISO standard "YYYY-M-d" format	442
		15.6 Convert Season2ToDate to Date format and convert to ISO standard "YYYY-M-d" format	423
		15.7 Convert to ISO Standard Date (or only Month) and remove default timestamp	442
16	Season2Time	16.1 Remove trailing and leading spaces	0
		16.2 Convert to lowercase	455
		16.3 Split into several columns based on the separator ";"	455
17	Season3Date	17.1 Remove trailing and leading spaces including consecutive spaces	2
		17.2 Create new columns From and To	78
		17.3 Remove trailing and leading spaces Season3FromDate	77
		17.4 Remove trailing and leading spaces Season3ToDate	76
		17.5 Convert Season3FromDate to Date format and convert to ISO standard "YYYY-M-d"	73
		17.6 Convert Season3ToDate to Date format and convert to ISO standard "YYYY-M-d" format	71
		17.7 Convert to ISO Standard Date (or only Month) and remove default timestamp	73
18	Season3Time	18.1 Remove trailing and leading spaces	0
		18.2 Convert to lowercase	75
		18.3 Split into several columns based on the separator ";"	75
19	Season4Date	19.1 Remove trailing and leading spaces	0
		19.2 Create new columns From and To	5

		19.3 Remove trailing and leading spaces Season3FromDate	5
		19.4 Remove trailing and leading spaces Season3ToDate	5
		19.5 Convert Season3FromDate to Date format and convert to ISO standard "YYYY-M-d"	5
		19.6 Convert Season3ToDate to Date format and convert to ISO standard "YYYY-M-d"	5
		19.7 Convert to ISO Standard Date and remove default timestamp	5
20	Season4Time	20.1 Remove trailing and leading spaces	0
		20.2 Convert to lowercase	5
		20.3 Split into several columns based on the separator ";"	5
21	x	21.1 Convert to Numbers	8778
22	y	22.1 Convert to Numbers	8778
23	Location	23.1 Remove trailing and leading spaces including consecutive spaces	0
24	Credit	24.1 Remove trailing and leading spaces	0
		24.2 Replace values to NULL if other than "Y" or "N"	0
25	WIC	25.1 Remove trailing and leading spaces	0
		25.2 Replace values to NULL if other than "Y" or "N"	0
26	WICCash	26.1 Remove trailing and leading spaces	0
		26.2 Replace values to NULL if other than "Y" or "N"	0
27	SFMNP	27.1 Remove trailing and leading spaces	0
		27.2 Replace values to NULL if other than "Y" or "N"	0
28	SNAP	28.1 Remove trailing and leading spaces	0
		28.2 Replace values to NULL if other than "Y" or "N"	0
29	Organic	29.1 Remove trailing and leading spaces	0
		29.2 Replace values to NULL if other than "Y" or "N"	5011
30	Bakedgoods	30.1 Remove trailing and leading spaces	0
		30.2 Replace values to NULL if other than "Y" or "N"	2821
31	Cheese	31.1 Remove trailing and leading spaces	0
		31.2 Replace values to NULL if other than "Y" or "N"	2821
32	Crafts	32.1 Remove trailing and leading spaces	0
		24.2 Replace values to NULL if other than "Y" or "N"	2821
33	Flowers	33.1 Remove trailing and leading spaces	0
		33.2 Replace values to NULL if other than "Y" or "N"	2821
34	Eggs	34.1 Remove trailing and leading spaces	0
		34.2 Replace values to NULL if other than "Y" or "N"	2821
35	Seafood	35.1 Remove trailing and leading spaces	0
		35.2 Replace values to NULL if other than "Y" or "N"	2821
36	Herbs	36.1 Remove trailing and leading spaces	0
		36.2 Replace values to NULL if other than "Y" or "N"	2821

37	Vegetables	37.1 Remove trailing and leading spaces 37.2 Replace values to NULL if other than "Y" or "N"	0 2821
38	Honey	38.1 Remove trailing and leading spaces 38.2 Replace values to NULL if other than "Y" or "N"	0 2821
39	Jams	39.1 Remove trailing and leading spaces 39.2 Replace values to NULL if other than "Y" or "N"	0 2821
40	Maple	40.1 Remove trailing and leading spaces 40.2 Replace values to NULL if other than "Y" or "N"	0 2821
41	Meat	41.1 Remove trailing and leading spaces 41.2 Replace values to NULL if other than "Y" or "N"	0 2821
42	Nursery	42.1 Remove trailing and leading spaces 42.2 Replace values to NULL if other than "Y" or "N"	0 2821
43	Nuts	43.1 Remove trailing and leading spaces 43.2 Replace values to NULL if other than "Y" or "N"	0 2821
44	Plants	44.1 Remove trailing and leading spaces 44.2 Replace values to NULL if other than "Y" or "N"	0 2821
45	Poultry	45.1 Remove trailing and leading spaces 45.2 Replace values to NULL if other than "Y" or "N"	0 2821
46	Prepared	46.1 Remove trailing and leading spaces 46.2 Replace values to NULL if other than "Y" or "N"	0 2821
47	Soap	47.1 Remove trailing and leading spaces 47.2 Replace values to NULL if other than "Y" or "N"	0 2821
48	Trees	48.1 Remove trailing and leading spaces 48.2 Replace values to NULL if other than "Y" or "N"	0 2821
49	Wine	49.1 Remove trailing and leading spaces 49.2 Replace values to NULL if other than "Y" or "N"	0 2821
50	Coffee	50.1 Remove trailing and leading spaces 50.2 Replace values to NULL if other than "Y" or "N"	0 2821
51	Beans	51.1 Remove trailing and leading spaces 51.2 Replace values to NULL if other than "Y" or "N"	0 2821
52	Fruits	52.1 Remove trailing and leading spaces 52.2 Replace values to NULL if other than "Y" or "N"	0 2821
53	Grains	53.1 Remove trailing and leading spaces 53.2 Replace values to NULL if other than "Y" or "N"	0 2821
54	Juices	54.1 Remove trailing and leading spaces 54.2 Replace values to NULL if other than "Y" or "N"	0 2821
55	Mushrooms	55.1 Remove trailing and leading spaces 56.2 Replace values to NULL if other than "Y" or "N"	0 2821
56	PetFood	56.1 Remove trailing and leading spaces	0

		56.2 Replace values to NULL if other than "Y" or "N"	2821
57	Tofu	57.1 Remove trailing and leading spaces	0
		57.2 Replace values to NULL if other than "Y" or "N"	2821
58	WildHarvested	58.1 Remove trailing and leading spaces	0
		58.2 Replace values to NULL if other than "Y" or "N"	2821
59	updateTime	59.1 Remove trailing and leading spaces	0
		59.2 Convert to ISO Standard Date	8251

3. Developing a relational schema [15%]

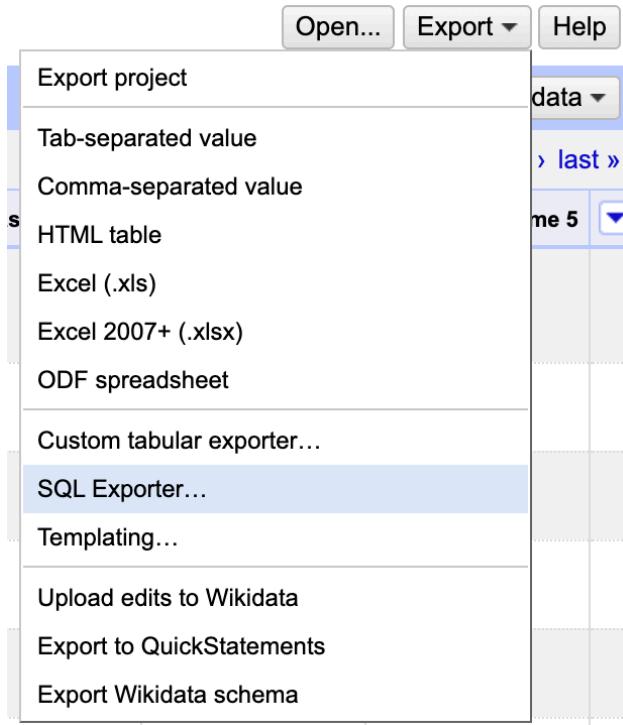
a. Identifying the appropriate integrity constraints [5%]

1. FMID: Integrity constraint check to ensure all FMIDs are unique with no blank values. To ensure this FMID is the primary key of the table STORE_DETAILS. FMID is not duplicated in the table STORE_DETAILS.
2. Market name: All stores should have a valid market name. Integrity constraint check to ensure that no FMID should be associated with NULL market name.
3. Product: Each farmers market should have at least one product available. Integrity constraint is to ensure no farmer market have all the products unavailable. There are 30 products in total provided in PRODUCTS table having PRODUCT_ID from 6 to 35.
4. Payment: Each farmers market should have at least one product method available Integrity constraint is to ensure no farmer market have all the payment options as unavailable I.e. marked as "N" or "UNKOWN". There are 5 payment options – Credit, WIC, WICCash, SFMNP and SNAP; provided in PRODUCTS table having PRODUCT_ID from 1 to 5.
5. Season From and End Date and Time: Every farmers market should have at least one time period meaning at least one season when it is open - both date/month and time.
6. Latitude and Longitude: As per geographical standards, Latitudes must range from -90 to 90 and Longitudes must range from -180 to 180. Integrity constraint check is to ensure that Latitude and Longitude column of STORE_DETAILS table fall between the relevant ranges only.
7. Zip: Integrity constraint check to ensure valid format for zip having numeric values of length 4, 5 or of the format comprising 5 numeric values I.e. XXXXX followed by a hyphen and then followed by a 4-digit numeric value I.e. XXXXX-XXXX. Only Puerto Rico has an exception of zip with length 3.
8. Combination of City, County and State: Integrity constraint to check the combination against actual data from US Gov data. 28889 rows from US Gov data was loaded into the table CITY_COUNTY_STATE_US_GOV. Columns from table STATES, COUNTIES, CITIES are compared with US GOV table to check for the data accuracy and validity.
9. Season From And To Dates: For every season I.e. Season1/2/3/4, integrity check to ensure Season From Date appears before Season End Date if the values represent date of the form "YYYY-MM-DD". Note that if the values represent Month then the ordering of the From and To Date doesn't matter.
10. Store timings: To ensure that Store timings are valid and meaningful, integrity constraint check to ensure all the starting 3 digits of the Time column of table STORE_SEASONAL_HOURS represent valid weekday values such as "mon" for Monday, "tue" for Tuesday and so on.

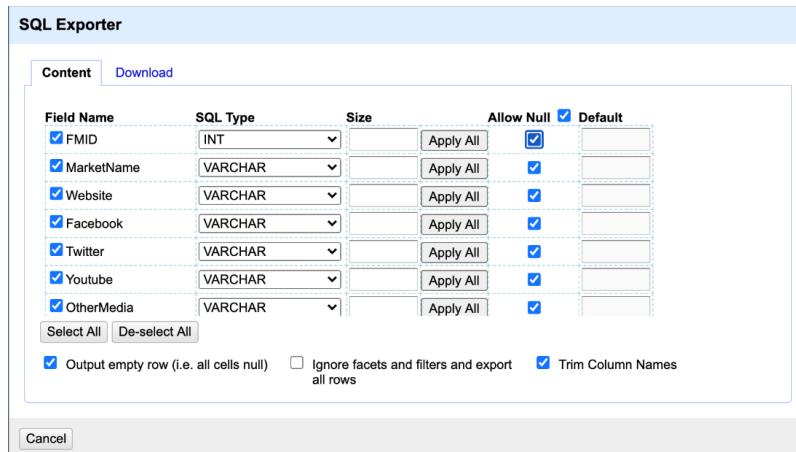
11. Website: Integrity constraint check to ensure all the non-blank website values have a valid format of the website that starts either with "http://" or "https://". Though we cannot validate the website, meaning if the URL is correct or invalid as such, but these integrity check to ensure the format is correct for the non-blank website.

b. Loading data into a database with proper schema [3%]

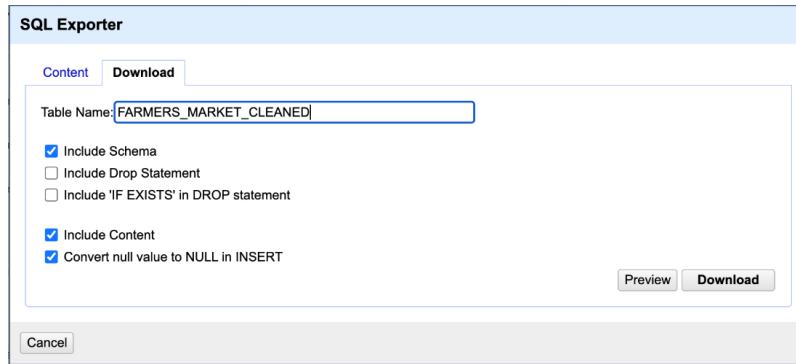
After the data has been cleaned in Open Refine, we used SQL Exporter functionality of Open Refine.



After selection SQL Exporter, select all columns which you would like to be exported, making sure all data type of columns is appropriate.



Then, go to Download tab, provide appropriate table name and download SQL file.



After getting SQL script, we loaded in MySQL database. We have hosted MySQL on Amazon Web Services Relational Database Service (RDS).

The screenshot shows the AWS RDS 'databases' section with the 'dream-data' database selected. The 'Summary' tab is active, displaying basic information: DB Identifier (dream-data), Role (Instance), CPU usage (2.17%), Current activity (4 Connections), Info (Available), Engine (MySQL Community), Class (db.t2.micro), Region & AZ (us-east-1f), and a 'Modify' and 'Actions' button. Below the summary, the 'Connectivity & security' tab is selected, showing detailed network configuration: Endpoint (dream-data.ctr52mmnzuzu.us-east-1.rds.amazonaws.com), Port (3306), Networking (Availability zone: us-east-1f, VPC: vpc-e4942b9e, Subnet group: default-vpc-e4942b9e, Subnets: subnet-23a2872c, subnet-3499fe00, subnet-0d49b1d4, subnet-3499fe0a, subnet-7991f1e, subnet-aef091f2), and Security (VPC security groups: default (sg-d3d1798) (active), Public accessibility: Yes, Certificate authority: rds-ca-2019, Certificate authority date: Aug 22nd, 2024).

The tables are in 3NF form to reduce the duplication of data, avoid data anomalies, ensure referential integrity, and simplify data management. This schema allows adaptability and scalability for using this data for developing any application:

- Farmers Market Details:**

States: All states have been stored.

Counties: All counties with foreign key from States table has been stored.

Cities: All cities with foreign key from Counties table has been stored.

Store Details: This table stores data related to Farmers Market – name, website, social media pages, address, zip, latitude, longitude, and foreign keys from States, Counties and City table.

- Farmers Market Product Availability:**

Departments: This table contains different departments of Farmers Markets. Departments have been formed by grouping similar product types.

Product Types: This table contains product types grouping similar products. This table contains foreign key from Departments table.

Products: All products including different payment types and products. This table foreign key from Product Types table.

Store Product Availability: This table stores all products' availability with foreign key from Store Details and Products table.

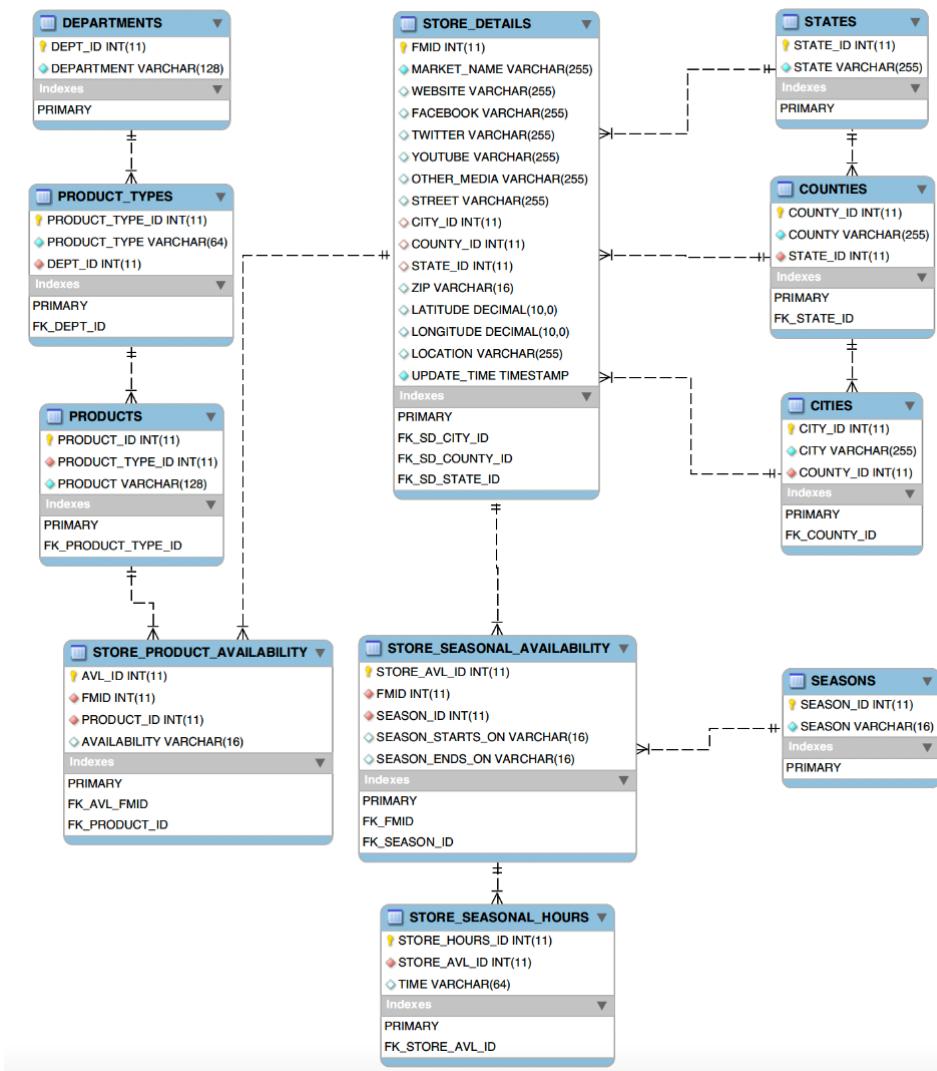
- **Farmers Market Availability:**

Seasons: Farmers Markets have availability based on seasons. This table stores all possible seasons.

Store Seasonal Availability: This table stores season start and end date for each market for each season. It has foreign key from Seasons and Store Details table.

Store Seasonal Hours: This table stores hours of store. This table references Store Seasonal Availability's primary key as foreign key.

The following ER diagram represents the relational schema for storing cleaned data:



Data Load Process – After downloading SQL file from Open Refine with table definition (DDL) and content (DML), we loaded in MySQL table (FARMERS_MARKET_CLEANED). We also enriched zip column by using Google API based on latitude and longitude data. This data was stored in different table (FARMERS_MARKET_CLEANED_ZIP).

The image shows two database tables side-by-side. The left table, titled 'FARMERS_MARKET_CLEANED', contains 54 columns, mostly of type VARCHAR(255), including fields like MARKETNAME, WEBSITE, FACEBOOK, TWITTER, YOUTUBE, OTHERMEDIA, STREET, CITY, COUNTY, STATE, ZIP, and various SEASON1 and SEASON2 time-related fields. The right table, titled 'FARMERS_MARKET_CLEANED_ZIP', contains two columns: FMID (INT(11)) and ZIP (VARCHAR(16)). Both tables have an 'Indexes' section at the bottom.

Based, on these two tables, we loaded relational tables by executing SQL queries.
i.e.

```
INSERT INTO TABLE (<COLUMN_1, COLUMN_2, .... COLUMN_N>)
SELECT COLUMN_1, COLUMN_2, .... COLUMN_N FROM FARMERS_MARKET_CLEANED;
```

You can refer to appendix section for detailed SQL queries for loading each relational table.

c. Writing queries to check the integrity constraints [7%]

1. FMID

```
1 •   SELECT FMID FROM STORE_DETAILS
2     GROUP BY FMID
3     HAVING COUNT(FMID) > 1;
```

The screenshot shows a database result grid with the following data:

FMID
NULL

Below the grid, there are navigation icons and a search bar labeled 'Filter Rows: Search'.

All FMID values are unique nonblank values in the table STORE_DETAILS.

2. Market Name

```

1 •   SELECT FMID,MARKET_NAME FROM STORE_DETAILS
2     WHERE MARKET_NAME IS NULL;

```

100% 27:2

Result Grid Filter Rows: Search Edit:

FMID	MARKET_NAME
NULL	NULL

All stores have a corresponding market name. No FMID found without a market name.

3. Product

```

1   SELECT COUNT(*) FROM
2   (SELECT SD.FMID,SD.MARKET_NAME FROM STORE_DETAILS SD
3   JOIN STORE_PRODUCT_AVAILABILITY PRD_AVL ON SD.FMID = PRD_AVL.FMID
4   JOIN PRODUCTS PRD ON PRD_AVL.PRODUCT_ID = PRD.PRODUCT_ID
5   JOIN PRODUCT_TYPES PRD_TYP ON PRD.PRODUCT_TYPE_ID = PRD_TYP.PRODUCT_TYPE_ID
6   JOIN DEPARTMENTS DEPT ON PRD_TYP.DEPARTMENT_ID = DEPT.DEPARTMENT_ID
7   WHERE DEPT.DEPARTMENT != 'Payment' AND PRD_AVL.AVAILABILITY = 'UNKNOWN'
8   GROUP BY SD.FMID,SD.MARKET_NAME HAVING COUNT(PRD_AVL.AVAILABILITY) = 30) A;

```

Result Grid Filter Rows: Export: Wrap Cell Content:

COUNT(*)
2784

There were 2784 market names out of 8806 market names found which had all its product as "UNKNOWN" meaning all their products were not known if available or not. Inner SELECT query will give the list of all such FMID and farmer market names that violate the integrity constraint.

4. Payment

```

1   SELECT SD.FMID,SD.MARKET_NAME
2   FROM STORE_DETAILS SD
3   JOIN STORE_PRODUCT_AVAILABILITY PRD_AVL
4   ON SD.FMID = PRD_AVL.FMID
5   JOIN PRODUCTS PRD
6   ON PRD_AVL.PRODUCT_ID = PRD.PRODUCT_ID
7   JOIN PRODUCT_TYPES PRD_TYP
8   ON PRD.PRODUCT_TYPE_ID = PRD_TYP.PRODUCT_TYPE_ID
9   JOIN DEPARTMENTS DEPT
10  ON PRD_TYP.DEPARTMENT_ID = DEPT.DEPARTMENT_ID
11  WHERE DEPT.DEPARTMENT = 'Payment' AND PRD_AVL.AVAILABILITY IN ('UNKNOWN','N')
12  GROUP BY SD.FMID,SD.MARKET_NAME
13  HAVING COUNT(PRD_AVL.AVAILABILITY) = 5) A;

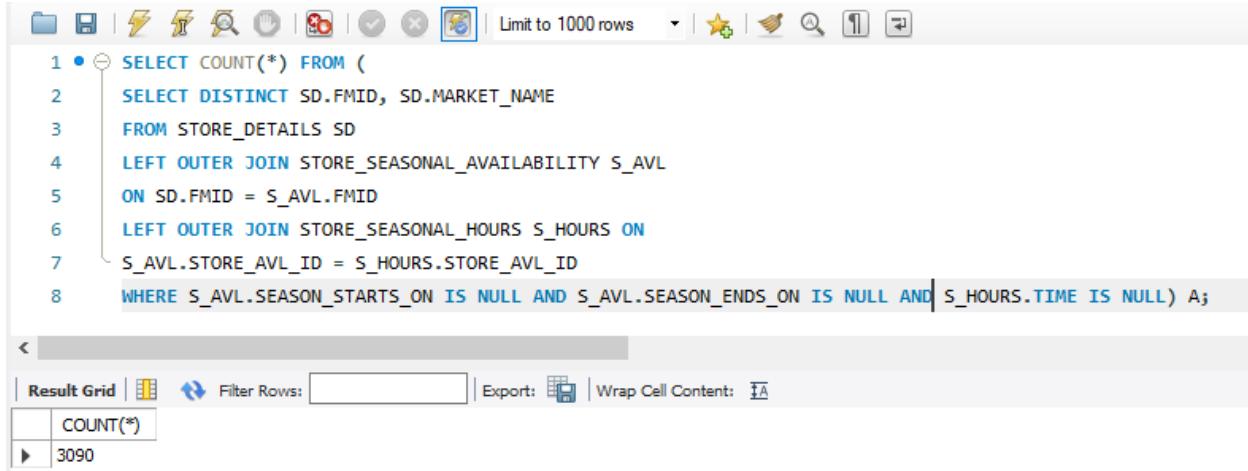
```

Result Grid Filter Rows: Export: Wrap Cell Content:

COUNT(*)
2854

There were 2854 market names found that didn't have any payment options available. Inner SELECT query provides the list of FMID and market names violating the integrity constraint check.

5. Season From and End Date and Time



```

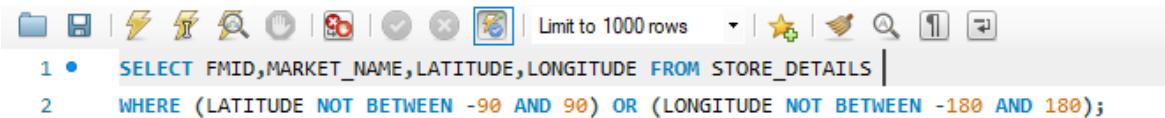
1 • 1 SELECT COUNT(*) FROM (
2   SELECT DISTINCT SD.FMID, SD.MARKET_NAME
3   FROM STORE_DETAILS SD
4   LEFT OUTER JOIN STORE_SEASONAL_AVAILABILITY S_AVL
5   ON SD.FMID = S_AVL.FMID
6   LEFT OUTER JOIN STORE_SEASONAL_HOURS S_HOURS ON
7   S_AVL.STORE_AVL_ID = S_HOURS.STORE_AVL_ID
8   WHERE S_AVL.SEASON_STARTS_ON IS NULL AND S_AVL.SEASON_ENDS_ON IS NULL AND S_HOURS.TIME IS NULL) A;
  
```

The screenshot shows a MySQL Workbench interface with a query editor and a result grid. The query counts FMIDs where season start and end dates are null, and associated hours have null times. The result grid shows one row with a count of 3090.

COUNT(*)
3090

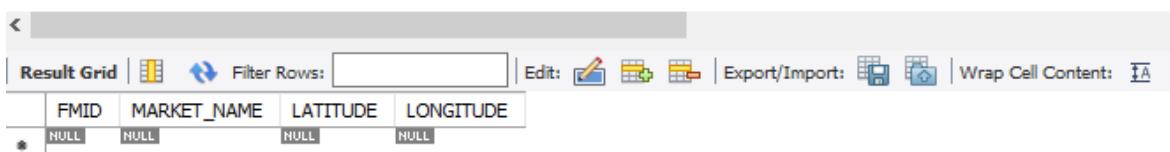
There were 3090 FMIDs with market names that don't have any information about Season From and End Date and Time. Inner SELECT query provides the list of FMID and Market Names violating the integrity constraint check.

6. Latitude and Longitude



```

1 • 1 SELECT FMID, MARKET_NAME, LATITUDE, LONGITUDE FROM STORE_DETAILS |
2 WHERE (LATITUDE NOT BETWEEN -90 AND 90) OR (LONGITUDE NOT BETWEEN -180 AND 180);
  
```



The screenshot shows a MySQL Workbench interface with a result grid. The query checks for coordinates outside valid ranges. The result grid shows one row with all columns set to NULL.

*	FMID	MARKET_NAME	LATITUDE	LONGITUDE
*	NULL	NULL	NULL	NULL

All Latitude and Longitude values were found to be valid.

7. Zip

Below query was used and it found 4 invalid zip codes. FMID 1012563 had nonnumeric zip as "726o1", FMID 1009876 had shorter zip as "33" of length 2, FMID 1005771 had longer zip as "513338" of length 6 and FMID 1006166 had longer zip as "550424" of length 6.

```

1   SELECT FMID,MARKET_NAME,ZIP
2   FROM (SELECT SD.FMID,SD.MARKET_NAME,SD.ZIP
3          FROM STORE_DETAILS SD
4         JOIN STATES ST
5        ON SD.STATE_ID = ST.STATE_ID
6       WHERE SD.ZIP IS NOT NULL AND ST.STATE != 'Puerto Rico') TEMP
7      WHERE ZIP NOT REGEXP '^[0-9][0-9][0-9][0-9]$'
8      AND ZIP NOT REGEXP '^[0-9][0-9][0-9][0-9]$'
9      AND ZIP NOT REGEXP '^([0-9][0-9][0-9][0-9][0-9]-[0-9][0-9][0-9][0-9]$'
10     UNION
11    SELECT FMID,MARKET_NAME,ZIP
12   FROM (SELECT SD.FMID,SD.MARKET_NAME,SD.ZIP
13          FROM STORE_DETAILS SD
14         JOIN STATES ST
15        ON SD.STATE_ID = ST.STATE_ID
16       WHERE SD.ZIP IS NOT NULL AND ST.STATE = 'Puerto Rico') TEMP
17      WHERE ZIP NOT REGEXP '^([0-9][0-9][0-9]$'
18      AND ZIP NOT REGEXP '^([0][0-9][0-9][0-9]$'
19      AND ZIP NOT REGEXP '^([0][0][0-9][0-9]$';

100% 43:18
Result Grid Filter Rows: Search Export:

```

FMID	MARKET_NAME	ZIP
1005771	Everly Farmers Market	513338
1006166	Buffalo Center Farmers Market	550424
1012563	Harrison Farmer's Market	72601
1009876	Bradford Farmers Market	33

8. Combination of City, County and State:

```

1 •  SELECT DISTINCT CT.CITY AS CITY,CNT.COUNTY AS COUNTY,ST.STATE AS STATE
2   FROM STATES ST
3   JOIN COUNTIES CNT
4   ON ST.STATE_ID = CNT.STATE_ID
5   JOIN CITIES CT
6   ON CNT.COUNTY_ID = CT.COUNTY_ID
7   WHERE NOT EXISTS (SELECT DISTINCT CITY,COUNTY,STATE
8                      FROM CITY_COUNTY_STATE_US_GOV);

```

100% 37:8

Result Grid Filter Rows: Search Export:

CITY	COUNTY	STATE
------	--------	-------

All the values of State, County and City were found to be valid and accurate

9. Season From And To Dates:

Season From and To dates have been validated in three steps as from and to values are not consistent:

1) Using Dates

```

1 •  SELECT SD.FMID,SD.MARKET_NAME,S_AVL.SEASON_STARTS_ON,S_AVL.SEASON_ENDS_ON
2   FROM STORE_DETAILS SD
3   JOIN STORE_SEASONAL_AVAILABILITY S_AVL
4   ON SD.FMID = S_AVL.FMID
5   WHERE STR_TO_DATE(S_AVL.SEASON_STARTS_ON, '%Y-%m-%d') > STR_TO_DATE(S_AVL.SEASON_ENDS_ON, '%Y-%m-%d');

```

100% 101:5

Result Grid Filter Rows: Search Export:

FMID	MARKET_NAME	SEASON_STARTS_ON	SEASON_ENDS_ON
1011959	Clark Park Farmer's Market	2016-10-1	2016-5-7
1018773	Sedalia Area Farmers' Market	2017-11-2	2017-4-27

There were 2 FMIDs found which failed the above integrity constraint check

2) Using Months

```

1 •  SELECT * FROM (SELECT SD.FMID,SD.MARKET_NAME,
2   CASE WHEN S_AVL_SEASON_STARTS_ON = 'January' THEN 1 WHEN S_AVL_SEASON_STARTS_ON = 'February' THEN 2 WHEN S_AVL_SEASON_STARTS_ON = 'March' THEN 3
3   WHEN S_AVL_SEASON_STARTS_ON = 'April' THEN 4 WHEN S_AVL_SEASON_STARTS_ON = 'May' THEN 5 WHEN S_AVL_SEASON_STARTS_ON = 'June' THEN 6
4   WHEN S_AVL_SEASON_STARTS_ON = 'July' THEN 7 WHEN S_AVL_SEASON_STARTS_ON = 'August' THEN 8 WHEN S_AVL_SEASON_STARTS_ON = 'September' THEN 9
5   WHEN S_AVL_SEASON_STARTS_ON = 'October' THEN 10 WHEN S_AVL_SEASON_STARTS_ON = 'November' THEN 11 WHEN S_AVL_SEASON_STARTS_ON = 'December' THEN 12
6   END AS SEASON_STARTS_ON_MON,
7   CASE WHEN S_AVL_SEASON_ENDS_ON = 'January' THEN 1 WHEN S_AVL_SEASON_ENDS_ON = 'February' THEN 2 WHEN S_AVL_SEASON_ENDS_ON = 'March' THEN 3
8   WHEN S_AVL_SEASON_ENDS_ON = 'April' THEN 4 WHEN S_AVL_SEASON_ENDS_ON = 'May' THEN 5 WHEN S_AVL_SEASON_ENDS_ON = 'June' THEN 6
9   WHEN S_AVL_SEASON_ENDS_ON = 'July' THEN 7 WHEN S_AVL_SEASON_ENDS_ON = 'August' THEN 8 WHEN S_AVL_SEASON_ENDS_ON = 'September' THEN 9
10  WHEN S_AVL_SEASON_ENDS_ON = 'October' THEN 10 WHEN S_AVL_SEASON_ENDS_ON = 'November' THEN 11 WHEN S_AVL_SEASON_ENDS_ON = 'December' THEN 12
11  END AS SEASON_ENDS_ON_MON,S_AVL_SEASON_STARTS_ON,S_AVL_SEASON_ENDS_ON
12  FROM STORE_DETAILS SD
13  JOIN STORE_SEASONAL_AVAILABILITY S_AVL ON SD.FMID = S_AVL.FMID
14  WHERE (S_AVL_SEASON_STARTS_ON IS NOT NULL AND S_AVL_SEASON_ENDS_ON IS NOT NULL)
15  AND (S_AVL_SEASON_STARTS_ON NOT REGEXP '^January|February|March|April|May|June|July|August|September|October|November|December$')
16  AND (S_AVL_SEASON_ENDS_ON REGEXP '^January|February|March|April|May|June|July|August|September|October|November|December$')) A
17  WHERE SEASON_STARTS_ON_MON > SEASON_ENDS_ON_MON

```

Result Grid Filter Rows: Search Export:

FMID	MARKET_NAME	SEASON_STARTS_ON_MON	SEASON_ENDS_ON_MON	SEASON_STARTS_ON	SEASON_ENDS_ON
100108	Ely Farmers Market	11	3	November	March
100148	Hill Country Farmers Market	11	4	November	April
100151	Marc Island Farmers Market	11	5	November	May
1002027	OceanSide Farmer's market at Lake Worth Beach	11	5	December	May
1002273	Green Spring Station Winter Farmers Market	12	5	September	April
1003898	LNU Farmers Market	9	4	September	May
1003900	Market in the Creek Farmers' Market	11	5	November	May
1004029	Certified SC Winter Farmers Market	11	9	November	March
100522	Manhattan Beach Farmers Market	9	5	September	May
1005379	Norwich Winter Farmers Market	11	4	November	April
1006039	Bronx Winter Farmers' Market	11	4	November	April
1006040	Winter Sun Farms Farmers Market	11	4	November	April
1006058	Winter Sun Farms Indoor Winter Market	12	4	December	April
100107	Ely Farmers Market	11	4	November	April
1001176	North Raleigh Farmers' Market	11	3	November	March
1002854	East Atlanta Farmers Market	11	4	November	April
1003000	Westfield Farmers Market	11	4	November	April
1004007	Marshall County Farmers Market	11	4	November	April
1004024	Farmer's Market @ Woburn's Spence Farm	11	5	November	May
1006078	Santa Fe Farmers' Market	12	3	December	March

Here no conflicting values found where season starting month is before season ending month. The print screen represents where starting month number is higher than ending month, but this is acceptable as a store can be open between November of current year to May of next year.

3) Using combination of dates and months

```

1 •  SELECT * FROM
2   (SELECT
3    SD.FMID,SD.MARKET_NAME,S_AVL_SEASON_STARTS_ON,S_AVL_SEASON_ENDS_ON
4   FROM STORE_DETAILS SD
5   JOIN STORE_SEASONAL_AVAILABILITY S_AVL ON SD.FMID = S_AVL.FMID
6   WHERE S_AVL_SEASON_STARTS_ON IS NOT NULL
7   AND S_AVL_SEASON_ENDS_ON IS NOT NULL
8   AND (S_AVL_SEASON_STARTS_ON NOT REGEXP '^^(January|February|March|April|May|June|July|August|September|October|November|December)$'
9   OR S_AVL_SEASON_ENDS_ON NOT REGEXP '^^(January|February|March|April|May|June|July|August|September|October|November|December)$'
10  AND (S_AVL_SEASON_STARTS_ON NOT REGEXP '^([0-9])'
11  OR S_AVL_SEASON_ENDS_ON NOT REGEXP '^([0-9])')) A;

```

Result Grid Filter Rows: Search Export:

FMID	MARKET_NAME	SEASON_STARTS_ON	SEASON_ENDS...
1001137	Hill Country Farmers Market Association	April	2011-11-4
1001138	Hill Country Farmers Market Association	April	2011-9-24
1001963	Garden Shop Nursery Farmers Market	October	ber
1002517	Fairhaven Winter Farmers' Market	October	ber
1002724	Seminole Heights Sunday Morning Ma...	October	ber
1002777	Main Street Calumet Market	October	ber
1001431	Marlboro County Farmers Market	October	ber
1002222	Alcona Farmers Market	October	ber
1004730	Lansing City Market	2012-7-1	December
► 1008935	water canyon farmers market	2013-1-1	End Date 12/31/1
1004248	Cadillac Area Farmers Market	2012-6-15	June
1004753	The Market at Bissell Gardens	2012-7-11	November
1000784	Weirton Farmers Market	2012-7-9	October
1000789	Weirton Events Center	2012-7-12	October
1002312	Highland Farmers Market	2012-6-2	October
1002313	South Pearl Street Farmers' Market	2012-5-20	October
1002713	Troostwood Youth Garden Market	2011-5-21	October
1003469	Acton-Bedxborough Farmers Market	2012-6-19	October
1003827	Ellsworth Farmers Market	2012-6-11	October
1004072	Downtown San Leandro Certified Farm...	2012-4-18	October
1004268	Emmett Farmers Market	2012-4-28	October
1005772	Silverton Farmers Market	2011-5-7	October
1007647	Andover Farmers Market	2012-5-1	October
1018884	Downtown Bloomington Association Fa...	2018-5-5	October
1011965	Wauconda Farmers Market	June 23	Sept. 8
1001722	Greenville Farmers Market	2012-6-2	September

Here we found 6 rows with invalid season end month "ber" and valid end date with "End Date" text.

10. Store Timings

```
1 •  SELECT DISTINCT SUBSTR(TIME,1,3) AS WEEK_DAY
2   FROM STORE_SEASONAL_HOURS;
```

WEEK_DAY
tue
sat
wed
thu
fri
mon
sun

All weekday values were found valid

12. Website

```
1 •  SELECT * FROM STORE_DETAILS
2   WHERE LOWER(SUBSTR(WEBBSITE,1,7)) NOT LIKE "http://"
3   AND LOWER(SUBSTR(WEBBSITE,1,8)) NOT LIKE "https://"
4   AND WEBSITE IS NOT NULL;
```

FMID	MARKET_NAME	WEBSITE	FACEBOOK	TWITTER	YOUTUBE	OTHER_MEDIA	STREET	CITY_ID	COUNTY_ID	STATE_ID	ZIP	LATITUDE	LONGITUDE	LOCATION	UPDATE_TIME
NULL	NULL	NULL	NULL	NULL	NULL	NULL	NULL	NULL	NULL	NULL	NULL	NULL	NULL	NULL	NULL

There we no nonblank Website values found that did not start with either “<http://>” or “<https://>”

4. Creating a workflow model [10%]

- Identifying the key inputs, outputs of your workflow along with the dependencies [3%]

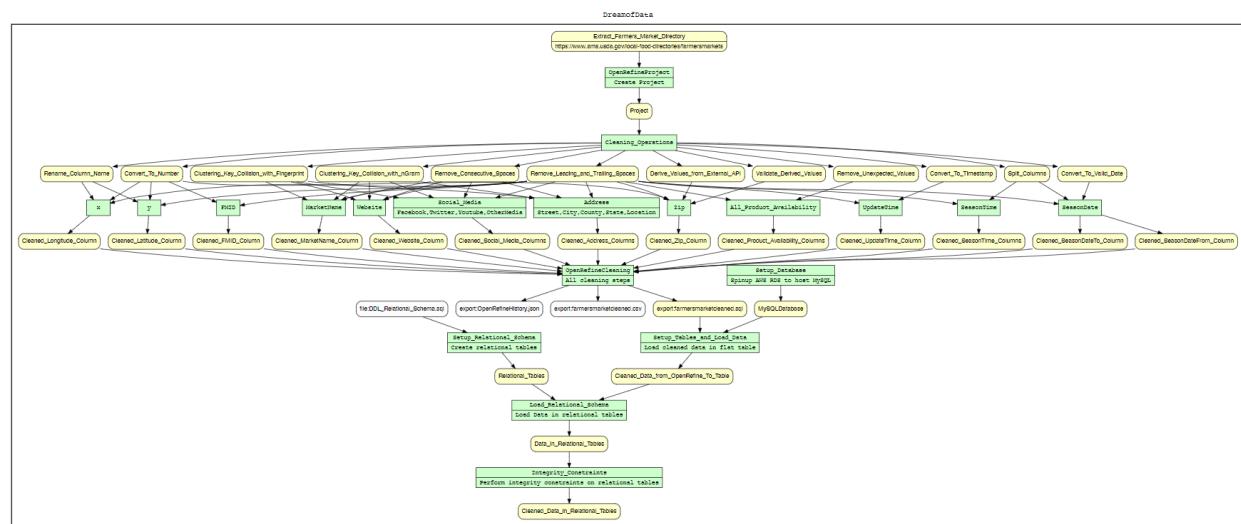
Step	Begin	In	Out	Params	URI
1	Extract Farmers Market Directory	Link to the data source	Farmers Market csv file		https://www.ams.usda.gov/local-food-directories/farmersmarkets
2	Create Open Refine Project	Farmers Market csv file	Open Refine Project		
3	Open Refine Project				

3.1	Cleaning Steps				
3.1.1	Remove Leading and Trailing Spaces	All Columns			
3.1.2	Remove consecutive spaces	Selected Columns			
3.1.3	Clustering				
3.1.3.1	Key-Collision with fingerprint	MarketName, Website, Facebook, Twitter, Youtube, OtherMedia, street, city, County, State, Location	MarketName , Website, Facebook, Twitter, Youtube, OtherMedia, street, city, County, State, Location		
3.1.3.2	Key-Collision with n-gram	MarketName, Website, Facebook, Twitter, Youtube, OtherMedia, street, city, County, State, Location	MarketName , Website, Facebook, Twitter, Youtube, OtherMedia, street, city, County, State, Location		
3.1.4	Convert to Number	FMID, x, y, zip	FMID, x, y, zip		
3.1.5	Rename Column Names	x, y	longitude, latitude		
3.1.6	Remove values except "Y" and "N" and replace with "UNKNOWN"	All products	All Products		
3.1.7	Convert to Timestamp	updateTime	updateTime		
3.1.8	Split Columns and converted to valid date formats	SeasonDates and SeasonTimes	Season{season_no}FromDate, Season{season_no}ToDate, , Season{season_no}Time{count}	Season{season _no}Date, seperator("to") , Season{season _no}Time	

3.1.9	Derive values using external API	Zip, x, y	GoogleAPI_ZIP	x,y	
3.1.10	Validate derived values	GoogleAPI_ZIP	Cleaned CSV File		
3.2	Export DDL and DML from Open Refine		SQL Files		
3.3	Setup Database instance		DB Instance		
3.4	Create Tables for Cleaned Data	SQL Files, DB Instance	Cleaned Data in Table from Open Refine		
3.5	Create Relational Schema and Tables		Tables, Schema		
3.6	Load data to relational schema from cleaned data	Schema, Cleaned Table Data as SQL	Cleaned Master Data in Relational Tables		
3.7	Perform and Identify Integrity Constraints on data which comply and doesn't comply with Integrity Constraints	Cleaned Master Data in Relational Tables	Cleaned Data for Application use and further analysis		

b. A visual representation of your overall workflow, e.g., using YesWorkflow [4%]

See attachment: Overall_YW_Diagram.pdf



- c. A visual representation of your OpenRefine workflow, e.g., using OR2YWTool [3%]

See attachment:

1. output_graph_all_openrefine_steps.pdf
2. output_graph_openrefine_steps_zip.pdf

5. Other factors [10%]

a. Further analysis/takeaways/challenges [5%] (In Progress)

Further Analysis

- After cleaning steps, identified Data which doesn't comply Integrity Constraints could be deleted as required
- Validation of Websites could be done using external API's to check for domain existence

Challenges:

- ZIP codes, all address fields majority of the data is blank and inaccurate
- Normalization is not ideal but could be improved

REFERENCES:

(Farmers Market, 2020)

<https://www.ams.usda.gov/local-food-directories/farmersmarkets>

<https://www.listplanit.com/list-of-categories-for-an-organized-grocery-list/>

https://developers.google.com/maps/documentation/geocoding/intro?hl=en_US

<http://try.yesworkflow.org/>

APPENDIX

Folder and File Details:

1. Project Report: This folder contains project report.
2. Source Data: This folder contains raw farmers market csv file used for cleaning.
3. SQL Scripts: This folder contains all SQL files used to create table, load tables and perform integrity constraints.
4. OpenRefine Cleaning: This folder contains json file representing all cleaning steps performed in OpenRefine.
5. Relational Schema Diagram: This folder contains relational schema diagram.
6. Cleaned Data: This folder contains cleaned farmers market csv file exported from OpenRefine.
7. YesWorkflow: This folder contains scripts used to generate YesWorkflow and diagrams for both – Overall cleaning process and OpenRefine cleaning steps.