

CS 513 – DATA CLEANING FINAL PROJECT

PROJECT TITLE	CS 513 SUMMER 2020 – DATA CLEANING PROJECT
TEAM	TEAM SPIRIT
COLLABORATORS	BALAJI SATHYAMURTHY (balajis2) GOWRI SHANKAR RAMANAN (gsr2)

PROJECT OBJECTIVE

The objective of this data cleaning project is to refine the publicly available dataset “US Farmer market” dataset available in the website to derive valuable business insights.

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

We used OpenRefine v3.3 to perform the initial cleaning of the raw dataset and then SQLite3 to perform integrity constraint checks and violation clean-ups on the dataset. The final cleaned dataset along with the original raw dataset was uploaded in the Box folder at <https://uofi.box.com/s/gkictao5xrnrpd3va3qwt6fkoxuv63gz> for review. They were further validated for special use-cases. The YesWorkflow model is used to annotate the various steps followed in the data cleaning process. We used the OpenRefine to YesWorkflow model toolkit (OR2YW v0.0.16) to render visual representations of workflows and provenance using Graph Viz.

PART 1: OVERVIEW AND INITIAL ASSESSMENT ON THE DATASET

1.1 DATASET OVERVIEW

- The US Farmers market dataset contains a total of 59 data elements.
- FMID is a numeric field with a unique identifier to identify each Farmers market. There seems to be no duplicate values in FMID.
- Market name is a text field containing the name of the farmer's market.
- Website, Facebook, Twitter, YouTube, and Other Media are all text fields containing either website or social media contacts such as twitter ID, channel name etc.,
- Street, City, County, State and Zip contains address information about the farmer's market.
- Season1Date, Season1Time, Season2Date, Season2Time, Season3Date, Season3Time, Season4Date, Season4Time contains various season date and time the respective Farmer's market are open during a year.

- Columns X and Y looks like to represent the geographic latitude and longitude of the of the farmer's market
- Column Location provides additional location details such as nearest landmark of the Farmer's market. This seems to be a text field.
- Credit is a single character field with either Y or N denoting whether the payment via credit card is accepted or not. On initial assessment it seems all the records have this value populated in the dataset.
- Columns WIC, WICcash, SFMNP, SNAP are all single character fields with either Y or N values denoting whether payment through food assistance programs is either accepted or not
- A list of 30 different indicator fields with respect to the characteristics of various food products sold in the farmers market. On initial assessment all these columns look good with no discrepancies.
- Updatetime denotes date and time when the data about farmer's market is captured

1.2 LIST OF DATA QUALITY ISSUES

- Market name seems to contain lot of text variations for same Farmer market name.
- Most of the values for the social media columns are blank
- County and city names contain variations in text for the same value.
- zip contains some invalid zip codes.
- The date format for season dates seems to be not consistent across the records. Some records contain only month name, and some contains only from date etc.,
- Most of the farmer's market have season1 and season 2 values and only very less records contain season three and four values.
- Columns X and Y names are not clear, and we assume it represents the geographic location of the latitude and longitude of the farmers market.

1.3 IS THE DATA CLEAN ENOUGH FOR THE USE CASES?

Though dataset have the above noted quality issues, we could use the dataset for the below list of use cases:

- To identify the list of Farmer's markets nearest based on the geographic location.
- To identify the list of Farmer's markets which are open based on seasonality.
- To identify list of Farmer's markets which accept credit card
- To identify list of Farmer's markets which accept food coupons such as WIC, WICcash etc.,
- To identify list of Farmer's markets based on the various sold products such as seafood, pet food, vegetables etc.,

PART 2: DATA CLEANING STEPS USING OpenRefine

The following steps have been performed to clean the data via OpenRefine.

1. FMID

The FMID field is converted to numeric field as shown below:

8812 rows

Show as: **rows** records Show: 5 10 25 50 rows

All	FMID	MarketName	Website	Facebook
1.	Facet	Market	https://sites.google.com/site/caledoniafarmersmarket/	https://www.facebook.com/D
2.	Edit cells	Transform...		
3.	Edit column	Common transforms		
4.	Transpose	Trim leading and trailing whitespace		
5.	Sort...	Collapse consecutive whitespace		
6.	View	Unescape HTML entities		
7.	Reconcile	Replace Smart quotes with ascii		
8.	1011100	12 South Farmers Mar	http://www.125thStreetFarmersMarket.com	Market
9.	1009845	125th Street Fresh Connect Farmers' Market	http://www.125thStreetFarmersMarket.com	Market
10.	1005586	12th & Brandywine Urban Farm	http://www.125thStreetFarmersMarket.com	Market

Context menu for cell 'MarketName' (row 1, column 3):

- Transform...
- Common transforms
 - Trim leading and trailing whitespace
 - Collapse consecutive whitespace
 - Unescape HTML entities
 - Replace Smart quotes with ascii
 - To titlecase
 - To uppercase
 - To lowercase
 - To number**
 - To date
 - To text
 - To null
 - To empty string

2. MARKET NAME

To canonicalize the same market name with different text variations, the text Facets and cluster feature available in OpenRefine is used as shown below:

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 239 clusters found

# Choices in Cluster	# Rows in Cluster	Average Length of Choices	Length Variance of Choices
2 — 4	2 — 34	14 — 71	0 — 2.5
2	2	2	2
3	2	2	2
2	5	2	2
2	2	2	2
2	2	2	2
2	2	2	2
2	2	2	2
2	2	2	2

Clusters:

- Bourbon Farmers' Market (1 rows)
- Montgomery Farmers Market (2 rows)
 - MONTGOMERY FARMERS MARKET (1 rows)
- St. Paul Farmers Market (1 rows)
 - St. Paul Farmers' Market (1 rows)
- Independence Farmers Market (4 rows)
 - Independence Farmers' Market (1 rows)
- Market in the Loop (1 rows)
 - The Market in the Loop (1 rows)
- Northville Farmers Market (1 rows)
 - Northville Farmers' Market (1 rows)
- Albany Downtown Farmers Market (1 rows)
 - Downtown Albany Farmers Market (1 rows)
- Hobart Farmers Market (1 rows)
 - Hobart Farmers' Market (1 rows)

Buttons:

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

3. WEBSITE, TWITTER, FACEBOOK, OTHER MEDIA and YOUTUBE

To canonicalize the social media columns text variations, the text Facets and cluster feature available in OpenRefine is used as shown below:

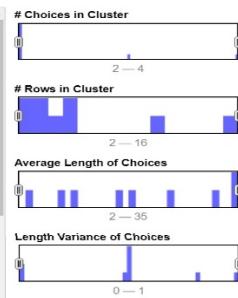
Cluster & Edit column "OtherMedia"

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 9 clusters found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
4	16	<ul style="list-style-type: none"> Instagram (10 rows) Instagram (3 rows) Instagram (2 rows) Instagram: (1 rows) 	<input checked="" type="checkbox"/>	Instagram
3	11	<ul style="list-style-type: none"> n/a (7 rows) N/A (3 rows) NA (1 rows) 	<input checked="" type="checkbox"/>	
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: visalifafarmersmarket (3 rows) visalifafarmersmarket (Instagram) (1 rows) 	<input checked="" type="checkbox"/>	instagram: visalifafarmersmarket
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/morgantownmkt/
2	3	<ul style="list-style-type: none"> smfms (Instagram) (2 rows) Instaface-SMfMS (1 rows) 	<input checked="" type="checkbox"/>	smfms (Instagram)

Export Clusters | Merge Selected & Re-Cluster | Merge Selected & Close | Close



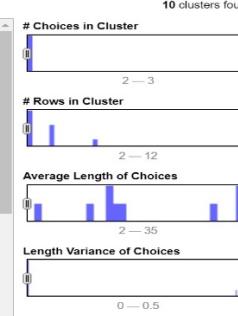
Cluster & Edit column "Twitter"

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 10 clusters found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
3	12	<ul style="list-style-type: none"> n/a (7 rows) N/A (3 rows) NA (2 rows) 	<input checked="" type="checkbox"/>	
2	3	<ul style="list-style-type: none"> @KSGrownMarket (2 rows) @ksgrownmarket (1 rows) 	<input checked="" type="checkbox"/>	@KSGrownMarket
2	3	<ul style="list-style-type: none"> @spotsyfarmmkt (2 rows) @SpotsyFarmMkt (1 rows) 	<input checked="" type="checkbox"/>	@spotsyfarmmkt
2	2	<ul style="list-style-type: none"> https://twitter.com/BosPublicMarket (1 rows) https://twitter.com/bospublicmarket (1 rows) 	<input checked="" type="checkbox"/>	https://twitter.com/BosPublicMarket
2	3	<ul style="list-style-type: none"> https://twitter.com/AthFarmersMkt (2 rows) https://twitter.com/AthFarmersMkt (1 rows) 	<input checked="" type="checkbox"/>	https://twitter.com/AthFarmersMkt
2	2	<ul style="list-style-type: none"> https://twitter.com/WilliamsburgFM (1 rows) 	<input checked="" type="checkbox"/>	https://twitter.com/WilliamsburgFM

Export Clusters | Merge Selected & Re-Cluster | Merge Selected & Close | Close



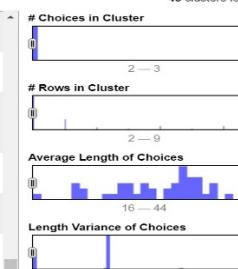
Cluster & Edit column "Website"

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: fingerprint 45 clusters found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
2	2	<ul style="list-style-type: none"> http://www.hornellhpg.com (1 rows) http://www.hornellhpg.com (1 rows) 	<input checked="" type="checkbox"/>	http://www.hornellhpg.com
2	2	<ul style="list-style-type: none"> http://www.greenmarketco-op.org (1 rows) http://www.greenmarketco-op.org (1 rows) 	<input checked="" type="checkbox"/>	http://www.greenmarketco-op.org
2	2	<ul style="list-style-type: none"> http://www.capitaldistrictfarmersmarket.org (1 rows) http://www.capitaldistrictfarmersmarket.org/ (1 rows) 	<input checked="" type="checkbox"/>	http://www.capitaldistrictfarmersm
2	5	<ul style="list-style-type: none"> http://foodinroot.com (3 rows) http://foodinroot.com/ (2 rows) 	<input checked="" type="checkbox"/>	http://foodinroot.com
2	3	<ul style="list-style-type: none"> http://prescottfarmersmarket.org (2 rows) http://prescottfarmersmarket.org/ (1 rows) 	<input checked="" type="checkbox"/>	http://prescottfarmersmarket.org
2	3	<ul style="list-style-type: none"> http://www.seela.org (2 rows) http://www.seela.org/ (1 rows) 	<input checked="" type="checkbox"/>	http://www.seela.org
2	2	<ul style="list-style-type: none"> http://www.saltairmarket.com (1 rows) http://www.saltairmarket.com/ (1 rows) 	<input checked="" type="checkbox"/>	http://www.saltairmarket.com

Export Clusters | Merge Selected & Re-Cluster | Merge Selected & Close | Close



4. CITY AND COUNTY NAME

To canonicalize the city and county name text variations, the text Facets and cluster feature available in OpenRefine is used as shown below:

Cluster & Edit column "city"

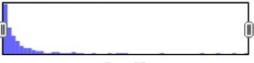
This feature helps you find groups of different cell values that might be alternative representations of the same thing. For example, the two strings "New York" and "new york" are very likely to refer to the same concept and just have capitalization differences, and "Gödel" and "Godel" probably refer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

Method: key collision Keying Function: fingerprint 467 clusters found

Count	Count	Choices	Action
2	2	• Bainbridge (1 rows) • bainbridge (1 rows)	<input checked="" type="checkbox"/> Bainbridge
2	3	• Bonita Springs (2 rows) • Bonita Springs (1 rows)	<input checked="" type="checkbox"/> Bonita Springs
2	5	• Rockville (4 rows) • Rockville (1 rows)	<input checked="" type="checkbox"/> Rockville
2	2	• Findlay (1 rows) • Findlay (1 rows)	<input checked="" type="checkbox"/> Findlay
2	4	• Livermore (3 rows) • Livermore (1 rows)	<input checked="" type="checkbox"/> Livermore
2	5	• Duluth (4 rows) • Duluth (1 rows)	<input checked="" type="checkbox"/> Duluth
2	2	• Fairmont (1 rows) • Fairmont (1 rows)	<input checked="" type="checkbox"/> Fairmont

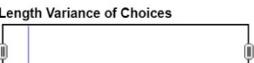
Choices in Cluster

2 — 4

Rows in Cluster

2 — 62

Average Length of Choices

3 — 23

Length Variance of Choices

0 — 4.97

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

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: fingerprint 28 clusters found

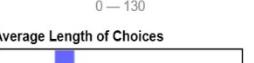
Count	Count	Choices	Action
2	2	• SONOMA (1 rows) • ELMORE (1 rows) • Elmore (1 rows)	<input checked="" type="checkbox"/> ELMORE
2	3	• Wilcox (2 rows) • WILCOX (1 rows)	<input checked="" type="checkbox"/> Wilcox
2	3	• LaGrange (2 rows) • Lagrange (1 rows)	<input checked="" type="checkbox"/> LaGrange
2	17	• San Luis Obispo (16 rows) • SAN LUIS OBISPO (1 rows)	<input checked="" type="checkbox"/> San Luis Obispo
2	3	• Inyo (2 rows) • INYO (1 rows)	<input checked="" type="checkbox"/> Inyo
2	27	• El Paso (26 rows) • EL Paso (1 rows)	<input checked="" type="checkbox"/> El Paso
2	54	• Monroe (53 rows) • MONROE (1 rows)	<input checked="" type="checkbox"/> Monroe

Choices in Cluster

2 — 3

Rows in Cluster

0 — 130

Average Length of Choices

3 — 15

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

5. GREL – VALID ZIP CODE

The regular expression is used to identify the valid zip code formats (XXXXXX and XXXXX-XXXX) and a new column "grel_valid_zip" is added to the dataset as shown below:

Add column based on column zip

New column name: `grel_valid_zip`

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

Expression: `if(value.match("^\\d{5}(?:\\d{4})?") != null,value,\"\")`

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

row	value	if(value.match("^\\d{5}(?:\\d{4})?") ...
1.	5828	
2.	null	
3.	29682	29682
4.	64759	64759
5.	10029	10029
6.	37204	37204
7.	10027	10027

OK Cancel

6. SEASON DATE – SPLIT AS FROM AND TO DATES

Python is used to identify the season dates in MM/DD/YYYY format and split them in two dates date columns as “From date” and “To Date” so that seasonality of the market can be identified easily. All the season date columns – season1, season2, season 3 and season 4 have been split into two date columns as shown below:

Add column based on column Season1Date

New column name: `Season1FromDate`

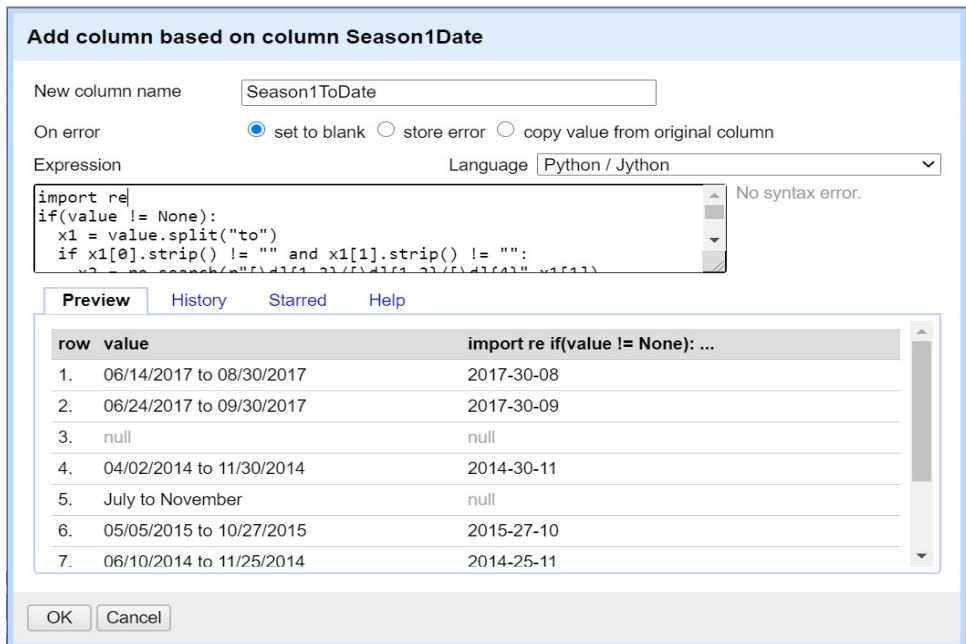
On error: set to blank store error copy value from original column

Expression: `x2 = re.search(r'(\d{1,2},\d{1,2},\d{4})-(\d{1,2},\d{1,2},\d{4})',x1[0])
if x2 != None:
 x3 = x1[0].split("/")
 return x3[2].strip() + "-" + x3[1].strip() + "-" +
x3[0].strip()`

Language: Python / Jython No syntax error.

row	value	import re if(value != None): ...
1.	06/14/2017 to 08/30/2017	2017-14-06
2.	06/24/2017 to 09/30/2017	2017-24-06
3.	null	null
4.	04/02/2014 to 11/30/2014	2014-02-04
5.	July to November	null
6.	05/05/2015 to 10/27/2015	2015-05-05
7.	06/10/2014 to 11/25/2014	2014-10-06

OK Cancel



7. COLUMNS X AND Y

The X and Y columns are converted to numeric field as shown below:

8812 rows

Show as: rows records Show: 5 10 25 50 rows

Extensions

« first < previous 1 - 1

Season4Date	Season4FromDa	Season4ToDate	Season4Time	x	y	Location	Credit	WIC	WICcash	SF
				-72.140335	Facet	Y	Y	N	Y	
				-81.73394	Text filter					
					Edit cells	Transform...				
						Common transforms				
						Trim leading and trailing whitespace				
						Collapse consecutive whitespace				
						Unescape HTML entities				
						Replace Smart quotes with ascii				
						To titlecase				
						To uppercase				
						To lowercase				
						To number				
						To date				
						To text				
						To null				
						To empty string				

8. UPDATE TIME

Update time column is converted to standard date column as shown below:

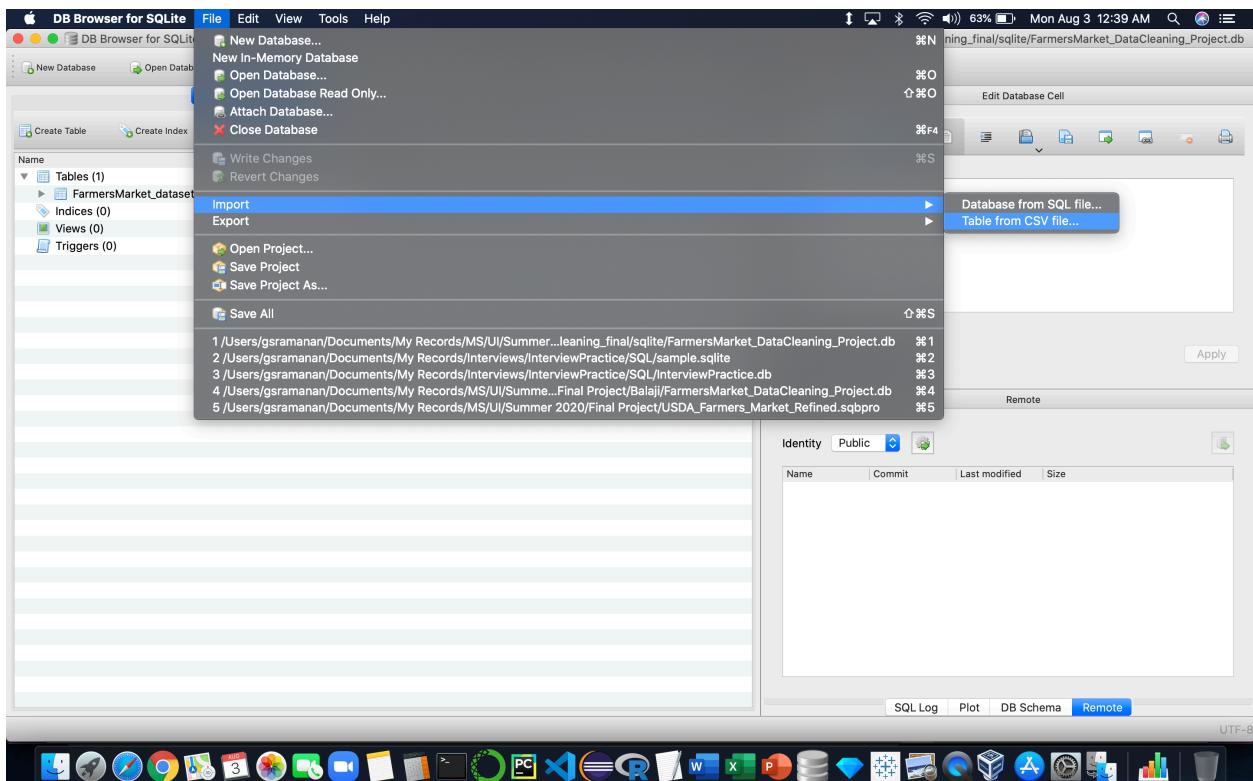
8812 rows Extensions: Wikidata

Show as: rows records Show: 5 10 25 50 rows

red	Soap	Trees	Wine	Coffee	Beans	Fruits	Grains	Juices	Mushrooms	PetFood	Tofu	WildHarvested	updateTime
Y	Y	N	Y	Y	Y	N	N	Y	Y	N	Facet	0/2017 22:43	
Y	N	N	N	N	Y	N	N	N	Transform...	Edit cells	1/2017 17:15		
									Trim leading and trailing whitespace	Common transforms	Edit column	13	
									Collapse consecutive whitespace		Transpose		
Y	N	N	N	N	Y	Unescape HTML entities			Fill down	Sort...	28/2014 9:49		
Y	N	N	N	N	Y	Replace Smart quotes with ascii			Blank down	View			
Y	N	N	N	N	Y	To titlecase			Split multi-valued cells...	Reconcile	2/2012 10:38		
Y	N	N	Y	N	Y	To uppercase			Join multi-valued cells...				
Y	N	Y	Y	N	Y	To lowercase			Cluster and edit...		5/1/2015 10:40		
N	N	N	N	N	N	To number			Replace				
						To date			N		4/7/2014 16:32		
						To text							
						To null					4/3/2014 15:43		
						To empty string							

PART 3: DEVELOPING A RELATIONAL SCHEMA

The refined data from OpenRefine is exported as csv file and imported into the below table in SQLite to perform further SQL data profiling and Integrity Check queries.

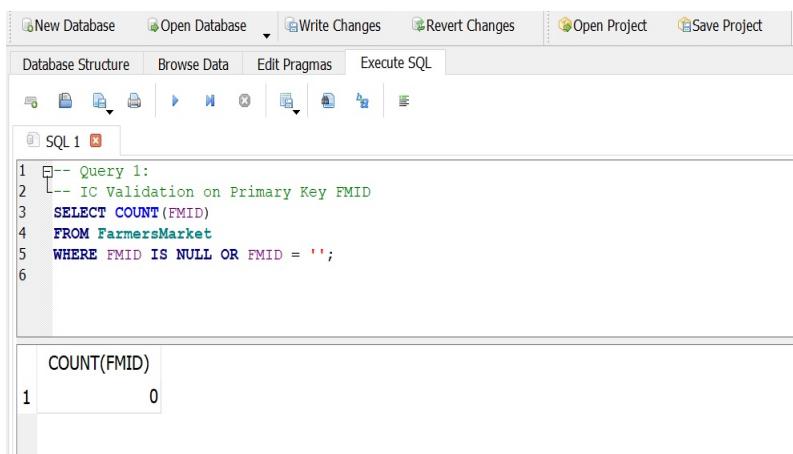


COLUMN NAME	DATA TYPE	Primary Key
FMID	INTEGER	Yes
MarketName	TEXT	No
Website	TEXT	No
Facebook	TEXT	No
Twitter	TEXT	No
Youtube	TEXT	No
OtherMedia	TEXT	No
street	TEXT	No
city	TEXT	No
County	TEXT	No
State	TEXT	No
zip	TEXT	No
grel_valid_zip	TEXT	No
Season1Date	TEXT	No
Season1FromDate	TEXT	No
Season1ToDate	TEXT	No
Season1Time	TEXT	No
Season2Date	TEXT	No
Season2FromDate	TEXT	No
Season2ToDate	TEXT	No
Season2Time	TEXT	No
Season3Date	TEXT	No
Season3FromDate	TEXT	No
Season3ToDate	TEXT	No
Season3Time	TEXT	No
Season4Date	TEXT	No
Season4FromDate	TEXT	No
Season4ToDate	TEXT	No
Season4Time	TEXT	No
x	REAL	No
y	REAL	No
Location	TEXT	No
Credit	TEXT	No
WIC	TEXT	No
WICcash	TEXT	No
SFMNP	TEXT	No
SNAP	TEXT	No
Organic	TEXT	No
Bakedgoods	TEXT	No
Cheese	TEXT	No
Crafts	TEXT	No

Flowers	TEXT	No
Eggs	TEXT	No
Seafood	TEXT	No
Herbs	TEXT	No
Vegetables	TEXT	No
Honey	TEXT	No
Jams	TEXT	No
Maple	TEXT	No
Meat	TEXT	No
Nursery	TEXT	No
Nuts	TEXT	No
Plants	TEXT	No
Poultry	TEXT	No
Prepared	TEXT	No
Soap	TEXT	No
Trees	TEXT	No
Wine	TEXT	No
Coffee	TEXT	No
Beans	TEXT	No
Fruits	TEXT	No
Grains	TEXT	No
Juices	TEXT	No
Mushrooms	TEXT	No
PetFood	TEXT	No
Tofu	TEXT	No
WildHarvested	TEXT	No
updateTime	TEXT	No

1. DUPLICATE FMID

On query analysis, there is no duplicate value for FMID field in the dataset.



The screenshot shows a SQLite database interface with the following details:

- Toolbar:** Includes buttons for New Database, Open Database, Write Changes, Revert Changes, Open Project, Save Project, Database Structure, Browse Data, Edit Pragmas, and Execute SQL.
- Query Window:** Labeled "SQL 1" with a "Run" button. The query is:

```

1 -- Query 1:
2 --- IC Validation on Primary Key FMID
3 SELECT COUNT(FMID)
4 FROM FarmersMarket
5 WHERE FMID IS NULL OR FMID = '';
6

```
- Results Window:** Shows the result of the query:

COUNT(FMID)
1 0

2. MARKETS WITH FMID OR MARKETNAME IS BLANK OR NULL

On analysis there seems to be no FMID or Market name having blank or null values in the dataset

The screenshot shows the SQLite Manager interface with a SQL tab open. The query is:

```
1 -- Query 2:
2 --- IC Validation based on Natural Unique Keys FMID and MarketName fields
3 SELECT COUNT(FMID)
4 FROM FarmersMarket
5 WHERE FMID IS NULL OR FMID = ''
6 OR MarketName IS NULL OR MarketName = '';
```

The results table shows:

COUNT(FMID)
1
0

3. CHECKING FOR SAME MARKET NAME WITHIN SAME CITY / STATE / ZIP BUT DIFFERENT FMID.

The below sql shows there are same market name within the same city and state but with different FMID.

The screenshot shows the SQLite Manager interface with a SQL tab open. The query is:

```
1 -- Query 3:
2 --- IC Validation based on Duplication check on MarketName column
3 SELECT MarketName, city, State, grel_valid_zip, COUNT(FMID)
4 FROM FarmersMarket
5 GROUP BY MarketName, city, State, grel_valid_zip
6 HAVING COUNT(FMID) > 1
7 ORDER BY COUNT(FMID) DESC;
```

The results table shows:

	MarketName	city	State	grel_valid_zip	COUNT(FMID)
1	Country Farm & Craft Market ...	Paso Robles	California	93446	3
2	Middlebury Farmers' Market	Middlebury	Vermont	NULL	3
3	Allen County Farmers' Market	IOLA	Kansas	66749	2
4	BOISE FARMERS MARKET	Boise City	Idaho	83702	2
5	Bigfork Farmers' Market ...	Bigfork	Montana	59911	2
6	Burlington Farmers Market	Burlington	Vermont	NULL	2

On further analysis, we identify around 70 markets to be removed due to the same FMID:

```

1 -- Query 4:
2 -- Identify ideal duplicate candidates to be removed based on duplicate MarketNames
3 WITH ValidMarkets AS
4 (SELECT MarketName, city, State, grel_valid_zip, MIN(FMID) AS ValidFMID
5 FROM FarmersMarket
6 GROUP BY MarketName, city, State, grel_valid_zip
7 HAVING COUNT(FMID) > 1
8 ORDER BY COUNT(FMID) DESC)
9
10 SELECT count(*) FROM FarmersMarket AS FM, ValidMarkets AS VM
11 WHERE FM.MarketName = VM.MarketName AND FM.city = VM.city AND FM.State = VM.state --AND FM.grel_valid_zip = V
12 AND FM.FMID <> VM.ValidFMID
13 ORDER BY FM.MarketName;
14
15

```

	count(*)
1	70

Deleting the identified duplicate market entries:

```

1 -- Query 5:
2 -- Remove Duplicates/ Invalid Entries which violates IC based on queries above
3 WITH ValidMarkets AS
4 (SELECT MarketName, city, State, grel_valid_zip, MIN(FMID) AS ValidFMID
5 FROM FarmersMarket
6 GROUP BY MarketName, city, State, grel_valid_zip
7 HAVING COUNT(FMID) > 1
8 ORDER BY COUNT(FMID) DESC)
9
10 DELETE FROM FarmersMarket WHERE FMID IN
11 (SELECT FM.FMID FROM FarmersMarket AS FM, ValidMarkets AS VM
12 WHERE FM.MarketName = VM.MarketName AND FM.city = VM.city AND FM.State = VM.state
13 --AND FM.zip = VM.zip
14 AND FM.FMID <> VM.ValidFMID
15 ORDER BY FM.MarketName)
16 ;
17 SELECT COUNT(*) FROM FarmersMarket;

```

	COUNT(*)
1	8742

4. MARKETS WITH VALID ZIP CODE FOR GEOGRAPHIC POSITIONING

For the use case to identify the nearby markets around a geographic position, the grel_valid_zip field is selecting accurate date on zip codes.

Database Structure Browse Data Edit Pragmas Execute SQL

SQL 1

```

1 --Query 6
2 --Check for records with valid zipcode
3 SELECT FMID,grel_valid_zip
4 FROM FarmersMarket
5 WHERE TRIM(grel_valid_zip) <> ""

```

	FMID	grel_valid_zip
1	1009364	29682
2	1010691	64759
3	1002454	10029
4	1011100	37204
5	1000845	10027

5. MARKETS WITH VALID SEASON DATES

For the use case to identify markets open for specific seasonality, the season from and to dates column fetch the appropriate values.

Database Structure Browse Data Edit Pragmas Execute SQL

SQL 1

```

1 -- Query 7:
2 --Check if the from date is greater than to date for the season dates
3 SELECT FMID,MarketName,
4       STRFTIME('%Y-%m-%d',Season1FromDate) AS SEASON1FROMDT,
5       STRFTIME('%Y-%m-%d',Season1ToDate)   AS SEASON1TODT
6 FROM FarmersMarket
7 WHERE DATE(Season1ToDate) < DATE(Season1FromDate)
8   AND Season1FromDate IS NOT NULL
9   AND Season1ToDate IS NOT NULL;
10

```

	FMID	MarketName	SEASON1FROMDT	SEASON1TODT
1	1000709	22nd Annual Highlands Busine...	2016-06-25	2016-05-11
2	1004950	3 French Hens French Country ...	2018-12-05	2018-10-13
3	1018656	31 & Main Farmers Market at ...	2017-11-06	2017-10-29
4	1012280	39 North Marketplace	2016-09-06	2016-08-18

Extracting the month and date part of the season dates:

Database Structure Browse Data Edit Pragmas Execute SQL

SQL 1

```

1 --Query 8
2 --Derive Month and Day Part of the season dates to determine the seasonality of the Farmers Market
3 SELECT FMID,
4       STRFTIME('%d', DATE(Season1FromDate)) AS season1fromdate_daypart,
5       STRFTIME('%m', DATE(Season1FromDate)) AS season1fromdate_monthpart,
6       STRFTIME('%d', DATE(Season1ToDate)) AS season1todate_daypart,
7       STRFTIME('%m', DATE(Season1ToDate)) AS season1todate_monthpart
8 FROM FarmersMarket
9
10 WHERE Season1ToDate IS NOT NULL

```

	FMID	season1fromdate_daypart	season1fromdate_monthpart	season1todate_daypart	season1todate_monthpart
1	1018261	14	06	30	08
2	1018318	24	06	30	09
3	1010691	04	02	30	11

6. MARKETS WITH CREDIT CARD PAYMENT

The data column for credit card payment is fetching the appropriate market names.

Database Structure Browse Data Edit Pragmas Execute SQL

SQL 1

```

1 --Query 9
2 --Check for MarketName which accept credit card payment
3 SELECT FMID,MarketName
4 FROM FarmersMarket
5 WHERE Credit = 'Y'

```

	FMID	MarketName
1	1018261	Caledonia Farmers Market ...
2	1018318	Stearns Homestead Farmers' ...
3	1009364	106 S. Main Street Farmers ...
4	1010691	10th Street Community Farmer...

7. MARKETS BASED ON VARIOUS SOLD PRODUCT CATEGORIES

The data column for various sold product categories is fetching the appropriate market names.

The screenshot shows the OpenRefine interface with the following components:

- Toolbar:** Database Structure, Browse Data, Edit Pragmas, Execute SQL.
- SQL Editor:** A panel titled "SQL 1" containing the following SQL code:


```

1 --Query10
2 --Check for MarketName which sells organic cheese
3 SELECT FMID,MarketName
4 FROM FarmersMarket
5 WHERE (Organic = 'Y' and Cheese = 'Y')

```
- Results Table:** A table showing the results of the query. The columns are FMID and MarketName. The data is as follows:

	FMID	MarketName
1	1018261	Caledonia Farmers Market ...
2	1011100	12 South Farmers Market
3	1009845	125th Street Fresh Connect ...
4	1008071	14&U Farmers' Market
5	1009994	18th Street Farmer's Market
6	1012158	21 Acres Farm Market

PART 4: CREATING A WORKFLOW MODEL

YesWorkflow model/ prototype has been used to annotate the data cleaning procedure followed, and to generate visual model representations using graphviz. The workflows have been created for both overall data cleaning process as well as the open refine tool specific operation history to represent the data provenance better.

4.1 Overall Workflow:

The binaries available at `yw-prototypes` repository (<https://github.com/yesworkflow-org/yw-prototypes>) have been employed to generate these workflows. The scripts (`Overall_Workflow.txt` and `OpenRefine_Workflow.txt`) containing special YesWorkflow (YW) comments will be interpreted by the prototype tool with simple commands, to render these graphical representations in DOT format (.gv).

`Overall_Workflow.txt` (Annotations as python script comments explaining the overall workflow)

```

@begin Overall_Data_Cleaning_Workflow @desc Overall Workflow for USDA Farmers Markets dataset cleaning project
@in FarmersMarket_dataset_Original_Dirty_CSV @uri file://data/FarmersMarket_dataset_Original_Dirty.csv

@begin CreateOpenRefineDataCleaningProject @desc Import the Raw Farmers Markets dataset and Create OpenRefine project for data cleaning
@in FarmersMarket_dataset_Original_Dirty_CSV @uri file://data/FarmersMarket_dataset_Original_Dirty.csv
@out USDA_Farmers_Market_OR_cleaning_project
@end CreateOpenRefineDataCleaningProject

@begin CaptureOpenRefineProvenance @desc Export the provenance/ operations history of data cleaning procedure from OpenRefine project in the
JSON format
@in USDA_Farmers_Market_OR_cleaning_project
@out Open_Refine_History_json @uri file://openrefine/Open_Refine_History.json
@end CaptureOpenRefineProvenance

@begin ExportOpenRefineCleanedDataSet @desc Export the cleaned data set from open refine in csv format
@in USDA_Farmers_Market_OR_cleaning_project
@out FarmersMarket_dataset_OR_Refined_CSV @uri file://openrefine/FarmersMarket_dataset_OR_Refined.csv
@end ExportOpenRefineCleanedDataSet

@begin CreateSQLiteDBSchema @desc Import the OR refined data set into SQLite and create a DB with table schema based on clean data set
@in FarmersMarket_dataset_OR_Refined_CSV @uri file://openrefine/FarmersMarket_dataset_OR_Refined.csv
@out FarmersMarket_DataCleaning_Project_DB @uri file://sqlite/FarmersMarket_DataCleaning_Project.db
@end CreateSQLiteDBSchema

```

```

@begin IntegrityCheckWithSQLite @desc Use sqlite to check integrity of data and record checks in Queries.txt
@in FarmersMarket_DataCleaning_Project_DB @uri file://sqlite/FarmersMarket_DataCleaning_Project.db
@out FarmersMarket_Integrity_Checks @uri file://sqlite/Queries.txt
@end IntegrityCheckWithSQLite

@begin CleanWithSQLite @desc Use sqlite to clean-up IC Violations and remove any duplicates or invalid entries
@in FarmersMarket_DataCleaning_Project_DB @uri file://sqlite/FarmersMarket_DataCleaning_Project.db
@out FarmersMarketFinalCleanedSQLite_DB @uri file://sqlite/FarmersMarket_DataCleaning_Project.db
@end CleanWithSQLite

@begin ExportSQLiteCleanedDataSet @desc Export the final clean data set suitable for required use cases from sqlite in csv format
@in FarmersMarketFinalCleanedSQLite_DB @uri file://sqlite/FarmersMarket_DataCleaning_Project.db
@out FarmersMarket_dataset_Final_Refined_CSV @uri file://data/FarmersMarket_dataset_Final_Refined.csv
@end ExportSQLiteCleanedDataSet

@begin QueryDBForUsecaseFitnessValidation @desc Use sqlite to retrieve in-sights for the required use cases
@in FarmersMarketFinalCleanedSQLite_DB @uri file://sqlite/FarmersMarket_DataCleaning_Project.db
@out FarmersMarket_Fitness_Validations @uri file://sqlite/Queries.txt
@end QueryDBForUsecaseFitnessValidation

@out FarmersMarket_dataset_Final_Refined_CSV @uri file://data/FarmersMarket_dataset_Final_Refined.csv
@end Overall_Data_Cleaning_Workflow
"""

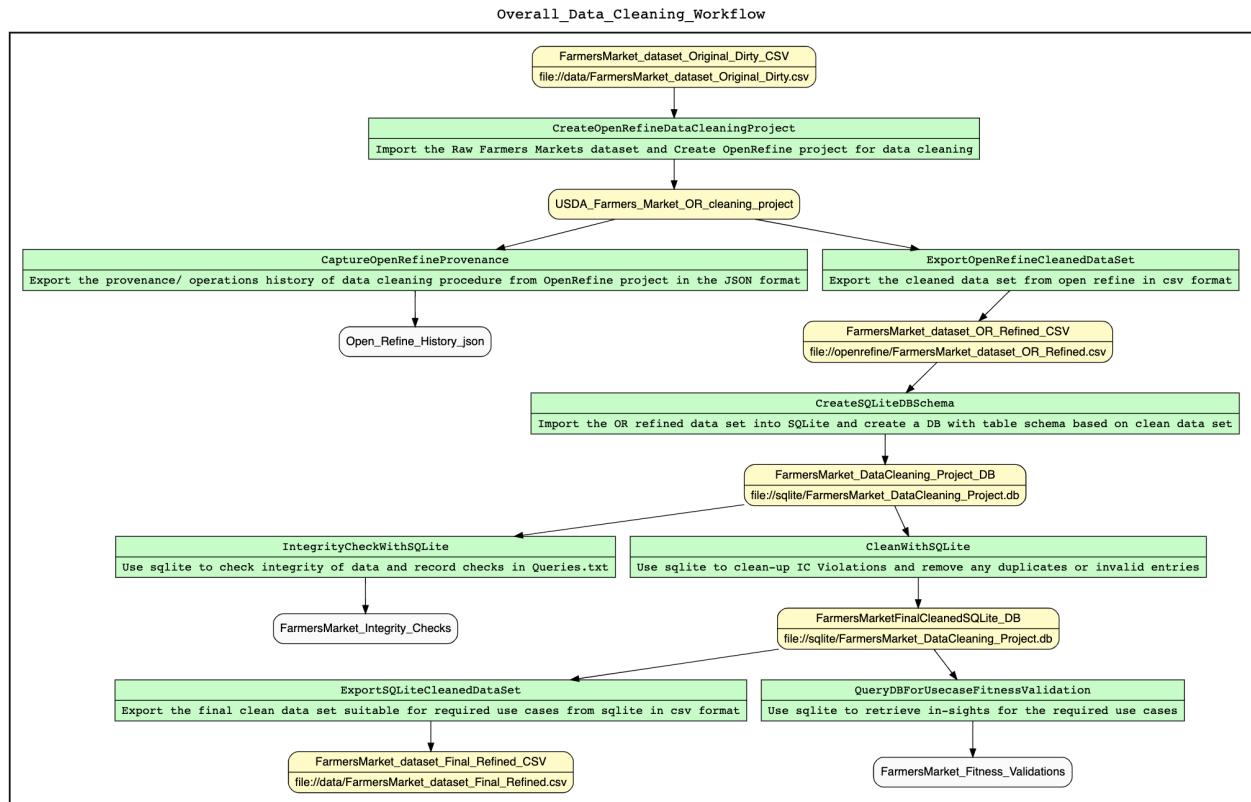
```

The annotations were converted into graphviz representations in DOT format (.gv file) using the following simple commands:

```

$ alias yw='java -jar ~/bin/yesworkflow-0.2.0-jar-with-dependencies.jar'
$ yw graph Overall_Workflow.txt > Overall_Workflow.gv
$ dot -Tpdf Overall_Workflow.gv -o Overall_Workflow.pdf

```



The above overall workflow explains clearly that the data set was cleaned initially with Open Refine to remove data anomalies and then checked for integrity and cleaned up the violations with deeper analysis using SQLite database before producing the final refined dataset in CSV format.

4.2 Open Refine Workflow:

The operations history/ provenance collected from Open Refine tool after the completion of necessary cleaning, is exported in JSON format. This JSON file is further converted into YesWorkflow models and visual representations using OR2YW tool available at <https://pypi.org/project/or2ywtool/>. This tool aims to provide an auto-parsing method from the Open Refine Operation Recipe JSON file to YesWorkflow model which classifies the operations into two levels: schema level and column level.

The Workflow annotations embedded in python script comments are generated using the following simple command:

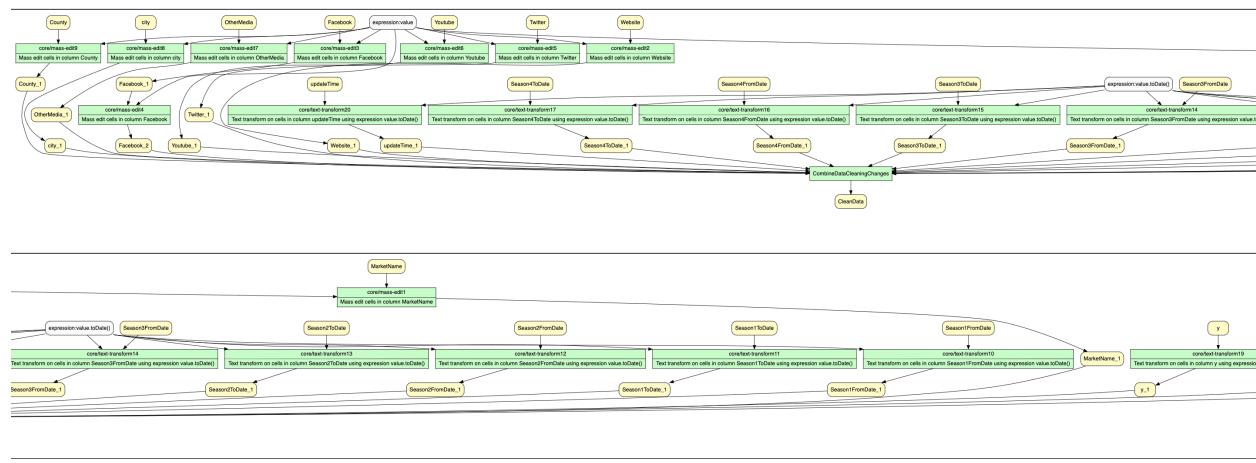
```
$ or2yw -i openrefine/Open_Refine_History.json -o workflows/Open_Refine_Workflow.txt -t parallel
File Open_Refine_Workflow.txt generated.
```

The graphical output is produced in DOT format (.gv file) using the following simple command:

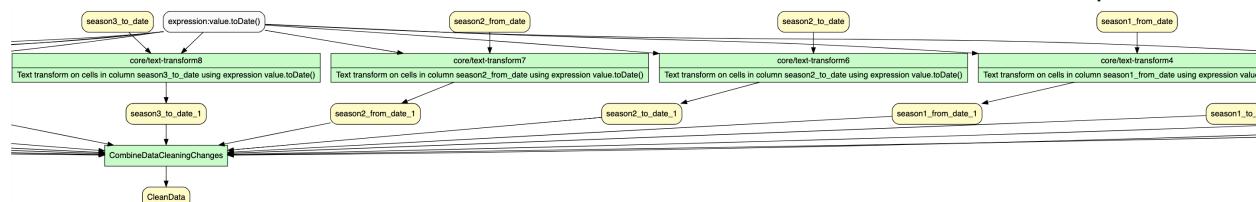
```
$ or2yw -i openrefine/Open_Refine_History.json -o workflows/Open_Refine_Workflow.gv -ot gv -t
parallel
java found: java
File Open_Refine_Workflow.gv generated.
```

The output gv file has then been tested with [GraphViz Online](#) for proper rendering of workflow image. The workflow image has been broken down into multiple chunks as the original output is too large to fit into this report with actual size.

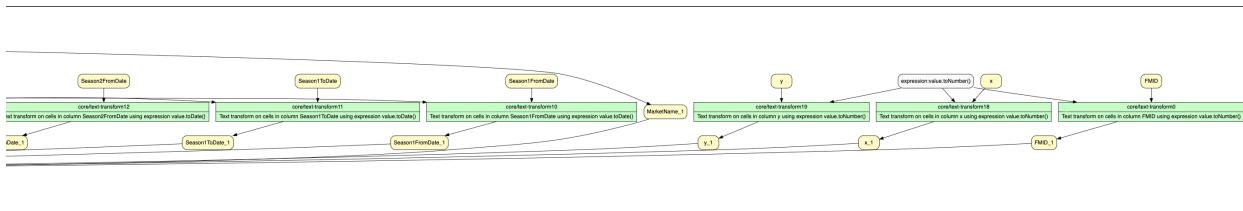
You can see the mass edit normalization of data performed with Open Refine Text Facet and clustering feature below for MarketName, City, County and other text columns.



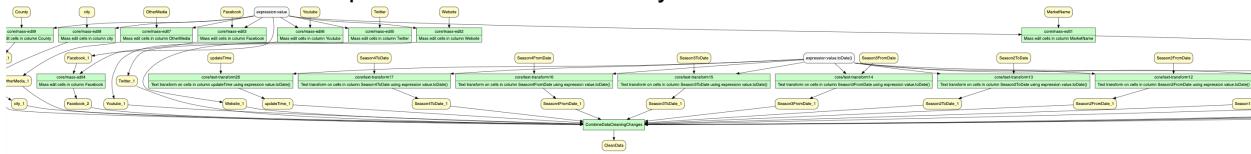
The date fields were also normalized to ISO date format with the transformation feature of Open Refine.



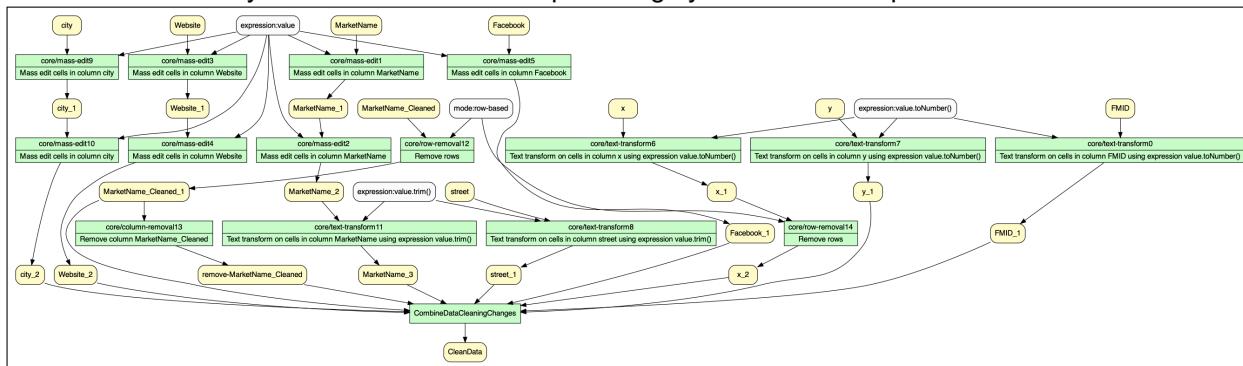
All data columns were trimmed to remove leading and lagging spaces, and numeric data columns were appropriately transformed into numeric fields using transformation features of Open Refine. Invalid data was identified using GREL expressions and cleaned up using filter features in Open Refine.



All normalized and cleaned up data columns were finally combined into a clean data set.



The final data set was further explored for opportunities to identify duplicate rows based on Market Name and other associated columns which could potentially be the same Farmers Market entity record. The refined data set was finally exported from Open Refine in the CSV format. This refined data set was fed into SQLite to identify and eliminate more complex integrity violations as explained in above sections.



CONCLUSION:

The USDA Farmers Markets Dataset is overall a well-structured dataset except that data is incomplete and inconsistent on several columns which required many parts of data to be cleaned and reformatted to improve consistency and integrity. There are few columns in the dataset which becomes completely unusable such as “updateTime”, and “Season1StartDate”, which is very inconsistent and contains either “Month” or “DD/MM/YYYY” in variable text format, and there were overlapping seasonal dates observed, which makes it unusable specifically for use-cases related to seasonal date specific search queries. Lot of anomalies have been observed with respect to invalid date ranges with Season_To_Date being less than the From_Date and those have been cleaned up. With these cleaned up season1 start and end dates, we can search for farmers market for a specific seasonality more accurately.

Duplicate entries of Market information with different FMIDs but with similar Market Name profile information, address etc. have also been identified as Integrity violations through queries done through SQLite database. These potential duplicates have been cleaned up for more accurate representations.

Overall, the dataset was cleaned appropriately for the hypothetical use cases. Zip code data has also been validated for accuracy and the discrepancies were eliminated, thus making it usable for comparison between the number of farmers markets in each city or states and to use with geographic positioning related use cases.

Contributions:

Tasks	Contributors
Raw Data Set Profiling & Analysis	Gowri Shankar Ramanan & Balaji Sathyamurthy
Clean-up using Open Refine & capturing operations history	Balaji Sathyamurthy
Integrity Checks with SQLite	Balaji Sathyamurthy
Clean-up IC violations using SQLite	Gowri Shankar Ramanan & Balaji Sathyamurthy
Create Overall Workflow using YesWorkflow	Gowri Shankar Ramanan
Generate Open Refine Workflow with history using OR2YW tool	Gowri Shankar Ramanan
Create Final Project Report	Gowri Shankar Ramanan & Balaji Sathyamurthy
Submit the artifacts of final project	Gowri Shankar Ramanan