# CS513-Project Phase 2

**By: Bhushan Bathani(bbath2@illinois.edu) and Zoheb Satta(satta2@illinois.edu)**

--------------------------------------------------------------------------------------------------------------

## Introduction

Farmers Markets dataset contains the data of all the markets , their timings, location, inventory, payment methods etc. This data is useful to locate any nearby market which provides customers what they want. e.g. Search for a market which is in 10 miles of radius and provides vegetables and seafood and is currently open.

We used OpenRefine, YesWorkflow and SQLLite tools to clean up these data and extract the meaningful data.

## Use Case:

As for the Use cases, we came up with 3 use cases representing:

**U0:** The use case that we could have implemented without any data cleaning on the data set was an app that could give you the names of the restaurants in your city

**U1:** The use case we decided to clean our dataset for was: an app that can match consumers with specialty stores by product in their city of choice
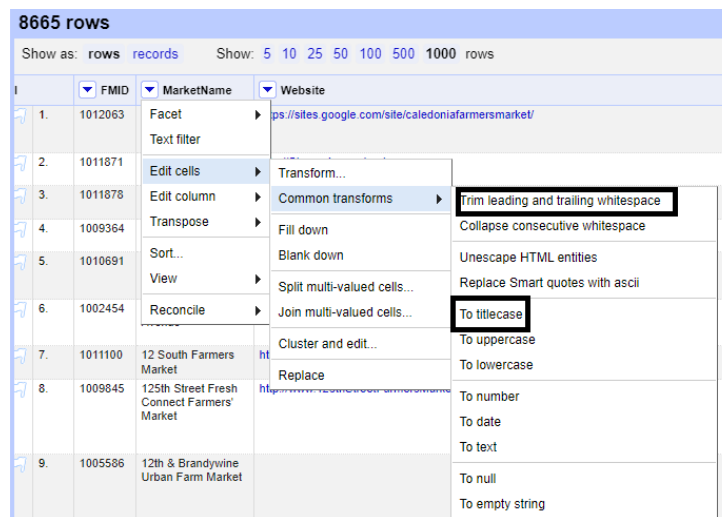
**U2:** The use case that could never be implemented is an app that sends you a link to all of a restaurant's social media when you are within 1 mile of the location. This is because all of the columns involving social media are significantly lacking in data as can be seen in these facets highlighting the null/blank values
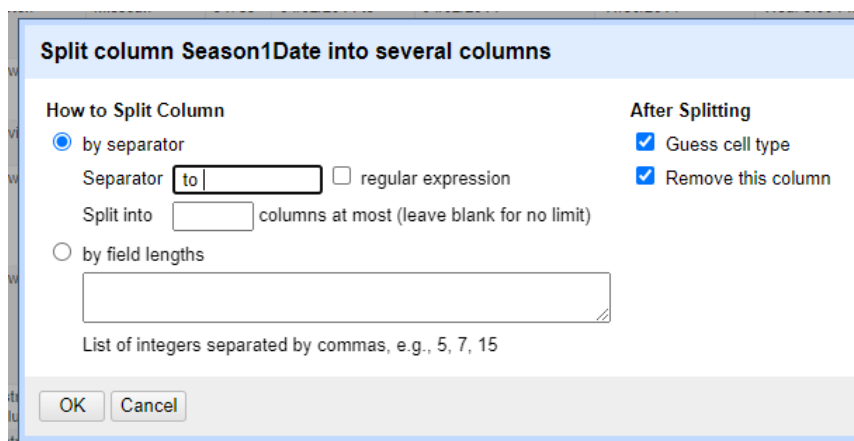


## Data Cleaning Workflow:

## OpenRefine

All the columns in the data were trimmed and the whitespace and MarketName, Website, Facebook, Twitter, Youtube, OtherMedia, Street, City, County and State columns were converted to TileCase. This was done so that these operations would not need to be done on the client side by the application.



Season1Date and Season1Time were then each split into multiple columns so that it would be easier to compare start and end dates/times later once it is dumped into SQLLite. We split Seaseon1Date, Season2Date, Season3Date and Season4Date with " to " and splitted Season1Time, Season2Time, Season3Time and Season4Time with ";" to create multiple time slots which created 8 different columns for each Time column.

X and Y columns were normalized to 4 decimal places for consistency's sake



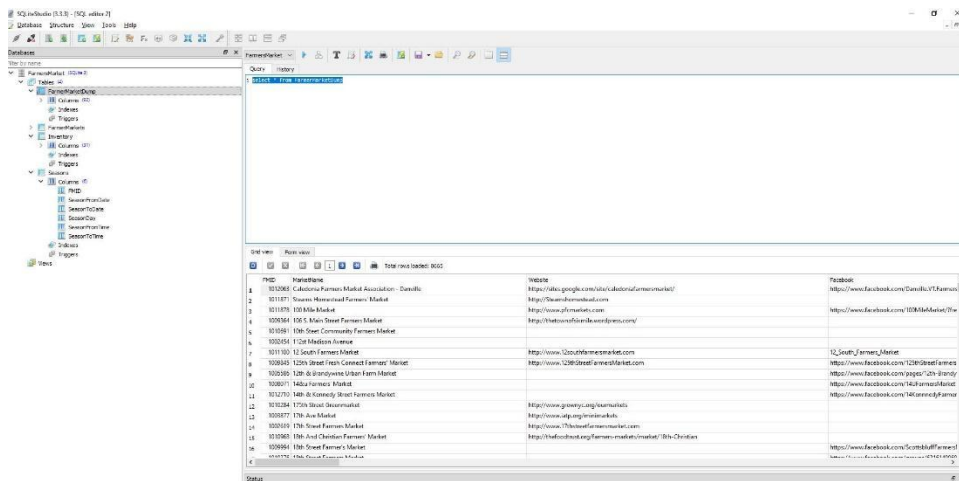All the remaining columns were converted from "Y" / "N" to true/false Boolean values, and any "-" values were converted to Null. This was so the use case could just directly compare the results to see if a product was available rather than needing to convert to Boolean in the client. Here we had to decide between choosing Null and defaulting to false, but OpenRefine defaults null values to false

## SQLLite

Once we finish the initial data cleaning in OpenRefine, we extracted the clean data from OpenRefine and load the data into SQLite (FarmerMarketDump). This table would represent the entirety of the cleaned data.



## Integrity Constraint Check

We checked Integrity Constraint violation by grouping by FMID which is unique as below. But we found multiple FMID data for the same Market under same state and city which we cleaned by by deleting duplicates as below.

**FarmersMarket**

Query | History

```
1 select count(*),FMID from FarmerMarketDump group by FMID having count(*)>1;
2
3 |
4
5
```

Grid view | Form view

Total rows loaded: 0

| count(*) | FMID |
|---|---|

---

**FarmersMarket**

Query | History

```
1 select count(*),MarketName,city,state from FarmerMarketDump group by MarketName,city,state having count(*)>1
2
3
```

Grid view | Form view

Total rows loaded: 64

| | count(*) | MarketName | city | state |
|---|---|---|---|---|
| 1 | 2 | Adams County Farmers' Market Association (pa) | Gettysburg | Pennsylvania |
| 2 | 2 | Bigfork Farmers' Market Cooperative | Bigfork | Montana |
| 3 | 2 | Bluegrass Farmers' Market, Inc. | Lexington | Kentucky |
| 4 | 2 | Brattleboro Area Farmers' Market | Brattleboro | Vermont |
| 5 | 3 | Bushwick Farmers' Market | Brooklyn | New York |
| 6 | 2 | Canton Farmers' Market | Canton | Ohio |
| 7 | 2 | Cape Ann Farmers' Market | Gloucester | Massachusetts |
| 8 | 3 | Capital City Farmers Market | Montpelier | Vermont |
| 9 | 2 | Charles Town Farmers Market | Charles Town | West Virginia |
| 10 | 2 | Charlestown Farmers Market | Lake Charles | Louisiana |
| 11 | 2 | City Of Grayling Farm Market | Grayling | Michigan |
| 12 | 2 | City Of Union Farmers' Market | Union | Missouri |
| 13 | 2 | Cleveland Farmers Market | Cleveland | Mississippi |
| 14 | 2 | Colorado Farm And Art Market | Colorado ... | Colorado |
| 15 | 3 | Columbia Farmers Market | Columbia | Missouri |
| 16 | 2 | Copley Creekside Farmers Market | Copley | Ohio |
| 17 | 3 | Country Farm & Craft Market Paso Park | Paso Robles | California |
| 18 | 4 | Crescent City Farmers Market | New Orleans | Louisiana |
| 19 | 2 | Cresco Farmers Market | Cresco | Iowa |

```
1 delete from FarmerMarketDump where rowid in (select max(rowid) from FarmerMarketDump group by MarketName,city,state having count(*)>1 )
2
3
```

After clean up, we can see that it does not consist of any duplicate rows.

FarmersMarket ∨

Query History

```
1
2 select count(*),MarketName,city,state from FarmerMarketDump group by MarketName,city,state having count(*)>1
3
4
```

Grid view   Form view

Total rows loaded: 0

| count(*) | MarketName | city | state |
| --- | --- | --- | --- |

**Split Data in Respected Tables and Data Quality Changes**

Then we created 3 separate tables to organize the data. **FarmerMarkets** holds the info about the Market itself, including the Name, and Social Networking values, as well as its geographic location. **Seasons** stores the open/close dates and times for seasonal markets. **Inventory** holds Boolean data for what products a market does or does not provide

```
FarmersMarket  ▼    ▶  🔒  T  🖉  ✂  🖨   🖺  🖫 ▼ 📂

Query   History

1  CREATE TABLE FarmerMarkets (
2      FMID        BIGINT   PRIMARY KEY,
3      MarketName  VARCHAR,
4      Website     VARCHAR,
5      Facebook    VARCHAR,
6      Twitter     VARCHAR,
7      Youtube     VARCHAR,
8      OtherMedia  VARCHAR,
9      street      VARCHAR,
10     city        VARCHAR,
11     County      VARCHAR,
12     State       VARCHAR,
13     zip         NUMERIC,
14     x           DECIMAL,
15     y           DECIMAL,
16     Location    VARCHAR,
17     Credit      BOOLEAN,
18     WIC         BOOLEAN,
19     WICcash     BOOLEAN,
20     SFMNP       BOOLEAN,
21     SNAP        BOOLEAN,
22     updateTime  DATETIME
23 );
24
25
26
27 CREATE TABLE Seasons (
28     FMID            BIGINT  REFERENCES FarmerMarkets (FMID),
29     SeasonFromDate  DATE,
30     SeasonToDate    DATE,
31     SeasonDay       VARCHAR,
32     SeasonFromTime  VARCHAR,
33     SeasonToTime    VARCHAR
34 );
35
36
37 CREATE TABLE Inventory (
38     FMID            BIGINT  REFERENCES FarmerMarkets (FMID),
39     Organic     BOOLEAN,
40     Bakedgoods  BOOLEAN,
41     Cheese      BOOLEAN,
42     Crafts      BOOLEAN,
43     Flowers     BOOLEAN,
44     Eggs        BOOLEAN,
45     Seafood     BOOLEAN,
46     Herbs       BOOLEAN,
47     Vegetables  BOOLEAN,
48     Honey       BOOLEAN,
49     Jams        BOOLEAN,
50     Maple       BOOLEAN,
51     Meat        BOOLEAN,
52     Nursery     BOOLEAN,
53     Nuts        BOOLEAN,
54     Plants      BOOLEAN,
55     Poultry     BOOLEAN,
56     Prepared    BOOLEAN,
57     Soap        BOOLEAN,
58     Trees       BOOLEAN,
59     Wine        BOOLEAN,
60     Coffee      BOOLEAN
```

The Data was split and distributed into their respective tables to organize the data so our app would be able to search along all 3 axes (By Store, By Time, By Products) with ease.

```
1  INSERT INTO FarmerMarkets (
2                          FMID,
3                          MarketName,
4                          Website,
5                          Facebook,
6                          Twitter,
7                          Youtube,
8                          OtherMedia,
9                          street,
10                         city,
11                         County,
12                         State,
13                         zip,
14                         x,
15                         y,
16                         Location,
17                         Credit,
18                         WIC,
19                         WICcash,
20                         SFMNP,
21                         SNAP,
22                         updateTime
23                    )
24             select FMID,
25                    MarketName,
26                    Website,
27                    Facebook,
28                    Twitter,
29                    Youtube,
30                    OtherMedia,
31                    street,
32                    city,
33                    County,
34                    State,
35                    zip,
36                    x,
37                    y,
38                    Location,
39                    Credit,
40                    WIC,
41                    WICcash,
42                    SFMNP,
43                    SNAP,
44                    updateTime from FarmerMarketDump
```

```sql
insert into Inventory
(
FMID ,
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)  select FMID ,
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 from FarmerMarketDump
```

```sql

update FarmerMarketDump set [Season1Time 1]= replace([Season1Time 1],' ','');
update FarmerMarketDump set [Season1Time 2]= replace([Season1Time 2],' ','');
update FarmerMarketDump set [Season1Time 3]= replace([Season1Time 3],' ','');
update FarmerMarketDump set [Season1Time 4]= replace([Season1Time 4],' ','');
update FarmerMarketDump set [Season1Time 5]= replace([Season1Time 5],' ','');
update FarmerMarketDump set [Season1Time 6]= replace([Season1Time 6],' ','');
update FarmerMarketDump set [Season1Time 7]= replace([Season1Time 7],' ','');
update FarmerMarketDump set [Season1Time 8]= replace([Season1Time 8],' ','');

update FarmerMarketDump set [Season2Time 1]= replace([Season2Time 1],' ','');
update FarmerMarketDump set [Season2Time 2]= replace([Season2Time 2],' ','');
update FarmerMarketDump set [Season2Time 3]= replace([Season2Time 3],' ','');
update FarmerMarketDump set [Season2Time 4]= replace([Season2Time 4],' ','');
update FarmerMarketDump set [Season2Time 5]= replace([Season2Time 5],' ','');
update FarmerMarketDump set [Season2Time 6]= replace([Season2Time 6],' ','');
update FarmerMarketDump set [Season2Time 7]= replace([Season2Time 7],' ','');
update FarmerMarketDump set [Season2Time 8]= replace([Season2Time 8],' ','');

update FarmerMarketDump set [Season3Time 1]= replace([Season3Time 1],' ','');
update FarmerMarketDump set [Season3Time 2]= replace([Season3Time 2],' ','');
update FarmerMarketDump set [Season3Time 3]= replace([Season3Time 3],' ','');
update FarmerMarketDump set [Season3Time 4]= replace([Season3Time 4],' ','');
update FarmerMarketDump set [Season3Time 5]= replace([Season3Time 5],' ','');
update FarmerMarketDump set [Season3Time 6]= replace([Season3Time 6],' ','');
update FarmerMarketDump set [Season3Time 7]= replace([Season3Time 7],' ','');
update FarmerMarketDump set [Season3Time 8]= replace([Season3Time 8],' ','');

update FarmerMarketDump set [Season4Time 1]= replace([Season4Time 1],' ','');
update FarmerMarketDump set [Season4Time 2]= replace([Season4Time 2],' ','');
update FarmerMarketDump set [Season4Time 3]= replace([Season4Time 3],' ','');
update FarmerMarketDump set [Season4Time 4]= replace([Season4Time 4],' ','');
update FarmerMarketDump set [Season4Time 5]= replace([Season4Time 5],' ','');
update FarmerMarketDump set [Season4Time 6]= replace([Season4Time 6],' ','');
update FarmerMarketDump set [Season4Time 7]= replace([Season4Time 7],' ','');
update FarmerMarketDump set [Season4Time 8]= replace([Season4Time 8],' ','');

--split time and insert into Seasons table

insert into Seasons
select FMID,Season1FromDate,Season1ToDate,substr([Season1Time 1],0,instr([Season1Time 1],':')),
substr([Season1Time 1],instr([Season1Time 1],':')+1,instr([Season1Time 1],'-')-5),
substr([Season1Time 1],instr([Season1Time 1],'-')+1,length([Season1Time 1]))
from FarmerMarketDump;




insert into Seasons
select FMID,Season1FromDate,Season1ToDate,substr([Season1Time 2],0,instr([Season1Time 2],':')),
substr([Season1Time 2],instr([Season1Time 2],':')+1,instr([Season1Time 2],'-')-5),
substr([Season1Time 2],instr([Season1Time 2],'-')+1,length([Season1Time 2]))
from FarmerMarketDump;
```
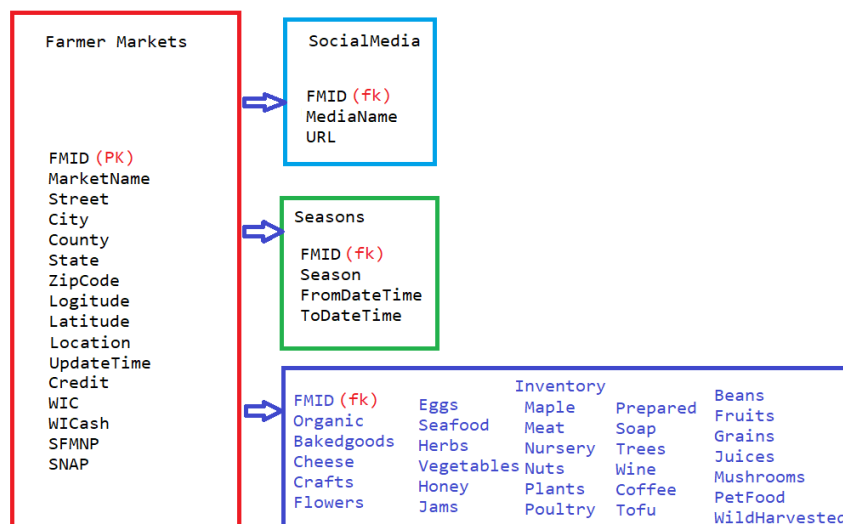
Once the data was separated, we generated the data for the use case. For example: "What are the nearby stores in New York City that accept credit cards and are open today that sells Organics and Vegetables"



There were some steps from our initial Phase 1 that did not go as planned. For example: unlike our ER diagram in phase one, we decided against creating a SocialMedia table as there wasn't enough available data for it to be valuable to our Use Case.
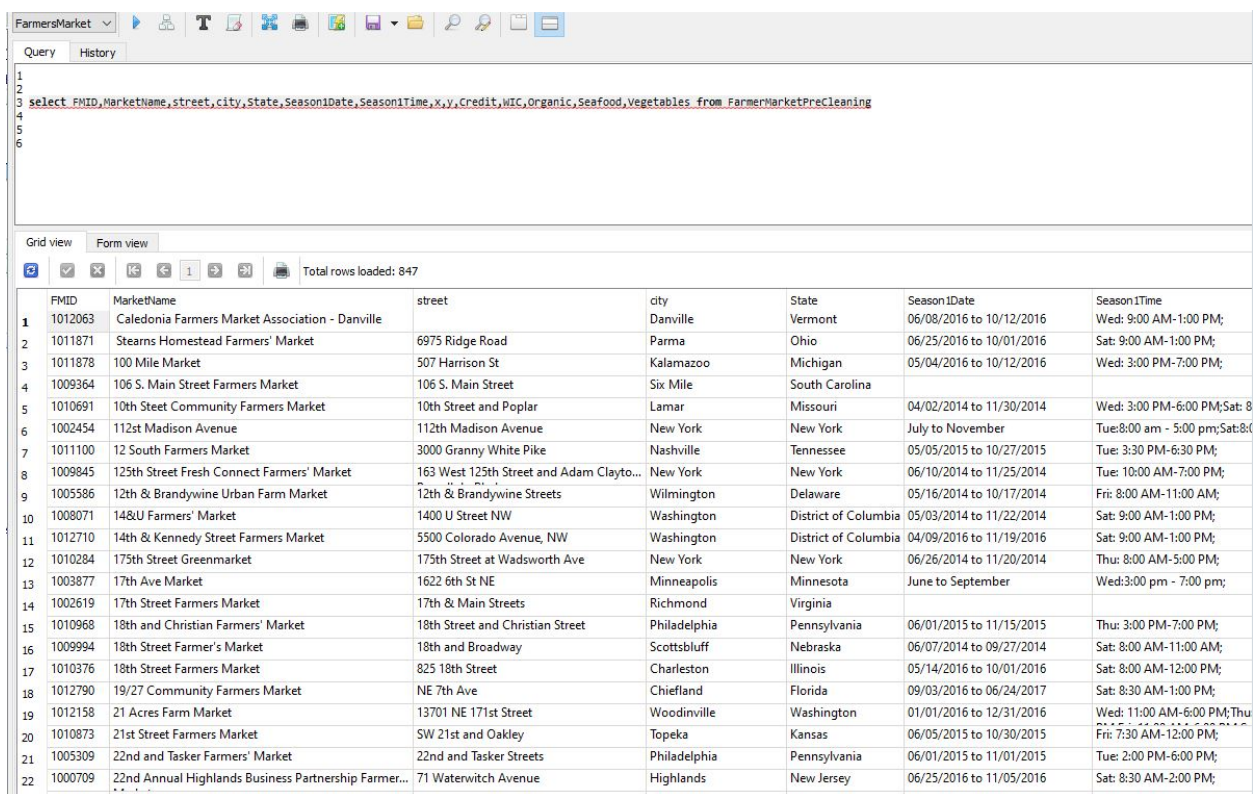


Similarly, we considered using data log however, so we used OpenRefine and SQLLite for the majority of our data cleaning because it was the most efficient tool out of the ones we've been

taught for the level of data cleaning that was required for our chosen use case. Similarly, SQLite provided a simple way to organize and retrieve the data for the use case.

## Cleaned Data:

Here you can see just some of the rows and columns that were improved by our data cleaning

**Before:**



**After:**

FarmersMarket ⌄ | Query | History

```
1
2
3 select FMID,MarketName,street,city,State,Season1FromDate,
4    Season1ToDate,
5    [Season1Time 1],
6    [Season1Time 2],
7    [Season1Time 3],
8    [Season1Time 4],
9    [Season1Time 5],
10   [Season1Time 6],
```

Grid view | Form view

Total rows loaded: 8665

| | FMID | MarketName | street | city | State | Season1FromD | Season1ToDate | Season1Time 1 | Season1Time 2 | Season1Time 3 |
|---|---|---|---|---|---|---|---|---|---|---|
| 16 | 1009994 | 18th Street Farmer's Market | 18th And Broadway | Scottsbluff | Nebraska | 06/07/2014 | 09/27/2014 | Sat:8:00AM-11:00AM | | |
| 17 | 1010376 | 18th Street Farmers Market | 825 18th Street | Charleston | Illinois | 05/14/2016 | 10/01/2016 | Sat:8:00AM-12:00PM | | |
| 18 | 1012790 | 19/27 Community Farmers Market | Ne 7th Ave | Chiefland | Florida | 09/03/2016 | 06/24/2017 | Sat:8:30AM-1:00PM | | |
| 19 | 1012158 | 21 Acres Farm Market | 13701 Ne 171st Street | Woodinville | Washington | 01/01/2016 | 12/31/2016 | Wed:11:00AM-6:00PM | Thu:11:00AM-6:00PM | Fri:11:00AM-6:00PM |
| 20 | 1010873 | 21st Street Farmers Market | Sw 21st And Oakley | Topeka | Kansas | 06/05/2015 | 10/30/2015 | Fri:7:30AM-12:00PM | | |
| 21 | 1005309 | 22nd And Tasker Farmers' Market | 22nd And Tasker Streets | Philadelphia | Pennsylvania | 06/01/2015 | 11/01/2015 | Tue:2:00PM-6:00PM | | |
| 22 | 1000709 | 22nd Annual Highlands Business ... | 71 Waterwitch Avenue | Highlands | New Jersey | 06/25/2016 | 11/05/2016 | Sat:8:30AM-2:00PM | | |
| 23 | 1011881 | 25th Street Market - North Logan... | 475 East 2500 North | North Logan | Utah | 05/07/2016 | 10/15/2016 | Sun:9:00AM-1:00PM | | |
| 24 | 1010966 | 26th And Allegheny | 26th Street And W ... | Philadelphia | Pennsylvania | 06/01/2015 | 11/15/2015 | Wed:1:00PM-5:00PM | | |
| 25 | 1005299 | 29th And Wharton Farmers' ... | 29th And Wharton Streets | Philadelphia | Pennsylvania | 06/01/2015 | 11/15/2015 | Tue:2:00PM-6:00PM | | |
| 26 | 1010994 | 2nd Street Farmers' Market | 194 Second Street | Amherst | Virginia | 05/05/2016 | 09/01/2016 | Thu:3:30PM-6:30PM | | |
| 27 | 1009959 | 2nd Street Market - Five Rivers ... | 600 E. 2nd Street | Dayton | Ohio | 01/01/2016 | 12/31/2016 | Thu:11:00AM-3:00PM | Fri:11:00AM-3:00PM | Sat:8:00AM-3:00PM |
| 28 | 1004950 | 3 French Hens French Country ... | 123 W. Illinois Ave. | Morris | Illinois | 05/10/2014 | 10/11/2014 | Sat:8:00AM-2:00PM | | |
| 29 | 1010775 | 30a Farmers' Market | Rosmary Beach Town ... | Rosemary Beach | Florida | 01/18/2015 | 01/05/2020 | Sun:9:00AM-1:00PM | | |
| 30 | 1012342 | 31 & Main Farmers Market At ... | 1928 Pennington Road | Ewing | New Jersey | 06/12/2016 | 10/30/2016 | Sun:10:00AM-2:00PM | | |
| 31 | 1005636 | 32nd Street/waverly Farmers ... | E. 32nd & Barclay Street | Baltimore | Maryland | 01/01/2013 | 12/31/2013 | Sat:7:00AM-12:00PM | | |
| 32 | 1005310 | 33rd And Diamond Farmers' ... | N 33rd And Diamond ... | Philadelphia | Pennsylvania | 06/01/2015 | 11/01/2015 | Thu:2:00PM-6:00PM | | |
| 33 | 1012784 | 38th & Meridian Farmers Market | 3808 N Meridian St | Indianapolis | Indiana | 06/02/2016 | 09/29/2016 | Thu:4:00PM-6:30PM | | |

## Supplementary Materials:

- **Workflow:**
  1. WorkFlowLinear.pdf
  2. WrokFlowParallel.pdf (please download and zoom)
  3. WorkFlowLinear.gv
  4. WorkFlowLinear.yw
  5. yw.properties

Commands used to generate above files are below.

```
or2yw -i OpenRefineOperations.json  -ot pdf -o WorkFlowLinear.pdf
or2yw -i OpenRefineOperations.json  -ot gv -o WorkFlowLinear.gv
or2yw -i OpenRefineOperations.json  -ot yw -o WorkFlowLinear.yw
or2yw -i OpenRefineOperations.json  -ot pdf -o WorkFlowParallel.pdf -t parallel
```

- **OpenRefine:**
  OpernRefineOperations.json

- **Queries:**
  queries.txt
- **Datasets:**
  1. Clean:
  2. Dirty:

## Conclusions:

One of the problems that we encountered was splitting the dates, using OpenRefine, it created a lot of columns because some rows had multiple date/times in one value so we decided to split the date into 2 columns (to and from) and listed the times split by ";" Then in SQLite we split that even further.  Some of the date data was also difficult to convert as they were just months rather than day/month/year. Next step we can make this application more intelligent by taking the user's current timezone and use it to check current open markets based on the geolocation of the user from where he is sending the request.

**Repository :** https://github.com/bkbathani/CS513_FinalProject

**References :**

https://pypi.org/project/or2ywtool/

## Team Member Contribution

Initially both Bhushan and Zoheb researched through all datasets and chose farmers market dataset to use for the final project considering team size and timeline. Zoheb worked through the data to identify the use case and tried to study data, meanwhile Bhushan worked on yesworkflow, openrefine and SQLLite setup and research. Both worked together later to go through the files and apply changes and scripted the data cleaning part on zoom call.