# CS 513 Theory & Practice of Data Cleaning

# Final Project: Data Cleaning Project

Author: Varun Sharma

Email :Varun4@illinois.edu

# Table of Contents

# 1. Introduction

In this report farmers market data is analyzed and cleansed to solve some common use cases related to farmer market information.

This report is done as part of Individual submission of data cleaning project done by varun sharma at UIUC.

The tools used in the project are OpenRefine for clustering and refining data, SQLite for relational database, Trifacta Data Wrangler and YesWorkflow for provenance.

# 2. Overview and Initial Assessment

## i. Dataset

I am using data from the United States Department of Agriculture. Maintained by the Agricultural Marketing Service, the Directory is designed to provide customers with convenient access to information about farmers market listings to include: market locations, directions, operating times, product offerings, accepted forms of payment, and more.

URL for the data set:

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

## ii. Proposed Use Cases

Proposed use case for this dataset is to provide easy access to consumers to clean Farmer market data where they can easily browse trusted data to figure out details of food markets in the area. Consumers can have more details about the farmer markets.

This data can be fed to other web application or IA engines/automated assistant like Alexa or google voice to make them smarter. This use case is not part of the project but worth mentioning because questions like which markets have organic food zip? 2. Which market are open during night time of 8-11pm? can be answered easily using the clean data set used in this project.

## iii. Data Cleaning Goals

The main goal for this project is to clean data based on data quality dimensions using farmer's market dataset.

I. **Completeness**: There are lot of columns with missing information. For example youtube column has 8558 null values which lowers the completeness score for the data.

II. **Consistency:** Data set has consistency issues like values NA,N.A,_ are used for missing values. Also, city values like ST. Louis and ST Louis exist in the dataset.

III.   **Uniqueness**: Data set scores high on unique as there are no duplicates in the data set.

Though there are multiple records for same store but their location are different in country.

IV.   **Validity:** There are some validity problems in data set like Facebook column has URL and

some places it doesn't have typical URL format.

V.   **Accuracy:** Accuracy % drops because of inconsistent formats used for the fields for example:

VI.   Update time field =201 has only year mentioned in it in one record and most of records have

VII.   **Timeliness:** Data is downloaded from the website directly and cleansed which has some

delay but it's up to date as its maintained by the United States Department of Agriculture.

Following whitepaper was consulted for this:
https://www.whitepapers.em360tech.com/wp-

content/files_mf/1407250286DAMAUKDQDimensionsWhitePaperR37.pdf

## *iv.      High Level Data Quality Issues:*

As discussed above data set has multiple issues. Following are the major high-level issues in the data set:
   I.      Lots of leading and trailing whitespaces in column values.
   II.      Markets with missing information like youtube
   III.      Standardization for date/time formats.
   IV.      Standard values for city names in the data set.

## *v.      Data Cleaning Tools and Approach*

OpenRefine is used in this project for common data cleaning steps like removing white spaces, line breaks and doing basic transformations. Clustering capability helped to achieve data consistency goal for data quality. Clustering in openrefine was used for columns like county, cities:

york" are very likely to refer to the same concept and just have capitalization differences, and "Godel" and "Godel" prob
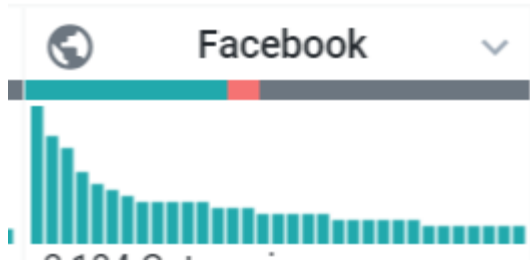
Method [ key collision ▼ ]          Keying Function [ fingerprint ▼ ]

| Cluster Size | Row Count | Values in Cluster | Merge? | New Cell Value |
|---|---|---|---|---|
| 2 | 10 | • ST. LOUIS (8 rows)<br>• ST LOUIS (2 rows) | ☐ | ST. LOUIS |
| 2 | 2 | • ST AUGUSTINE (1 rows)<br>• ST. AUGUSTINE (1 rows) | ☐ | ST AUGUSTINE |
| 2 | 2 | • LAND O LAKES (1 rows)<br>• LAND O' LAKES (1 rows) | ☐ | LAND O LAKES |
| 2 | 2 | • WHEELING (1 rows)<br>• WHEELING, (1 rows) | ☐ | WHEELING |

## Cluster & Edit column "city_clean"

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

Method [ key collision ▼ ]     Keying Function [ ngram-fingerprint ▼ ]     Ngram Size [ 2 ]

| Cluster Size | Row Count | Values in Cluster | Merge? | New Cell Value |
|---|---|---|---|---|
| 2 | 3 | • - (2 rows)<br>• O (1 rows) | ☐ | - |
| 2 | 2 | • LE ROY (1 rows)<br>• LEROY (1 rows) | ☐ | LE ROY |
| 2 | 5 | • NORTHPORT (4 rows)<br>• NORTH PORT (1 rows) | ☐ | NORTHPORT |
| 2 | 4 | • LA CROSSE (3 rows) | ☐ | LA CROSSE |

After OpenRefine step Trifecta Data Wrangler is used to check data quality. I felt trifacta data wrangler's visual representation is better than open refine. Wrangler gives capability to profile data column wise and understand data better.

Each column data pattern can be seen in the UI
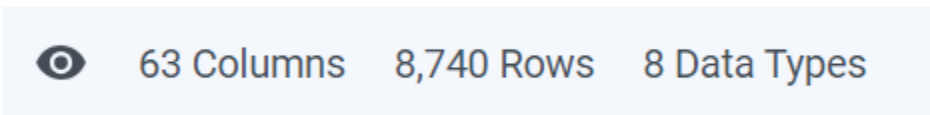
Red here marks invalid values or mismatched values.

I have used YesWorkflow to create a graphical representation of the data cleaning process of the farmers market data. With http://try.yesworkflow.org/

I have developed one workflow for steps in OpenRefine.

## 2.5 Data file and cleaning summary

Farmers market data, and this a summary of some of the columns and cleaning method adopted

Summary of the file:



| Column Name | Description | Cleaning Operations |
|---|---|---|
| FMID | ID farmer Market | Data is clean no operations on this column |
| MarketName | Name of the Market | Open Refine:<br><br>1.Trim white spaces and collapse consecutive whitespaces<br><br>2.Cluster and edit using key collision and finger print for data standardization. |
| Website | Website URL | Open Refine: |

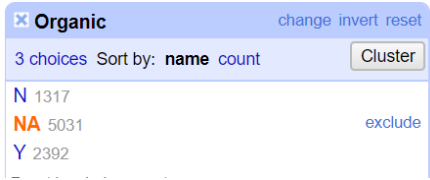| | | |
|---|---|---|
| | | 1.Trim white spaces and collapse consecutive whitespaces<br><br>2.Cluster and edit using key collision and finger print for data standardization.<br><br>Example:<br><br>http://www.seela.org<br><br>http://www.seela.org/<br><br>Standardized value:<br><br>http://www.seela.org<br><br><span style="color:red">Trifacta:</span><br><br>Check mismatch of column value formats |
| Facebook | Facebook link | Open Refine:<br><br>1.Trim white spaces and collapse consecutive whitespaces<br><br>2.Cluster and edit using key collision and finger print for data standardization.<br><br><span style="color:red">Trifacta:</span><br><br>Check mismatch of column value formats |
| Twitter | Twitter Username | Open Refine:<br><br>1.Trim white spaces and collapse consecutive whitespaces<br><br>2.Cluster and edit using key collision and finger print for data standardization<br><br>Not available information is converted to space or null in open refine.<br><br>n/a(7 rows)<br>N/A(2 rows) |

| | | |
|---|---|---|
| | | Converted to space |
| YouTube | Youtube Channel | Open Refine: 1.Trim white spaces and collapse consecutive whitespaces 2.Cluster and edit using key collision and finger print for data standardization This info is missing in the file so standardized to NA |
| OtherMedia | Other social media accounts | Open Refine: 1.Trim white spaces and collapse consecutive whitespaces 2.Cluster and edit using key collision and ngram fingerprint data standardization Ngram size 2 |
| Street | Street address for store | Open Refine: 1.Trim white spaces and collapse consecutive whitespaces 2.Cluster and edit using key collision and finger print for data standardization |
| City | City location for store | Open Refine: 1.Trim white spaces and collapse consecutive whitespaces 2.Upper case for city name 3.Cluster and edit using key collision and finger print for data standardization |

|  |  | Created city clean column |
| --- | --- | --- |
| County | Name of county | Open Refine:<br><br>1.Trim white spaces and collapse consecutive whitespaces<br><br>2.Cluster and edit using key collision and finger print for data standardization |
| State | State location for farmers market | Open Refine:<br><br>1.Trim white spaces and collapse consecutive whitespaces<br><br>2.Cluster and edit using key collision and finger print for data standardization |
| Zip | Zip code for farmers market | Open Refine:<br><br>1.Trim white spaces and collapse consecutive whitespaces<br><br>No clustering is required for this column |
| Season1Date | First season date for maket | Open Refine:<br><br>1.Trim white spaces and collapse consecutive whitespaces<br><br>2.Cluster and edit using key collision and finger print for data standardization |
| Season1Time | First season time for market | Trim white spaces and replace new lines in OpenRefine |

| | | |
|---|---|---|
| Season2Date | Second season date for market | Open Refine:<br><br>1.Trim white spaces and collapse consecutive whitespaces<br><br>2.Cluster and edit using key collision and finger print for data standardization |
| Season2Time | Second season time for market | Open Refine:<br><br>1.Trim white spaces and collapse consecutive whitespaces<br><br>2.Cluster and edit using key collision and finger print for data standardization |
| Season3Date | Third season date for market | Open Refine:<br><br>1.Trim white spaces and collapse consecutive whitespaces<br><br>2.Cluster and edit using key collision and finger print for data standardization |
| Season3Time | Third season time for market | Open Refine:<br><br>1.Trim white spaces and collapse consecutive whitespaces<br><br>2.Cluster and edit using key collision and finger print for data standardization |
| Season4Date | Fourth season date for market | Open Refine:<br><br>1.Trim white spaces and collapse consecutive whitespaces<br><br>2.Cluster and edit using key collision and finger print for data standardization |

| | | |
|---|---|---|
| | | |
| Season4Time | Fourth season time for market | Open Refine:<br><br>1.Trim white spaces and collapse consecutive whitespaces<br><br>2.Cluster and edit using key collision and finger print for data standardization |
| Location | Location of the market | Data was clean for this |
| Organic | Flag to indicate if store sells organic products or not | Open Refine:<br><br>1.Trim white spaces and collapse consecutive whitespaces<br><br>Changed facet – to NA to identify if FLAG is NA i.e. Not Available.<br><br> |
| updateTime | Time the info was updated | Open Refine:<br><br>1.Trim white spaces and collapse consecutive whitespaces |

## 3.  Data Cleaning with Open Refine and Trifacta tool set

Open Refine tool was used to do data cleaning. Common transforms were used to remove white spaces and to collapse consecutive white spaces. Clustering was done to do data standardization for fields.

Trifecta data wrangler was used for data wrangling and cleanup individual fields.

Sample screenshot for the openrefine operations:



Here is the sample extract for open refine operations done for data quality:

**Extract Operation History**

Extract and save parts of your operation history as JSON that you can apply to this or other projects in the future.

☑ Create column city_clean at index 9 based on column ci ty using expression grel:value

☑ Text transform on cells in column city_clean using expre ssion value.toUppercase()

☑ Text transform on cells in column city_clean using expre ssion value.trim()

☑ Text transform on cells in column city_clean using expre ssion value.replace(/\s+/,' ')

☑ Mass edit cells in column city_clean

☑ Mass edit cells in column city_clean

☑ Create column clean_county at index 11 based on colu mn County using expression grel:value

☑ Rename column clean_county to county_clean

☑ Mass edit cells in column county_clean

☑ Text transform on cells in column FMID using expressio n value.trim()

☑ Text transform on cells in column FMID using expressio n value.replace(/\s+/,' ')

☑ Text transform on cells in column MarketName using ex pression value.trim()

[ Select All ] [ Unselect All ]

```json
[
  {
    "op": "core/column-addition",
    "description": "Create column city_clean at i
    "engineConfig": {
      "mode": "row-based",
      "facets": []
    },
    "newColumnName": "city_clean",
    "columnInsertIndex": 9,
    "baseColumnName": "city",
    "expression": "grel:value",
    "onError": "set-to-blank"
  },
  {
    "op": "core/text-transform",
    "description": "Text transform on cells in co
    "engineConfig": {
      "mode": "row-based",
      "facets": []
    },
    "columnName": "city_clean",
    "expression": "value.toUppercase()",
    "onError": "keep-original",
    "repeat": false,
    "repeatCount": 10
  },
  {
```

More details are provided as file in the final deliverables.

Trifacta data wrangler steps:



Some Facebook column value don't have URL. They are represented in different column by flag called Not_facebookURL to identify such kind of columns.

Data Profiling in Trifacta Wrangler tool:



## 4. Relation Database Schema

### I.    Schema

Sqlite3 was used to load the data file exported from the trifacta data wrangler for next step in the project. Screenshot of the schema:

```
sqlite> .schema
CREATE TABLE Farmers_Market(
  "FMID" Integer,
  "MarketName" TEXT,
  "Website" TEXT,
  "ismismatched_Website" TEXT,
  "Facebook" TEXT,
  "Not_FacebookURL" TEXT,
  "Twitter" TEXT,
  "Youtube" TEXT,
  "OtherMedia" TEXT,
  "street" TEXT,
  "city" TEXT,
  "city_clean" TEXT,
  "County" TEXT,
  "county_clean" TEXT,
  "State" TEXT,
  "zip" TEXT,
  "Season1Date" TEXT,
  "Season1Time" TEXT,
  "Season2Date" TEXT,
  "Season2Time" TEXT,
  "Season3Date" TEXT,
  "Season3Time" TEXT,
  "Season4Date" TEXT,
  "Season4Time" TEXT,
  "x" TEXT,
  "y" TEXT,
  "Location" TEXT,
  "Credit" TEXT,
  "WIC" TEXT,
 "WICcash" TEXT,
  "SFMNP" TEXT,
  "SNAP" TEXT,
  "Organic" TEXT,
  "Bakedgoods" TEXT,
  "Cheese" TEXT,
  "Crafts" TEXT,
  "Flowers" TEXT,
  "Eggs" TEXT,
  "Seafood" TEXT,
```

```
  "Herbs" TEXT,
  "Vegetables" TEXT,
  "Honey" TEXT,
  "Jams" TEXT,
  "Maple" TEXT,
  "Meat" TEXT,
  "Nursery" TEXT,
  "Nuts" TEXT,
  "Plants" TEXT,
  "Poultry" TEXT,
  "Prepared" TEXT,
  "Soap" TEXT,
  "Trees" TEXT,
  "Wine" TEXT,
  "Coffee" TEXT,
  "Beans" TEXT,
  "Fruits" TEXT,
  "Grains" TEXT,
  "Juices" TEXT,
  "Mushrooms" TEXT,
  "PetFood" TEXT,
  "Tofu" TEXT,
  "WildHarvested" TEXT,
  "updateTime" DATE
);
```

## II.    Integrity Constraints

Following integrity constraints where identified and described here:

### Data Integrity Constraints:
Basic data sanity check:

✓ Check number of rows loaded into the sqlite3. Please note If the table already exists, the sqlite3 tool uses all the rows, including the first row, in the CSV file as the actual data to import. Therefore while matching data rows I have deleted the first row from the file as its metadata and not the actual data.

**Following IC are written in denial form:**
✓ Every market has a website URL.
✓ URL format is consistent for farmer's market
✓ Every record has city and county information.

## 5. Workflow Model

I have used YesWorkflow to create a provenance representation of the data cleaning process of the farmers market data. I have used http://try.yesworkflow.org/ to create the workflow with the web-based version which allows us to see our changes as they are made.

I have developed one workflow focusing on the two main steps used to clean the data, OpenRefine and Trifacta. Within those two main steps, there are numerous sub-steps to show the exact steps taken to clean the steps.

Here is the flow diagram for open refine:



Trifacta:

## CleanWithTrifacta



Recipe example:

'''

@begin CleanWithOpenRefine @desc OpenRefine Workflow for farmersmarkets dataset

@in farmersmarkets.csv @uri file://data/farmersmarkets.csv

@begin LoadToOpenRefine @desc Upload farmersmarkets.csv data to OpenRefine

@in farmersmarkets.csv @uri file://data/farmersmarkets.cv

@out spreadsheet

@end CleanWithOpenRefine


@begin ColumnsToClean @desc Columns that are cleaned

@in spreadsheet

@out MarketName

@out Website

@out Facebook

@out Twitter

@out Youtube

@out OtherMedia

@out street

@out city

@out County

@out State

@out zip

@out updatedTime

@end ColumnsToClean


@begin TrimSpaces @desc Trim leading and trailing white spaces, Collapse white spaces

@in MarketName

@in Website

@in Facebook

@in Twitter

@in Youtube

@in OtherMedia

@in street

@in city

@in county

@in State

@in zip

@in UpdatedTime

@out MarketName_trim

@out Website_trim

@out Facebook_trim

@out Twitter_trim

@out Youtube_trim

@out OtherMedia_trim

@out city_trim

@out county_trim

@out State_trim

@out zip_trim

@in Organic

@out UpdatedTime_trim

@end TrimSpaces

@begin FormatTime @desc Change 'Month DD YYYY' to 'MM/DD/YYYY'

@in updatedTime

@end FormatTime


@begin ClusterValues @desc Cluster Similar values

@in MarketName_trim

@in city_trim

@in Website_trim

@in county_trim

@in Facebook_trim

@in city_trim

@in State_trim

@in Youtube_trim

@in Twitter_trim

@in OtherMedia_trim

@end ClusterValues


@begin UpperCase @desc Uppercase Values

@in city_trim

@end DUpperCase


@begin ReplaceNull @desc Replace Null,- by NA Not Available

@in Organic

@end ReplaceNull

@out farmersmarkets-or.csv @uri file://data/farmersmarkets-or.csv

@end CleanWithOpenRefine

'''

'''

@begin CleanWithTrifacta @desc Trifacta Workflow for farmersmarkets dataset

@in farmersmarkets.csv @uri file://data/farmersmarkets.csv


  @begin LoadToTrifacta @desc Upload farmersmarkets_Openrefine.csv data to Trifacta

  @in farmersmarkets_openrefine.csv @uri file://data/farmersmarketsopenrefine.cv

  @out spreadsheet

  @end CleanWithTrifacta


  @begin ColumnsToClean @desc Columns that are cleaned

  @in spreadsheet

  @out Facebook

  @out Website

  @out updateTime

  @end ColumnsToClean


  @begin ismismatched @desc new mismatchformat field

  @in Facebook

  @in Website

  @end ismismatched

@begin FormatTimeDefault @desc Change 'YYYY' format values to '1-1-YYYY 1:00:00AM'

  @in updateTime

  @out updateTime_1

  @end FormatTimeDefault


@begin FormatTime @desc Change 'Month DD YYYY' to mm-dd-yy hh:mm:ss'

  @in updateTime_1


  @end FormatTime


@out farmersmarkets-or.csv @uri file://data/farmersmarkets-or.csv

@end CleanWithOpenRefine

'''

# Appendix A – Trifacta Recipe

```
New Step          Recipe          ×

                                  •••

1  Create ismismatched_Facebook from
   ISMISMATCHED(Facebook, ['Url'])

2  Create ismismatched_Website from
   ISMISMATCHED(Website, ['Url'])

3  Replace matches of `{start}{digit}{4}{end}` from
   updateTime with '1\/1\/2013 1:00:00 AM'

4  Set updateTime to IFMISMATCHED($col,
   ['Datetime','mm-dd-yy
   hh:mm:ss','mm*dd*yyyy*hh:MM:SSa'], NULL())

5  Rename ismismatched_Facebook to
   'Not_FacebookURL'
```

## Appendix B – SQLite Logical Integrity Checking Script and Output

1.Every market has a website URL

Not all farmer's market have website

*select count(*) from farmers_market where Website="" ;*

3490

2.IC check website validity

In negate form finds the non compliant number

*select count(\*) from farmers_market where not_FacebookURL="TRUE";*

564

3.Every record has city and county information

*sqlite> select count(\*) from farmers_market where city="" OR county="";*

523

# Appendix C – YesWorkflow Script

This is the sample workflow text, more details on the github page[2]

'''

@begin CleanWithOpenRefine @desc OpenRefine Workflow for farmersmarkets dataset

@in farmersmarkets.csv @uri file://data/farmersmarkets.csv


  @begin LoadToOpenRefine @desc Upload farmersmarkets.csv data to OpenRefine

  @in farmersmarkets.csv @uri file://data/farmersmarkets.cv

  @out spreadsheet

  @end CleanWithOpenRefine


  @begin ColumnsToClean @desc Columns that are cleaned

  @in spreadsheet

  @out MarketName

  @out Website

  @out Facebook

  @out Twitter

  @out Youtube

  @out OtherMedia

  @out street

  @out city

  @out County

  @out State

  @out zip

  @out updatedTime

  @end ColumnsToClean


  @begin TrimSpaces @desc Trim leading and trailing white space

  @in MarketName

@in Website

@in Facebook

@in Twitter

@in Youtube

@in OtherMedia

@in street

@in city

@in county

@in State

@in zip

@in UpdatedTime

@out MarketName_trim

@out Website_trim

@out Youtube_trim

@out Season1Date_trim

@end TrimSpaces


@begin CollapseSpaces @desc Collapse consecutive white spaces

@in MarketName

@in Facebook

@in OtherMeida

@in street

@in city

@in Season1Date

@out MarketName_o

@out city_

@out Facebook_

@out Season1Date_

@end CollapseSpaces


@begin FormatTime @desc Change 'Month DD YYYY' to 'MM/DD/YYYY'

@in updatedTime

@end FormatTime


@begin ClusterValues @desc Cluster Similar values

@in MarketName_

@in city_

@end ClusterValues


@begin DeleteNAorNONE @desc Delete na or none values

@in Facebook_

@in Twitter_

@in Youtube_

@end DeleteNAorNONE


@begin SplitInto2Columns @desc Split date into Start Date and End Date

@in Season1Date_

@out Season1StartDate

@out Season1EndDate

@end SPlitInto2Columns


@begin CleanFormatZip @desc Clean any string and reformat to NNNNN or NNNNN-NNNN

@in zip

@end CleanFormatZip

@begin ToDateType @desc Change data type to date

@in Season1StartDate

@in Season1EndDate

@end ToDateType


@out farmersmarkets-or.csv @uri file://data/farmersmarkets-or.csv

@end CleanWithOpenRefine

'''

# Appendix D – References and Github URL link

1. https://www.ams.usda.gov/local-food-directories/farmersmarkets
2. https://github.com/varunuiuc/DataCleaning

1. https://www.ams.usda.gov/local-food-directories/farmersmarkets
2. https://github.com/varunuiuc/DataCleaning