



CS513: Data Cleaning Project Report

ABSTRACT

Project report on data cleaning performed on New York Public Library's crowd-sourced historical menus dataset to make the dataset usable for data analytics purpose

TEAM

Rohit Bansal (rbansal3@illinois.edu)
Parul Mainwal (mainwal2@illinois.edu)

Table of Contents

Introduction	3
Initial Assessment	3
Structure Summary:	5
Quality Issues	5
Use Cases of the Dataset	6
Data Cleaning Goals	6
Data Cleaning Process.....	7
Step1: Python Script.....	7
Step2: OpenRefine	8
File: Menu.csv	8
File: Dish.csv.....	16
File: MenuItem.csv.....	18
File: MenuPage.csv	18
Step 3: Integrity Checks using SQLite database constraints	18
Database Model.....	18
Tables Creation and Data Import.....	19
Integrity Constraints check	20
Process Workflow (Provenance via YesWorkflow)	22
Summary Workflow	23
Key Inputs/Outputs.....	23
Dependencies.....	23
Diagram.....	23
Data Cleanup Detailed Workflow	23
SQLite Detailed Workflow.....	25
Use Case Implementation	25
Use Case 1: Top 10 Breakfast Dishes	25
Use Case 2: Dishes Served in Hotel Eastman	26
Use Case 3: Five Oldest Dishes.....	27
Challenges	28
Ambiguity	28
OpenRefine Efficiency	28

Results.....	29
Appendix	31
Dataset	31
Tools/Software Used.....	31
Project Deliverables	31

Introduction

This document is a project report for the end-to-end data cleaning process performed on New York Public Library's crowd-sourced historical menus dataset. It captures purpose of the data cleaning, details of the steps performed for data cleaning and our findings during the process. This document also highlights the provenance of the steps performed using a tool called "Yes Workflow". The comparison of the raw data and clean data is published via tableau dashboard [here](#). The cleaned data can be accessed by downloading the dashboard and data source from the same link.

The New York Public Library's restaurant menu collection holds data about menus and dishes from 1840 to the present. This is a crowdsourced dataset collected through spreadsheets and APIs. Since the data is crowdsourced and collected via various means, the data quality is very poor. The objective of this project is to apply various data cleaning techniques and analyze their effectiveness on such a large data set. After performing data cleaning, it is expected that the data will be usable for various data analytics purpose effectively.

Initial Assessment

The New York Public Library's restaurant menu collection is one of the largest in the world. It holds about 45,000 menus dating from the 1840s to the present. It has almost half a million dish items listed in the database along with almost 20 thousand menus of various eateries. We downloaded the data in comma separated values format. As a result, we got 4 files:

- Dish.csv

This file has name of the dish along with its lowest price and highest price. It also has some popularity related data where it informs when the dish first appeared in a menu and when it last appeared and in how many menus it appeared. Here is the list of columns:

- id
- name
- description
- menus_appeared
- times_appeared
- first_appeared
- last_appeared
- lowest_price
- highest_price

- Menu.csv

This file has details about the menus. The details include the name, place, occasion, event, location, venue, sponsor of the menu. It also has some other information such as number of pages in the menu, number of dishes in the menu, language of the menu, currency of payment and some other descriptive fields. Here is the list of columns:

- Id
- Name
- Sponsor
- Event
- Venue
- Place
- physical_description
- occasion
- notes
- call_number
- keywords
- language
- date
- location
- location_type
- currency
- currency_symbol
- status
- page_count
- dish_count

- MenuPage.csv

This file refers to menu file and provides details of menu pages. It provides the image of the menu page, height and width of the page and page number. It also provides the menu id from the menu file so that page could be associated to specific menu. Here is the list of columns:

- id
- menu_id
- page_number
- image_id
- full_height
- full_width
- uuid

- MenuItem.csv

This file refers to the menu, dish and menu page files and provides the information related to each menu item. For each menu item, it provides the dish id from dish.csv and menu page id from menupage.csv. It also provides other details of menu items such as price, creation type, and the location coordinate of the place. Here is the list of columns:

- Id
- menu_page_id
- price

- high_price
- dish_id
- created_at
- updated_at
- xpos
- ypos

Structure Summary:

Features	Menu	Dish	MenuPage	MenuItem
Row count	17547	423400	66937	1332726
Attribute count	20	9	7	9
Unique column	id	id	id	id
Refers			Menu-> Id	MenuPage -> Id Dish-> Id
Description	File with record for each menu and its details	File with records of dishes	File with menu page information along with reference to menu	File with menu item details along with reference to dish and menu page

Quality Issues

We analyzed each of the data files and identified below data quality issues:

- Leading/trailing spaces
- Additional spaces
- Inconsistent case
- Blank values for certain fields
- Inconsistent representation of missing information
- Special Characters like #%?\()\[]
- Spelling mistakes
- In consistent format for values
- Some fields enclosed in double quotes.
- Inconsistent data -abbreviations/full form present for same field.
- Inconsistent date format for date type field.

- Key values mismatch between associated tables.
- Negative/Incorrect values assigned e.g. 0 for price
- Data integrity issues e.g. Created Date>Updated Date.
- Poor formatting
- Ambiguity

Use Cases of the Dataset

As highlighted above, the data had many quality issues. It had to be cleaned in order to make it usable for any use case. Cleaned data could be used for data analytics purpose. Some of the use cases that can be applied on the cleansed data are:

- Answering specific questions for researchers such as when did a dish first appeared on the menu, how has the price/location of a dish changed over the years.
- What dishes are contained in a Menu along with their price and other details.
- The change in demand of a dish over a period and possible contributing factors.
- Correlation between the demand of a dish with its price/location/other factors.

Because of the multiple issues with the data, the dataset was not fit enough for being used in any application. After cleaning the data, we will implement following use case on the cleaned data for demo:

- Use Case 1: Top 10 Breakfast Dishes
- Use Case 2: Dishes Served in Hotel Eastman
- Use Case 3: Five Oldest Dishes

Data Cleaning Goals

After analysis of the dataset, we decided to clean the data to a point that it can be loaded into the database and some basic data analysis could be performed. The goal of our cleaning process is to have it clean enough to be able to support basic user queries like menu structure, popularity of a dish over time etc. We also decided to make the clean data available via public tableau dashboards.

We **did not intend to remove the ambiguity** in the data. We did not want to remove ambiguity by making assumptions. For example, if we found an item as omelette 1 and omelette 2 then we did not try to merge it to “omelette”. The reason behind that was that these two could be a variant of omelette with different ingredient and sometime different size of servings. Assuming them to represent one dish could hamper the data integrity. So, we left those cases as it is. We only cleaned data when it had following issues:

- Leading/trailing spaces
- Additional spaces
- Inconsistent case
- Blank values for certain fields
- Inconsistent representation of missing information
- Special Characters like #%?\`()[]
- Spelling mistakes
- In consistent format for values

- Some fields enclosed in double quotes.
- Inconsistent data -abbreviations/full form present for same field.
- Inconsistent date format for date type field.
- Key values mismatch between associated tables.
- Negative/Incorrect values assigned e.g. 0 for price
- Data integrity issues e.g. Created Date>Updated Date.
- Poor formatting

Data Cleaning Process

Step1: Python Script

We developed a python script in order to get rid of few of the common issues with various fields. These issues included:

- Leading/trailing spaces
- Additional spaces
- Inconsistent case
- Blank values for certain fields
- Inconsistent representation of missing information
- Special Characters like #%?\(\[]

The reason for going for python script instead of Open Refine was that these issues were present with multiple columns in different files. With parallel processing architecture using map-reduce, these cleanups could be done very efficiently in one pass via python script. Using the script, we cleaned following items in one pass:

- File: Menu.csv (12 columns)
 - Name
 - Sponsor
 - Event
 - Venue
 - Place
 - Physical_Description
 - Occasion
 - Notes
 - Date
 - Location
 - Status
 - Currency_Symbol
- File: Dish.csv (1 column)
 - Name

It is noticeable that menu.csv contains 17547 rows while dish.csv contains 423400 rows. Thus, with the script $(12 \times 17547) + (1 \times 423400) = (210564 + 423400) = 633964$ items were cleaned in one step within

few seconds. This turned out to be most efficient than any other method we considered including open refine transformations.

Step2: OpenRefine

We used OpenRefine to further clean the data. Menu, MenuItem and Dish file was cleaned using OpenRefine. MenuPage did not have any issues which could be removed by clustering of facet. Hence it was not considered for cleaning with OpenRefine. We used clustering and facet feature of open refine extensively. We considered following methods for clustering:

- Method: Key-Collision, Function: Fingerprint
- Method: Key-Collision, Function: n-gram fingerprint (n=2)
- Method: Key-Collision, Function: metaphone3
- Method: Key-Collision, Function: cologne-phonetic
- Method: Nearest Neighbor, Function: PPM (Radius: 1, Block Chars: 6)
- Method: Nearest Neighbor, Function: Levenshtein (Radius: 1, Block Chars: 6)

Not all the method returned good candidates and sometimes the method did not return any candidate. We picked and chose only good cluster to perform data cleaning while ignore the other clusters. Here is the list of detailed steps performed on each of the files:

File: Menu.csv

Step 1: Converted Date Field to Date

Step 2: Analyzed Date field outliers and based on that removed outliers from Date field

Step 3: Used Clustering from open refine on column "Sponsor". Following methods were considered:

- a) Method: Key-Collision, Function: Fingerprint
- b) Method: Key-Collision, Function: n-gram fingerprint (n=2)
- c) Method: Key-Collision, Function: metaphone3
- d) Method: Key-Collision, Function: cologne-phonetic
- e) Method: Nearest Neighbor, Function: PPM (Radius: 1, Block Chars: 6)
- f) Method: Nearest Neighbor, Function: Levenshtein (Radius: 1, Block Chars: 6)
- g) Group updates


Screenshot:

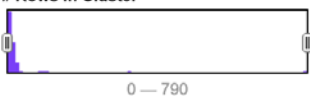
This feature helps you find groups of different cell values that might be alternative representations of the same thing. For example, the two strings "New York" and "new york" are very likely to refer to the same concept and just have capitalization differences, and "Gödel" and "Godel" probably refer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

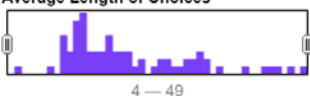
Method key collision Keying Function ngram-fingerprint Ngram Size 2 74 clusters found

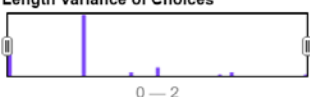
2	6	<ul style="list-style-type: none">A.H. Meyer Rathskeller (3 rows)A.H.Meyer Rathskeller (3 rows)	<input type="checkbox"/>	A.H. Meyer Rathskeller
2	13	<ul style="list-style-type: none">Nippon Yusen Kaisha - S.S.Kasuga (9 rows)Nippon Yusen Kaisha - S.S. Kasuga (4 rows) Browse this cluster	<input type="checkbox"/>	Nippon Yusen Kaisha - S.S.Kas
2	8	<ul style="list-style-type: none">Laurel In The Pines (4 rows)Laurel-In-The-Pines (4 rows)	<input type="checkbox"/>	Laurel In The Pines
2	2	<ul style="list-style-type: none">S. S. President Wilson (1 rows)S.S. President Wilson (1 rows)	<input type="checkbox"/>	S. S. President Wilson
2	2	<ul style="list-style-type: none">Bellevue - Stratford (1 rows)Bellevue-Stratford (1 rows)	<input type="checkbox"/>	Bellevue - Stratford
2	7	<ul style="list-style-type: none">United Air Lines (4 rows)United Airlines (3 rows)	<input type="checkbox"/>	United Air Lines
2	6	<ul style="list-style-type: none">Flat Iron Restaurant And Cafe (3 rows)Flatiron Restaurant And Cafe (3 rows)	<input type="checkbox"/>	Flat Iron Restaurant And Cafe
2	18	<ul style="list-style-type: none">Cafeteria Lunch (15 rows)Cafeteria.Lunch (3 rows)	<input type="checkbox"/>	Cafeteria Lunch

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

Choices in Cluster


Rows in Cluster


Average Length of Choices


Length Variance of Choices


Step 4: Used Clustering from open refine on column "Event". Following methods were used:

- Method: Key-Collision, Function: Fingerprint
- Method: Key-Collision, Function: n-gram fingerprint (n=2)
- Method: Key-Collision, Function: metaphone3
- Method: Key-Collision, Function: cologne-phonetic
- Method: Nearest Neighbor, Function: PPM (Radius: 1, Block Chars: 6)
- Method: Nearest Neighbor, Function: Levenshtein (Radius: 1, Block Chars: 6)
- Group updates

Screenshot:

Cluster & Edit column "event"

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 47 clusters found

2	12	<ul style="list-style-type: none">New Year'S Day Dinner (7 rows)New Years Day Dinner (5 rows)	<input checked="" type="checkbox"/>	New Year'S Day Dinner
2	6	<ul style="list-style-type: none">Eighth Annual Banquet (5 rows)Eighth Annual Banquet' (1 rows)	<input checked="" type="checkbox"/>	Eighth Annual Banquet
2	10	<ul style="list-style-type: none">Daily Menu, Breakfast (9 rows)Daily Breakfast Menu (1 rows)	<input checked="" type="checkbox"/>	Daily Menu, Breakfast
2	21	<ul style="list-style-type: none">Daily Menu, Luncheon (20 rows)Daily Luncheon Menu (1 rows)	<input checked="" type="checkbox"/>	Daily Menu, Luncheon
2	547	<ul style="list-style-type: none">Luncheon (543 rows)Luncheon; (4 rows)	<input checked="" type="checkbox"/>	Luncheon
2	5	<ul style="list-style-type: none">Inspection Trip, Breakfast (3 rows)Inspection Trip - Breakfast (2 rows)	<input checked="" type="checkbox"/>	Inspection Trip, Breakfast
2	11	<ul style="list-style-type: none">Mittagessen - Dinner (9 rows)Mittagessen Dinner (2 rows) Browse this cluster	<input checked="" type="checkbox"/>	Mittagessen - Dinner
2	2	<ul style="list-style-type: none">Lunch-Dinner (1 rows)Lunch;Dinner (1 rows)	<input checked="" type="checkbox"/>	Lunch-Dinner

Choices in Cluster
2 — 4

Rows in Cluster
0 — 2200

Average Length of Choices
4 — 81

Length Variance of Choices
0 — 1

Step 5: Used Clustering from open refine on column "Place". Following methods were used:

- Method: Key-Collision, Function: Fingerprint
- Method: Key-Collision, Function: n-gram fingerprint (n=2)
- Method: Key-Collision, Function: metaphone3
- Method: Key-Collision, Function: cologne-phonetic
- Method: Nearest Neighbor, Function: PPM (Radius: 1, Block Chars: 6)
- Method: Nearest Neighbor, Function: Levenshtein (Radius: 1, Block Chars: 6)
- Group updates

Screenshot:

Cluster & Edit column "place"

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 **nearest neighbor** Distance Function **PPM** Radius **1.0** Block Chars **6** **122 clusters found**

		<ul style="list-style-type: none">Bleecker,Thompson,And Sullivan Streets,Ny (1 rows)		
2	2	<ul style="list-style-type: none">Birmingham, Al (1 rows)Birmingham (1 rows)	<input checked="" type="checkbox"/>	Birmingham, Al
2	5	<ul style="list-style-type: none">Grand Pacific Hotel, Chicago, Il (4 rows)Grand Pacific Hotel,Chicago,IlI (1 rows)	<input type="checkbox"/>	Grand Pacific Hotel, Chicago, Il
2	3	<ul style="list-style-type: none">Broadway & 10Th Street, New York, Ny (2 rows)Broadway And 10Th Street, New York, Ny (1 rows)	<input checked="" type="checkbox"/>	Broadway & 10Th Street, New Y
2	2	<ul style="list-style-type: none">Hotel Jermyn, Scranton, Pa (1 rows)The Hotel Jermyn, Scranton, Pa (1 rows)	<input checked="" type="checkbox"/>	Hotel Jermyn, Scranton, Pa
2	2	<ul style="list-style-type: none">Portland, Oregon (1 rows)Portland,The, Portland, Oregon (1 rows) Browse this cluster	<input checked="" type="checkbox"/>	Portland, Oregon
2	2	<ul style="list-style-type: none">Hotel Vendome,Ny (1 rows)Hotel Vendome (1 rows)	<input checked="" type="checkbox"/>	Hotel Vendome,Ny
2	4	<ul style="list-style-type: none">Murray Hill Hotel,Ny (3 rows)	<input type="checkbox"/>	Murray Hill Hotel,Ny

Choices in Cluster**# Rows in Cluster****Average Length of Choices****Length Variance of Choices**

Step 6: Used Clustering from open refine on column "Venue". Following methods were used:

- Method: Key-Collision, Function: Fingerprint
- Method: Key-Collision, Function: n-gram fingerprint (n=2)
- Method: Key-Collision, Function: metaphone3
- Method: Key-Collision, Function: cologne-phonetic
- Method: Nearest Neighbor, Function: PPM (Radius: 1, Block Chars: 6)
- Method: Nearest Neighbor, Function: Levenshtein (Radius: 1, Block Chars: 6)
- Group updates

Screenshot:

Cluster & Edit column "venue"

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 metaphone3 3 clusters found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
3	6	<ul style="list-style-type: none"> Other Privately Hosted Dinner Party (4 rows) Other Privately Hosted Banquet (1 rows) Other Privately Hosted Luncheon (1 rows) 	<input type="checkbox"/>	Other Privately Hosted Dinner P
2	3	<ul style="list-style-type: none"> Other Private Hosts (2 rows) Other Private Hostess (1 rows) 	<input type="checkbox"/>	Other Private Hosts
2	33	<ul style="list-style-type: none"> Other Private Party (26 rows) Other Private Party By Single Host (7 rows) 	<input type="checkbox"/>	Other Private Party

Choices in Cluster

2 — 3

Rows in Cluster

3 — 33

Average Length of Choices

20 — 32

Length Variance of Choices

1 — 7.5

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

Step 7: Used Clustering from open refine on column "Occasion". Following methods were used:

- Method: Key-Collision, Function: Fingerprint
- Method: Key-Collision, Function: n-gram fingerprint (n=2)
- Method: Key-Collision, Function: metaphone3
- Method: Key-Collision, Function: cologne-phonetic
- Method: Nearest Neighbor, Function: PPM (Radius: 1, Block Chars: 6)
- Method: Nearest Neighbor, Function: Levenshtein (Radius: 1, Block Chars: 6)
- Group updates

Screenshot:

Cluster & Edit column "occasion"

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 metaphone3 12 clusters found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
6	576	<ul style="list-style-type: none">Complimentary/Testimonial (569 rows)Complimentary (2 rows)Complimentary Dinner (2 rows)Complimentary Banquet For Montgomery Guards (1 rows)Complimentary To Commodore C. We. Bliss (1 rows)Complimentary/Testimonial,Anniversary (1 rows)	<input type="checkbox"/>	Complimentary/Testimonial
4	4	<ul style="list-style-type: none">Other Commercial (1 rows)Other, Commercial Convention Dinner (1 rows)Other,Commercial Anniversary (1 rows)Other,Commercial Opening (1 rows)	<input type="checkbox"/>	Other Commercial
2	2	<ul style="list-style-type: none">Other Private Party (1 rows)Other Private Party Ferdinard De Lesseps And Party (1 rows)	<input type="checkbox"/>	Other Private Party
2	2	<ul style="list-style-type: none">Other,University Anniversary (1 rows)Other,University Paper Annual Dinner (1 rows)	<input type="checkbox"/>	Other,University Anniversary
2	2	<ul style="list-style-type: none">Anniversary,Fifteenth Annual Convention (1 rows)AnniversaryFifth Of Founding (1 rows)	<input type="checkbox"/>	Anniversary,Fifteenth Annual Cc

Choices in Cluster

Rows in Cluster

Average Length of Choices

Length Variance of Choices

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

Step 8: Used Clustering from open refine on column "Location". Following methods were used:

- Method: Key-Collision, Function: Fingerprint
- Method: Key-Collision, Function: n-gram fingerprint (n=2)
- Method: Key-Collision, Function: metaphone3
- Method: Key-Collision, Function: cologne-phonetic
- Method: Nearest Neighbor, Function: PPM (Radius: 1, Block Chars: 6)
- Method: Nearest Neighbor, Function: Levenshtein (Radius: 1, Block Chars: 6)
- Group updates

Screenshot:

Cluster & Edit column "location"

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: cologne-phonetic 75 clusters found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
4	16	<ul style="list-style-type: none"> Hotel Dennis (9 rows) Hotel Du Musee (4 rows) Hotel Athens (2 rows) Hotel Adams (1 rows) Browse this cluster	<input type="checkbox"/>	Hotel Dennis
3	3	<ul style="list-style-type: none"> Diva (1 rows) Djawah (1 rows) The Ivy (1 rows) 	<input type="checkbox"/>	Diva
3	4	<ul style="list-style-type: none"> Coffee House (2 rows) Chevy'S (1 rows) Coffee Etc (1 rows) 	<input type="checkbox"/>	Coffee House
3	7	<ul style="list-style-type: none"> 21 Club (4 rows) Club (2 rows) Club 31 (1 rows) 	<input type="checkbox"/>	21 Club
3	4	<ul style="list-style-type: none"> LutÄ"ce (2 rows) Le Duc (1 rows) Lite House (1 rows) 	<input type="checkbox"/>	LutÄ"ce
3	6	<ul style="list-style-type: none"> La Choza (4 rows) 	<input type="checkbox"/>	La Choza

Choices in Cluster

Rows in Cluster

Average Length of Choices

Length Variance of Choices

Step 9: Used Clustering from open refine on column "Name". Following methods were used:

- Method: Key-Collision, Function: Fingerprint
- Method: Key-Collision, Function: n-gram fingerprint (n=2)
- Method: Key-Collision, Function: metaphone3
- Method: Key-Collision, Function: cologne-phonetic
- Method: Nearest Neighbor, Function: PPM (Radius: 1, Block Chars: 6)
- Method: Nearest Neighbor, Function: Levenshtein (Radius: 1, Block Chars: 6)
- Group updates

Screenshot:

Cluster & Edit column "name"

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 metaphone3 14 clusters found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
4	6	<ul style="list-style-type: none"> Ritz Carlton (2 rows) Ritz Carlton Hotel (2 rows) Ritz Carlton New York (1 rows) Ritz-Carlton Restaurant (1 rows) 	<input type="checkbox"/>	Ritz Carlton
3	4	<ul style="list-style-type: none"> Long Island Business College (2 rows) Long Island Business College Hall (1 rows) Long Island Business College South (1 rows) 	<input type="checkbox"/>	Long Island Business College
3	5	<ul style="list-style-type: none"> Great Northern Hotel (3 rows) Great Northern (1 rows) Great Northern Railway (1 rows) 	<input type="checkbox"/>	Great Northern Hotel
3	7	<ul style="list-style-type: none"> Pennsylvania Lunch (5 rows) Pennsylvania Lines (1 rows) Pennsylvania Lines: West Of Pittsburgh (1 rows) 	<input type="checkbox"/>	Pennsylvania Lunch
3	565	<ul style="list-style-type: none"> Waldorf Astoria (559 rows) Waldorf-Astoria: South Gate (5 rows) Waldorf-Astoria: Grand Ball Room (1 rows) 	<input type="checkbox"/>	Waldorf Astoria
2	2	<ul style="list-style-type: none"> Horn & Hardart Company (1 rows) 	<input type="checkbox"/>	Horn & Hardart Company

Choices in Cluster

Rows in Cluster

Average Length of Choices

Length Variance of Choices

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

Step 10: Split the physical_description column. It resulted in 7 new columns:

physical_description 1

physical_description 2

physical_description 3

physical_description 4

physical_description 5

physical_description 6

physical_description 7

We updated the column physical_description 1 to physical description type and then created a new column Physical_description_additional by merging following columns:

physical_description 2

physical_description 3

physical_description 4

Screenshot:

Add column based on column physical_description 2

New column name

Physical_Description_Additional

☒ set to blank

☐ store error

☐ copy value from original column

Expression

Language

General Refine Expression Language (GREL) ▼

```
if(or(isBlank(value),isNull(value)) , "",value+" - "+
(if(or(isBlank(cells["physical_description
3"].value),isBlank(cells["physical_description
3"].value)), "",cells["physical_description 3"].value))+" - "+
```

No syntax error.

Preview

History

Starred

Help

row	value	if(or(isBlank(value),isNull(value)) , "",value+" - "+ (if(or(isBlank(cells["physical_description 3"].value),isBlank(cells["physical_description 3"].value)), "",cells["physical_description 3"].value))+" - "+ (if(or(isBlank(cells["physical_description 4"].value),isBlank(cells["physical_description 4"].value)), "",cells["physical_description 4"].value)))
1.	4.75X7.5	4.75X7.5 - -
2.	Illus	Illus - Col - 7.0X9.0
3.	Illu	Illu - Col - 5.5X8.0

OK

Cancel

File: Dish.csv

Step 1: We found even after cleaning with python script, there were few different kinds of special characters were present in the file. So, we transform data to remove unwanted characters. Formula:

```
value.replace(/[>:<%#@!\\\()\\[]?\"\\-\\*\\.\\+]/, " ").replace(/\\s+/, " ").trim()
```

Step 2: Used Clustering from open refine on column "Name". Following methods were used:

- Method: Key-Collision, Function: Fingerprint
- Method: Key-Collision, Function: n-gram fingerprint (n=2)
- Method: Key-Collision, Function: metaphone3
- Method: Key-Collision, Function: cologne-phonetic
- Method: Nearest Neighbor, Function: PPM (Radius: 1, Block Chars: 6)
- Method: Nearest Neighbor, Function: Levenshtein (Radius: 1, Block Chars: 6)
- Group updates

Screenshot:

Cluster & Edit column "name"

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 ...](#)

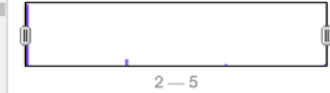
Method **key collision**

Keying Function **fingerprint**

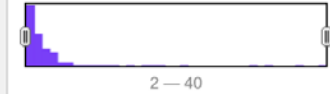
396 clusters found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
5	7	<ul style="list-style-type: none">Canadian Cheese Imported (2 rows)Imported Canadian Cheese (2 rows)Canadian Imported Cheese (1 rows)Cheese Canadian Imported (1 rows)Imported Cheese Canadian (1 rows)	<input type="checkbox"/>	Canadian Cheese Imported
5	13	<ul style="list-style-type: none">Iced Tea Or Coffee (6 rows)Iced Tea Or Iced Coffee (3 rows)Iced Coffee Or Tea (2 rows)Coffee Or Tea Iced (1 rows)Iced Coffee Or Iced Tea (1 rows)	<input type="checkbox"/>	Iced Tea Or Coffee
5	17	<ul style="list-style-type: none">Imported Swiss Cheese (6 rows)Swiss Cheese Imported (5 rows)Cheese Imported Swiss (2 rows)Imported Cheese Swiss (2 rows)Swiss Imported Cheese (2 rows) Browse this cluster	<input type="checkbox"/>	Imported Swiss Cheese
5	13	<ul style="list-style-type: none">Club Soda C & C Imported (5 rows)Imported Club Soda C & C (4 rows)Club Soda C & C Imported (2 rows)C & C Club Soda Imported (1 rows)Club Soda Imported C & C (1 rows)	<input type="checkbox"/>	Club Soda C & C Imported

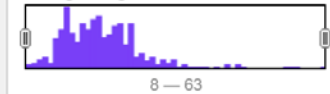
Choices in Cluster



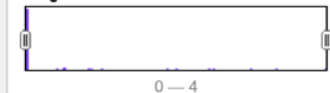
Rows in Cluster



Average Length of Choices



Length Variance of Choices



Select All Unselect All

Merge Selected & Re-Cluster

Merge Selected & Close

Close

Cluster & Edit column "name"

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 **nearest neighbor**

Distance Function **levenshtein**

Radius **1.0**

Block Chars **6**

20642 clusters filtered from 21970 total

3	3	<ul style="list-style-type: none">Antedeluvian (1 rows)Antediluvian (1 rows)Antidiluvian (1 rows)	<input checked="" type="checkbox"/>	Antedeluvian
3	3	<ul style="list-style-type: none">Schwepp'S Imported Ginger Ale Pint 50 (1 rows)Schwepp'S Imported Gingerale Pint (1 rows)Schwepp'S Imported Gingerale Pint (1 rows)	<input type="checkbox"/>	Schwepp'S Imported Ginger
3	3	<ul style="list-style-type: none">Tomates Genoise (1 rows)Tomates GÄnoise (1 rows)Tomatoes Genoise (1 rows) Browse this cluster	<input type="checkbox"/>	Tomates Genoise
3	3	<ul style="list-style-type: none">Guinness Export Stout (1 rows)Guinnes Export Stout (1 rows)Guinness Export Stout (1 rows)	<input type="checkbox"/>	Guinness Export Stout
3	5	<ul style="list-style-type: none">Manhattan Salad (3 rows)Manhatam Salad (1 rows)Mountain Salad (1 rows)	<input type="checkbox"/>	Manhattan Salad
3	5	<ul style="list-style-type: none">Orange Pekoe (3 rows)Orange PekoÅ« (1 rows)Orange PÄ«koÅ« (1 rows)	<input type="checkbox"/>	Orange Pekoe
3	5	<ul style="list-style-type: none">Consomme Florentine (3 rows)	<input type="checkbox"/>	Consomme Florentine

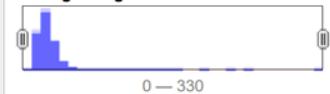
Choices in Cluster



Rows in Cluster



Average Length of Choices



Length Variance of Choices



Select All Unselect All

Merge Selected & Re-Cluster

Merge Selected & Close

Close

File: MenuItem.csv

Step 12: Transform created_at field to date by creating a new field created_date

Step 13: Transformed updated_at field to date by creating a new field updated_date.

(Note: both created_at and created_date field exists in the csv file. Similarly, both updated_at and updated_date field exists in the database.)

Step 14: Transform xpos field to number

Step 15: Transformed ypos field to number

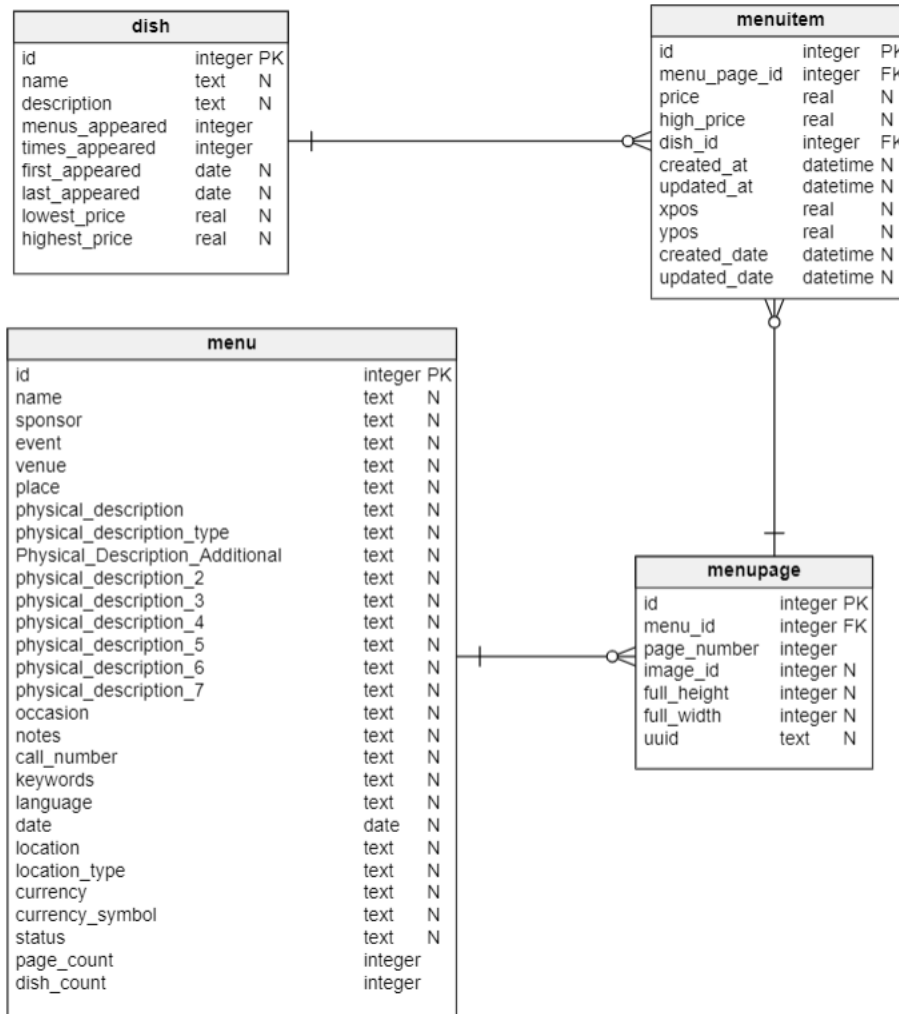
File: MenuPage.csv

No transformations were performed for this csv in OpenRefine.

Step 3: Integrity Checks using SQLite database constraints

Database Model

Based on the cleaned data from step 2, we created a database model which included the logical integrity constraints of 'Primary Key', 'Foreign Key' and 'Not Null'. The database model also defined the one-to-one and one-to-many relationships between the four tables. Building the schema based on this model ensured that only the records satisfying all the defined constraints are loaded into the tables and the dirty data is discarded. Here is the ER diagram of our DB model:



Tables Creation and Data Import

A database schema 'nyp1' was created using SQLiteStudio. We created 4 tables and imported data into these tables from the cleaned csv files. As a result, following 4 tables were created:

- Dish
- Menu
- Menupage
- MenuItem

Screenshot (for data import)

← Import data

Data source to import from

Data source type: CSV

Options

Input file: Jb/513FinalProject-master/Step2_OpenRefine/Dish/513_Data_cleaning_Dish.csv

Text encoding: UTF-8

☒ Ignore errors

Data source options

☒ First line represents CSV column names

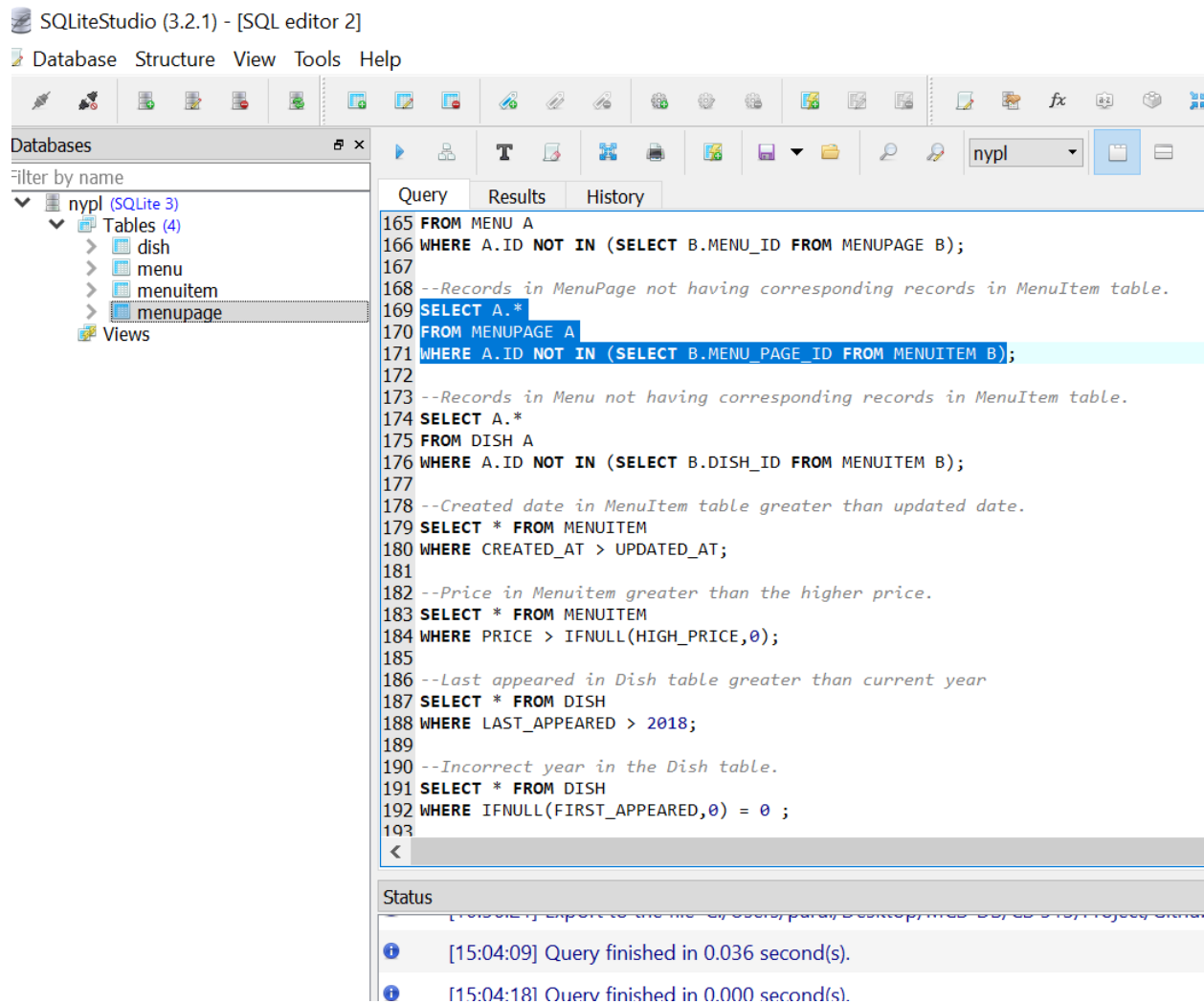
Field separator: , (comma)

☐ NULL values:

Cancel Finish

Integrity Constraints check

We created several queries to check the logical integrity constraints on the tables. All the scripts for the schema/table creation and integrity constraints check are listed under the 'SQL_scripts'.sql file. SQLiteStudio was used to run the queries and obtain results.



Below is a brief description of the various checks done on the data.

Key Data Duplication Check

This step checked presence of any duplicate data in the table based on ID field. Since, we applied primary key constraint on ID field during the data load, the duplicate data was filtered out already from the tables. Our final data in the tables returned no duplicate records.

Foreign Key Constraint Check

This step checked the foreign key constraints violation on the tables. That is if any table that used primary key of another table, had data that not matched with the primary key value in the source table. We applied this constraint during the data load itself for Menu, Dish and MenuItem tables. We had to disable the foreign key constraints only on the MenuItem table during the data load as it would have otherwise resulted in loss of too many rows of data and we wouldn't get enough data for further analysis.

Not Null Check

We applied not null constraint on certain table columns to ensure that only those records were loaded into the table that had values in these columns. This resulted in discarding irrelevant records which would not be useful for further analysis. This constraint was applied on all the tables during the data load itself.

Other Checks

We applied following integrity checks to identify and correct other integrity constraint violations:

- Created date in Menuitem table greater than updated date.
- Price in Menuitem greater than the higher price.

We found that 1091 records were in violation of these constraints hence got deleted.

We checked for further integrity violations as well but for these cases we did not perform any corrective action. Corrective action is often context dependent and can be different for different use cases. Thus, we refrained from taking corrective actions. Rather, we have just highlighted the violations via sql queries. Here are the constraint descriptions along with record count returned after running the queries:

- Last appeared in Dish table greater than current year (2018)
Result: 179 records returned.
- Incorrect year in the Dish table.
Result: 55287 records returned.
- First_appeared greater than last_appeared in Dish table.
Result: 6 records returned.
- Lowest price greater than highest price in Dish table.
Result: 0 records returned.
- Count of menus_appeared in dish table not matching count obtained from menuitem and menupage tables.
Result: 2375 records returned.
- Mismatch of page count between menu and menupage tables.
Result: 147 records returned.

Process Workflow (Provenance via YesWorkflow)

We created overall process workflow using online YesWorkflow tool (link in appendix section) for provenance. This workflow provides a birds-eye view of the overall process used for this data cleaning exercise.

We also created two separate detailed workflows for SQLite and Data Cleanup Operations. Each of these operations involved multiple steps and these workflows provide the details about these operations.

Summary Workflow

Key Inputs/Outputs

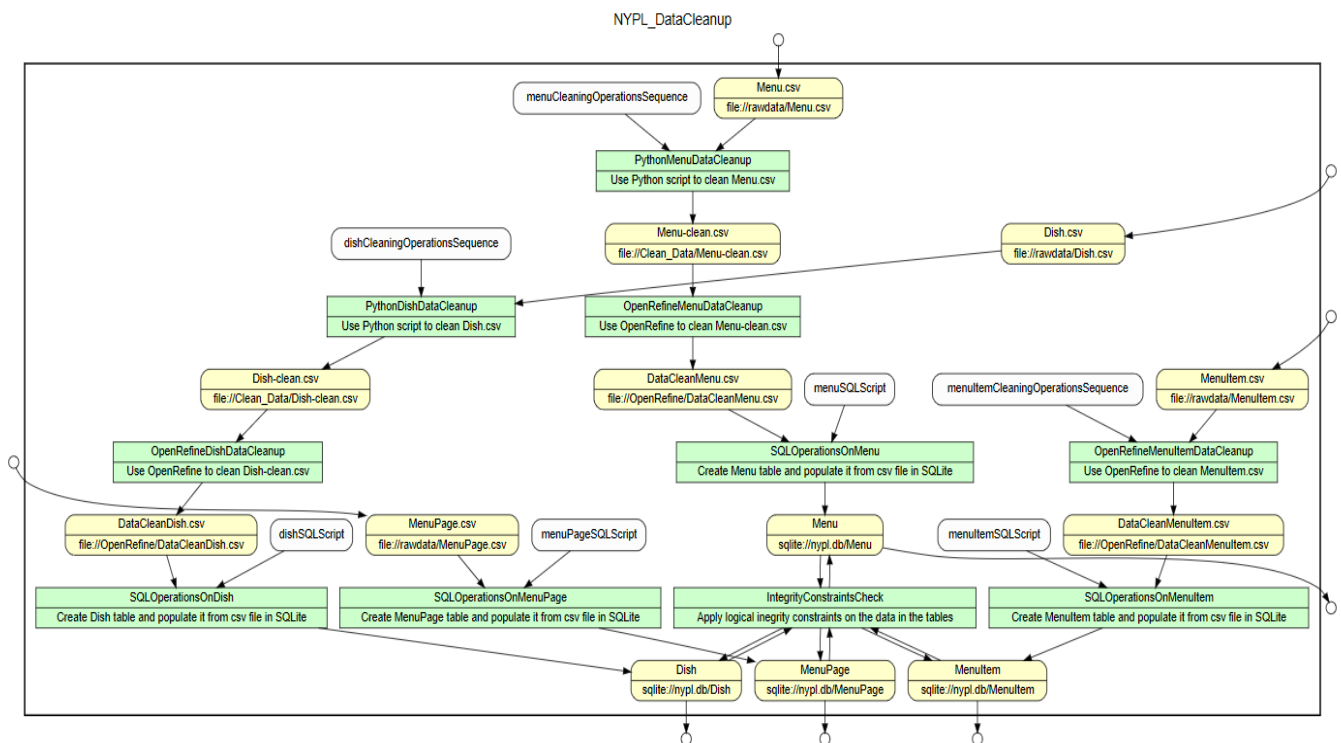
The key inputs to the summary workflow process are the four input data files namely, menu.csv, dish.csv, menuitem.csv, menupage.csv along with the data cleaning operations comprising of OpenRefine step, Python script step, and SQLite step.

Dependencies

The complete workflow includes several steps that are interdependent on each other. The data cleaning operation done with OpenRefine and Python script is dependent on the raw data files. OpenRefine step is dependent on python script output for dish.csv and menu.csv. The SQL operations of schema creation and integrity constraints is dependent on the files cleansed through open refine for dish.csv, menu.csv and menuitem.csv.

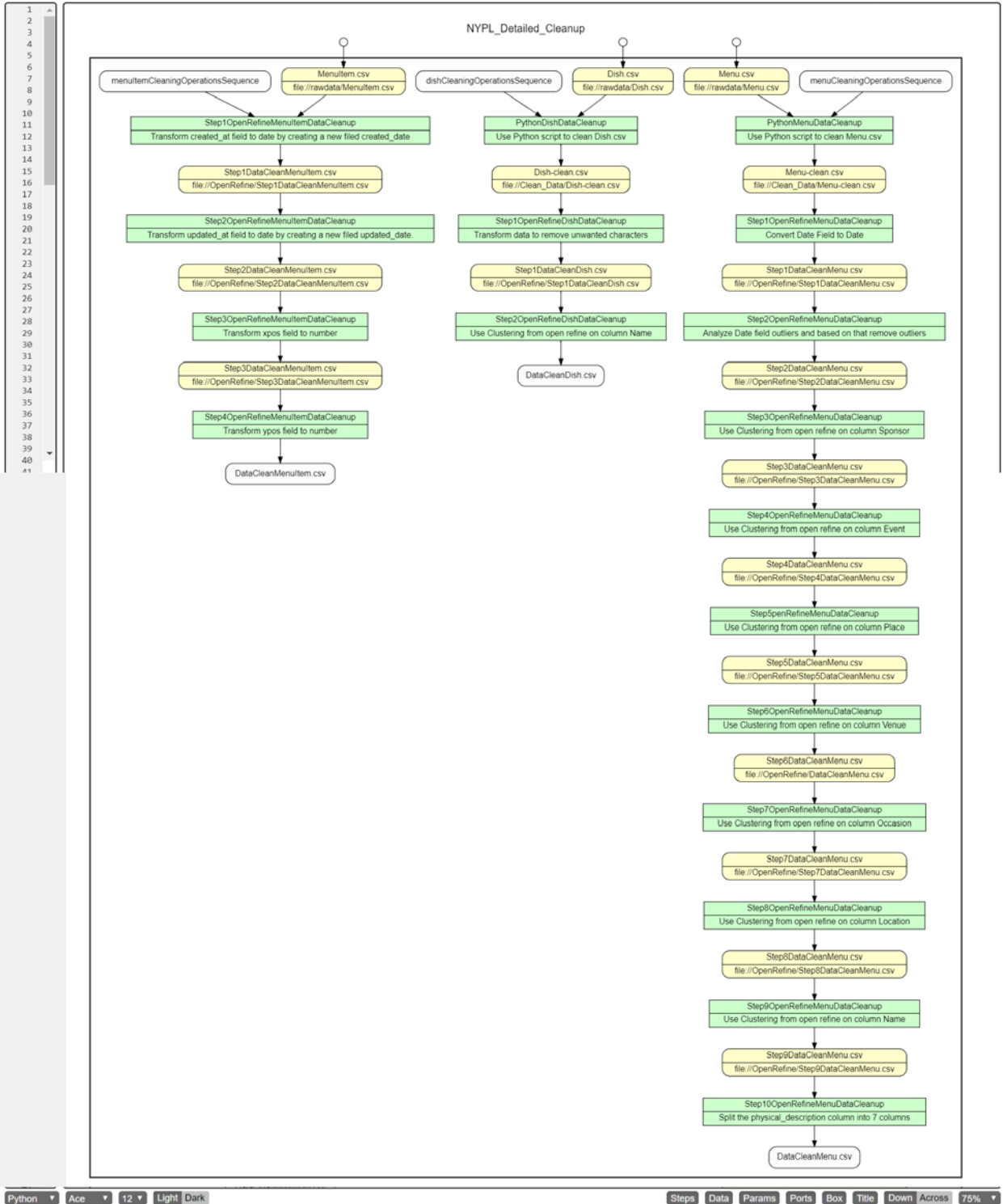
Diagram

Below is a snapshot of the Overall Workflow Diagram:



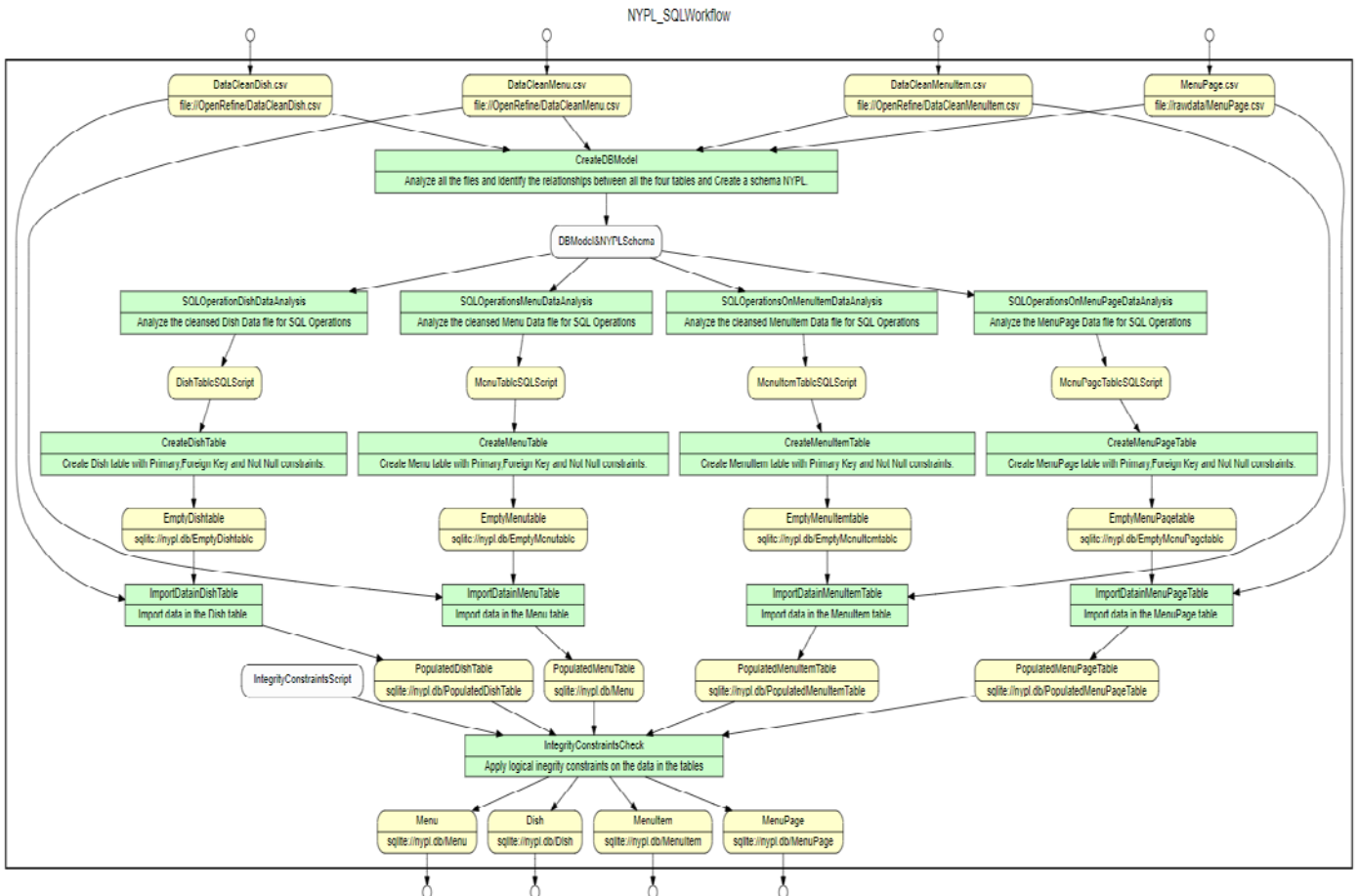
Data Cleanup Detailed Workflow

The input here are the raw data files and the output are the cleansed Menu, MenuItem and Dish data files. Since, MenuPage.csv was already clean, it was not part of this step. The multiple sub steps created partially cleaned files which have been named accordingly in the diagram.



SQLite Detailed Workflow

The input here is the MenuPage file and other data files cleaned via OpenRefine and Python and the output are the 4 tables with clean data after all the integrity constraints violation checks.



Use Case Implementation

After we were done with all the above steps, we implemented three use cases that provide meaningful data analysis to the users. This step was to check the usability of finally cleaned data. Here are the details of the use cases conducted on the final data available in the database:

Use Case 1: Top 10 Breakfast Dishes

The goal here is to identify 10 most popular Breakfast items in New York since 2008 based on the data available. Here is the query and the output:

The screenshot shows a database query tool interface. On the left, a sidebar lists the database schema: 'me' (SQLite 3), 'Tables (4)', and a list of tables: 'dish', 'menu', 'menuitem', and 'menupage'. The main area is divided into 'Query' and 'History' tabs. The 'Query' tab contains a SQL query. Below the query, there are 'Grid view' and 'Form view' tabs. The 'Grid view' tab is active, showing a table of results. The table has 10 rows and 9 columns: event, place, name, menu, times, status, first a, and last ar. The results show the top 10 most popular breakfast items in New York since 2008, with 'Tea' and 'Potatoes Mashed' being the most popular.

```

1 --Use Case 1
2 --Top 10 most popular Breakfast items in New York since 2008.
3 select distinct a.event, a.place, d.name, menus_appeared, times_appeared, status, first_appeared, last_appeared
4 from menu a, menupage b, menuitem c, dish d
5 where a.id=b.menu_id
6 and b.id=c.menu_page_id
7 and c.dish_id=d.id
8 and event = 'Breakfast'
9 and first_appeared < > 1
10 and last_appeared > 2008 and last_appeared < 2018
11 and place= 'New York'
12 order by menus_appeared desc limit 10;

```

	event	place	name	menu	times	status	first a	last ar
1	Breakfast	New York	Tea	4159	4769	Complete	1858	2012
2	Breakfast	New York	Potatoes Mashed	2583	2670	Complete	1852	2012
3	Breakfast	New York	Milk	1894	2290	Complete	1890	2012
4	Breakfast	New York	Oranges	1423	1479	Complete	1851	2012
5	Breakfast	New York	American Cheese	1372	1398	Complete	1865	2012
6	Breakfast	New York	Assorted Cakes	1331	1355	Complete	1880	2012
7	Breakfast	New York	Potatoes Baked	1293	1326	Complete	1859	2012
8	Breakfast	New York	Apples	1215	1239	Complete	1851	2012
9	Breakfast	New York	Watermelon	674	792	Complete	1881	2012
10	Breakfast	New York	Plums	399	503	Complete	1889	2012

The output indicates that tea and mashed potatoes are the most popular dishes for breakfast since 2008 which seems reasonable.

Use Case 2: Dishes Served in Hotel Eastman

The goal here is to output to the user the dishes served on Hotel Eastman's Lunch menu. Following is the query and output:

QueryHistory

```
1 --Use Case 3
2 --5 of the oldest dishes on the Menu (first appeared in 1851)
3 select * from dish where
4 first_appeared= (select min(first_appeared) from dish where first_appeared not in ( 0,1)) order by id desc limit 5
```

Grid viewForm view

✓

✕

↺

↻

1

➡

⌂

Total rows loaded: 5

id	name	descri	menu	times	first a	last a	lowest	highes
1 493867	Claret Chateau Lafitte		1	1	1851	1851	1	2
2 493866	Claret Chateau La Rose		1	1	1851	1851	0.75	1.5
3 493865	Porter And Ale Guinness'S Xx Dublin Brown Stout		1	1	1851	1851	0.25	0.5
4 493864	Porter And Ale Byass'S London Porter		1	1	1851	1851	0.25	0.5
5 493863	Porter And Ale Tennant'S Scotch Ale		1	1	1851	1851	0.25	0.5

The result shows us few oldest wines which makes sense.

Challenges

Ambiguity

One of the key challenges that we faced was the ambiguity in the data. We did not intend to remove the ambiguity in the data as that would have required us to make assumptions. For example, if we found an item as omelette 1 and omelette 2 then we did not try to merge it to “omelette”. The reason behind it was that these two could be a variant of omelette with different ingredient and sometimes different size of servings. Assuming them to represent one dish could hamper the data integrity. Thus, there was a tradeoff between ambiguity and integrity. So, we left those cases as it is

OpenRefine Efficiency

Open refine was not able to load large datasets such as dish.csv. We had to change following configurations in refine.ini

```
REFINE_MAX_FORM_CONTENT_SIZE=99999999
```

```
REFINE_MEMORY=4096M
```

Same can be achieved by following command: openrefine /m 4096M.

Further, for some of the operations, especially facet creation on large dataset, open refine was not able to perform optimally. We had to reduce the data size by applying a filter and then we were able to

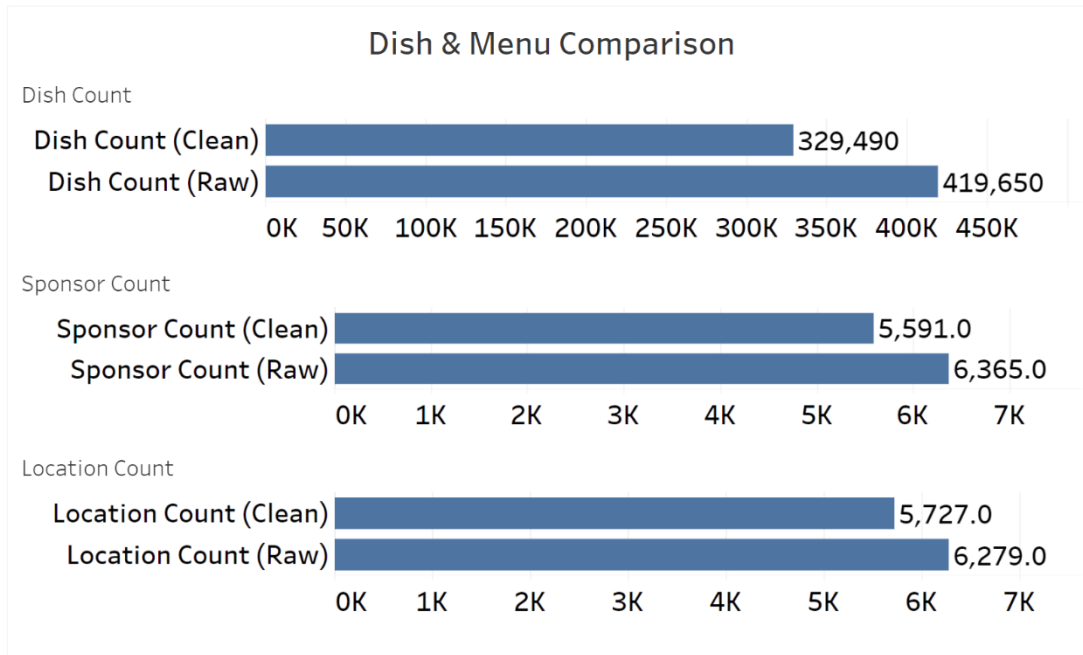
create the facet. We also found that it is more time-efficient to apply clustering directly instead of first creating facet and then applying the clustering.

Results

After cleaning the data, we compared the raw data with the cleaned data. We used tableau to visualize the comparison. The comparison can be accessed [here](#).

Comparison			
Raw Data		Clean Data	
Number of Dishes		Number of Dishes	
	419,650		329,490
Number of Sponsors		Number of Sponsors	
	6,365		5,591
Number of Menus		Number of Menus	
	19,816		17,517
Avg Dish Price		Avg Dish Price	
	12.84		17.05

We also compared some other fields in the dish and menu files. The comparison can be accessed [here](#).



Also, from the same link user can access both the raw data and the clean data via tableau dashboards. These tableau dashboards are available on public.tableau.com. This way clean data can be accessed for further data analytics purpose.

We can see that a there is a huge impact of the cleaning exercise. There were several fields with dirty data and were cleaned as part of this exercise. For dish names alone almost 90K dish names were found to be either duplicate or had some sort of spelling mistakes!!

Appendix

Dataset

New York Public Library Menus data

Tools/Software Used

Software	Version/Url	Usage
Python	3.6.0	Scripting Language used for Data Cleaning
OpenRefine	2.6-rc.2	Data Cleaning Tool
SQLiteStudio	3.2.1	UI tool for working with SQLite
YesWorkflow	http://try.yesworkflow.org/	Provenance Tool
Tableau	2018.1.2	Visualization Representation
Vetabelo	https://my.vertabelo.com/	Database Modeler
GitHub	https://github.com/rohitbansal83/513FinalProject	Source control repository
Slack	https://uiuc-mcsds.slack.com/	Communication channel used for project work
zoom	https://www.zoom.us/	Video Conferencing tool for coordination

Project Deliverables

All the projects deliverables including the raw dataset, sql scripts, cleansed table data, Yesworkflow scripts can be accessed at this [link](#).