

# CS 513: Theory and Practice of Data Cleaning

Final Project: Data Cleaning NYPL Menus  
(Team 52)

## [Data Cleaning Project Phase II](#)

[Improved Phase - 1 report](#)

[Target Use Case U1](#)

[Data Cleaning Performed](#)

[Introduction](#)

[Using OpenRefine](#)

[Menu.csv \(Total number of records = 17,546\)](#)

[Dish.csv \(Total number of records = 423397\)](#)

[MenuItem.csv \(Total number of records = 104,586\)](#)

[MenuPage.csv \(Total number of records = 66,938\)](#)

[Using Python](#)

[Using SQLite3](#)

[Table name menu](#)

[Table menu\\_item](#)

[Table menu\\_pages](#)

[Table dish](#)

[Cleaned records report](#)

[Summary table changes](#)

[Rationale For design](#)

[ER Diagram](#)

[Overview Of the cleaning process](#)

[Workflow Diagrams](#)

[Summary / Lesson Learned](#)

[Acknowledgements](#)

[User Stories and Contributions](#)

[Project Team members](#)

[Github Details](#)

[FilePath Details](#)

# Data Cleaning Project Phase II

## Improved Phase - 1 report

### Target Use Case U1

We wanted to clean the data to find out historical information about the menu and dishes which were served. For example, what were the menu items which were served for Dinner at New York back in the 1960s? And how menu items were priced and what were the dishes served during that event? What were the locations where these dishes were served?

## Data Cleaning Performed

### Introduction

In this section, we have described how the data cleaning was performed using various technologies and tools. We used the dataset NYPL menus. We used Openrefine primarily to clean most of the data, we also used Python to clean up special characters from strings using Regex. And finally we used SQLite3 database and queries to fix integrity constraint violations. The [Overview Of the cleaning process](#) section shows how the cleaning process was done for the use case mentioned above.

### 1) Using OpenRefine

We performed data cleaning on all the four csv files of NYPL. All the 4 csv files were uploaded into the Openrefine to do the cleaning. The cleaning operations were captured in the JSON files.

#### **Menu.csv (Total number of records = 17,546)**

Below is the list of columns for the Menu.csv file.

- Please be advised that we performed a lot of steps (100+) so it is not possible to list down all of them in this document. Many steps were repeatedly performed to clean data on various columns.
- Below are some of the updates we made on the columns of the csv -

◆ **Id**

- Converted the column values to number

◆ **name**

- Removed the white spaces (trailing spaces, consecutive white spaces, etc.)
- Converted the column values to Uppercase.
- We used Regex on GREL to remove the special characters.
- We used clustering to group the data values into multiple groups. The below methods were used to cluster the data (Key collision and Nearest Neighbor).

◆ **sponsor**

- Removed the white spaces (trailing spaces, consecutive white spaces, etc.)
- Converted the column values to Uppercase.
- We used Regex on GREL to remove the special characters.

◆ **Event**

- Removed the white spaces (trailing spaces, consecutive white spaces, etc.)
- Converted the column values to Uppercase.
- We used Regex on GREL to remove the special characters.
- We used clustering to group the data values into multiple groups. The below methods were used to cluster the data (Key collision).

◆ **Venue**

- Removed the white spaces (trailing spaces, consecutive white spaces, etc.)
- Converted the column values to Uppercase.
- We used Regex on GREL to remove the special characters.
- We used clustering to group the data values into multiple groups. The below methods were used to cluster the data (Key collision and Nearest Neighbor).

◆ **place**

- Removed the white spaces (trailing spaces, consecutive white spaces, etc.)
- Converted the column values to Uppercase.
- We used Regex on GREL to remove the special characters.

- We used clustering to group the data values into multiple groups. The below methods were used to cluster the data (Key collision and Nearest Neighbor).

#### ◆ **physical\_description**

- Removed the white spaces (trailing spaces, consecutive white spaces, etc.)
- Converted the column values to Uppercase.
- We used Regex on GREL to remove the special characters.

#### ◆ **Occasion**

- Removed the white spaces (trailing spaces, consecutive white spaces, etc.)
- Converted the column values to Uppercase.
- We used Regex on GREL to remove the special characters.
- We used clustering to group the data values into multiple groups. The below methods were used to cluster the data (Key collision and Nearest Neighbor).

#### ◆ **Notes**

- Removed the white spaces (trailing spaces, consecutive white spaces, etc.)
- Converted the column values to Uppercase.
- We used Regex on GREL to remove the special characters.

#### ◆ **Call\_number**

- Converted the column values to number

#### ◆ **keywords,language,location\_type**

- Removed the columns

#### ◆ **date**

- Converted the column values to date

#### ◆ **location**

- Removed the white spaces (trailing spaces, consecutive white spaces, etc.)
- Converted the column values to Uppercase.
- We used Regex on GREL to remove the special characters.

#### ◆ **currency**

- Removed the white spaces (trailing spaces, consecutive white spaces, etc.)
- Converted the column values to Uppercase.

#### ◆ **currency\_symbol**

- Removed the white spaces (trailing spaces, consecutive white spaces, etc.)
- Converted the column values to Uppercase.

◆ **status**

- Removed the white spaces (trailing spaces, consecutive white spaces, etc.)
- Converted the column values to Uppercase.

◆ **page\_count**

- Converted the column values to number

◆ **dish\_count**

- Converted the column values to date

**Dish.csv (Total number of records = 423397)**

Data was mostly clean. Most columns were numeric and some cleanup has been done. The name column has cleanup in terms of clustering of data and cleanup of special characters.

→ Id

- ◆ Transformed 423397 cells in column id to Number

→ Menus\_appeared

- ◆ Transformed 423397 cells in column menus\_appeared to Number

→ Times\_appeared

- ◆ Transformed 423397 cells in column times\_appeared to Number

→ First\_appeared

- ◆ Transformed 423397 cells in column first\_appeared to Number

→ Last\_appeared

- ◆ Transformed 423397 cells in column last\_appeared to Number

→ Lowest\_price

- ◆ Transformed 394297 cells in column lowest\_price to Number

→ Highest\_price

- ◆ Transformed 394297 cells in column highest\_price to Number

→ Name

- ◆ Trimmed 9045 cells in column name to remove leading and trailing spaces

→ Trimmed 6415 cells in column name to remove consecutive white spaces

- ◆ Transformed 406726 cells in column name to have Upper Case

→ Created a column facet\_name\_facet and applied RegEx to update 423397 rows with below:

- ◆ `value.replace('[', ' ')`  
`.replace(']', ' ')`  
`.replace('{', ' ')`  
`.replace('}', ' ')`  
`.replace('"', ' ')`  
`.replace(/[,:;]$/, ' ')`  
`.replace('?', ' ')`  
`.replace(/\//, ' ')`  
`.replace('*', ' ')`

- ◆ `.replace('','')`
- ◆ Trimmed 8550 cells in column name to remove leading and trailing spaces
- ◆ Trimmed 11302 cells in column name to remove consecutive white spaces
- Did iterative clustering with name\_facet column using "key Collision/fingerprint" to remove the clusters till there were none left. Below are the results in multiple iterations:
  - ◆ Mass edit 8284 cells in column name\_facet
  - ◆ Mass edit 7666 cells in column name\_facet
  - ◆ Mass edit 7152 cells in column name\_facet
  - ◆ Mass edit 6848 cells in column name\_facet
  - ◆ Mass edit 6741 cells in column name\_facet
  - ◆ Mass edit 6343 cells in column name\_facet
  - ◆ Mass edit 5959 cells in column name\_facet
  - ◆ Mass edit 5957 cells in column name\_facet
  - ◆ Mass edit 5993 cells in column name\_facet
  - ◆ Mass edit 5984 cells in column name\_facet
  - ◆ Mass edit 5995 cells in column name\_facet
  - ◆ Mass edit 5983 cells in column name\_facet
  - ◆ Mass edit 5671 cells in column name\_facet
- Removed the original column name
- Renamed the column "name\_facet" to "name"

#### **Menuitem.csv(Total number of records = 104,586)**

- Data was mostly clean. Most of the columns were numeric so not much cleaning was needed. Some generic transformations were done.
  - ◆ **id**
    - Converted the column values to number
  - ◆ **menu\_page\_id**
    - Converted the column values to number
  - ◆ **price**
    - Converted the column values to number
  - ◆ **high\_price**
    - Converted the column values to number
  - ◆ **dish\_id**
    - Converted the column values to number
  - ◆ **created\_at**
    - Converted the column values to Date
  - ◆ **updated\_at**
    - Converted the column values to Date

◆ **xpos**

- Converted the column values to number

◆ **ypos**

- Converted the column values to number

**MenuPage.csv(Total number of records = 66,938)**

→ Data was mostly clean. Most of the columns were numeric so much cleaning was needed.

- **id**
  - Converted the column values to number
- **menu\_id**
  - Converted the column values to number
- **page\_number**
  - Converted the column values to number
- **image\_id**
  - Converted the column values to number
- **full\_height**
  - Converted the column values to number
- **full\_width**
  - Converted the column values to number
- **uuid**
  - Removed the white spaces (trailing spaces, consecutive white spaces, etc.)
  - Converted the column values to Uppercase.
  - We used Regex on GREL to remove the special characters.

## 2) Using Python

We used the **pandas** library to primarily clean the Menu.csv. We used regex on the following columns using pandas dataframe -

- a. name
- b. sponsor
- c. venue
- d. event
- e. place
- f. Location

Also, we fixed the date column using Pandas and then we filtered out the future dates.

One challenge we faced was that there was no one size fits all regex string to work with. We had to write different regexes to remove different sets of characters. Here is a code snippet of cleaning the text columns using regex.

```
# Removing white trailing spaces
df[column_name] = df[column_name].str.strip()

# Removing double quotes
df[column_name] = df[column_name].str.replace('"', '', regex=True)
df[column_name] = df[column_name].str.replace('"""', '', regex=True)

# Removing Square brackets
df[column_name] = df[column_name].str.replace('[', '', regex=True)

df[column_name] = df[column_name].str.replace(']', '', regex=True)

# cleaning text (?)
df[column_name] = df[column_name].str.replace('\(?\)', '', regex=True)

# cleaning text (;$)
df[column_name] = df[column_name].str.replace('\;$', '', regex=True)

# cleaning text ?
df[column_name] = df[column_name].str.replace('?', '', regex=True)

# cleaning text ;
df[column_name] = df[column_name].str.replace(';', '', regex=True)

# cleaning text ()
df[column_name] = df[column_name].str.replace('\(\)', '', regex=True)

# cleaning text (
df[column_name] = df[column_name].str.replace('\(', '', regex=True)

# cleaning text )
df[column_name] = df[column_name].str.replace('\)', '', regex=True)
```



### 3) Using SQLite3

After the Openrefine and Python cleaning is completed, we loaded the data into a SQLite3 database with 4 tables. In the section [ER Diagram](#), the relationship diagram captures how the tables are connected via the foreign key relationship.

Below are the summary of changes we did -

- Table name menu
  - Verify null value in id column.
    - Result - No rows were found
  - Check duplicate ids in the id column.
    - Result - No rows were found
  - check count of menu rows where name is null and place is null.
    - Result - 6202 rows were found
    - Action - We deleted such rows.
  - check how many event rows have blank or null values
    - Result - 3387 rows were retrieved.
    - Action - We updated such rows with the value event = Unknown
  - Check how many name values are null or blank
    - Result - 7558 rows were retrieved.
    - Action - We updated such rows with the value name = Unknown
  - Check how many place values are null or blank
    - Result - 3195 rows were retrieved.
    - Action - We updated such place rows with the value = Unknown
  - Check how many location values are null or blank
    - Result - 27 rows were retrieved.
    - Action - We updated such location rows with the value = Unknown
- Table menu\_item
  - Verify null value in id column.
    - Result - No rows were found
  - Check duplicate ids in the id column.
    - Result - No rows were found
  - Check if there are any Integrity Constraint violations with the dish table.
    - Result - 3 rows were found
    - Action - We deleted such rows.
  - Check if there are any Integrity Constraint violations with the menu page table.
    - Result - 0 rows were retrieved.

- Table menu\_pages
  - Verify null value in id column.
    - Result - No rows were found
  - Check duplicate ids in the id column.
    - Result - No rows were found
  - We checked for empty row values for all columns and updated them to 0.  
Below are the columns -
    - page\_number = 1202 rows updated with 0.
    - full\_height = 329 rows were updated with 0.
    - full\_width = 329 rows were updated with 0.
  - We also checked for negative values in all numeric columns.
    - No rows found.
  
- Table dish
  - Verify null value in id column.
    - Result - No rows were found
  - Check duplicate ids in the id column.
    - Result - No rows were found
  - We searched for negative values in times\_appeared column-
    - Result - 21 rows were retrieved.
    - Action - Updated them to 0.
  - Checked for the first\_appeared and last\_appeared columns where values are less than 100 and greater than 2021.
    - Result - 55503 rows found.
    - Action - updated those values to a blank value.
  - Checked for blank values in the lowest\_price and highest\_price values
    - Result - 29100 rows found.
    - Action - updated those values to 0.

## Cleaned records report

### Summary table changes

Below is the report of the number of records updated.

Table Name	Column Name	Change log
Menu.csv	place	<ul style="list-style-type: none"><li>Removed white spaces</li><li>Removed unwanted commas, extra characters like [, {, }, ] using regex.</li><li>Removed clustered the similar names into one using openrefine.</li></ul>
	id	<ul style="list-style-type: none"><li>Text transform on 17545 cells in column id: value.toNumber()</li></ul>
	name	<ul style="list-style-type: none"><li>Text transform on 9 cells in column name: value.trim()</li><li>Text transform on 1 cells in column name: value.replace(/\s+/, '')</li><li>Text transform on 3196 cells in column name: value.toUpperCase()</li><li>Create new column name_facet based on column name by filling 3197 rows with grel:value</li></ul>
	sponsor	<ul style="list-style-type: none"><li>Text transform on 14 cells in column sponsor: value.trim()</li><li>Text transform on 127 cells in column sponsor: value.replace(/\s+/, '')</li><li>Text transform on 8535 cells in column sponsor: value.toUpperCase()</li><li>Create new column sponsor_facet based on column sponsor by filling 15984 rows with grel:value</li><li>Remove column sponsor</li><li>Rename column sponsor_facet to sponsor</li></ul>
	venue	<ul style="list-style-type: none"><li>Text transform on 0 cells in column venue: value.trim()</li><li>Text transform on 0 cells in column venue: value.replace(/\s+/, '')</li><li>Text transform on 9 cells in column venue: value.toUpperCase()</li><li>Create new column venue_facet based on column venue by filling 8131 rows with grel:value</li><li>Mass edit 2251 cells in column venue_facet</li><li>Remove column venue</li></ul>

		<ul style="list-style-type: none"> <li>Rename column venue_facet to venue</li> </ul>
	physical_description	<ul style="list-style-type: none"> <li>Text transform on 384 cells in column physical_description: value.trim()</li> <li>Text transform on 38 cells in column physical_description: value.replace(/s+/, ' ')</li> <li>Text transform on 7371 cells in column physical_description: value.toUpperCase()</li> <li>Create new column physical_description_facet based on column physical_description by filling 14768 rows with grel:value</li> </ul>
	notes	<ul style="list-style-type: none"> <li>Text transform on 125 cells in column notes: value.trim()</li> <li>Text transform on 195 cells in column notes: value.replace(/s+/, ' ')</li> <li>Text transform on 3873 cells in column notes: value.toUpperCase()</li> <li>Create new column notes_facet based on column notes by filling 10613 rows with grel:value</li> </ul>
	call_number	<ul style="list-style-type: none"> <li>Text transform on 1 cells in column call_number: value.toNumber()</li> </ul>
	keywords	<ul style="list-style-type: none"> <li>Text transform on 0 cells in column keywords: value.trim()</li> <li>Text transform on 0 cells in column keywords: value.replace(/s+/, ' ')</li> <li>Text transform on 0 cells in column keywords: value.toUpperCase()</li> </ul>
	language	<ul style="list-style-type: none"> <li>Text transform on 0 cells in column language: value.trim()</li> <li>Text transform on 0 cells in column language: value.replace(/s+/, ' ')</li> <li>Text transform on 0 cells in column language: value.toUpperCase()</li> </ul>
	date	<ul style="list-style-type: none"> <li>Text transform on 16959 cells in column date: value.toDate()</li> </ul>
	location	<ul style="list-style-type: none"> <li>Text transform on 14 cells in column location: value.trim()</li> <li>Text transform on 555 cells in column location: value.replace(/s+/, ' ')</li> <li>Text transform on 17434 cells in column location: value.toUpperCase()</li> </ul>

		<ul style="list-style-type: none"> <li>• Create new column location_facet based on column location by filling 17545 rows with grel:value</li> </ul>
	location_type	<ul style="list-style-type: none"> <li>• Remove column location_type</li> </ul>
	currency	<ul style="list-style-type: none"> <li>• Text transform on 0 cells in column currency: value.trim()</li> <li>• Text transform on 0 cells in column currency: value.replace(/\s+/, ' ')</li> <li>• Text transform on 6456 cells in column currency: value.toUpperCase()</li> <li>• Create new column currency_facet based on column currency by filling 6456 rows with grel:value</li> <li>• Remove column currency</li> <li>• Rename column currency_facet to currency</li> </ul>
	currency_symbol	<ul style="list-style-type: none"> <li>• Text transform on 4 cells in column currency_symbol: value.trim()</li> <li>• Text transform on 0 cells in column currency_symbol: value.replace(/\s+/, ' ')</li> <li>• Text transform on 169 cells in column currency_symbol: value.toUpperCase()</li> </ul>
	status	<ul style="list-style-type: none"> <li>• Text transform on 0 cells in column status: value.trim()</li> <li>• Text transform on 0 cells in column status: value.replace(/\s+/, ' ')</li> <li>• Text transform on 17545 cells in column status: value.toUpperCase()</li> </ul>
	page_count	<ul style="list-style-type: none"> <li>• Text transform on 17545 cells in column page_count: value.toNumber()</li> </ul>
	dish_count	<ul style="list-style-type: none"> <li>• Text transform on 17545 cells in column dish_count: value.toNumber()</li> </ul>
	occasion	<ul style="list-style-type: none"> <li>• Text transform on 0 cells in column occasion: value.trim()</li> <li>• Text transform on 3 cells in column occasion: value.replace(/\s+/, ' ')</li> <li>• Text transform on 17 cells in column occasion: value.toUpperCase()</li> <li>• Create new column occasion_facet based on column occasion by filling 3803 rows with grel:value</li> <li>• Mass edit 2779 cells in column occasion_facet</li> <li>• Mass edit 608 cells in column occasion_facet</li> </ul>

		<ul style="list-style-type: none"> <li>• Mass edit 8 cells in column occasion_facet</li> <li>• Remove column occasion</li> <li>• Rename column occasion_facet to occasion</li> </ul>
	event	<ul style="list-style-type: none"> <li>• Text transform on 3 cells in column event: value.trim()</li> <li>• Text transform on 6 cells in column event: value.replace(/\s+/, ' ')</li> <li>• Text transform on 729 cells in column event: value.toUpperCase()</li> <li>• Mass edit 5314 cells in column event - Cluster using "Key Collision"</li> <li>• Text transform on 0 cells in column event: value.trim()</li> <li>• Text transform on 0 cells in column event: value.replace(/\s+/, ' ')</li> <li>• Create new column event_facet based on column event by 8154 rows with grel:value: grel:value.replace('[', ' ')\n.replace(']', ' ')\n.replace('{', ' ')\n.replace('}', ' ')\n.replace('"', ' ')\n.replace(/[,;]\$/, ' ')\n.replace('?', ' ')\n.replace(/\ /, ' ')\n.replace(')', ' ')\n.replace('(', ' ')\n.replace('*', ' ')\n.replace('\\', ' ')</li> <li>• Text transform on 124 cells in column event_facet: value.trim()</li> <li>• Text transform on 89 cells in column event_facet : value.replace(/\s+/, ' ')</li> <li>• Text transform on 0 cells in column event: value.toUpperCase()</li> <li>• Remove column event</li> <li>• Rename column event_facet to event</li> </ul>
Dish.csv	id	Transformed 423397 cells in column id to Number
	menus_appeared	Transformed 423397 cells in column menus_appeared to Number
	times_appeared	Transformed 423397 cells in column times_appeared to Number
	first_appeared	Transformed 423397 cells in column first_appeared to Number
	last_appeared	Transformed 423397 cells in column last_appeared to Number
	lowest_price	Transformed 394297 cells in column lowest_price to Number

	highest_price	Transformed 394297 cells in column highest_price to Number
	name	<ul style="list-style-type: none"> <li>Trimmed 9045 cells in column name to remove leading and trailing spaces</li> <li>Trimmed 6415 cells in column name to remove consecutive white spaces</li> <li>Transformed 406726 cells in column name to have Upper Case</li> <li>Created a column facet name_facet and applied RegEx to update 423397 rows with below:</li> <li>value.replace('[', ' ' )          .replace(']', ' ' )          .replace('{', ' ' )          .replace('}', ' ' )          .replace('"' , ' ' )          .replace(/[,;]\$/, ' ' )          .replace('?', ' ' )          .replace(/\//, ' ' )          .replace('*', ' ' )          .replace('"' , ' ' )</li> <li>Trimmed 8550 cells in column name to remove leading and trailing spaces</li> <li>Trimmed 11302 cells in column name to remove consecutive white spaces</li> <li>Did iterative clustering with name_facet column using "key Collision/fingerprint" to remove the clusters till there were none left. Below are the results in multiple iterations:             <ul style="list-style-type: none"> <li>Mass edit 8284 cells in column name_facet</li> <li>Mass edit 7666 cells in column name_facet</li> <li>Mass edit 7152 cells in column name_facet</li> <li>Mass edit 6848 cells in column name_facet</li> <li>Mass edit 6741 cells in column name_facet</li> <li>Mass edit 6343 cells in column name_facet</li> <li>Mass edit 5959 cells in column name_facet</li> <li>Mass edit 5957 cells in column name_facet</li> <li>Mass edit 5993 cells in column name_facet</li> <li>Mass edit 5984 cells in column name_facet</li> <li>Mass edit 5995 cells in column name_facet</li> <li>Mass edit 5983 cells in column name_facet</li> <li>Mass edit 5671 cells in column name_facet</li> </ul> </li> <li>Removed the original column name</li> <li>Renamed the column "name_facet" to "name"</li> </ul>
MenuItem.csv		<ul style="list-style-type: none"> <li>Not much to clean.</li> <li>Most of the data was clean and columns were mostly numeric.</li> <li>Text transform on 1332726 cells in column menu_page_id: value.toNumber()</li> </ul>

		<ul style="list-style-type: none"> <li>• Text transform on 886810 cells in column price: value.toNumber()</li> <li>• Text transform on 91905 cells in column high_price: value.toNumber()</li> <li>• Text transform on 1332485 cells in column dish_id: value.toNumber()</li> <li>• Text transform on 0 cells in column created_at: value.toDate()</li> <li>• Text transform on 1332726 cells in column xpos: value.toNumber()</li> </ul>
MenuPage.csv	id	<ul style="list-style-type: none"> <li>• Text transform on 66937 cells in column id: value.toNumber()</li> </ul>
	menu_id	<ul style="list-style-type: none"> <li>• Text transform on 66937 cells in column menu_id: value.toNumber()</li> </ul>
	page_number	<ul style="list-style-type: none"> <li>• Text transform on 65735 cells in column page_number: value.toNumber()</li> </ul>
	image_id	<ul style="list-style-type: none"> <li>• Text transform on 66914 cells in column image_id: value.toNumber()</li> </ul>
	full_height	<ul style="list-style-type: none"> <li>• Text transform on 66608 cells in column full_height: value.toNumber()</li> </ul>
	full_width	<ul style="list-style-type: none"> <li>• Text transform on 66608 cells in column full_width: value.toNumber()</li> </ul>
	uuid	<ul style="list-style-type: none"> <li>• Text transform on 66936 cells in column uuid: value.toUpperCase()</li> <li>• Text transform on 0 cells in column uuid: value.trim()</li> <li>• Text transform on 0 cells in column uuid: value.replace(/s+/, '')</li> </ul>



## Rationale For design

Each cleaning step above played an important role for the U1 use case. Let's revisit the U1

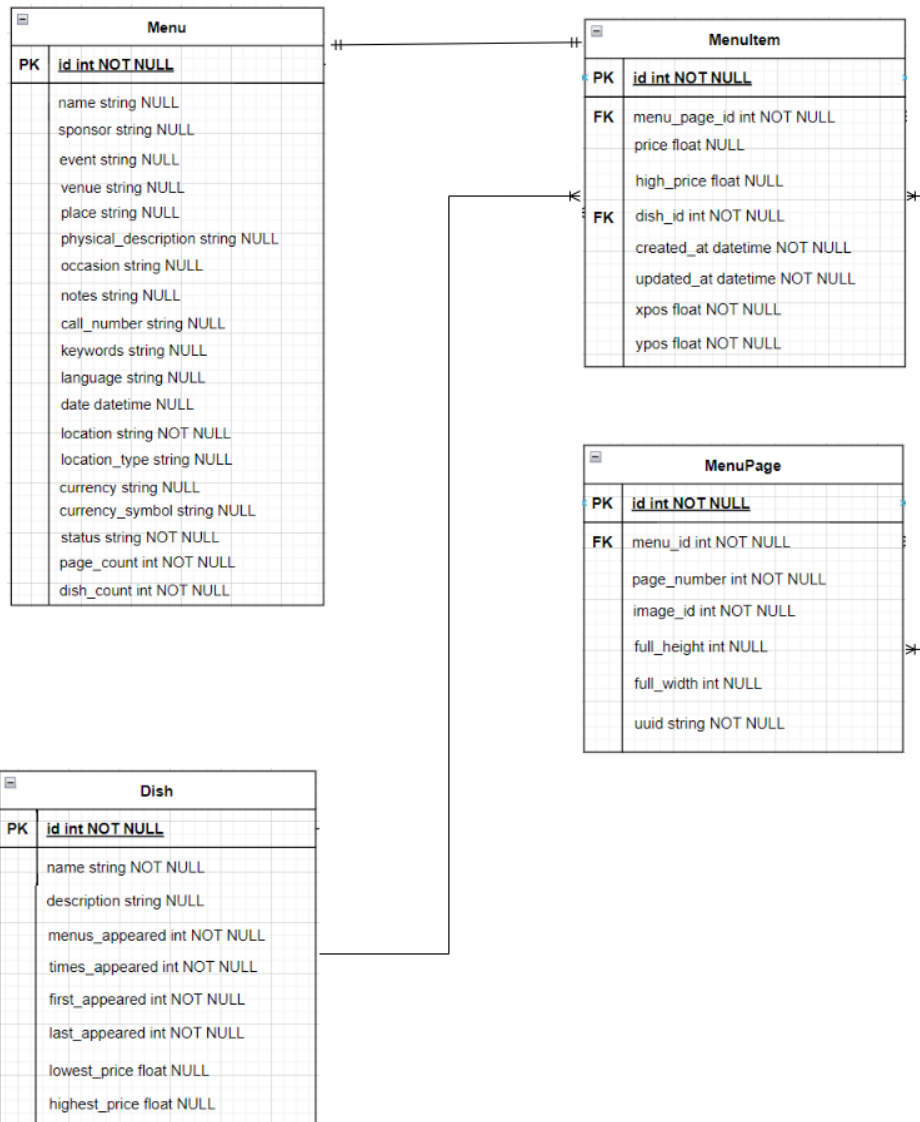
- *“What were the menu items which were served for Dinner at New York between any certain time-frame(ex - 1950 - 1960)? And how menu items were priced and what were the dishes served during that event? What were the locations where these dishes were served?”*

The data cleaning processes applied were relevant because -

1. We used clustering to clean up the name, location, places, events column to serve the purpose of the above question.
2. We removed white spaces to make the data consistent across columns.
3. We removed the unwanted special characters to clean up the text values.
4. We used the data type conversion to convert columns (date) so that date queries can be performed easily.
5. The Sqlite3 database helped us remove the integrity constraints violations because then we were able to join two tables without worrying about missing or unwanted data. For example, we cleaned the Menu\_item table by checking if any dish\_id exists which are not valid. See [here](#).

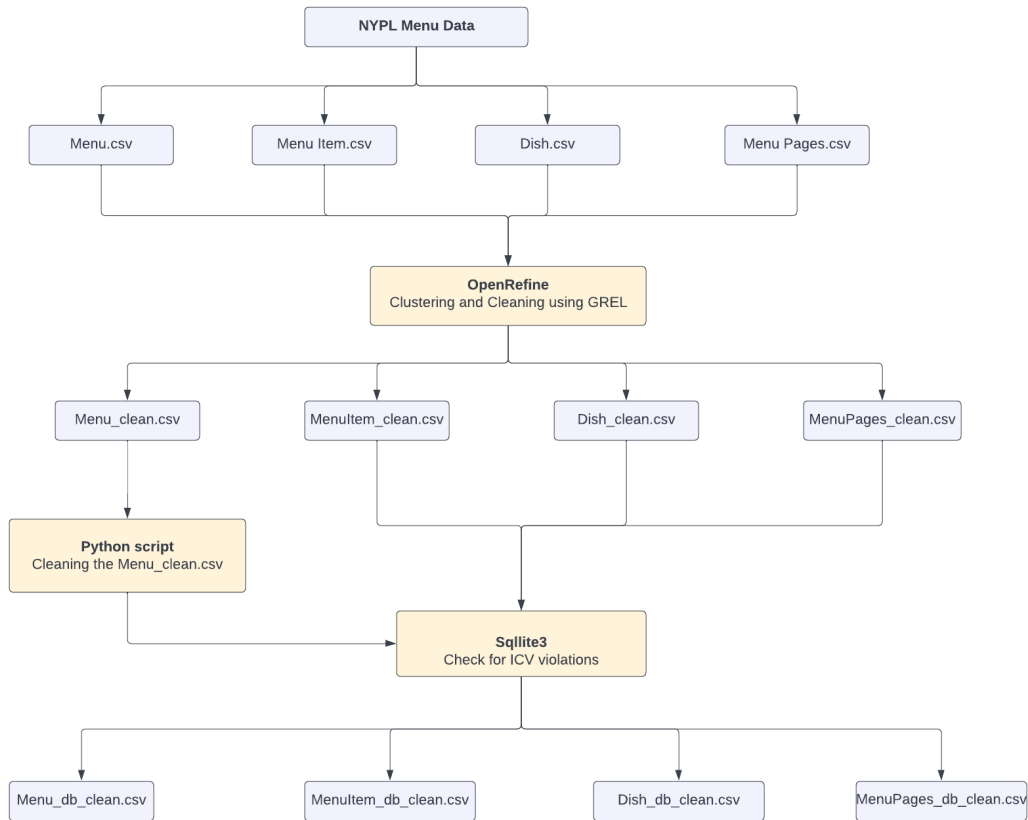
## ER Diagram

Here is the ER diagram from the SQLite3 database. Available under the folder **/sqlit3/Docs & Images/**



## Overview Of the cleaning process

Below is a diagram of how the overview of the cleaning process looks like



## Workflow Diagrams

Once all the csv files were cleaned we used the YesWorkflow to generate the diagrams.

We have provided the YesWorkflow diagrams for the Openrefine logs and python scripts.

A Python based [OR2YWTool](#) application was used to generate the YW and PDF files for the Workflow.

Commands used for YesWorkflow artifact generation:

- set PYTHONUTF8=1 (before each of below step to account for special UTF-8 characters)
- or2yw -i menu.json -o Serial\_menu.yw
- or2yw -i menu.json -o Parellel\_menu.yw -t parallel
- or2yw -i menu.json -o menu.pdf -ot pdf

We also used LucidChart and [Draw.io](#) to generate the above charts.

We will provide the paths of those diagrams belows.

## Summary / Lesson Learned

We enjoyed working on this project a lot and data cleaning is fun to be a part of. We ran the query on the dirty dataset and on the cleaned dataset.

Here is the query -

```
SELECT distinct d.name, m.event, m.date, m.place, m.location, mi.price, mp.image_id
FROM dish d, menu_item mi, menu m, menu_page mp
WHERE d.id = mi.dish_id
AND mi.menu_page_id = mp.id
AND mp.menu_id = m.id
AND m.place = 'NEW YORK, NY'
AND m.date > '1900-01-01'
order by date desc;
```

Here are the results -

**Before** -We retrieved 28,302 rows from the dataset.

**After** - Once the cleaning was done, 53,946 rows.

As we cleaned the dataset, all places named “New York, NY “ were cleaned and we replaced different places with this above string.

We also learned a few things along the way -

1. Openrefine gets very slow when we load more than 100K+ records, which is not too many in industry standard. I would prefer to use Pandas at that point. However, we then have to write our own code.
2. Openrefine is good for clustering using different Algorithms. But I prefer using the nearest neighbor algorithm much better.
3. Trying to figure out a common Regex for all columns was very difficult. There were so many unwanted characters in so many different forms. For example, some columns had brackets like this “(…)” and some had “[...]”.
4. There were so many erroneous values while dealing with Date columns, e.g. - Future dates.
5. We learnt that cleaning improves the data quality. This was observed by the different number of rows retrieved for the same query on querying the uncleaned dataset and cleaned dataset.

## Acknowledgements

- Our sincere thanks to Prof. Bertram Ludaescher and all the TAs in CS 513 - Theory and Practice Data Cleaning course for making the course engaging and helping us with all the queries.
- Also special thanks to the open source community of Python to show us endless Regex examples.
- Thanks to our project team members for working late night despite being in different Timezones and making this project successful.
- A heartfelt appreciation to University of Illinois - Urbana Champaign for providing students with endless tools to become successful in their courses. Here are some of the tool examples such as box, sharepoint, Google Drive, zoom, slack, CampusWire and many more.

## User Stories and Contributions

### 1. Epic - Cleaning the files using OpenRefine

- a. Sanitize Menu.csv - Kavithaa Suresh Kumar
- b. Sanitize Dish.csv - Sudipto Sarkar
- c. Sanitize MenuItem.csv - Arnab Kar Sarkar
- d. Sanitize MenuPages.csv - Kavithaa Suresh Kumar

### 2. Epic - Cleaning the files using PythonScript

- a. Clean all text columns using Pandas - Arnab Kar Sarkar

### 3. Epic - SQLite3

- a. Sanitize Menu.csv - Kavithaa Suresh Kumar
- b. Sanitize Dish.csv - Sudipto Sarkar
- c. Sanitize MenuItem.csv - Arnab Kar Sarkar
- d. Sanitize MenuPages.csv - Kavithaa Suresh Kumar

### 4. Epic - Documentation and Github setup

- a. Write different sections - Kavithaa Suresh Kumar , Arnab Kar Sarkar  
Sudipto Sarkar

### 5. Epic - Yes Workflow and Other Diagrams

- a. OpenRefine - Sudipto Sarkar

- b. FlowChart - Arnab Kar Sarkar
- c. ER Diagram - Kavithaa Suresh Kumar

## Project Team members

Name	NetId	Email
Kavithaa Suresh Kumar	ks64	ks64@illinois.edu
Arnab KarSarkar	arnabk2	arnabk2@illinois.edu
Sudipto Sarkar	sudipto2	sudipto2@illinois.edu

## Github Details

Github link - <https://github.com/arnabksarkar/CS513-DataCleaning-Project>

## FilePath Details

- OpenRefine details
  - Logs - **/OpenRefine/Logs/json**
  - YesWorkflow docs - **/OpenRefine/YWDocs**
- Python details
  - Code - **/Python/CodeFiles**
  - YesWorkflow - **/Python/YWDocs**
- SQLite3 details
  - Queries - **/sqlite3/Queries.sql**
  - SQL Workflow - **/sqlite3/Docs & Images**
- Documentations
  - Final Project PDF document - **/Docs/Final Project**
  - First Draft PDF Document - **/Docs/First Draft**