

# Final Project

## End-To-End Data Cleaning Workflow – NYPL Data

by

Sudhir Behani (sbehani2@illinois.edu)

Bollam Raja Shekar

## Table of Contents

<b>Introduction and Overview .....</b>	<b>3</b>
1. Dish.csv .....	3
2. Menu.csv.....	3
3. MenuItem.csv .....	3
4. MenuPage.csv .....	3
5. Missing Data .....	3
Other Information .....	3
<b>Initial Assessment of the dataset .....</b>	<b>4</b>
<b>Data Cleaning methods and Process .....</b>	<b>5</b>
Open Refine .....	5
<b>Data Cleaning Results .....</b>	<b>8</b>
Relation Database Schema .....	8
ER Diagrams .....	8
Integrity Constraint Checks .....	9
Workflow Model.....	9
Or2yw overall workflow .....	9
Or2yw open Refine csv .....	9
Overview of changes.....	10
<b>Conclusion and Future Work .....</b>	<b>11</b>

## Introduction and Overview

### 1. Dish.csv

Information about all the dishes from all the menus transcribed by the project are stored in data.csv. In this data file, a dish is represented by a row of values. Columns identify attributes of a dish. One of these attributes is an identifier, which identifies the dish. However, the identity of a dish appears to be based on the exact form of the string labeled "name." Thus, dishes with variant orthographic forms of their names, e.g. "half chicken", "Half Chicken" and "chicken [half]") are treated as separate entries with different identifiers.

### 2. Menu.csv

Information about the menus as physical objects, including historical information about their origins, uses, and formats are stored in menu.csv

### 3. MenuItem.csv

This is the largest data file. A "MenuItem" represents a single instance of a dish appearing somewhere on a menu page image. "MenuItem.csv" is useful as a mapping between multiple other data files/tables.

### 4. MenuPage.csv

Information about individual pages of the menus is stored in "Menu.csv". Pages are modeled here as digital images produced as a result of digitization by the NYPL. Menus often have multiple pages.

### 5. Missing Data

We note the presence of missing data points in the column "missing values." However, there are also strings in the present data that point to missing information (e.g. "unknown" or "?"), which do not formally appear as null values.

### Other Information

The data in these files comes from several different sources. We reflect this in the "generated by" category of the tables below. Some of the data is supplied by "volunteer transcribers," by which we mean people who have participated in the project through the What's On the Menu? site. Some of the information is generated by "web application." This means that some of this information was automatically created as the database supporting the application was constructed and populated (e.g., various ids); some is created as the web application runs (e.g. timestamps as data values are updated). Finally, a lot of information is generated from "NYPL metadata." This metadata comes from many places and reflects the long history of the project and the many parts of New York Public Library involved in it. Much of the data "supplied by NYPL metadata" in the menu spreadsheet is from the catalog cards made by Frank E. Buttolph in the early twentieth century.

## **Initial Assessment of the dataset**

# Data Cleaning methods and Process

## Open Refine Tasks

### 1. Menu

Activities performed:

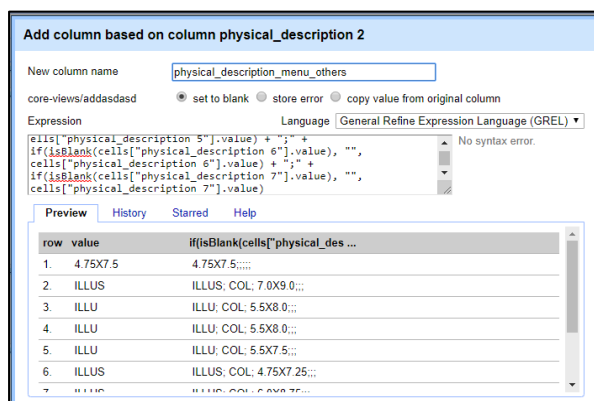
- Trim leading and trailing white spaces
- Collapse consecutive white spaces
- Convert columns to number type
- Convert columns date to format YYYY-MM-DD
- Remove special characters “()[]{} \?#!%” from the text columns. Also, remove characters at the end of the string like semicolon, comma and hyphen “,-”

Note: Can't remove (semicolon) character in the middle of string as it signifies

multiple options- example multiple sponsors, event, venue, place etc

- Remove continuous hyphen characters from the text columns
- Convert all values in columns with string data type to upper case so that problem in clustering similar values is reduced considerably
- After removing special characters again perform white spaces removal for each the string data type columns
- Clusters columns based on method key collision and keying function as fingerprint. Accepted the default values for new cell value. Merge cluster and then perform ngram-fingerprint using ngram size as 2. Metaphore3 and cologne-phonetic were mostly creating clusters with somewhat different names too so were not used.
- Physical\_description seems to contain information about the menu hard copy format and its size. Menu could be in different forms – accordion fold, book, booklet, broadside, card, folder, tri-folder, two cards joined by ribbons etc

Split the physical\_description column with separator as semicolon. It is split into 7 columns. Rename the first column as “physical\_description\_menu\_type” and join other columns as rename it as “physical\_description\_menu\_others”



- Remove semilon at the end of the string using expression

## Not updated columns:

Columns which are not updated: id, name, keywords, language, location\_type, currency, currency\_symbol

### 1. Dish

Activities performed:

- Trim and Collapse white spaces in “name” column.
- Convert column “name” to title case.
- Make a facet and perform the cluster operation using the “key-collision” method and “fingerprint” as shown below.

**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** 30780 clusters found

Method	Keying Function	Clusters Found
key collision	fingerprint	30780

26 35

- Poached Eggs, 2 On Toast (1 rows)
- [eggs] Poached On Toast (2) (1 rows)
- Imported Ginger Ale, C. & C. (6 rows)
- C & C Imported Ginger Ale (3 rows)
- Ginger Ale, Imported C & C (3 rows)
- "C. & C." Ginger Ale, Imported (1 rows)
- C. & C. Ginger Ale, Imported (1 rows)
- C. & C. Imported Ginger Ale (1 rows)
- Ginger Ale (c. & C.), Imported (1 rows)
- Ginger Ale (imported) C. & C. (1 rows)
- Ginger Ale Imported, C. & C. (1 rows)
- Ginger Ale Imported, C. & C. (1 rows)
- Ginger Ale, C. & C. (imported) (1 rows)
- Ginger Ale, C. & C., Imported (1 rows)
- Ginger Ale, Imported C. & C. (1 rows)
- Ginger Ale, Imported, C & C (1 rows)
- Ginger Ale, Imported, C. & C. (1 rows)
- Imported (c. & C.) Ginger Ale (1 rows)
- Imported C. & C. Ginger Ale (1 rows)
- Imported Ginger Ale (c. & C) (1 rows)
- Imported Ginger Ale (c. & C.) (1 rows)
- Imported Ginger Ale C & C (1 rows)
- Imported Ginger Ale C. & C. (1 rows)
- Imported Ginger Ale, C & C (1 rows)
- Imported Ginger Ale, C. & C. (1 rows)
- Imported Ginger Ale, C. & C. (1 rows)

# Choices in Cluster: 2 — 27

# Rows in Cluster: 2 — 49

Average Length of Choices: 0 — 670

Length Variance of Choices: 0 — 34

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

- Repeat step c with “ngram-fingerprint” as shown below.

**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: **ngram-fingerprint** Ngram Size: 2 6159 clusters found

Method	Keying Function	Ngram Size	Clusters Found
key collision	ngram-fingerprint	2	6159

5 30

- Hors D'oeuvres Varies (16 rows)
- Hors D'oeuvres Varie (5 rows)
- Hors-d'oeuvres Varies (6 rows)
- Hors D' Oeuvres Varies (1 rows)
- Hors D'oeuvres-varies (1 rows)

4 4

- Commodore Coffee With Cream (p. P.) (1 rows)
- Commodore Coffee With Cream (p.P.) (1 rows)
- Commodore Coffee With Cream(p. P.) (1 rows)
- Commodore Coffee With Cream(p.P.) (1 rows)

4 5

- Home Made Sausage (2 rows)
- Home Made Sausages (1 rows)
- Home-made Sausages (1 rows)
- Homemade Sausage (1 rows)

4 10

- Blue Point Oysters On Half Shell (5 rows)
- Bluepoint Oysters On Half Shell (3 rows)
- Blue Point Oysters On Half-shell (1 rows)
- Blue Point Oysters-on Half Shell (1 rows)

4 4

- N.Y.N.H. & H. R. R. Marine Disc. No 1 (1 rows)
- N.Y.N.H. & H.R.R. Marine Disc. No 1 (1 rows)
- N.Y.N.H.&h.R.R. Marine Disc. No. 1 (1 rows)
- N.Y.N.H.&h.R.R.Marine Disc. No 1 (1 rows)

# Choices in Cluster: 2 — 14

# Rows in Cluster: 2 — 48

Average Length of Choices: 0 — 510

Length Variance of Choices: 0 — 12

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

- Remove “description” column, since most of the entries are empty.

## 2. MenuPage

- Convert columns to number type

There are no other columns cleaned

## 3. MenuItem

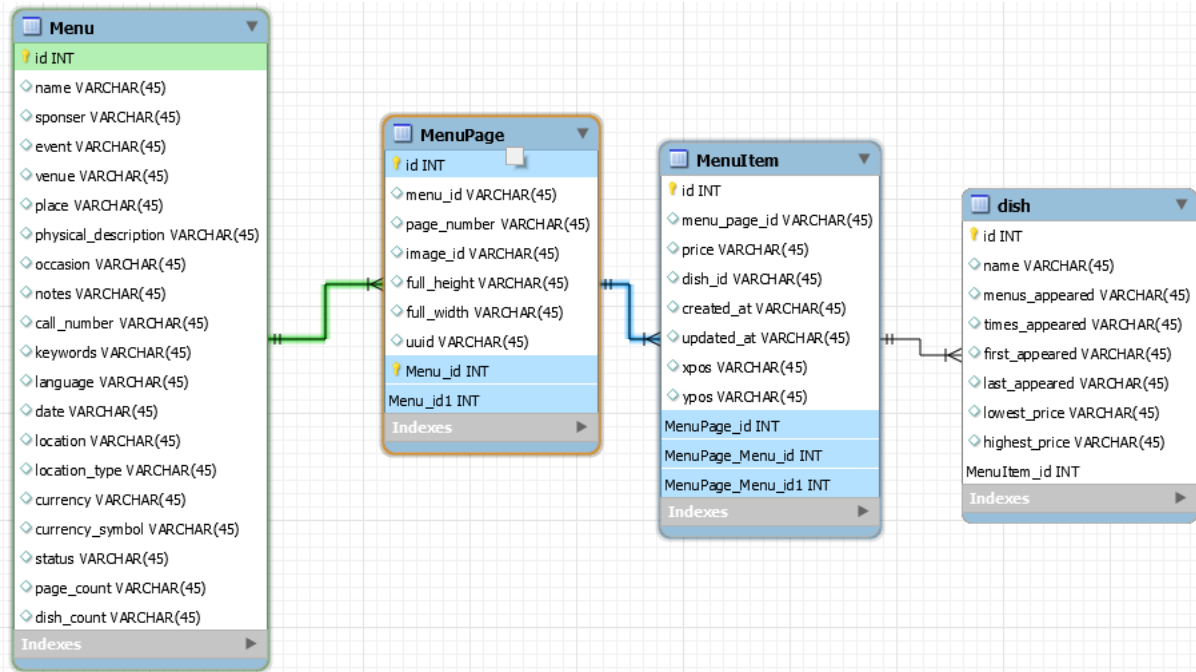
Following data cleaning steps were performed on the Dish dataset.

- "high\_price" column is removed
- "created\_at" is transposed into "toDate()"
- "updated\_at" is transposed into "toDate()"

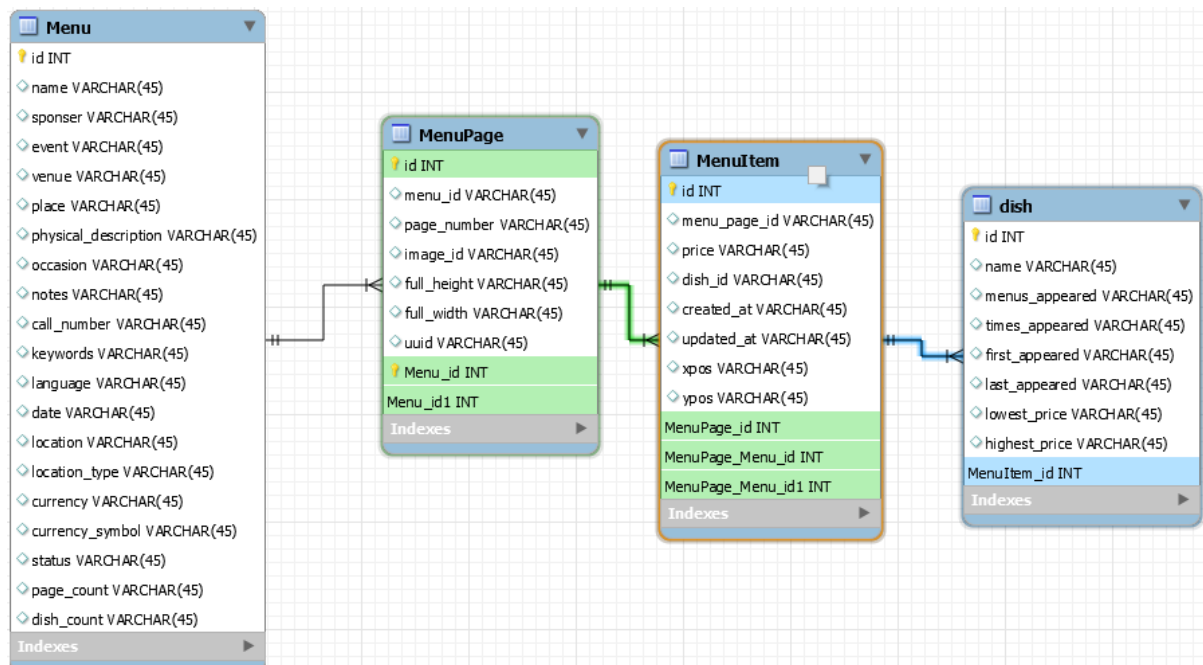
# Data Cleaning Results

## Relation Database Schema

### ER Diagrams



<<<Add text >>>

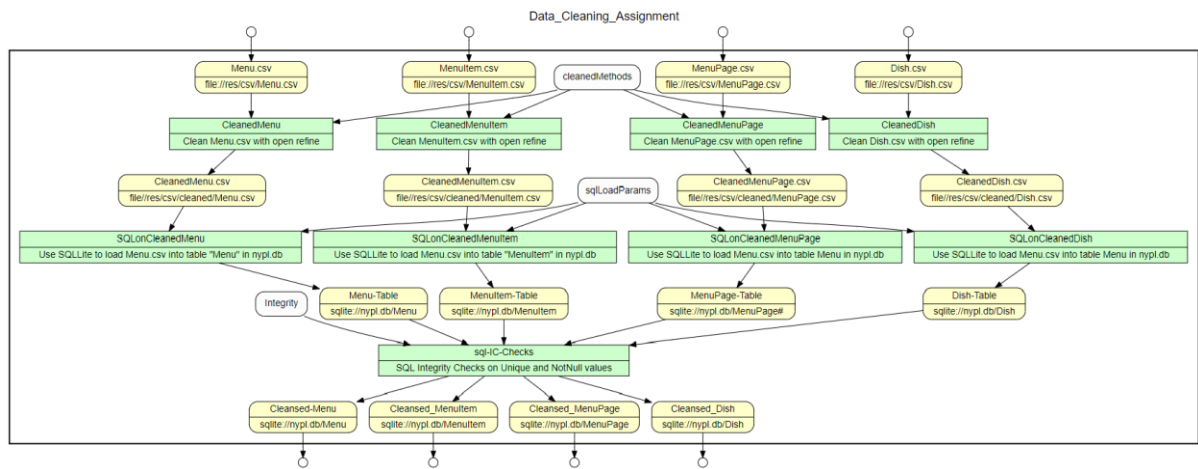




## Integrity Constraint Checks

### Workflow Model

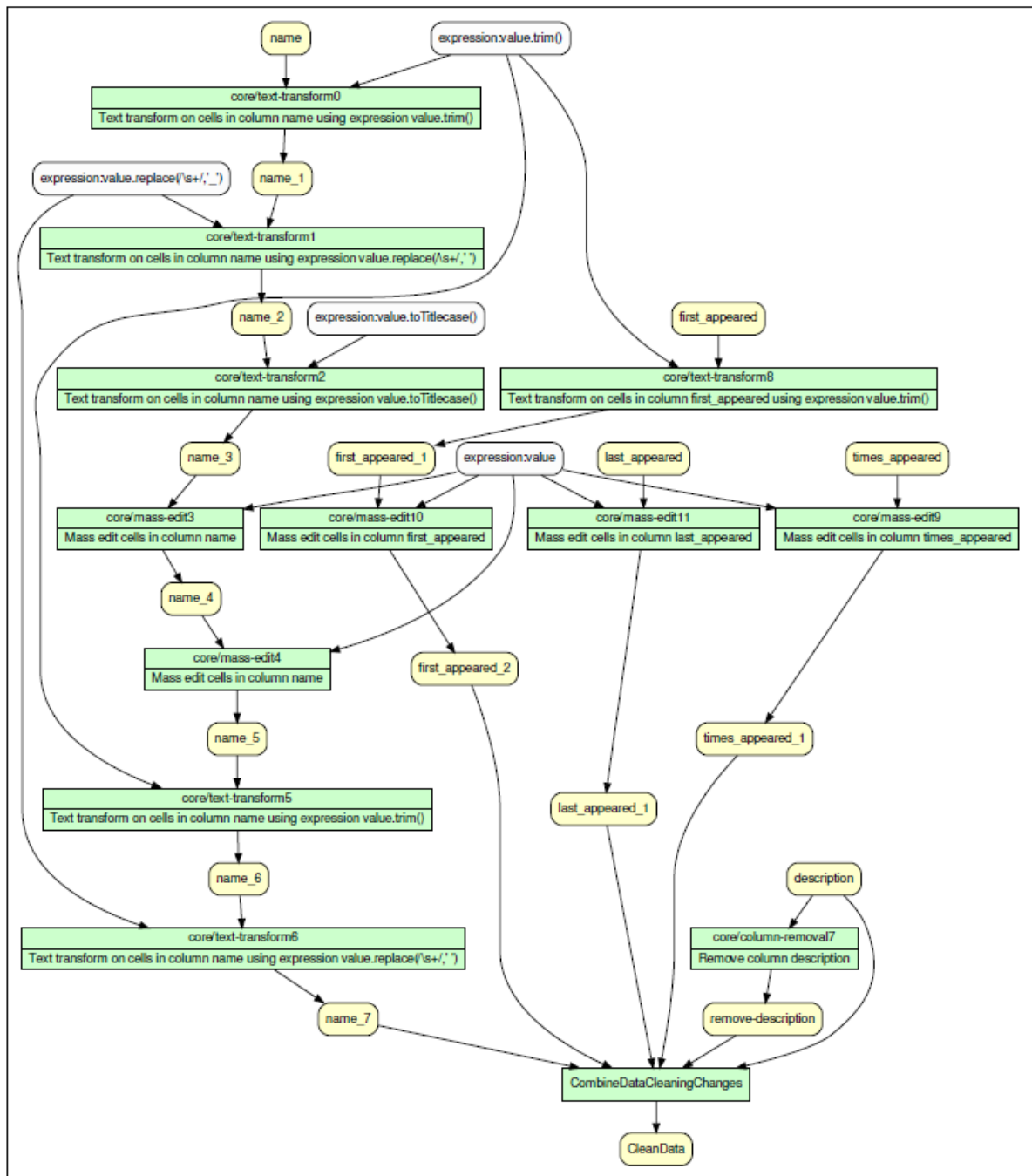
#### Or2yw overall workflow



#### Or2yw open Refine csv

Following graphs are generated using or2yw tool based on the recipes generated from the openrefine too.

1. Dish.csv



## 2. Menu.csv

### Overview of changes

## Conclusion and Future Work

Data cleaning course helped us to understand the importance of having a cleaned dataset. NYPL being a crowd sourced dataset is full of issues and needs good amount of cleaning activities. Majority of the cleaning works were required for string-based columns in the dataset files. As textual information could be written in different ways by different people.

Main outcome of our cleaning project is to create a cleaned dataset. Though this is a huge task in hand given the time and effort required to bring data into a completely consistent state so that this could be useful for further tasks like machine learning or natural language processing. We also tried our best to create a workflow visualization so that its easier for clients to get visual details about the cleaning activities performed.

Important takeaways from this cleaning activity is to first get an understanding about the importance of cleaned dataset. This was understood once we investigated the NYPL dataset and cleaned data was user readable and one could develop better understanding of the information. A simple removal of leading and trailing whitespaces and removal of consecutive white spaces bring so much consistency to the dataset. It is utmost importance to create a cleaned dataset. From NYPL dataset it was also understood that crowd source dataset needs a lot of cleaning activities.

There were number of challenges ....

Problems or Challenges:

Additional steps:

Certain foreign letters which needs to be take care by studying the words in more detail. These foreign characters seems to occur in menu table under columns – sponsor and location.

Future additional tasks:

For Menu we could use status column having two values “Complete” and “Under Review” and perform classification task. We could try to predict which features in menu table led to the value Complete and Under Review.

Appendix A

Link to project deliverables: