

# CS513: Theory & Practice of Data Cleaning

Suruchi Sinha – [suruchi3@illinois.edu](mailto:suruchi3@illinois.edu)

Pradeep Sharma – [pbs6@illinois.edu](mailto:pbs6@illinois.edu)

Upul Gunasena – [ug2@illinois.edu](mailto:ug2@illinois.edu)

# CS513: THEORY & PRACTICE OF DATA CLEANING

**Github:** <https://github.com/upulindika/CS-513-Team75>

**Abstract-** This report describes a data cleaning workflow using the New York Public Library Rare Books Division historical menus dataset has an example to demonstrate various cleaning techniques. This project will use the following opensource software and tools to clean and organize the dataset: OpenRefine, SQLite, and YesWorkflow, Creately for creating ER diagrams

**Index Terms-** Data Cleaning, Provenance, OpenRefine, YesWorkflow, SQLite

## 1. INTRODUCTION

We are going to select the New York public library historical menus dataset. The data set NY menus is a collection of more than 45000 historical menus from the 1850s to the 2000s. Most data were organized by Frank E Buttolph [5] around 1900-1921. The dataset was digitized in 2011 via the “What’s on Menu?” project and so far, 17500 historical menus have been digitized. We will try to clean this data using the tools and the techniques we have learned in the CS513 Theory and Practice of Data cleaning.

## 2. DATA SET INTRODUCTION OVERVIEW AND INITIAL ASSESSMENT

### 2.1 THE DATASET

The New York Public Library Rare Book Division holds over 45,000 historical menus. About half of these were collected and curated by Frank E. Buttolph between 1900 and 1921. The menus date from the 1850s to the present and include menus from the restaurant, railroad, and steamship companies, as well as a range of other organizations. Beginning in 2011, menus from the NYPL’s collection were digitized and transcribed with the help of thousands of volunteers. Through the NYPL’s What’s on the Menu? [4] the project, volunteers looked at digitized copies of the menus and typed in the many pieces of information included on each one, such as restaurant names, locations, dishes, prices, and dates.

The detailed fields information about the fields in the file is described as below:

#### MenuItem.csv

Id: unique id for this menu  
menu\_page\_id: id of the menu page that this item appears.  
price: the price of the smallest amount of the item  
high\_price: the price of the max amount of item  
dish\_id: menu item referring dish id  
created\_at: the date when it was created  
updated\_at: last updated  
xpos: x-axis position of the item on the scanned image  
ypos: y-axis position of the item on the scanned image

#### MenuPage.csv

Id: id of the menu page  
menu\_id: id for menu  
page\_number: page number of menus  
image\_id: unique id of the image  
full\_height: height of the menu  
full\_width: width of menu  
uuid: uniqueid of page/image

### **Dish.csv**

Id: id of dish

name: Name of the dish

description: Description of the dish

menus\_appeared: dish appearing on a number of menus

times\_appeared: number of times this dish appears in additional section also

first\_appeared: year this dish first appeared

last\_appeared: year this dish last appeared

lowest\_price: lowest price of the dish

highest\_price: the highest price of the dish

### **Menu.csv**

Id: id for this menu

name: name on menu/restaurant or blank

sponsor: sponsor is the name of the restaurant:

event: name of the meal or the event the menu was created for

venue: location where the food is served

place: includes city/state/country/address or name of venue

physical\_description: paper stock, dimension, color, design of the menu

occasion: special occasion, holiday, daily or blank

notes: additional details about the menu

call\_number: number within NYPL collection

keywords: keywords on menu

language: language the menu is printed in

date: date menu was collected

location: menu used location

location\_type: type of location

currency: money type charged

currency\_symbol: symbol for currency

status: digitization of the menu (complete/under review)

page\_count: number of pages on the menu

dish\_count: number of dishes on the menu

## **2.2 USE CASES**

**U1:** We can clean the data to answer the below use cases.

- How the food preferences changed over time based on the event, years, and location in the menu? After cleaning, we can ask this question, and the solution can be visualized as a trend analysis.
- Is there any specific type of dish whose price change has been greater than or less than the average change over time?
- How has the median price of dishes changed over time? Are there particular types of dishes whose price changes have been greater than or less than the average change over time?

**U0:** As the data set is really messy for any practical use case, but still we can use the data to gather some statistics which can be done without cleaning the data. Some of those use cases can be as below.

- We can get the information about the popular dishes based on how many times it appeared in menus using the dish dataset.
- What is the max/min of height and width of the Menu structure of previous years using the menu Page dataset?
- Using the help of Menu Page and Menu item together we can get information about a particular menu item as to when and where it appeared on a menu.

**U2:** What we cannot answer even after data cleaning and steers our thoughts towards the arguments that data cleaning should be done with a purpose in mind.

- How we would predict anything about a dish's price based on its name or description as it requires modeling
- There's been some work on how the words used in advertisements for potato chips are reflective of their price: is that also true of the words used in the name of the food?
- Based on French or Italian words in the name, does the price of the dish become expensive?

The dataset is quite messy and needs to be organized. For example, the date columns seem full of repeat or missing dates. However, some files are cleaner than others. We are not choosing to clean MenuPage.csv, as it is clean enough. As because we are getting towards the final outcome as trend analysis, we want to remove the data with outliers. We observed that most of the data has been recorded between 1850 to 2014 so we will remove the data for the rest of the dates.

### 3. DATA CLEANING STEPS

Based on the initial view of the data in the four files we can do the following sort of cleaning for all the files separately and then we can use SQL lite to create tables and put constraints:

#### **Menu.csv:**

Generally, for any data cleaning activity for the file below things can be done which will make data suitable for conveying information and making it more presentable.

- The sponsor column should trim leading and trailing white spaces and replace consecutive spaces.
- Convert columns to uppercase
- Remove special characters
- Replace; with space when it is inside the text of column value and may create separate columns with values that might need to be updated in ER. Saw the value in the physical description
- Convert similar values of the important columns and making in a standardized format using facets
- Make standard date formats
- Remove null columns after the processing

To be more specific we observed the below data issues in the sponsor column with special characters leading and trailing white spaces.

The screenshot shows the OpenRefine interface with a 'Custom text transform on column sponsor' dialog box open. The 'Expression' field contains the formula `value.replace(/[!@#(){}%*]/, '')`. The 'Language' is set to 'General Refine Expression'. Below the expression field is a 'Preview' tab showing a table of data rows. The table has two columns: 'row' and 'value'. The 'value' column shows the result of the text transformation, where special characters have been removed from the original text. The 'On error' section has three radio buttons: 'keep original' (selected), 'set to blank', and 'store error'. There is also a checkbox for 'Re-transform up to 10 times until no change'.

row	value
51.	NATIONAL VERBANDES DEUTSCH AMERIKANISCHER JOURNALISTEN UND SCHRIFTSTELLER (UNUSUAL)
52.	COLUMBIA RESTAURANT
53.	MRS CLUFF
54.	BAILY CATERING CO.THE
55.	NORTHERN STEAMSHIP CO.
56.	CITIZENS' STEAMBOAT COMPANY

As we see many data items having multiple variations of the same fields, we will be using OpenRefine different clustering facets to standardize sponsor for display purposes.

OpenRefine

Facet / Filter

Refresh

— sponsor

110 choices

Sort by:

BERGEN"" 5

CONTE BIANCAMAN

EX LIBRIS"" 1

HAMBURG"" 2

LET ER BUCK"" 1

PARIS"" 1

VICTORIA LUISE" 1

CHEZ HANSI"" 1

EATING THE ITALIAN

THE MANOR"" 4

2TH REGIMENT INF

4TH REGIMENT AR

Cluster & Edit column "sponsor"

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

163 clusters found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
8	24	<ul style="list-style-type: none"> <li>RED STAR LINE - ANTWERPEN - NY (7 rows)</li> <li>RED STAR LINE - ANTWERPEN NY (6 rows)</li> <li>RED STAR LINE - ANTWERPEN -NY (5 rows)</li> <li>RED STAR LINE -ANTWERPEN -NY (2 rows)</li> <li>RED STAR LINE -ANTWERPEN - NY (1 rows)</li> <li>RED STAR LINE - ANTWERPEN - NY (1 rows)</li> <li>RED STAR LINE - ANTWERPEN -NY (1 rows)</li> <li>RED STAR LINE- ANTWERPEN NY (1 rows)</li> </ul>	<input checked="" type="checkbox"/>	RED STAR LINE - ANTWERPE
5	667	<ul style="list-style-type: none"> <li>NORDEUTSCHER LLOYD BREMEN (632 rows)</li> <li>NORDEUTSCHER LLOYD - BREMEN (31 rows)</li> <li>NORDEUTSCHER LLOYD, BREMEN (2 rows)</li> <li>BREMEN NORDEUTSCHER LLOYD (1 rows)</li> <li>NORDEUTSCHER LLOYD -BREMEN (1 rows)</li> </ul>	<input checked="" type="checkbox"/>	NORDEUTSCHER LLOYD BI
3	29	<ul style="list-style-type: none"> <li>NORDEUTSCHER LLOYD BREMEN (24 rows)</li> <li>NORDEUTSCHER LLOYD, BREMEN (3 rows)</li> <li>NORDEUTSCHER LLOYD - BREMEN (2 rows)</li> </ul>	<input checked="" type="checkbox"/>	NORDEUTSCHER LLOYD BRI
3	29	<ul style="list-style-type: none"> <li>GRAMERCY PARK HOTEL (19 rows)</li> <li>HOTEL GRAMERCY PARK (9 rows)</li> <li>GRAMERCY PARK HOTEL HOTEL GRAMERCY PARK (1 rows)</li> </ul>	<input checked="" type="checkbox"/>	GRAMERCY PARK HOTEL

Select All

Unselect All

Export Clusters

Merge Selected & Re-Cluster

Merge Selected & Close

Close

# Choices in Cluster

2 — 8

# Rows in Cluster

0 — 700

Average Length of Choices

3 — 93

Length Variance of Choices

0 — 9.43

We are going to do similar processing for columns: event, venue, name, place, location, occasion. For the date column, we will remove outliers as below.

OpenRefine Menu

Permalink

Open...

Facet / Filter

Undo / Redo 47 / 47

Refresh

Reset All

Remove All

— date

change reset

1828-03-26 02:00:00 — 2028-03-25 14:00:00

☒ Time
 ☐ Non-Time
 ☒ Blank
 ☐ Error

16961

0

586

0

17544 matching rows (17547 total)

Show as: rows records

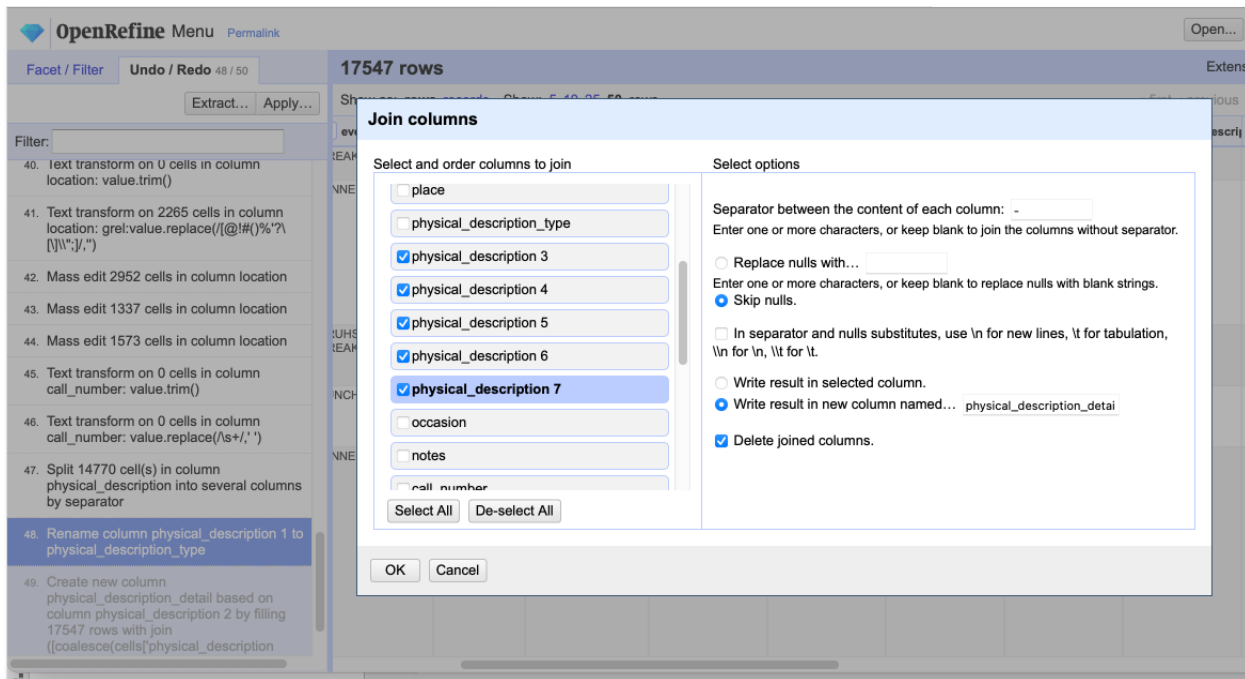
Show: 5 10 25 50 rows

« first < previous 1

iption	occasion	notes	call_number	keywords	language	date	location	location_type
	EASTER		1900-2822			1900-04-15T00:00:00Z	HOTEL EASTMAN	
7.0X9.0;	EASTER	WEDGEWOOD BLUE CARD; WHITE EMBOSSED GREEK KEY BORDER; "EASTER SUNDAY" EMBOSSED IN WHITE; VIOLET COLORED SPRAY OF FLOWERS IN UPPER LEFT CORNER;	1900-2825			1900-04-15T00:00:00Z	REPUBLICAN HOUSE	
.5X8.0;		MENU IN GERMAN AND ENGLISH; ILLUS, STEAMSHIP AND SAILING VESSEL;	1900-2827			1900-04-16T00:00:00Z	NORDEUTSCHER LLOYD BREMEN	
.5X8.0;		MENU IN GERMAN AND ENGLISH; ILLUS, HARBOR SCENE WITH SAILING VESSEL;	1900-2828			1900-04-16T00:00:00Z	NORDEUTSCHER LLOYD BREMEN	
; 5.5X7.5;		MENU IN GERMAN AND ENGLISH; ILLUS, HARBOR SCENE WITH ROCKS AND LIGHTHOUSE; STEAMSHIP AND SAILING VESSELS; CONCERT PROGRAM; DATES: ON GERMAN SIDE OF MENU "MONTAG, DEN 16 APRIL 1900"; ON ENGLISH SIDE OF MENU "MONDAY, APRIL 15TH, 1900";	1900-2829			1900-04-16T00:00:00Z	NORDEUTSCHER LLOYD BREMEN	

July 2021

For the physical description column, we will split it into multiple columns to deduce more meaningful information about the menu sizes and structures.



We will not do any kind of cleaning activities for the below columns excepts uppercases and removing spaces making a suitable for viewing in the trend analysis tying to our use case U1:

ID, Keywords, Language, Status

No changes are needed for page count and dish count.

### MenuPage.csv:

looks to be clean. As per the initial set up looks like it might not need cleaning or a bit of cleaning.

### MenuItem.csv:

Convert all the dates to a standard format for columns *created\_at* and *updated\_at*.

Check for the special characters and remove if any.

### Dish.csv:

For columns *name* and *description* using key collision to cluster values and not nearest neighbor because it was computationally expensive. Check for the special characters and remove them while cleaning.

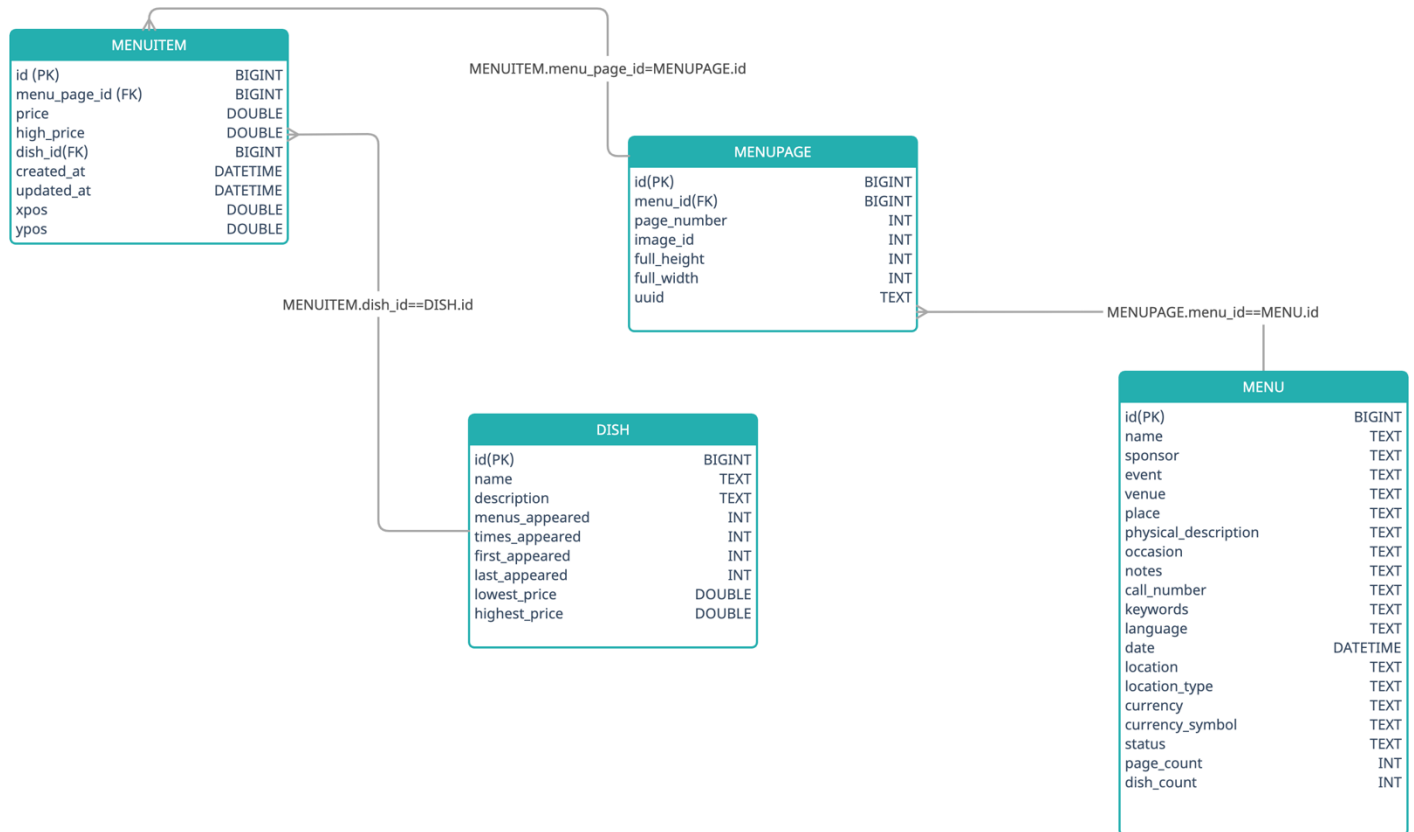
After we clean the data, we will run through the integrity constraints to check the data for the fields in the table.

We will use SQLite [3] for this. We can do a second round of cleanup for the data missing from OpenRefine [1].

## 4. SCHEMA

### 4.1 ER DIAGRAM

Below is the ER diagram of the files.



### 4.2 INTEGRITY CONSTRAINT

The checks that we will perform on the data for all the files will be as below:

1. Not null checks
2. Dates within a range
3. Primary keys are not null and empty as well as unique
4. Null foreign keys will be deleted
5. Low price should be below the highest price
6. None of the columns should have only NULL as the value
7. Update date should be greater than the created date.

## 5. CONCLUSION

We will discuss our observations and provide more information in phase II.

## CONTRIBUTIONS

Cleaning raw data files using OpenRefine: Suruchi Sinha, Pradeep Sharma

Creating tables and integrity constraints checks: Suruchi Sinha, Upul Gunasena, Pradeep Sharma

YesWorkflow diagrams: Pradeep Sharma, Upul Gunasena

Github uploadings: Suruchi Sinha, Pradeep Sharma, Upul Gunasena

## ACKNOWLEDGMENT

We would like to thank Prof. Bertram Ludaescher and TAs for their guidance to work on this project.

## REFERENCES

- [1] “Openrefine: A free, open-source, powerful tool for working with messy data,” <http://openrefine.org/>.
- [2] B. L. et al, “Yesworkflow,” <https://github.com/yesworkflow-org/>.
- [3] “Sqlite,” <https://www.sqlite.org/>.
- [4] N. Y. P. Library, “What’s on the menu?” <http://menus.nypl.org/>
- [5] “Frank e. buttolph,” [https://en.wikipedia.org/wiki/Frank E.Buttolph](https://en.wikipedia.org/wiki/Frank_E.Buttolph).