

CS513: Theory & Practice of Data Cleaning

NYPL menus dataset Final Project

Thomas Kalnik and Pradosh Kumar Subudhi

Our team chooses the following Phase-1 option:

- [X] (A) No change to Phase-1 report
- [] (B) Improved (or extended) Phase-1 report

1. [40 Pts] **Data cleaning performed (with OpenRefine, Python, and/or other tools)**
 - a. [20 Pts] Identify and describe all data cleaning steps you have performed.

For the data cleaning portion of this assignment, we focused our efforts on using OpenRefine, SQLite and Python (with Pandas and the re Regex module). Some of the challenges we encountered with this assignment related to, extraneous character values (such as semicolons, quotation marks, square brackets and various other extraneous characters), erroneous date values (such as dates with year values that are well in the future, or almost 1000 years in the past) and multiple values contained in the same columns.

The data cleaning steps we performed were as follows:

For the event, venue, place and date columns, we used Python and Regex to sanitize the data values (removing extraneous characters and bringing the dates to within a reasonable range) in an appropriate way. The data output was then saved to a new csv file and later merged with the output of the data cleaning that was completed using OpenRefine and SQL.

For the other columns in the dataset, namely the id, name, sponsor, physical_description, occasion, notes, call_number, keywords, language, location, location_type, currency, status, page_count and dish_count columns we used OpenRefine and SQL to clean and sanitize the data.

Once most of the cleaning was completed using Python and OpenRefine, SQLite Studio was used to load the data into different tables and do additional data validation (Integrity Constraint Violation) and cleaning where possible. As most of the data were cleaned using OpenRefine, the scope for cleaning was minimal. Below is the detail of the ER Diagrams, Integrity Constraint Violation Validation and several Cleaning steps performed using SQLite.

- b. [20 Pts] For each data cleaning step you have performed, explain its rationale. Was it required to support the use case U1? If not, explain why those steps were still useful.

Each data cleaning step had an important role in the overall process of correcting the data. Referring to our original use case described in the Final Project Proposal, the U1 use case we described was:

U1 - Data is Fit-for-Purpose:

A consulting group is evaluating menu options and offerings as a service to a new fine dining restaurant set to open in New York City. The consulting group needs a clean dataset with no duplicates to conduct their research, but some minor discrepancies in the data are acceptable. Data Cleaning is valuable in this context because in the present form (raw data, i.e. D) the dataset contains blank values, inconsistency in naming convention and data types and duplicate values that can be merged. Completing these steps will improve the quality of the data enough that it is a valuable resource to a consulting firm to make recommendations to a new fine dining establishment on what type of menus they should offer.

The data cleaning steps applied were relevant or not relevant to this use case in the following ways:

- *Clustering was used to group different data values, which allows for the removal of duplicate values in the dataset. As described in the U1 use case, this is a necessary prerequisite for the U1 use case (namely. A consulting group evaluating menu options and providing recommendations to a fine dining establishment client.*
- *Removing trailing and leading white spaces has little relevance to the U1 use case, a consulting group could still use the data without this transformation step, however, readability and usability are improved by this step.*
- *Removing extraneous character values like “ , () { } * [] ‘ . ? ; has some relevance to the U1 use case. While the consulting group could theoretically make use of data containing these characters, readability is improved through the removal and this transformation step also allows for better clustering of the data values, which removes additional duplicates and improves the quality and usability of the dataset.*
- *Data Type conversion had some relevance to the U1 use case. While generally not directly applicable to a consultant evaluating menu options to make recommendations to a restaurant business, data type is sometimes relevant to the readability and accuracy of the data. For example, dates and number values can contain higher precision when stored with the correct data type and this precision can allow for more accurate data value storage in addition to better interpretation and utilization of the data.*

2. [20 Pts] Document data quality changes

Data Cleaning With OpenRefine:

All the 4 dirty datasets were loaded into the OpenRefine and cleaning operation was performed. The cleaning operation is captured in the json file that is generated by OpenRefine.

1) Menu.csv

Number of records – 17, 547

a. Please note that 90+ steps were performed on the dataset and it was the most dirty dataset of all. Some steps were repeated and the data set was cleaned. So, it was difficult to capture all the numbers in the process.

b. Generic : All the columns where numeric fields were present were converted to number type. Date fields were converted to date type and text types to text.

i. name :

1. Removed leading and trailing spaces.
2. Collapsed consecutive white spaces.
3. Converted the values to title case.
4. Clustering was used to group different data values and appropriate action was taken on the values. All these below methods were used to cluster the data.

a. Method : key collision

- i. fingerprint
- ii. ngram-fingerprint
- iii. metaphone-3
- iv. cologne-phonetic
- v. Daitch-Mokotoff
- vi. Beider-Morse

b. Nearest neighbor

- i. levenshtein
- ii. ppm

5. Used GREL to remove values like “ , () { } * [] ‘ . ? ;

6. All these above steps were repeated multiple times as each clustering and cleaning opened up an opportunity to apply the transformation again.

ii. sponsor :

1. Removed leading and trailing spaces.
2. Collapsed consecutive white spaces.
3. Converted the values to title case.
4. Clustering was used to group different data values and appropriate action were taken on the values. All these below methods were used to cluster the data.

a. Method : key collision

- i. fingerprint
- ii. ngram-fingerprint
- iii. metaphone-3
- iv. cologne-phonetic
- v. Daitch-Mokotoff
- vi. Beider-Morse

b. Nearest neighbor

- i. levenshtein

ii. ppm

5. Used GREL to remove values like “ , () { } * [] ‘ . ? \ \ ;

6. All these above steps were repeated multiple times as each clustering and cleaning open up an opportunity to apply the transformation again.

iii. occasion :

1. Removed leading and trailing spaces.

2. Collapsed consecutive white spaces.

3. Converted the values to title case.

4. Used GREL to remove values like “ , () { } * [] ‘ . ? \ \ ;

5. Clustering was used to group different data values and appropriate action was taken on the values. All these below methods were used to cluster the data.

a. Method : key collision

i. fingerprint

ii. ngram-fingerprint

iii. metaphone-3

iv. cologne-phonetic

v. Daitch-Mokotoff

vi. Beider-Morse

b. Nearest neighbor

i. levenshtein

ii. ppm

6. The opportunity to clean the data based on the clustering was very helpful. Data like ‘Anniv’, “ANNIVERSARYERSARY;” etc. were clubbed together as ‘Anniversary’. This was one of the examples from the list of values that were cleaned.

7. Used GREL to remove values like “ , () { } * [] ‘ . ? \ \

8. All these above steps were repeated multiple times as each clustering and cleaning opened up an opportunity to apply the transformation again.

iv. event, venue :

1. Removed leading and trailing spaces.

2. Collapsed consecutive white spaces.

3. Converted the values to title case.

4. Used GREL to remove values like “ , () { } * [] ‘ . ? \ \

5. Clustering was used to group different data values and appropriate action was taken on the values. All these below methods were used to cluster the data.

a. Method : key collision

i. fingerprint

ii. ngram-fingerprint

iii. metaphone-3

iv. cologne-phonetic

v. Daitch-Mokotoff

vi. Beider-Morse

b. Nearest neighbor

i. levenshtein

ii. ppm

6. Used GREL to remove values like “ , () { } * [] ‘ . ? \ \

7. All these above steps were repeated multiple times as each clustering and cleaning opened up an opportunity to apply the transformation again.

v. place

1. Removed leading and trailing spaces.

2. Collapsed consecutive white spaces.

3. Used GREL to remove values like “ , () { } * [] ‘ . ? \ \

4. Clustering was used to group different data values and appropriate action was taken on the values. All these below methods were used to cluster the data.

a. Method : key collision

i. fingerprint

ii. ngram-fingerprint

iii. metaphone-3

iv. cologne-phonetic

v. Daitch-Mokotoff

vi. Beider-Morse

b. Nearest neighbor

i. levenshtein

ii. ppm

5. Used GREL to remove values like “ , () { } * [] ‘ . ? \ \

6. All these above steps were repeated multiple times as each clustering and cleaning opened up an opportunity to apply the transformation again.

vi. physical_description

1. Removed leading and trailing spaces.

2. Collapsed consecutive white spaces.

3. Used GREL to remove values like “ , () { } * [] ‘ . ? \ \

4. Clustering was used to group different data values and appropriate action was taken on the values. All these below methods were used to cluster the data.

5. The opportunity to clean the data was minimal.

6. The cleaning was not necessary and could have been left as it is.

vii. occasion

1. Removed leading and trailing spaces.

2. Collapsed consecutive white spaces.

3. Used GREL to remove values like “ , () { } * [] ‘ . ? \ \

4. Clustering was used to group different data values and appropriate action was taken on the values. All these below methods were used to cluster the data.

a. Method : key collision

i. fingerprint

ii. ngram-fingerprint

iii. metaphone-3

iv. cologne-phonetic

- v. Daitch-Mokotoff
 - vi. Beider-Morse
- b. Nearest neighbor
 - i. levenshtein
 - ii. ppm
- 5. Used GREL to remove values like “ , () { } * [] ‘ . ? \ \
- 6. All these above steps were repeated multiple times as each clustering and cleaning open up an opportunity to apply the transformation again

viii. notes

- 1. Removed leading and trailing spaces.
- 2. Collapsed consecutive white spaces.
- 3. Used GREL to remove values like “ , () { } * [] ‘ . ? \ \
- 4. Few clusters were available for cleaning considering the nature of data available.
- 5. Opportunity was less.

ix. call_number

- 1. Few clusters were available for cleaning considering the nature of data available.
- 2. Opportunity was less.

x. keywords, language, location, location_type

- 1. The columns were removed because of lack of data in them.

xi. date

- 1. Converted the column to Date type.
- 2. Formatted the data to ISO standard - toString(toDate(value),"yyyy-MM-dd")

xii. currency

- 1. Removed leading and trailing spaces.
- 2. Collapsed consecutive white spaces.
- 3. Converted the values to title case.

xiii. status

- 1. Removed leading and trailing spaces.
- 2. Collapsed consecutive white spaces.
- 3. Converted the values to title case.
- 4.

xiv. page_count, dish_count

- 1. No scope for transformation other than changing the column type to number type.

xv. currency_symbol

- 1. No Cleaning.

2) Dish.csv

Number of records – 428,082

a. Please note that because of the size of the dataset, the OpenRefine tool was not able to apply any Facets. Because of this, the mess in the data was not possible to identify.

b. With careful observation, many dirty data were cleaned.

c. Generic : All the columns where numeric fields were present were converted to number type. Date fields were converted to date type and text types to text.

i. name :

1. Removed leading and trailing spaces. – Rows Affected - 6582

2. Collapsed consecutive white spaces. - Rows Affected - 9288

3. Converted the values to title case.

4. Clustering was not possible.

5. Used GREL to remove values like “ , () { } * [] ‘ . ? -

a. Rows Affected – 7009, 2568, 5 1801, 18, 309, 33015 (different types of characters)

ii. id, menus_appeared, times_appeared, lowest_price, highest_price :

1. Removed leading and trailing spaces.

2. Collapsed consecutive white spaces.

3. Convert the values to numbers.

iii. first_appeared, last_appeared

1. Date transformations.

3) MenuItem.csv

Number of records – 1,334,779

a. Please note that because of the size of the dataset 118MB, the OpenRefine tool was not able to apply any Facets.

b. Data was mostly clean.

c. Some generic transformations applied.

d. Text transform on 1334779 cells in column menu_page_id: value.toNumber()

e. Text transform on 888520 cells in column price: value.toNumber()

f. Text transform on 91979 cells in column high_price: value.toNumber()

g. Text transform on 1334538 cells in column dish_id: value.toNumber()

h. Text transform on 1334779 cells in column xpos: value.toNumber()

i. Text transform on 1334779 cells in column ypos: value.toNumber()

4) MenuPage.csv

Number of records – 66,937

a. Data was mostly clean.

b. Few generic transformation was applied.

c. Text transform on 66937 cells in column id: value.toNumber()

d. Text transform on 66937 cells in column menu_id: value.toNumber()

e. Text transform on 65735 cells in column page_number: value.toNumber()

f. Text transform on 66914 cells in column image_id: value.toNumber()

g. Text transform on 66608 cells in column full_height: value.toNumber()

- h. Text transform on 66608 cells in column full_width: value.toNumber()
- i. Text transform on 1 cells in column uuid: value.toLowerCase()

a. [10 Pts] Quantify the results of your efforts (e.g., by providing a summary table of changes: Which columns changed? How many cells (per column) have changed? Etc.

Summary Table

File	Column	Quality Change Measure
Menu.csv	id	Text transform on 676 cells in column sponsor: <code>grel:value.replace(^?/,").replace(^(\)/,").replace(^"/,")</code>
		Text transform on 474 cells in column sponsor: <code>grel:value.replace(\\ \\/,").replace(^"/,").replace(^([\\)/,")</code>
		Text transform on 8275 cells in column sponsor: <code>value.toTitlecase()</code>

		<p>Clustering :</p> <p>key collision -</p> <p>fingerprint</p> <p>Mass edit 3937 cells in column sponsor</p> <p>ngram-fingerprint</p> <p>Mass edit 1787 cells in column sponsor</p> <p>metaphone-3</p> <p>Mass edit 524 cells in column sponsor Undo</p> <p>cologne-phonetic</p> <p>Mass edit 1676 cells in column sponsor</p> <p>Daitch-Mokotoff</p> <p>Mass edit 163 cells in column sponsor Undo</p> <p>Beider-Morse</p> <p>Mass edit 23 cells in column sponsor Undo</p> <p>nearest-neighbor -</p> <p>levenshtein</p> <p>Mass edit 1617 cells in column sponsor Undo</p> <p>ppm</p> <p>0</p>
	name	<p>Text transform on 28 cells in column name: <code>grel:value.replace(/^?/,").replace(/^(\\)/,").replace(/^"/,")</code></p> <p>Text transform on 139 cells in column name: <code>grel:value.replace(/^\\ \\/,").replace(/^"/,").replace(/^(\\)/,").replace(/^*/,")</code></p>

		<p>Clustering :</p> <p>key collision -</p> <p>fingerprint</p> <p>Mass edit 556 cells in column sponsor</p> <p>ngram-fingerprint</p> <p>Mass edit 622 cells in column sponsor</p> <p>metaphone-3</p> <p>Mass edit 762 cells in column sponsor Undo</p> <p>cologne-phonetic</p> <p>Mass edit 4 cells in column sponsor</p> <p>Daitch-Mokotoff</p> <p>Mass edit 12 cells in column sponsor Undo</p> <p>Beider-Morse</p> <p>Mass edit 0 cells in column sponsor Undo</p> <p>nearest-neighbor -</p> <p>levenshtein</p> <p>Mass edit 7 cells in column sponsor Undo</p> <p>ppm</p> <p>60</p>
	event	Text transform on 9393 cells in column event: value.toString()
		Text transform on 99 cells in column event: grel:value.replace(/?/,").replace(/(\\)/,").replace(/"/,")
		Text transform on 90 cells in column event: grel:value.replace(/\\V/,").replace(/"/,").replace(/["]/,")
		Text transform on 7826 cells in column event: value.toTitlecase()

		<p>Clustering :</p> <p>key collision -</p> <p>fingerprint</p> <p>Mass edit 4027 cells in column sponsor</p> <p>ngram-fingerprint</p> <p>Mass edit 2325 cells in column sponsor</p> <p>metaphone-3</p> <p>Mass edit 3057 cells in column sponsor Undo</p> <p>cologne-phonetic</p> <p>Mass edit 2297 cells in column sponsor</p> <p>Daitch-Mokotoff</p> <p>Mass edit 98 cells in column sponsor Undo</p> <p>Beider-Morse</p> <p>Mass edit 2 cells in column sponsor Undo</p> <p>nearest-neighbor -</p> <p>levenshtein</p> <p>Mass edit 407 cells in column sponsor Undo</p> <p>ppm</p> <p>0</p>
	venue	Text transform on 9449 cells in column venue: value.toString()
		Text transform on 8096 cells in column venue: value.toTitlecase

		<p>Clustering: fingerprint Mass edit 2126 cells in column sponsor ngram-fingerprint Mass edit 21 cells in column sponsor metaphone-3 Mass edit 5044 cells in column sponsor Undo cologne-phonetic Mass edit 0 cells in column sponsor Daitch-Mokotoff Mass edit 505 cells in column sponsor Undo Beider-Morse Mass edit 0 cells in column sponsor Undo</p> <p>nearest-neighbor - levenshtein Mass edit 0 cells in column sponsor Undo ppm Mass edit 7 cells in column sponsor Undo</p>
	physical_description	Text transform on 2777 cells in column physical_description: value.toString()
		Text transform on 38 cells in column physical_description: value.replace(/s+/, '')

		<p>Clustering: fingerprint Mass edit 1514 cells in column sponsor ngram-fingerprint Mass edit 4575 cells in column sponsor metaphone-3 Mass edit 0 cells in column sponsor Undo cologne-phonetic Mass edit 2 cells in column sponsor Daitch-Mokotoff Mass edit 0 cells in column sponsor Undo Beider-Morse Mass edit 0 cells in column sponsor Undo</p> <p>nearest-neighbor - levenshtein Mass edit 0 cells in column sponsor Undo ppm Mass edit 0 cells in column sponsor Undo</p>
	Occasion	Text transform on 13744 cells in column occasion: value.toString()
		Text transform on 508 cells in column sponsor: grel:value.replace(^?/,").replace(^(\)/,").replace(^"/,")
		Text transform on 474 cells in column sponsor: grel:value.replace(^\\ /,"").replace(^"/,").replace(^[\]/,")
		Text transform on 2954 cells in column occasion: grel:value.replace(^\\ /,"").replace(^"/,").replace(^[\]/,").replace(/:/,").replace(^*/,")

		<p>Clustering ; fingerprint Mass edit 969 cells in column sponsor ngram-fingerprint Mass edit 272 cells in column sponsor metaphone-3 Mass edit 1393 cells in column sponsor cologne-phonetic Mass edit 693 cells in column sponsor Daitch-Mokotoff Mass edit 215 cells in column sponsor Beider-Morse Mass edit 0 cells in column sponsor</p> <p>nearest-neighbor - levenshtein Mass edit 281 cells in column sponsor ppm Mass edit 716 cells in column sponsor</p> <p>Text transform on 3752 cells in column occasion: value.toTitlecase()</p>
	notes	<p>Text transform on 195 cells in column notes: value.replace(/s+/, ' ')</p> <p>Text transform on 9460 cells in column notes: value.toTitlecase()</p>

		<p>Clustering:</p> <p>fingerprint</p> <p>Mass edit 1355 cells in column sponsor</p> <p>ngram-fingerprint</p> <p>Mass edit 320 cells in column sponsor</p> <p>metaphone-3</p> <p>Mass edit 0 cells in column sponsor</p> <p>cologne-phonetic</p> <p>Mass edit 268 cells in column sponsor</p> <p>Daitch-Mokotoff</p> <p>Mass edit 0 cells in column sponsor</p> <p>Beider-Morse</p> <p>Mass edit 0 cells in column sponsor</p> <p>nearest-neighbor -</p> <p>levenshtein</p> <p>Mass edit 0 cells in column sponsor</p> <p>ppm</p> <p>Mass edit 0 cells in column sponsor</p>
		Text transform on 851 cells in column sponsor: <code>grel:value.replace(^?/,").replace(^(\)/,").replace(^"/,")</code>
		Text transform on 267 cells in column notes: <code>grel:value.replace(^\\ \\/,").replace(^"/,").replace(^[\\]/,")</code>
		Text transform on 6034 cells in column notes: <code>grel:value.replace(^;\$/,")</code>

	call_number	<p>Clustering :</p> <p>Mass edit 0 cells in column sponsor ngram-fingerprint</p> <p>Mass edit 0 cells in column sponsor metaphone-3</p> <p>Mass edit 0 cells in column sponsor cologne-phonetic</p> <p>Mass edit 0 cells in column sponsor Daitch-Mokotoff</p> <p>Mass edit 0 cells in column sponsor Beider-Morse</p> <p>Mass edit 0 cells in column sponsor</p> <p>nearest-neighbor - levenshtein</p> <p>Mass edit 0 cells in column sponsor ppm</p> <p>Mass edit 0 cells in column sponsor</p>
	date	<p>Text transform on 16958 cells in column date: value.toDate()</p> <p>Text transform on 16958 cells in column date: grel:toString(toDate(value),"yyyy-MM-dd")</p>
	status	Text transform on 17547 cells in column status: value.toTitlecase()
	page_count, dish_count	Text transform on 17547 value.toNumber()
Menupage.csv	id	Text transform on 66937 cells in column id: value.toNumber()
	menu_id	Text transform on 66937 cells in column menu_id: value.toNumber()
	page_number	Text transform on 65735 cells in column page_number: value.toNumber()

	image_id	Text transform on 66914 cells in column image_id: value.toNumber()
	full_height	Text transform on 66608 cells in column full_height: value.toNumber()
	full_width	Text transform on 66608 cells in column full_width: value.toNumber()
	uuid	Text transform on 1 cells in column uuid: value.toLowerCase()
Menuitem.csv	menu_page_id	Text transform on 1334779 cells in column menu_page_id: value.toNumber()
	price	Text transform on 888520 cells in column price: value.toNumber()
	high_price	Text transform on 91979 cells in column high_price: value.toNumber()
	dish_id	Text transform on 1334538 cells in column dish_id: value.toNumber()
	xpos	Text transform on 1334779 cells in column xpos: value.toNumber()
	ypos	Text transform on 1334779 cells in column ypos: value.toNumber()
Dish.csv	name	<p>Removed leading and trailing spaces. – Rows Affected - 6582</p> <p>Collapsed consecutive white spaces. - Rows Affected - 9288</p> <p>Converted the values to titlecase.</p> <p>Clustering was not possible.</p> <p>Used GREL to remove values like “ , () { } * [] ‘ . ? -</p> <p>Rows Affected – 7009, 2568, 5 1801, 18, 309, 33015 (different types of characters)</p>

	Id, menus_appeared, times_appeared, lowest_price, highest_price	Removed leading and trailing spaces. Collapsed consecutive white spaces. Converted the values to numbers.
	first_appeared, last_appeared	Date transformations. (mostly all records)

b. [10 Pts] Demonstrate that data quality has been improved, e.g., by devising ICV queries and showing the difference between “before” and “after” (cleaning).

Please see the ICV queries and summary statistics provided below. Please also see the provided queries.txt file in the attached zip archive.

Integrity Constraint Check :

1) Criteria: id shouldn't have duplicate or NULL Values on menu table.

a) Table : menu

Query:

```
SELECT id, COUNT(id) AS cnt_id
FROM menu
GROUP BY id
HAVING COUNT(id) > 1;
```

Result : 0 Records

Action : No action.

b) NULL Value Check:

Query:

```
SELECT * FROM menu
WHERE id IS NULL;
```

Result: 0 such records.

Action : No action.

2) Criteria: id shouldn't have duplicate or NULL Values on menuitem table.

a) Table : menuitem

Query:

```
SELECT id, COUNT(id) AS cnt_id
FROM menu
GROUP BY id
HAVING COUNT(id) > 1;
```

Result : 0 Records

Action : No action.

b) NULL Value Check:

Query:

```
SELECT * FROM menuitem
```

WHERE id IS NULL;

Result: 0 such records.

Action : No action.

3) Criteria: id shouldn't have duplicate or NULL Values on dish table.

a) Table : menuitem

Query:

```
SELECT id, COUNT(id) AS cnt_id
FROM dish
GROUP BY id
HAVING COUNT(id) > 1;
```

Result : 0 Records

Action : No action.

b) NULL Value Check:

Query:

```
SELECT * FROM dish
WHERE id IS NULL;
```

Result: 0 such records.

Action : No action.

4) Criteria: id shouldn't have duplicate or NULL Values on menupage table.

a) Table : menuitem

Query:

```
SELECT id, COUNT(id) AS cnt_id
FROM menupage
GROUP BY id
HAVING COUNT(id) > 1;
```

Result : 0 Records

Action : No action.

b) NULL Value Check:

Query:

```
SELECT * FROM menupage
WHERE id IS NULL;
```

Result: 0 such records.

Action : No action.

Cleaning Using SQL:

1) Table : menu

a. Identify and delete irrelevant records.

Query:

```
SELECT COUNT(*)
```

```
FROM menu
```

```
WHERE name = ''
```

```
AND sponsor= '';
```

Result: 1619

Observation : 1680 records where both name and sponsor are blank don't have most of its other columns not populated. So, all such records were deleted since they don't hold enough data for those records to be relevant.

Action : Delete such records.

Query:

```
DELETE FROM menu
```

```
WHERE name = ''
```

```
AND sponsor= '';
```

Result: 1619 records were deleted.

b. Default columns where most of the data are missing so that they can be used for categorizing.

Query:

```
SELECT COUNT(*)
```

```
FROM MENU
```

```
WHERE NAME="";
```

Result: 12729

Action : Default such records with value as 'Unavailable'.

Query:

```
UPDATE menu
```

```
SET name = 'Unavailable'
```

```
WHERE name = '';
```

Result: 12729 records were updated.

c. Default name and sponsor column where most of the data are like (Restaurant Andor Location Not Given./ Restaurant Name Andor Location Not Given/ Not Given) so that they can be used for categorizing

Query:

```
SELECT COUNT(*) FROM menu
```

```
WHERE name LIKE '%NOT GIVEN%'
```

Result: 138

Action : Default such records with value as 'Unavailable'.

Query:
UPDATE menu
SET name = 'Unavailable'
WHERE name LIKE '%NOT GIVEN%';
Result: 138 records were updated.

Query:
SELECT COUNT(*) FROM menu
WHERE sponsor LIKE '%NOT GIVEN%';
Result: 216

Query:
UPDATE menu
SET sponsor = 'Unavailable'
WHERE sponsor LIKE '%NOT GIVEN%';
Result: 216 records were updated.

2) Table : dish

a. Identify and delete irrelevant records.

Query:
SELECT count(*) FROM dish
WHERE times_appeared = "";
Result: 0

Query:
SELECT count(*) FROM dish
WHERE menus_appeared = "";
Result: 0

Action : None. No such records were found.

b. Observation : 1680 records where both name and sponsor are blank don't have most of its other columns not populated. So, all such records were deleted since they don't hold enough data for those records to be relevant.

Query:
SELECT COUNT(*) FROM menu
WHERE name = ""
AND sponsor = "";
Result: 1680 records found
Action : Delete such records.

Query:

```
DELETE FROM menu
```

```
WHERE name = ''
```

```
AND sponsor= '';
```

Result: 1680 records were deleted.

c. Last appeared_date had years beyond 2021.

Query:

```
SELECT COUNT(*) FROM dish
```

```
WHERE last_appeared > '2021';
```

Result: 180

Action : Delete such records.

```
DELETE FROM dish
```

```
WHERE last_appeared > '2021';
```

Result: 180 records were deleted.

d. Last appeared_date had years beyond 2021.

Query:

```
SELECT COUNT(*) FROM dish
```

```
WHERE last_appeared > '2021';
```

Result: 180

Action : Delete such records.

```
DELETE FROM dish
```

```
WHERE last_appeared > '2021';
```

Result: 180 records were deleted.

e. highest_price column contains a lot of spaces. Update them to 0.

Query :

```
SELECT * FROM dish
```

```
WHERE highest_price = ''
```

Result: 29,098

Action : Default the spaces to 0.

Query :

```
UPDATE dish
```

```
SET highest_price = 0
```

```
WHERE highest_price = '';
```

Result: 29,098 records updated.

f. Check the same for lowest_price.

Query :

```
SELECT * FROM dish
```

```
WHERE lowest_price = ''
```


Result: 0
Action : None.

Query :
UPDATE dish
SET lowest_price = 0
WHERE highest_price = "";
Result: 446,259 records updated.

3) Table : menuitem

a. The high_price column contains a lot of spaces. Update them to 0.

Query :
SELECT * FROM menuitem
WHERE high_price = '
Result: 1,242,800
Action : Default the spaces to 0.

Query :
UPDATE menuitem
SET high_price = 0
WHERE high_price = '
Result: 1,242,800 records updated.

b. The price column contains a lot of spaces. Update them to 0.

Query :
SELECT * FROM menuitem
WHERE price = '
Result: 446,259
Action : Default the spaces to 0.

Query :
UPDATE menuitem
SET price = 0
WHERE price = '
Result: 446,259 records updated.

4) Table : menupage
No relevant data found to be cleaned as such.

3. [20 Pts] Create a workflow model

a. [10 Pts] A visual representation of your overall (or “outer”) workflow W_o , e.g., using a suitable tool such as YesWorkflow. At a minimum, you should identify key inputs, outputs, and steps of the workflow, along with dependencies between these. Key phases and steps of your data cleaning project may include, e.g., data profiling, data loading, data cleaning, IC violation checks, etc. Explain the design of W_o and why you’ve chosen the tools that you have in your overall workflow.

Please see the provided Flow Diagrams and YW Files in the Flow Diagrams & YW Files and Python Code & python YW diagrams folders of the zip attachment.

b. [10 Pts] A detailed (possibly visual) representation of your “inner” data cleaning workflow W_i (e.g., if you’ve used OpenRefine, you can use the OR2YW tool).

Please see the provided Flow Diagrams and YW Files in the Flow Diagrams & YW Files and Python Code & python YW diagrams folders of the zip attachment.

4. [10 Pts] **Conclusions & Summary.** Please provide a concise summary and conclusions of your project, including lessons learned. If you haven’t done so earlier in the report, you should also summarize the contributions of each team member here (for teams with ≥ 2 members).

This project involved the use of OpenRefine, SQL and Python in cleaning the NY Public Library Menus dataset.

This project was beneficial in that it allowed us to gain deeper experience with the tools and workflows presented in the class, specifically, OpenRefine, SQL and Python Pandas. These tools are powerful assistants in the Data Analyst’s arsenal and allow for fast and efficient manipulation of large, messy datasets that contain problems such as missing features, extraneous characters, and issues related to data type and dates. Proficiency in these tools is an invaluable skill set for any aspiring or seasoned Data Analyst, Data Engineer, Data Scientist or Software Developer.

The contributions of the teammates were as follows:

Pradosh led the data cleaning efforts with OpenRefine, YesWorkflow and SQL and ICV and contributed to writing the report.

Thomas was involved in the generation of the SQL Database Schema Diagram, YesWorkflow Diagrams wrote the python code used for the data cleaning component of the project and took the lead in compiling and writing the report.

[10 Pts] Submitting the report and supplementary materials Submit a single ZIP file with the final project report (including both Part-1 and Part-2) and the following supplementary materials:

a. Name of the report should include team number, e.g.: team1001-final-report.pdf

Team number: cs-513-theory-t-nrq

b. Supplementary materials:

i. Conceptual model / database schema: Provide these (e.g., your ER-diagram and/or your SQL schema) as separate files

Please see the provided SQL Database Schema as a separate attachment.

ii. If you used OpenRefine for data cleaning: Operation History: A copy of the OpenRefine operation history (copy-paste it into a json file named OpenRefineHistory.json) If you used other tools, e.g., Trifacta Wrangler or Python programs: include those scripts or programs and any other auxiliary files required.

Please see the provided OpenRefineHistory.json and parse.py python files as separate attachments.

iii. Queries: A copy of the queries written in SQL or Datalog that you used to profile the dataset and/or check integrity constraints (copy-paste them into a text file named Queries.txt or Queries.sql).

Please see the provided Queries.sql file.

iv. Outer and Inner Workflow Models: For both workflow models, provide the necessary files. For example, if you've used YesWorkflow for defining the model, then please include the following: the YW annotations files (OverallWorkflow.yw), and the YW-generated Graphviz/DOT file (OverallWorkflow.gv). If you've used another diagramming tool, then include your workflow diagrams as separate files as well.

Please see the provided Flow Diagrams and YW Files in the Flow Diagrams & YW Files and Python Code & python YW diagrams folders of the zip attachment.

v. Raw and Cleaned Datasets: Please do not include the datasets in the ZIP file! Rather, upload the raw and cleaned datasets to a Box folder and share the link in a plain text file (DataLinks.txt). Make sure the Box folder link is open to the public (or at least to Illinois.edu-registered users)!

Please see the provided DataLinks.txt file as a separate attachment.