

BIG DATA PROJECT REPORT - Citywide Payroll Data (Fiscal Year)

RUI WANG RW3068, YUTING PENG YP2212, and XIXUAN WANG XW2579*, New York University, US

[Github link](#)

ACM Reference Format:

Rui Wang rw3068, Yuting Peng yp2212, and Xixuan Wang xw2579. 2021. BIG DATA PROJECT REPORT - Citywide Payroll Data (Fiscal Year). 1, 1 (December 2021), 21 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 INTRODUCTION

Data is extremely important in a variety of fields, which plays an unshakable role in making informed decisions, finding solutions to problems, and making people's lives much easier. However, the quality of data is further more essential than the data itself, and there are numbers of cases that awful data quality leads to severe consequences such as substantial money waste and fatal accidents. Therefore, data cleaning is deemed to be an indispensable pre-process before actually making the best use of those data. It's hence reasonable to ask: what is data cleaning?

Data cleaning is the process of eliminating or fixing incorrect, corrupted, incorrectly formatted, duplicate, or incomplete data within a dataset, which has significant meaning in increasing the overall productivity and allowing for the highest quality information in decision makings.

However, data cleaning is usually so difficult and time-consuming for several reasons. First, a variety of data formats must be dealt with. Second, there are so many different data quality issues as databases try to fit a complex world into a simple abstraction. Finally, cleaning the data requires domain knowledge.

Fortunately, openclean saves us a lot of work. Openclean is a Python library for data profiling and data cleaning, which will be discussed further in the later Related Work section. Openclean, together with spark, a unified engine for large-scale data analytics, supported most of our work in part one and part two of this project respectively.

2 PROBLEM FORMULATION

The project has two part. The first part is to profile and clean data for the Citywide Payroll Data (Fiscal Year), and the second part is to apply the clean method to a volume of datasets, testing the scale ability of our clean method.

2.1 project part one

2.1.1 *profile data.*

We mainly used OpenClean for data profile, for instance, from the basis use of `profiles.stats()`, `profiles.minmax()`, and `profiles.column().get()` to more deep methods helping us to find the dirty data like the similarity library and create

*All authors contributed equally to this research.

Authors' address: Rui Wang rw3068; Yuting Peng yp2212; Xixuan Wang xw2579, New York University, US.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2021 Association for Computing Machinery.

Manuscript submitted to ACM

cluster function. The formulation of this part of work is as follow:

functions	what we found from it
profiles.stats()	overview the attribute of dataset
profiles.minmax()	see the numerical data's range
profiles.column().get('topValues')	see the top 10 values for each attribute and count their frequency
different kinds of cluster	to check the spelling error and no consistent data

For different kind of data, we design different methods to find the error or inconsistent data.

Attribute	The method we use
Agency Name	distinct statistic and Key-collision cluster
Last and First Name	uppercase
Mid Init	distinct and invalid character check
Agency Start Date	invalid date check and impossible date check
Work Location Borough	distinct statistic and key-collision cluster
Title Description	distinct statistic and KNN clustering
other number values	negative dealing

by doing the profile using the above method, we find a lot of data problems like data defected, data inconsistent, data error, and extra, the below just show a small part of them:

- (1) Empty values in 'Last Name', 'First Name'
- (2) characters in 'Mid Init'
- (3) 'Agency Name' value : police department, board of correction
- (4) 'Agency Start Date' value : 12/31/9999, 10/16/2049
- (5) 'OT hours', 'Regular Hours' negatives:
- (6) 'Work Location Borough' : MANHATTAN manhattan....
- (7) 'Title Description': many typos

2.1.2 clean data. We design clean method for each type of data. and the main function is as below:

functions	what we found from it
Agency Name	uppercase and specific change
Last and First Name	uppercase
Mid Init	using ASCII to select the A to Z character
Agency Start Date	judge the invalid days for each month and year, delete the future time
Work Location Borough	uppercase and empty data to 'UNKNOWN'
Title Description	uppercase and specific change
other number values	change negative value to null

2.2 project part two

2.2.1 scale method and analysis.

We used spark to clean the scale datasets and analysis our clean method in the scale datasets. To better focus on different attribute, for each column, we find 10 datasets from the NYC open data source having the same attribute and apply our clean method to those datasets, while refine the clean method to make the cleaning be more comprehensive

and precise.

refine method for each attribute

functions	the refine method we added
Agency Name	space eliminated, inconsistent refine and uncompleted data supplement
Mid Init	using spark
Agency Start Date	to date function
Work Location Borough	space eliminated
Title Description	using spark

2.2.2 data references.

For the String attribute: Agency Name, Work Location Borough and Title description, we used the final cleaned 10 datasets for each attribute to create a reference data. Since the 10 datasets volume is large enough and have a lot of duplicate, which can double make sure our final data are clean, we UNION all the data through the 10 datasets and get the result of the reference data.

3 RELATED WORK

3.1 Openclean

Openclean is a Python library for data profiling and data cleaning, whose goal is to bring together data cleaning tools in a single environment that is easy and intuitive to use for a data scientist, allowing users to compose and execute cleaning pipelines that are built using a variety of different tools. Moreover, it is flexible and extensible to allow easy integration of new functionality.

The following code presents where and how we mainly used the openclean library. As shown below, openclean provides numbers of APIs which significantly reduced the complexity of operations, making us be able to concentrate and focus more on the data cleaning core logic.

```

1 from openclean.data.source.socrata import Socrata
2 dataset = Socrata().dataset("k397-673e")
3
4 from openclean.pipeline import stream
5 ds_full = stream(datafile)
6
7 from openclean.profiling.column import DefaultColumnProfiler
8 profiles = ds.profile(default_profiler=DefaultColumnProfiler)
9 profiles.stats()
10
11 for col in ds.columns:
12     p = profiles.column(col)
13     print("{}'{}_({})'.format(col, p['datatypes']['distinct'].most_common(1)[0][0]))
14
15 profiles.minmax('Fiscal_Year')
16 profiles.column('Agency_Name').get('topValues')
17
18 from openclean.cluster.knn import knn_clusters
19 from openclean.function.similarity.base import SimilarityConstraint
20 from openclean.function.similarity.text import LevenshteinDistance
21 from openclean.function.similarity.text import HammingDistance
22 from openclean.function.similarity.text import JaroSimilarity
23 from openclean.function.similarity.text import StringSimilarityFunction

```

```
24 from openclean.function.value.threshold import GreaterThan
```

3.2 Clustering

The concept of clustering is about using statistic methods for processing data by organizing items into groups, or clusters, on the basis of how closely associated they are. In this project, we mainly used two methods related to clustering including the KNN clustering method and the key collision method.

KNN is short for K-nearest neighbors, and it is one of the most popular learning algorithms. The purpose of it is to use a database in which the data points are separated into several classes to predict the classification of a new data point. In our project, we don't need to make predictions for new data points, so we actually used KNN to cluster the data of existing columns so that we are able to find out similar data fields. For example, by utilizing the KNN clustering method, we are able to distinguish "ADM MANAGER-NON-MGR" and "ADM MANAGER-NON-MGRL", or "DIRECTOR OF MANAGEMENT PLANNING" and "DIRECTOR OF MANAGEMENT PLANNING SS", etc. This really helps use to make judgement on which data fields are incorrectly formatted or needed to be unified to be the same data field.

```
1 Cluster 1
2   NON-TEACHING ADJUNCT III (x 3419)
3   NON-TEACHING ADJUNCT I (x 8407)
4   NON-TEACHING ADJUNCT V (x 1047)
5   NON-TEACHING ADJUNCT IV (x 1279)
6   NON-TEACHING ADJUNCT II (x 3774)
7 Cluster 2
8   EDUCATIONAL ADMINISTRATOR UFT (x 58)
9   EDUCATIONAL ADMINISTRATOR (x 1056)
10  EDUCATIONAL ADMINISTRATOR CSA (x 6281)
11 Cluster 3
12  SUPERVISOR II (x 1503)
13  SUPERVISOR I (x 3390)
14  SUPERVISOR III (x 490)
15 Cluster 4
16  ADM MANAGER-NON-MGR (x 39)
17  ADM MANAGER-NON-MGRL (x 5406)
18 Cluster 5
19  INTELLIGENCE RESEARCH MANAGER (x 6)
20  INTELLIGENCE RESEARCH MANAGER-PD (x 39)
21 Cluster 6
22  HEALTH SERVICES MANAGER NON MANAGERIAL LEVEL II (x 92)
23  HEALTH SERVICES MANAGER NON MANAGERIAL LEVEL I (x 301)
24 Cluster 7
25  DIRECTOR OF MANAGEMENT PLANNING (x 3)
26  DIRECTOR OF MANAGEMENT PLANNING SS (x 23)
27  ....
```

As mentioned above, KNN clustering helps us be able to discover similar data fields. They are considered similar due to one or more character differences. The other clustering method - key collision, is also used to help us find out similar data fields, whereas they are considered similar because of special character differences or the words within two data fields are in different orders. For example, " *ADM SCHOOL SECURITY MANAGER-U" and "ADM SCHOOL SECURITY MANAGER-U" are basically the same except for the leading "*" character. Also, "ASSISTANT DEPUTY

COMMISSIONER" and "DEPUTY ASSISTANT COMMISSIONER" are quite similar with the only difference of word order. We can assert that "DEPUTY ASSISTANT" and "ASSISTANT DEPUTY" both refer to the same thing.

```

1 Cluster 1
2 *ADM SCHOOL SECURITY MANAGER-U (x 16)
3 ADM SCHOOL SECURITY MANAGER-U (x 3)
4
5 Cluster 2
6 *ADMIN SCHL SECUR MGR-MGL (x 4)
7 ADMIN SCHL SECUR MGR-MGL (x 1)
8
9 Cluster 3
10 *ADMINISTRATIVE ATTORNEY (x 24)
11 ADMINISTRATIVE ATTORNEY (x 5)
12 ...
13 Cluster 19
14 ASSISTANT DEPUTY COMMISSIONER (x 9)
15 DEPUTY ASSISTANT COMMISSIONER (x 326)
16
17 Cluster 20
18 *CUSTODIAL ASSISTANT (x 147)
19 CUSTODIAL ASSISTANT (x 1340)
20 .....

```

3.3 Apache Spark

Apache Spark is an open-source, distributed processing system used for big data workloads, which utilizes in-memory caching and optimized query execution for fast analytic queries against data of any size. It provides development APIs in Java, Scala, Python and R, and in the part two of our project, we used the APIs in Python named Pyspark. For example:

```

1 %spark.pyspark
2 work_location_borough1 = dataset1.map(lambda x : (x[7], 1)).countByKey()
3 for key, value in work_location_borough1.items():
4     print(key, value)
5
6 work_location_borough2 = dataset2.map(lambda x : (x[2], 1)).countByKey()
7 for key, value in work_location_borough2.items():
8     print(key, value)
9
10 %spark.pyspark
11 work_location_borough3 = dataset3.map(lambda x : (x[4], 1)).countByKey()
12 for key, value in work_location_borough3.items():
13     print(key, value)
14 ... ..

```

By using Pyspark, we can easily obtain much faster speed to even concurrently find out the frequencies of each data field of each dataset. Even if our datasets become extremely large, Apache Spark or Pyspark still performs quite well in dealing with them with as fast speed as possible.

4 METHODS, ARCHITECTURE AND DESIGN

4.1 'Mid Name Initial' Attribute

Initial is an attribute that everyone has reached a consensus: only capital letters are permitted to appear. Therefore, the first thing we did was to find out the distinct values of this column. Since there should be at most 26 possible values, it is quite quick to locate the problems. Our results for part one are shown as follow:

initial	frequency
null	1,596,166
A	348,499
M	330,479
...	...
Q	2,181
.	131
-	86
1	71
x	39
2	30
...	...

Among all values above, null value has the greatest frequency. It means that most people don't have a middle name or they just don't want to write it down. What's more, there are some weird characters, including numbers, dashes, parentheses etc. We need to do judgements and turn them to null. Last but not least, there is a 'x' value. It's a tiny problem but it can be deadly. We need all the letters to be uppercase, just to be on the safe side.

4.2 'Title Description' Attribute

It's a little tricky when we are dealing with those columns which are stored as strings. First, a couple of typos exist such as 'RADIO AND TELEVISION OPERATOR'. Second, the same data may be written as several different forms. For example, 'SERGEANT D/A SPECIAL ASSIGNMENT' and 'SERGEANT-D/A SPECIAL ASSIGNMENT', 'SECRETARY TO COMMISSIONER' and 'SECRETARY OF COMMISSIONER', 'CHAUFFEUR ATTENDANT' and 'CHAUFFEUR-ATTENDANT'... It's difficult to make a decision. We are not sure which one is the correct answer, because we don't have reference data. Under such circumstances, we calculate the frequency of values in each similar pairs and choose the value with the highest frequency as the correct one. Furthermore, we find some cultural issues which are really interesting.

title description	frequency
AGENCY ATTORNEY INTERN	110
AGENCY ATTORNEY INTERNE	698

It's strange that the frequency of word 'INTERNE' is much higher than 'INTERN' which we thought would be the only right spelling of this word. As people from remote China, we are surprised to find 'INTERNE' is also right as a NORTH AMERICAN idiom. This makes the cleaning strategy more complicated. More factors should be taken into account.

4.2.1 *Original Method in Part1.*

In part1, since the dataset is relatively small, we used `reduceByKey()` on each value. Then, we did clustering to find similar strings by the methods we discussed above(3.2). The results are clear:

```

1 Cluster 28
2 SECRETARY TO THE CHAIR (x 1)
3 SECRETARY TO THE CHAIRMAN (x 11)
4 ...
5 Cluster 30
6 SECRETARY TO DEPARTMENT (x 7)
7 SECRETARY OF DEPARTMENT (x 1)
8 ...
9 Cluster 34
10 SUPV OF HOUSING EXTERMINATORS (x 5)
11 SUPV OF HOUSING EXTERMINATOR (x 23)
12 ...
13 Cluster 8
14 AGENCY ATTORNEY (x 8921)
15 *AGENCY ATTORNEY (x 9)
16 ...

```

Then, We clean data manually according to the rules we discussed above.

4.2.2 *Apply to new datasets.*

When we apply our methods to new datasets, we find some new issues. Some values are truncated suddenly: 'PROGRAM OFFICER (DEPT FOR THE'. It means details are lost. We need to find the correct information and fix them. This phenomenon always happen across a whole single dataset. It's impossible to do completion on data by do clustering on the same dataset(each value lost its details more or less). If we can get detail from another dataset, then the problem can be solved. In conclusion, we hope the datasets can complement each other.

4.2.3 *Data cleaning on scale.*

We find ten datasets that have a 'title' attribute. We find it is difficult to work with openclean and spark together in hpc. Therefore, we need to write our own clusters.

First, we read all the 'title' data.

```

1 %spark.pyspark
2 # title description
3
4 # Read All the data
5 # NYC Civil Service Titles 1
6 ds1 = sc.textFile('hdfs://user/CS-GY-6513/project_data/data-cityofnewyork-us.nzjr-3966.csv').mapPartitions(lambda x :
7     reader(x)).map(lambda x: [x[1], 1]).filter(lambda line : len(line)>1 and 'Title_Description' not in line)
8 # Civil Service Titles Subject to Investigation 1
9 ds2 = sc.textFile('hdfs://user/CS-GY-6513/project_data/data-cityofnewyork-us.h7fs-cqma.csv').mapPartitions(lambda x :
10     reader(x)).map(lambda x: [x[1], 1]).filter(lambda line : len(line)>1 and 'TITLE_DESCRIPTION' not in line)
11 ...
12 ...
13 # Appeals Closed In 2016 7
14 ds9 = sc.textFile('hdfs://user/CS-GY-6513/project_data/data-cityofnewyork-us.wbtw-zkex.csv').mapPartitions(lambda x :
15     reader(x)).map(lambda x: [x[7], 1]).filter(lambda line : len(line)>1 and 'Title' not in line)
16 # Civil Service List (Terminated) 7

```

```

365 14 ds10 = sc.textFile('hdfs://user/CS-GY-6513/project_data/data-cityofnewyork-us.qu8g-sxqf.csv').mapPartitions(lambda x :
366 reader(x)).map(lambda x: [x[7], 1]).filter(lambda line : len(line)>1 and 'List_Title_Desc' not in line)

```

union

```

369
370 1 all = ds1.union(ds2).union(ds3).union(ds4).union(ds5).union(ds6).union(ds7).union(ds8).union(ds9).union(ds10)
371 2 all.count()
372 3 amount: 2672736

```

Currently, we have got more than 2.5 million rows of title data. Next, we use spark to do mapreduce.

```

375
376 1 %spark.pyspark
377 2 title = all.reduceByKey(add)
378 3 title.count()
379 4 amount: 2818

```

We are excited to find the data amount is reduced to a very significant extent(2672736->2818). Now, it's time to do clustering. We think Levenshtein distance will help us to achieve the best results because our samples are all strings.

```

383
384 1 def LevenshteinDistance(s1, s2):
385 2     d = [[x for x in range(len(s1)+1)] for _ in range(len(s2)+1)]
386 3     for y in range(1, len(s2)+1):
387 4         d[y][0] = d[y-1][0] + 1
388 5     for x in range(1, len(s1)+1):
389 6         for y in range(1, len(s2)+1):
390 7             if s1[x-1] == s2[y-1]:
391 8                 d[y][x] = d[y-1][x-1]
392 9             else:
393 10                 substate = d[y-1][x-1] + 1
394 11                 add = d[y][x-1] + 1
395 12                 delete = d[y-1][x] + 1
396 13                 d[y][x] = min(add, substate, delete)
397 14     return d[-1][-1]

```

Among all the cluster algorithms, we think Affinity Propagation is the most suitable one, because we don't know how many groups we may need at the very beginning. We calculate the Levenshtein Distance between each string in the title sets. We use it as our similarity matrix.

```

402 1 %spark.pyspark
403 2 import numpy as np
404 3 import sklearn.cluster
405 4 title = title.collect()
406 5 similarity = -1*np.array([[LevenshteinDistance(t1,t2) for t1 in title] for t2 in title])
407 6 print(similarity)
408 7 affprop = sklearn.cluster.AffinityPropagation(affinity="precomputed", damping=0.5)
409 8 affprop.fit(similarity)
410 9 [[ 0 -38 -23 ... -28 -26 -39]
411 10 [-38 0 -39 ... -46 -42 -44]
412 11 [-23 -39 0 ... -19 -19 -40]
413 12 ...
414 13 [-28 -46 -19 ... 0 -16 -35]
415 14 [-26 -42 -19 ... -16 0 -33]
416 15 [-39 -44 -40 ... -35 -33 0]]

```

We still can't figure out how to use spark to do machine learning, so we have to ask sklearn for help. Although it takes some time to run the model, the results are nice. We do find the datasets can help each other! The data that lost details are always be clustered into the same group with the data with full details. So finally we know how to correct the dirty data!

data that lost detail	complement data
PROGRAM OFFICER (DEPT FOR THE	PROGRAM OFFICER (DEPARTMENT FOR THE AGING)
GENERAL SUPERINTENDENT (SANITA	GENERAL SUPERINTENDENT (SANITATION)
AREA SUPERVISOR (HIGHWAY MAINT	AREA SUPERVISOR (HIGHWAY MAINTENANCE)
SECRETARY OF THE DEPARTMENT (D	SECRETARY OF THE DEPARTMENT (DBS)
PUBLIC HEALTH NURSE (SCHOOL HE	PUBLIC HEALTH NURSE (SCHOOL HEALTH)
SUPERVISOR (WATER & SEWER SYST	SUPERVISOR (WATER AND SEWER SYSTEMS)
CONSULTANT (PUBLIC HEALTH-SOCI	CONSULTANT (PUBLIC HEALTH SOCIAL WORK)
MULTIPLE DWELLING SPECIALIST (MULTIPLE DWELLING SPECIALIST (BUILDINGS)
...	...

Still, we have to admit that we can not fix all the problems, for example, we can't find the corresponding value of 'DEPUTY COMMISSIONER (SPECIAL S', 'ASST COMMISSIONER FOR PUBLIC &', 'DIRECTOR OF INFORMATION AND RE' and so on. But we believe if the dataset is large enough, it is possible to solve all the problems.

4.3 'Agency Start Date' Attribute

4.3.1 Original Method in Part1.

In the NYCPayroll dataset, we used the following steps to do cleaning on the date column. Finding out the minimum and the maximum of dates can help us get a picture of the data and narrow the range. Furthermore, since it's an attribute that represents an event happened before, it cannot be a date in the future. And Of course, we should filter impossible dates as well.

- (1) find the min date and the max date
- (2) check whether there is a date in the future
- (3) check impossible date, E.g. 2021-13-01, 2021-2-30
- (4) for each abnormal value, do cleaning manually

```

1 agency_start_dates = ds.distinct('Agency_Start_Date')
2 for rank, val in enumerate(agency_start_dates.most_common()):
3     dt, freq = val
4     if dt == '':
5         continue
6     if int(dt.split('/')[1]) > 2021:
7         print(dt)
8     elif int(dt.split('/')[1]) > 31 or int(dt.split('/')[1]) < 1:
9         print(dt)
10    elif int(dt.split('/')[1]) > 30 and int(dt.split('/')[0]) in [4,6,9,11]:
11        print(dt)
12    elif int(dt.split('/')[1]) > 28 and int(dt.split('/')[0]) == 2 and int(dt.split('/')[1]) % 4 != 0:
13        print(dt)
14    elif int(dt.split('/')[1]) > 29 and int(dt.split('/')[0]) == 2 and int(dt.split('/')[1]) % 4 == 0:
15        print(dt)
16    elif int(dt.split('/')[0]) > 12 or int(dt.split('/')[0]) < 1:
17        print(dt)

```

Using the openclean package and the code above, we are able to find out abnormal value:

```
12/31/9999
10/16/2049
```

4.3.2 Apply to new datasets.

There is no doubt that this method works well on other datasets, because date is a kind of attribute which is relatively stable. But If we take efficiency into account, this method may not that good. For each date, we determine whether the day is valid, whether the month is valid and whether the year is valid. It's really time-consuming. What's more, we have to clean the dirty data manually like this:

```
1 def agency_start_date(x):
2     if x == '12/31/9999':
3         return ''
4     elif x == '10/16/2049':
5         return ''
6     else:
7         return x
8 ds = update(ds, columns='Agency_Start_Date', func = agency_start_date)
```

It is lucky we just meet with two problematic dates in this NYCPayroll dataset. What if there are a lot? We need a general and efficent method to do this.

So, we decide to use Spark. We are happy to find several amazing convenience Spark brings us during this process. If we use Spark to read data from a .csv file. The date attribute will be parsed as a String, which is not what we want. So we need to turn them back into date type. Spark offers us the to_date() method to do this. We use the following code to do this(just an example).

```
1 test = spark.createDataFrame([('1-30-2021',)], ['date'])
2 test.select(to_date(test.date, 'MM-dd-yyyy').alias('date')).show()
```

By observing the results, we will get some interesting facts

(p.s.:2020 is a leap year):

string	format	to_date()
01-32-2021	'MM-dd-yyyy'	null
02-29-2021	'MM-dd-yyyy'	null
02-29-2020	'MM-dd-yyyy'	datetime.date(2020,2,29)
02/29/2020	'MM-dd-yyyy'	null
02/29/2020	'MM/dd/yyyy'	datetime.date(2020,2,29)
02/29-2020	'MM/dd-yyyy'	datetime.date(2020,2,29)
2-29-2020	'MM-dd-yyyy'	datetime.date(2020,2,29)

Amazing! Spark help us do date checking automatically when we call the to_date() function. All we need to do is to match a suitable format String to the data, then all data will be edited to 'yyyy-MM-dd'.

Thanks to the great job by Spark, the only thing we need to do next is to check whether there exists future dates. To achieve the goal, we can use the sparksql:

```
1 spark.sql("select_*_from_ds_where_date_>_current_date()")
```

4.3.3 Refined method.

Through all the above analysis, we can get the final improved cleaning plan.

- (1) read raw data with pyspark
- (2) figure out what the date format is
- (3) turn string to date type
- (4) use sparksql to find whether future dates exist
- (5) do cleaning according to requirements(drop null records or leave them alone)

Following is the code, take NYCPayroll dataset for instance:

```

1 %spark.pyspark
2 ds = spark.read.csv('hdfs://user/CS-GY-6513/project_data/data-cityofnewyork-us.k397-673e.csv', header=True,
3 inferSchema=True)
4 ds = ds.withColumn('Agency_Start_Date', to_date(ds['Agency_Start_Date'], 'MM/dd/yyyy'))
5 #ds.printSchema()
6 #ds.show()
7 ds.createOrReplaceTempView('ds')
8 spark.sql("select_*_from_ds_where_`Agency_Start_Date`>current_date()").show()

```

4.4 'Agency Name' Attribute

4.4.1 Original Method in Part1.

In the NYCPayroll dataset, we used the following steps to do cleaning on date column. We tried different kinds of knn clusters to find data problems. Specifically, We use LevenshteinDistance, HammingDistance, and JaroSimilarity

- (1) Cluster for Agency Name used kNN cluster and the Levenshtein distance as the similarity measure.
- (2) Cluster for Agency Name used key collision and multiple threads (4) to generate value keys in parallel.
- (3) We change all the agency name to the upper case to make the data consistence
- (4) According to the outcome of the cluster, we change the same agency name to the same description

KNN cluster by using Levenshtein distance

```

1 minsize = 2
2 clusters = knn_clusters(
3     values=agency_names,
4     sim=SimilarityConstraint(func=JaroSimilarity(), pred=GreaterThan(0.9)),
5     minsize=minsize
6 )
7
8 def print_k_clusters(clusters, k=5):
9     clusters = sorted(clusters, key=lambda x: len(x), reverse=True)
10    val_count = sum([len(c) for c in clusters])
11    print('Total_number_of_clusters_is_{0}with_{1}_values'.format(len(clusters), val_count))
12    for i in range(min(k, len(clusters))):
13        print('\nCluster_{0}'.format(i + 1))
14        for key, cnt in clusters[i].items():
15            if key == '':
16                key = ''
17            print(f'_{key}_{cnt}')
18
19 print_k_clusters(clusters,13)

```

573 key collision with 4 threads

```

574 1 from openclean.cluster.key import key_collision
575 2 minsize = 2
576 3
577 4 clusters = key_collision(values=agency_names, minsize=minsize, threads=4)
578 5
579 6 print('{}_clusters_of_size_{}_or_greater'.format(len(clusters), minsize))
580 7 print_k_clusters(clusters)

```

582 due to the problem we find, we design the clean function as follow:

```

583 1 def agency_name(x):
584 2     if x != '':
585 3         x = x.upper()
586 4     if x == 'BOARD_OF_CORRECTIONS':
587 5         x = 'BOARD_OF_CORRECTION'
588 6     return x

```

590 However, if the dataset becomes larger and larger, we will face the problem that the dataset can not fit in to the memory
591 for us to clean it. So for the part two of the project, we change the way by using pyasppark to find the dataset problems
592 and clean the column of the Agency Name.

594 4.4.2 Apply to New Datasets in Part 2.

595 In order to make sure if the techniques for cleaning this kind of data columns are general and can be used to generate
596 reference data, we need to find 10 datasets that include 'Agency Name' columns. the dataset for the Agency Name is as
597 below:

```

599 1 dataset1 = sc.textFile("hdfs:/user/CS-GY-6513/project_data/data-cityofnewyork-us.ye3c-m4ga.csv").mapPartitions(lambda
600 x: reader(x))
601 2 dataset2 = sc.textFile("hdfs:/user/CS-GY-6513/project_data/data-cityofnewyork-us.423i-ukqr.csv").mapPartitions(lambda
602 x: reader(x))
603 3 dataset3 = sc.textFile("hdfs:/user/CS-GY-6513/project_data/data-cityofnewyork-us.gzfs-3h4m.csv").mapPartitions(lambda
604 x: reader(x))
605 4 dataset4 = sc.textFile("hdfs:/user/CS-GY-6513/project_data/data-cityofnewyork-us.mdcw-n682.csv").mapPartitions(lambda
606 x: reader(x))
607 5 dataset5 = sc.textFile("hdfs:/user/CS-GY-6513/project_data/data-cityofnewyork-us.k397-673e.csv").mapPartitions(lambda
608 x: reader(x))
609 6 dataset6 = sc.textFile("hdfs:/user/CS-GY-6513/project_data/data-cityofnewyork-us.mwzb-yiwb.csv").mapPartitions(lambda
610 x: reader(x))
611 7 dataset7 = sc.textFile("hdfs:/user/CS-GY-6513/project_data/data-cityofnewyork-us.a9md-ynri.csv").mapPartitions(lambda
612 x: reader(x))
613 8 dataset8 = sc.textFile("hdfs:/user/CS-GY-6513/project_data/data-cityofnewyork-us.vx8i-nprf.csv").mapPartitions(lambda
614 x: reader(x))
615 9 dataset9 = sc.textFile("hdfs:/user/CS-GY-6513/project_data/data-cityofnewyork-us.xrwg-eczf.csv").mapPartitions(lambda
616 x: reader(x))
617 10 dataset10 = sc.textFile("hdfs:/user/CS-GY-6513/project_data/data-cityofnewyork-us.qbjq-atxv.csv").mapPartitions(
618 lambda x: reader(x))

```

619 In the first part, the main clean method for the Agency Name is:

```

620 1 def clean(x):
621 2     x = toUpperCase(x)
622 3     return x
623 4 def toUpperCase(x):

```

624 Manuscript submitted to ACM

```

5     if x != '':
6         return x.upper();
7     else:
8         return "_"

```

we apply it to the 10 dataset and get the 1st version of the clean data for each dataset. In order to do analysis and compare the result with the refine way, we named the 1st version of the clean dataset as dataset old.

the spark funtion for the clean is as below:

```

1 agency_name_old1 = dataset1.map(lambda x: (clean(x[3]),1)).countByKey()
2 for key, value in agency_name1.items():
3     print(key, value)

```

4.4.3 refine the part1 clean method for new datasets.

there are some new problems in the new datasets

- (1) some of the dataset have a specific length for the attribute, if the Agency Name length is small, there will be a space in the String to make the same Agency Name count more than twice when we do analysis
- (2) some of the dataset have a specific length for the attribute, if the Agency Name length is large, some part of the String will be lost, and the information is not complete
- (3) there will be different way to represent a same Agency Name for example, the following two tables have the same type of data but with completely different present way.

1st dataset	2nd dataset
BRONX COMMUNITY BOARD 2	COMMUNITY BOARD NO.4-MANHA
BRONX COMMUNITY BOARD 3	COMMUNITY BOARD NO.9-QUEEN
BRONX COMMUNITY BOARD 4	COMMUNITY BOARD NO.17-BROO
BRONX COMMUNITY BOARD 5	COMMUNITY BOARD NO.8-BRONX
BRONX COMMUNITY BOARD 6	COMMUNITY BOARD NO.5-MANHA
BRONX COMMUNITY BOARD 7	COMMUNITY BOARD NO.10-BROO
BRONX COMMUNITY BOARD 8	COMMUNITY BOARD NO.11-QUEE
BRONX COMMUNITY BOARD 9	COMMUNITY BOARD NO.2-QUEEN
BRONX COMMUNITY BOARD 10	COMMUNITY BOARD NO.7-QUEEN
BRONX COMMUNITY BOARD 11	COMMUNITY BOARD NO.12-MANH
BRONX COMMUNITY BOARD 12	COMMUNITY BOARD NO.4-BROOK
BRONX COMMUNITY BOARD 1	COMMUNITY BOARD NO.6-BROOK

We check lots of other datasets and find that the first version to present it is more frequently, and the second version with incomplete district name. Therefore, to make the data more consistent, we decided to change the 2nd version to the 1st one.

According to the new problems we found, we design a new clean method to refine the original one

- (1) we delete the space at the front and end of each String

```

1     def elimiateBlank(x):
2         x = x.strip()
3         return x;

```

- (2) we design a identical function to make two types of representations become the same

```

677 1 def identical(x):
678 2     if len(x.split('_')) == 3 and x.split('_')[0] == 'COMMUNITY' and x.split('_')[1] == 'BOARD':
679 3         number = x.split('.')[1].split('-')[0]
680 4         reg = x.split('-')[1]
681 5         print(len(reg))
682 6         reg = reg[0:4]
683 7         print(reg)
684 8         if reg == 'QUEE':
685 9             c_reg = 'QUEENS'
686 10            x = c_reg + '_' + 'COMMUNITY_BOARD_' + number
687 11            if reg == 'BROO':
688 12                c_reg = 'BROOKLYN'
689 13                x = c_reg + '_' + 'COMMUNITY_BOARD_' + number
690 14            if reg == 'BRON':
691 15                c_reg = 'BRONX'
692 16                x = c_reg + '_' + 'COMMUNITY_BOARD_' + number
693 17            if reg == 'MANH':
694 18                c_reg = 'MANHATTAN'
695 19                x = c_reg + '_' + 'COMMUNITY_BOARD_' + number
696 20            if reg == 'RICH':
697 21                c_reg = 'RICHMOND'
698 22                x = c_reg + '_' + 'COMMUNITY_BOARD_' + number
699 23        return x

```

- (3) we manually complete the Agency Name for the specific rows to let the dataset become more complete, then we can get the references data

```

702 1 def completeTheDatas(x):
703 2     if x == 'DEPARTMENT_OF_CITYWIDE_ADM':
704 3         x = 'DEPARTMENT_OF_CITYWIDE_ADMINISTRATIVE_SERVICES'
705 4     if x == 'HRA/DEPARTMENT_OF_SOCIAL_S':
706 5         x = 'HRA/DEPARTMENT_OF_SOCIAL_SERVICES'
707 6     if x == 'CUNY_QUEENSBOROUGH_COMMUNI':
708 7         #extra...
709 8     return x

```

the new clean method is as below

```

712 1 def clean_new(x):
713 2     x = toUpperCase(x)
714 3     x = elimiateBlank(x)
715 4     x = identical(x)
716 5     return x

```

Then, we apply the new clean method to the dataset and do the celan through spark to get the 2nd version of the clean data. For each dataset, we named it as dataset new

```

721 1 agency_name_new1 = dataset1.map(lambda x: (clean_new(x[3]),1)).countByKey()
722 2 for key, value in agency_name_new.items():
723 3     print(key, value)

```

For some dataset, we needed manually update. We usee the clean method as below, through spark to get the 2nd version of the clean data. For each dataset, we named it as dataset d. And in the analysis part, we assumed the dataset d as the correct dataset.

```

1 def clean_for_dataset1(x):
2     x = toUpperCase(x)
3     x = elimiateBlank(x)
4     x = identical(x)
5     x = completeTheDatas(x)
6     return x

```

```

1 agency_name_d1 = dataset1.map(lambda x: (clean_for_dataset1(x[3]),1)).countByKey()
2 for key, value in agency_name_d1.items():
3     print(key, value)

```

4.4.4 Create the reference dataset for Agency Name.

by combining the ten datasets from above, we can get a Agency Name reference data

```

1 reference_data_agency_name = set(agency_name_d10).union(set(agency_name_d9).union(set(agency_name_d8).union(set(
    agency_name_d7).union(set(agency_name_d6).union(set(agency_name_d5).union(set(agency_name_d4).union(set(
    agency_name_d3).union(set(agency_name_d2).union(set(agency_name_d1))))))))))

```

4.5 'Work Location Borough' Attribute

4.5.1 Original Method in Part 1.

In the NYCPayroll dataset, we used the following steps to do cleaning on the date column. Finding out the most frequent values in column 'Work Location Borough' gave us an overall picture of this data column, letting us realize 'MANHATTAN', 'QUEENS', 'BROOKLYN' and 'BRONX' are among the top four most popular work location boroughs. In addition, getting a set of distinct values for this column made us know how many different work locations in total are there as well as their counts respectively. Finally, we used different kinds of knn cluster methods to find out more data issues and analyze them accordingly.

- (1) find out the frequencies of top values
- (2) find out the frequencies of all distinct values
- (3) cluster work location borough using key collision
- (4) change all the data fields to be upper cases
- (5) change empty values to be "UNKNOWN"

```

1 profiles.column('Work_Location_Borough').get('topValues')
2 wlb = ds.distinct('Work_Location_Borough')
3 for rank, val in enumerate(wlb.most_common()):
4     wlb, freq = val
5     print(f'{rank+1:<3}_{wlb}_{freq:>10,}')
6
7 def print_k_clusters(clusters, k=5):
8     clusters = sorted(clusters, key=lambda x: len(x), reverse=True)
9     val_count = sum([len(c) for c in clusters])
10    print('Total_number_of_clusters_is_{0}_with_{1}_values'.format(len(clusters), val_count))
11    for i in range(min(k, len(clusters))):
12        print('\nCluster_{0}'.format(i + 1))
13        for key, cnt in clusters[i].items():
14            if key == '':
15                key = ""
16            print(f'_{key}_{cnt}')

```

```

781 17
782 18 work_locations = ds.select('Work_Location_Borough').distinct()
783 19 clusters = KeyCollision(func=Fingerprint()).clusters(work_locations)
784 20
785 21 def work_location_borough(x):
786 22     if x == '':
787 23         return "UNKNOWN"
788 24     else:
789 25         return x.upper()
790 26
791 27 ds = update(ds, columns='Work_Location_Borough', func = work_location_borough)

```

By running the code above and utilizing the openclean package, we were able to find the following abnormal values:

```

792 BRONX
793
794 bronx
795
796 MANHATTAN
797
798 manhatton
799
800 QUEENS
801
802 queens
803
804 RICHMOND
805
806 Richmond
807
808 (null)

```

However, as mentioned before, if our dataset becomes increasingly larger, we will face the challenge that the dataset can't be fed into the limited memory. Therefore, we'll use spark/pyspark to continue our work on part two and do the data cleaning on columns which are similar to or the same as 'work location borough'.

4.5.2 Apply to New Datasets in Part 2.

In order to make sure if the techniques for cleaning this kind of data columns are general and can be used to generate reference data, we need to find as many datasets as possible that include 'borough' columns.

```

813 1 %spark.pyspark
814 2 dataset1 = sc.textFile("hdfs://user/CS-GY-6513/project_data/data-cityofnewyork-us.k397-673e.csv").mapPartitions(lambda
815 3     x: reader(x))
816 3 dataset2 = sc.textFile("hdfs://user/CS-GY-6513/project_data/data-cityofnewyork-us.ic3t-wcy2.csv").mapPartitions(lambda
817 4     x: reader(x))
818 4 dataset3 = sc.textFile("hdfs://user/CS-GY-6513/project_data/data-cityofnewyork-us.rbx6-tga4.csv").mapPartitions(lambda
819 5     x: reader(x))
820 5 dataset4 = sc.textFile("hdfs://user/CS-GY-6513/project_data/data-cityofnewyork-us.vz8c-29aj.csv").mapPartitions(lambda
821 6     x: reader(x))
822 6 dataset5 = sc.textFile("hdfs://user/CS-GY-6513/project_data/data-cityofnewyork-us.6khm-nrue.csv").mapPartitions(lambda
823 7     x: reader(x))
824 7 dataset6 = sc.textFile("hdfs://user/CS-GY-6513/project_data/data-cityofnewyork-us.avir-tzek.csv").mapPartitions(lambda
825 8     x: reader(x))
826 8 dataset7 = sc.textFile("hdfs://user/CS-GY-6513/project_data/data-cityofnewyork-us.xywu-7bv9.csv").mapPartitions(lambda
827 9     x: reader(x))
828 9 dataset8 = sc.textFile("hdfs://user/CS-GY-6513/project_data/data-cityofnewyork-us.3qem-6v3v.csv").mapPartitions(lambda
829 10    x: reader(x))
830 10 dataset9 = sc.textFile("hdfs://user/CS-GY-6513/project_data/data-cityofnewyork-us.ur7y-ziyb.csv").mapPartitions(lambda
831 11    x: reader(x))
832 11 dataset10 = sc.textFile("hdfs://user/CS-GY-6513/project_data/data-cityofnewyork-us.esmb-8zkm.csv").mapPartitions(

```



```

12
13 %spark.pyspark
14 work_location_borough1 = dataset1.map(lambda x : (x[7], 1)).countByKey()
15 for key, value in work_location_borough1.items():
16     print(key, value)
17
18 work_location_borough2 = dataset2.map(lambda x : (x[2], 1)).countByKey()
19 for key, value in work_location_borough2.items():
20     print(key, value)
21
22 %spark.pyspark
23 work_location_borough3 = dataset3.map(lambda x : (x[4], 1)).countByKey()
24 for key, value in work_location_borough3.items():
25     print(key, value)
26
27 ... ..

```

The above code locates ten datasets that contain columns related to New York boroughs. By using pyspark, we can easily obtain much faster speed to even concurrently find the frequencies of each data field of each dataset. After that, all the related issues become crystal clear merely by inspecting those data fields as well as their frequencies:

1. Characters of a location borough are not capitalized.
2. Empty string of a location borough.
3. Some white-spaces at the front or the tail of a location borough.

4.5.3 Refine the Techniques Used in Part 1 for New Datasets.

First, it is both reasonable to make all of the characters of a location borough to be lower cases or capitalize each character given a non-empty location borough, and we decided to choose the latter. Besides, there are multiple ways of dealing with empty string location boroughs, and we decided to change them to be "UNKNOWN", which is more reasonable because we actually don't know what data fields they are. Finally, in terms of the white-spaces, we decided to eliminate them with no doubt. Therefore, the refined cleaning techniques are as follows:

```

1 def clean(x):
2     x = x.strip()
3     if x != '':
4         return x.upper()
5     else:
6         return "UNKNOWN"

```

For example, we'll apply the techniques to the dataset seven. Before applying the techniques:

```

1 %spark.pyspark
2 work_location_borough7 = dataset7.map(lambda x : (x[1], 1)).countByKey()
3 for key, value in work_location_borough7.items():
4     print(key, value)
5
6 NYC Total 1
7 Bronx 1
8 Brooklyn 1
9 Manhattan 1
10 Queens 1
11 Staten Island 1

```

After applying the refined techniques:

```

1 %spark.pyspark
2 %spark.pyspark
3 temp = dataset7.map(lambda x: (clean(x[1]), 1)).countByKey()
4 for key, value in temp.items():
5     print(key, value)
6
7 NYC TOTAL 1
8 BRONX 1
9 BROOKLYN 1
10 MANHATTAN 1
11 QUEENS 1
12 STATEN ISLAND 1

```

It's now clear that the relevant issues with the dataset seven have been resolved.

5 RESULTS

For big data analysis, we need data to be as clean as possible to avoid the incorrect data record affecting our data research. So the outcome of the data clean results is extremely important. Since the above method we proposed cannot change correct data to the wrong data, we decided to use recall for the String part of the analysis. And for other data, we finally get a precise version to check the data problem and fix it to the right version.

5.1 The cleaning results of each attribtue

5.1.1 Mid Name initial.

It is easy to get the Right Mid Name initial by the above method, we can just pick all the characters form 'A' to 'Z' and make everyone is uppercase. It can have a reference data just be 'A,B,C,.....,Z'.

5.1.2 Agency Start Date.

To dealing with the Date attribute, we can find the illegal date and future date by sql. So for each dataset, we can set all the illegal date and futrue data to the Null value. The only error data can exist in the dataset is that some mistaken input with the right date in the past. It just a tiny part of the all datasets, and will not have any big influence in the following data analysis.

5.1.3 Agency Name.

We use the recall to check the percentage of the dirty data and the numbers of the data we clean

- (1) the whole number of data needed to be clean is: the precision data length to subtract the length of the intersection between the origin data and the precision data
 - (2) the original clean data volume is: the difference of the origin dataset and the 1st clean dataset
 - (3) the new clean data volume is: the difference of the origin dataset and the 2nd clean dataset
-

```

1 arr_origin = []
2 arr_old = []
3 arr_new = []
4 arr_pres = []
5 for key, value in agency_name1.items():
6     arr_origin.append(key)

```

Manuscript submitted to ACM

```

7 for key, value in agency_name_old1.items():
8     arr_old.append(key)
9 for key, value in agency_name_new1.items():
10    arr_new.append(key)
11 for key, value in agency_name_d1.items():
12    arr_pres.append(key)
13
14 clean_recall = len(set(arr_old).difference(set(arr_origin))) / (len(arr_pres) - len(set(arr_origin).intersection(set(
15    arr_pres))))
16 clean_new_recall = len(set(arr_new).difference(set(arr_origin))) / (len(arr_pres) - len(set(arr_origin).intersection(
17    set(arr_pres))))
18
19 print(clean_recall)
20 print(clean_new_recall)

```

The result of the anylasis in the 10 dataset is as below:

dataset	clean recall	clean new recall
dataset1	0.0	0.6129
dataset2	1.0	1.0
dataset3	1.0	1.0
dataset4	1.0	1.0
dataset5	1.0	1.0
dataset6	1.0	1.0
dataset7	1.0	1.0
dataset8	0.5	0.5
dataset9	1.0	1.0
dataset10	0.0238	1.0

As we can see from the above table, the accuracy of the first clean method is lower than the new clean method for some datasets, which may have the problem of space in the end and front. By using the new clean method, we can get cleaner data, but there still will have some data inconsistent that need us to manually correct. And then we can use the final data to generate the Agency Name reference data.

5.1.4 Work Location Borough.

Inherently the refined techniques have merely one more operation than the original methods for cleaning the 'work location borough' column - eliminating the leading and trailing white-spaces. This is because when we did data cleaning on this column in part one, there was no such issues within that dataset. When we manage to experiment with more datasets in part two, such issues become clear.

Therefore, it becomes quite simple in this case to measure their effectiveness. The number of operations in the original techniques are $m = 2$; The number of operations in the refined techniques are $n = 3$; Suppose at most 3 operations are needed to finish data cleaning with regard to columns which are similar to Work Location Borough, we can see that $m / n = 2/3$, and thus $2/3$ should be the overall performance evaluation of our original data cleaning techniques.

5.2 The reference data for some attributes

5.2.1 Agency Name.

The final reference data generated has 504 lines.

```

1 0
2 1      DEPARTMENT OF YOUTH AND COMMUNITY DEVELOPMENT
3 2      CITY UNIVERSITY CONSTRUCTION FUND
4 3      OFFICE OF COLLECTIVE BARGAINING
5 4      SANITATION
6 ..      ...
7 501     PORT AUTHORITY OF NEW YORK & NEW JERSEY
8 502     SMALL BUSINESS SERVICES
9 503     DEPT OF PARKS & RECREATION
10 504    DISABILITIES, OFFICE FOR PEOPLE WITH DEVELOPME...
11
12 [505 rows x 1 columns]
13 FINISHED

```

5.2.2 Work Location Borough Reference Data Generation.

The final reference data generated are as follows:

```

1 BROOKLYN
2 MANHATTAN
3 BRONX
4 RICHMOND
5 QUEENS
6 WASHINGTON DC
7 ULSTER
8 WESTCHESTER
9 ALBANY
10 NASSAU
11 DELAWARE
12 SULLIVAN
13 ORANGE
14 SCHOHARIE
15 DUTCHESS
16 GREENE
17 PUTNAM
18 STATEN ISLAND

```

6 REFERENCES

<https://data.cityofnewyork.us/City-Government/Civil-List/ye3c-m4ga>
<https://data.cityofnewyork.us/City-Government/Local-Law-18-Pay-and-Demographics-Report-Agency-Re/423i-ukqr>
<https://data.cityofnewyork.us/City-Government/Agency-Spending-by-Budget-Function/gzfs-3h4m>
<https://data.cityofnewyork.us/City-Government/Greenbook/mdcw-n682>
<https://data.cityofnewyork.us/City-Government/Expense-Budget/mwzb-yiwb>
<https://data.cityofnewyork.us/City-Government/Citywide-Payroll-Data-Fiscal-Year-/k397-673e>
<https://data.cityofnewyork.us/City-Government/Civil-Service-List-Active-/vx8i-nprf>
<https://data.cityofnewyork.us/City-Government/Civil-Service-List-Certification/a9md-ynri>

1041 <https://data.cityofnewyork.us/City-Government/SCOUT-CORE/xrwg-eczf>
1042 <https://data.cityofnewyork.us/City-Government/FY2018-PMMR-Spending-and-Budget/qbjq-atxv>
1043 <https://data.cityofnewyork.us/City-Government/Citywide-Payroll-Data-Fiscal-Year-/k397-673e>
1044 <https://data.cityofnewyork.us/Housing-Development/DOB-Job-Application-Filings/ic3t-wcy2>
1045 <https://data.cityofnewyork.us/Housing-Development/DOB-NOW-Build-Approved-Permits/rbx6-tga4>
1046 <https://data.cityofnewyork.us/Education/Borough-Enrollment-Offices/vz8c-29aj>
1047 <https://data.cityofnewyork.us/Social-Services/Demographics-by-Borough/6khm-nrue>
1048 <https://data.cityofnewyork.us/City-Government/New-York-City-Population-by-Borough-1950-2040/xywu-7bv9>
1049 <https://data.cityofnewyork.us/Social-Services/Buildings-by-Borough-and-Community-District/3qem-6v3v>
1050 <https://data.cityofnewyork.us/Social-Services/Associated-Address-by-Borough-and-Community-Distri/ur7y-ziyb>
1051 <https://data.cityofnewyork.us/Education/School-Based-Programs-by-Borough/esmb-8zkm>
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092