

Initial Cleaning:

1. Taking the original data (found in so_etl_code/inspections.csv) I utilized the classes CleanMapper1.class, CleanerReducer.class, and Clean.class (all located in etl_code/so_etl_code/Initial Cleaning) to clean the data. The new clean data is in data_ingest/so_data_ingest/inspections_clean2.csv.

Data cleaning took all the rows from the original dataset that had adhered to our desired schema. For example, boroughs must be a member of the set {Manhattan, Brooklyn, Bronx, Queens, Staten Island}, zip must be a number, critical must be a member of the set {Critical, Not Critical, Not Applicable}, etc.

Second Cleaning:

I cleaned the data further by conducting text formatting and creating a binary column. For the text formatting, I made all records without a valid grade (those with a null, N, P, or Z value) to be normalized to a value of N. For creating a binary column, each record with a health score better than the mean (17, see data profiling section) got a value of 1, while all that had a worse record got a value of 0. The code cleaning was done by running the etl_code/so_etl_code/Second Cleaning/cleaner.jar file on the cleaned data set (etl_code/so_etl_code/Second Cleaning/inspections_clean2.csv) and the result is stored in etl_code/so_etl_code/Second Cleaning/inspections_clean3.csv. Sample running script: `hadoop jar cleaner.jar Clean inspections_clean2.csv` output. Below is a screenshot of what part of the new data looks like:

```
Brooklyn,11201,02/19/2020,Not Critical,9,A,40.693719525273,-73.985694530011,0
Brooklyn,11201,02/20/2020,Not Critical,2,A,40.703293906326,-73.992047363499,0
Brooklyn,11201,02/21/2023,Critical,10,A,40.690826232689,-73.98345224546,0
Brooklyn,11201,02/21/2023,Critical,11,A,40.687022995402,-73.993661083746,0
Brooklyn,11201,02/21/2023,Critical,12,A,40.687692441332,-73.989795612161,0
Brooklyn,11201,02/21/2023,Critical,12,A,40.70247842835,-73.988700542295,0
Brooklyn,11201,02/21/2023,Critical,13,A,40.690044411028,-73.98683839769,0
Brooklyn,11201,02/21/2023,Critical,13,A,40.690343671862,-73.987559525333,0
Brooklyn,11201,02/21/2023,Critical,29,N,40.694376195237,-73.992845391234,1
Brooklyn,11201,02/21/2023,Critical,32,N,40.688364794871,-73.98856592177,1
Brooklyn,11201,02/21/2023,Critical,36,N,40.697845490794,-73.991434937772,1
Brooklyn,11201,02/21/2023,Critical,9,A,40.694886823634,-73.994698913437,0
Brooklyn,11201,02/21/2023,Not Critical,,N,40.68985672822,-73.979720371884,0
Brooklyn,11201,02/21/2023,Not Critical,11,A,40.687022995402,-73.993661083746,0
Brooklyn,11201,02/21/2023,Not Critical,12,A,40.70247842835,-73.988700542295,0
Brooklyn,11201,02/21/2023,Not Critical,13,A,40.690343671862,-73.987559525333,0
Brooklyn,11201,02/21/2023,Not Critical,29,N,40.694376195237,-73.992845391234,1
Brooklyn,11201,02/21/2023,Not Critical,32,N,40.688364794871,-73.98856592177,1
Brooklyn,11201,02/21/2023,Not Critical,36,N,40.697845490794,-73.991434937772,1
Brooklyn,11201,02/21/2023,Not Critical,4,A,40.68985672822,-73.979720371884,0
Brooklyn,11201,02/21/2023,Not Critical,9,A,40.694886823634,-73.994698913437,0
Brooklyn,11201,02/22/2022,Critical,0,N,40.690826232689,-73.98345224546,0
Brooklyn,11201,02/22/2022,Critical,111,N,40.684714524862,-73.991836855705,1
Brooklyn,11201,02/22/2022,Not Critical,0,N,40.690826232689,-73.98345224546,0
Brooklyn,11201,02/22/2022,Not Critical,111,N,40.684714524862,-73.991836855705,1
Brooklyn,11201,02/22/2023,Critical,11,A,40.70157861295,-73.995527884559,0
Brooklyn,11201,02/22/2023,Critical,36,N,40.6917323815,-73.986282737649,1
```

Data Profiling:

* All profiling was done with MapReduce

Unclean data can be found at: `etl_code/so_etl_code/inspections.csv`

Clean data can be found at: `etl_code/so_etl_code/Second Cleaning/inspections_clean3.csv`

There were several stages of profiling. First, determining how many unique records there are.

1. This was done by simply counting the number of records (lines) present in our CSVs.
2. The original, unclean dataset had 206149 records. This was found by conducting a map-reduce line count on the original, uncleaned data. Sample script: “`hadoop jar CountRecs.jar CountRecsinspections.csv output`.” Here is a screenshot of this profiling that was done for homework 7:

```
rvs6736_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -cat hw7/CountRecs/output/part-r-00000
Total number of records in dataset:      206149
```

3. The clean dataset had 104265. We lost many records due to incorrect adherence to our desired schema. Here is a screenshot of this profiling that was done for homework 7:

```
rvs6736_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -cat hw7/CountRecs/output2/part-r-00000
Total number of records in dataset:      104265
```

Next, we had counting unique grade occurrences, unique critical violation categorizations occurrences, unique zip codes, and unique borough occurrences.

Here are the results from all of these profiling processes. The necessary jar files are all included in the `profiling_code/so_profiling_code` directory. Here are the results from homework 8:

```

rvs6736_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -cat hw8/UniqueCriticalCount/part-r-00000
Critical      52326
Not Applicable 5727
Not Critical   52454
rvs6736_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -cat hw8/UniqueGradeCount/part-r-00000
      46127
A      52218
B      4464
C      1699
N      2876
P      507
Z      2616
rvs6736_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -cat hw8/uniqueBoroughCounts
Bronx      10432
Brooklyn   29288
Manhattan  41403
Queens     25494
Staten Island 3890
rvs6736_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -cat hw8/UniqueZipcodesCount/part-r-00000
Total Unique Zipcodes: 219

```

I also wanted to know the mean score... it was 17.

```

rvs6736_nyu_edu@nyu-dataproc-m:~$ hdfs dfs -cat hw8/meanScore/part-r-00000
Average Score: 17

```

This can be reproduced by running the meanScore.jar with the MeanScore class on the clean data.