%sh

python -m pip install socrata-py

python -m pip install openclean

python -m pip install openclean
python -m pip install openclean-geo
python -m pip install geopy

Defaulting to user installation because normal site-packages is not writeable

Requirement already satisfied: socrata-py in ./.local/lib/python3.7/site-packages (1.0.10)

Requirement already satisfied: requests in /share/apps/peel/python/3.7.9/gcc/lib/python3.7/site-packages (from socrata-py) (2.25.1)

Requirement already satisfied: urllib3<1.27,>=1.21.1 in /share/apps/peel/python/3.7.9/gcc/lib/python3.7/site-packages (from requests->socrata-py) (1.26.4)

Requirement already satisfied: idna<3,>=2.5 in /share/apps/peel/python/3.7.9/gcc/lib/python3.7/site-packages (from requests->socrata-py) (2.10)

Requirement already satisfied: certifi>=2017.4.17 in /share/apps/peel/python/3.7.9/gcc/lib/python3.7/site-pac kages (from requests->socrata-py) (2020.12.5)

Requirement already satisfied: chardet<5,>=3.0.2 in /share/apps/peel/python/3.7.9/gcc/lib/python3.7/site-pack ages (from requests->socrata-py) (4.0.0)

WARNING: You are using pip version 20.1.1; however, version 21.3.1 is available.

You should consider upgrading via the '/share/apps/peel/python/3.7.9/gcc/bin/python -m pip install --upgrade pip' command.

Defaulting to user installation because normal site-packages is not writeable

 $Requirement\ already\ satisfied:\ openclean\ in\ /share/apps/peel/openclean/0.2.1/lib/python 3.7/site-packages\ (0.2.1/lib/python 3.7/site-packages)$

~ 4 \

Took 6 sec. Last updated by nbuser at December 12 2021, 5:11:10 AM.

```
%pyspark
historical_dob_df = spark.read.options(header='true').csv('/user/CS-GY-6513/project_data/data-cityofnewyork
dob_cert_occupancy_df = spark.read.options(header='true').csv('/user/CS-GY-6513/project_data/data-cityofnew
housing_litigations_df = spark.read.options(header='true').csv('/user/CS-GY-6513/project_data/data-cityofnew
housing_maintenance_code_complaints_df = spark.read.options(header='true').csv('/user/CS-GY-6513/project_data
housing_maintenance_code_violations_df = spark.read.options(header='true').csv('/user/CS-GY-6513/project_data
historical_dob_df.createOrReplaceTempView('historical_dob')
dob_cert_occupancy_df.createOrReplaceTempView('dob_cert_occupancy')
housing_litigations_df.createOrReplaceTempView('housing_litigations')
housing_maintenance_code_complaints_df.createOrReplaceTempView('housing_maintenance_code_complaints')
housing_maintenance_code_violations_df.createOrReplaceTempView('housing_maintenance_code_violations')
```

Took 1 sec. Last updated by nbuser at December 12 2021, 5:41:33 AM.

%pyspark
historical_dob_df.select('BIN', 'Number', 'Street', 'Postcode', 'BOROUGH').show()

#------+
BIN|Number|
Street|Postcode|BOROUGH|

++-	+-	+	+-	+
2118801	2960	WEBSTER AVENUE	10458	BRONX
2096812	100	DEKRUIF PLACE	10475	BRONX
2008604	1898	HARRISON AVENUE	10453	BRONX
2007652	1998	MORRIS AVENUE	10453	BRONX
2084155	565	WEST 235 STREET	10463	BRONX
2012264	606	EAST FORDHAM ROAD	10458	BRONX
[2103486]	730 C0	ONCOURSE VILLAGE	10451	BRONX I

```
|2000391|
            345
                        BROOK AVENUE
                                         10454
                                                 BRONX
|2011594|
           4487
                        THIRD AVENUE
                                         10457
                                                 BRONX |
|2001106|
            575
                       WALTON AVENUE
                                         10451
                                                 BRONX
|2119514|
                         BEACON LANE
             69|
                                         10473
                                                 BRONX
|200 <del>| 3</del>50 |
                      PISSUMNCE61|
                                                 BRONX
                                         10471
                                                 BRONX |
|2085168|
            4441
                          259 STREET
            4401
Langeage
                       DEMONITE DI ACEL
```

Took 0 sec. Last updated by anonymous at December 11 2021, 4:05:48 AM.

%pyspark SPARK JOB (http://hc02.nyu.cluster:39642/jobs/job?id=11) FINISHED historical_dob_df.count()

2428526

Took 6 sec. Last updated by anonymous at December 11 2021, 4:15:47 AM.

%pyspark SPARK JOB (http://hc02.nyu.cluster:39642/jobs/job?id=25) FINISHED historical_dob_df['Postcode'].isNull()).filter(historical_dob_df['Latitude'].isNull()).

761

Took 8 sec. Last updated by anonymous at December 11 2021, 4:56:04 AM.

%pyspark
historical_dob_distinct_building_df = historical_dob_df.select('BIN', 'Number', 'Street', 'Postcode', 'BOROU(

Took 1 sec. Last updated by anonymous at December 11 2021, 4:17:01 AM.

%pyspark
historical_dob_distinct_building_cnt = historical_dob_distinct_building_df.groupby('BIN').count()

Took 0 sec. Last updated by anonymous at December 11 2021, 4:42:38 AM.

%pyspark

■ SPARK JOB (http://hc02.nyu.cluster:39642/jobs/job?jd=14) FINISHED historical_dob_distinct_building_cnt['count'] > 1).show()

```
+----+
     BIN|count|
+----+
 | <u>1</u>01<u>58</u>27 |
)@B₁Permit Issuance
|1006620|
|1082303|
             5|
|1005754|
             4
|1024757|
            22
|4533191|
             2|
|2057310|
             3|
|3340850|
             2|
|1038572|
             2|
             2|
|4005275|
|5149496|
             2
|2023445|
             2
|2007940|
             3|
104040771
Took 11 sec. Last updated by anonymous at December 11 2021, 4:43:18 AM.
```

```
%pvspark
                                                                                SPARK JOB FINISHED
historical_dob_distinct_building_df.filter(historical_dob_distinct_building_df['BIN'] == 1024757).show()
+----+
                       Street|Postcode| BOROUGH|
    BIN|Number|
+----+
                      BROADWAY
                                10036 | MANHATTAN |
|1024757| 1590|
|1024757| 1588|
                      BROADWAY
                                10036 | MANHATTAN |
                SEVENTH AVENUE
|1024757|
          708
                                 10036 | MANHATTAN |
|1024757|
           708
                       7 AVENUE
                                 10036 | MANHATTAN |
|1024757|
             2|
                   TIME SQUARE
                                  10036 | MANHATTAN |
                    7TH AVENUE
                                  10036 | MANHATTAN |
|1024757|
           714
|1024757|
             2|
                  TIMES SQUARE
                                  10036 | MANHATTAN |
|1024757|
         1592
                      BROADWAY
                                  10036 | MANHATTAN |
|1024757|
           710|
                SEVENTH AVENUE
                                  10036 | MANHATTAN |
                      BROADWAY
                                  10036 | MANHATTAN |
|1024757| 1584|
|1024757|
           704
                      7 AVENUE
                                  10036 | MANHATTAN |
|1024757|
                      7 AVENUE
                                  10036 | MANHATTAN |
           714
|1024757|
          1580
                      BROADWAY
                                  10036 | MANHATTAN |
|1024757|
          1576
                      BROADWAY
                                  10036 | MANHATTAN |
Took 9 sec. Last updated by anonymous at December 11 2021, 4:44:55 AM.
```

```
■ SPARK JOB FINISHED 'StreetName', 'Zip', 'Bor
%pvspark
 housing_litigations_df.select('LitigationID', 'BuildingID', 'BIN', 'HouseNumber', 'StreetName',
|LitigationID|BuildingID|
                          BIN|HouseNumber|
                                                   StreetName | Zip | Boro | Latitude | Longitude |
+------
                                    784
                                             AMSTERDAM AVENUE | 10025 | 1 | 40.795592 | -73.969261 |
      299988
                 6037 | 1056401 |
       96958
                805376 | 1085178 |
                                     1305
                                                 THIRD AVENUE | 10021 | 1 | 40.771509 | -73.959304 |
      288066
                 43492 | 1063707 |
                                     524
                                              WEST 184 STREET | 10033 | 1 | 40.850261 | -73.930470 |
       15892
                116740 | 2003497 |
                                    1055|DR M L KING JR BO...|10452|
                                                                    2|40.834676|-73.930269|
```

1	246954	366832 3079587	280	SCHAEFER	STREET 11237	3 40.692319 -73.905219
	208468	265319 3109164	402	EAST 46	STREET 11203	3 40.652120 -73.934341
	129962	186156 3087868	2869	ATLANTIC	AVENUE 11207	3 40.677032 -73.888300
_	28053	428406 4113860	144-45	41	AVENUE 11355	4 40.760695 -73.820986
DOE	3₃₽æri	m•it₁lss⊌a	nce	CROTONA	AVENUE 10457	2 40.851232 -73.885242
	254002	46172 2012027	2110	ARTHUR	AVENUE 10457	2 40.849868 -73.891292
	197682	116660 2014626	2422	UNIVERSITY	AVENUE 10468	2 40.863401 -73.904512
	42324	111955 2025351	1235	STRATFORD	AVENUE 10472	2 40.830390 -73.875785
1	244158	117443 null	578 FRONT	VAN NEST	AVENUE 10460	2 40.842165 -73.868502
I	170344 l	225096 3244549	3076 l	CONFY TSI AND	ΔVFNIIF 11235	3 40 581118 ₋ 73 959892

Took 0 sec. Last updated by anonymous at December 10 2021, 4:17:21 AM.

%pyspark
housing_litigations_boro_df = spark.sql("select LitigationID, BuildingID, BIN, HouseNumber, StreetName, Zip
end as Borough, Latitude, Longitude from housing_litigations")

Took 0 sec. Last updated by nbuser at December 12 2021, 5:41:53 AM.

%pyspark ■ SPARK JOB (http://hc11.nyu.cluster:44088/jobs/job?id=14) FINISHED housing_litigations_boro_df.filter(housing_litigations_boro_df['Zip'].isNull()).filter(housing_l

63

Took 1 sec. Last updated by nbuser at December 12 2021, 5:42:55 AM.

Took 1 sec. Last updated by anonymous at December 10 2021, 4:32:58 AM.

```
■ SPARK JOB FINISHED
per', 'StreetName', 'Zip',
%pyspark
housing_litigations_boro_df.select('LitigationID', 'BuildingID', 'BIN', 'HouseNumber',
|LitigationID|BuildingID|
                                BIN|HouseNumber|
                                                                                          Borough | Latitude | Longitude |
                                                               StreetName| Zip|
       131047
                     66460 | 2012965 |
                                               786
                                                         EAST 182 STREET | 10460 |
                                                                                            Bronx | 40.848909 | -73.884165 |
       161956
                     27988 | 1063950 |
                                              1441
                                                     ST NICHOLAS AVENUE | 10033 |
                                                                                        Manhattan | 40.849991 | -73.933253 |
                                                                                            Bronx | 40.849329 | -73.890233 |
       364010
                    842702 | 2114115 |
                                              2111
                                                           HUGHES AVENUE | 10457 |
       258106
                     42827 | 1062752 |
                                               538
                                                         WEST 159 STREET | 10032 |
                                                                                       Manhattan | 40.834579 | -73.942245 |
         48315
                    466695 | 4451496 |
                                           150-10|
                                                                71 AVENUE | 11367 |
                                                                                           Queens | 40.729238 | -73.817700 |
                                                                                         Brooklyn | 40.629530 | -73.965065 |
        32415
                    243141 | 3179168 |
                                               800 l
                                                          EAST 12 STREET | 11230 |
       239944
                     59195 | 2013896 |
                                              2301
                                                          CRESTON AVENUE | 10468 |
                                                                                            Bronx | 40.858062 | -73.900771 |
                    346877 | 3029764 |
                                                         NOSTRAND AVENUE | 11216 |
                                                                                         Brooklyn | 40.678107 | -73.949732 |
       182307
                                               596
                    726025 | 5002483 |
                                                            CORSON AVENUE | 10301 | Staten Island | 40.637004 | -74.086736 |
        41646
                                               247
       238699
                    363795 | 3034421 |
                                               332
                                                           ROGERS AVENUE | 11225 |
                                                                                         Brooklyn | 40.665153 | -73.953735 |
       125219
                     75641 | 2067469 |
                                              4077
                                                             EDSON AVENUE | 10466 |
                                                                                            Bronx | 40.891863 | -73.844233 |
                                                                 1 AVENUE | 10003 |
                                                                                       Manhattan | 40.730517 | -73.983183 |
       313118
                       434 | 1006490 |
                                               215
       132792
                    352670 | 3044159 |
                                                52
                                                          PATCHEN AVENUE | 11221 |
                                                                                         Brooklyn | 40.690669 | -73.927618 |
                                               657|
                                                             LENOX AVENUE | 10037 |
                                                                                        Manhattan | 40.819016 | -73.937287 |
       241901
                     23670 | 1060119 |
```

```
%pyspark SPARK JOB (http://hc11.nyu.cluster:40268/jobs/job?id=69) FINISHED housing_litigations_boro_df.filter(housing_litigations_boro_df['Zip'] == 10463).show()
```

+-----+
|LitigationID|BuildingID| BIN|HouseNumber| StreetName| Zip| Borough| Latitude| Longitude|

	129965	54954 2084079	3235 CAMBRIDGE AVENUE 10463 Bronx 40.883997 -73.908573
	82592	118583 2083947	3660 WALDO AVENUE 10463 Bronx 40.886851 -73.904392
	111668	108597 2092463	3605 SEDGWICK AVENUE 10463 Bronx 40.881542 -73.896483
$\mathbf{O}(\mathbf{C})$	B Reri	m#t¤lss⊌ar	C20 KINGSBRIDGE TERRACE 10463 Bronx 40.872342 -73.903288
	373725	46543 2015948	3300 BAILEY AVENUE 10463 Bronx 40.879454 -73.901253
	106442	89811 2015702	2765 KINGSBRIDGE TERRACE 10463 Bronx 40.872517 -73.903252
	276716	89811 2015702	2765 KINGSBRIDGE TERRACE 10463 Bronx 40.872517 -73.903252
	231703	25387 1064584	99 MARBLE HILL AVENUE 10463 Manhattan 40.877059 -73.909248
	253339	89724 2083169	3120 KINGSBRIDGE AVENUE 10463 Bronx 40.880369 -73.906014
	298527	879536 2113577	2726 KINGSBRIDGE TERRACE 10463 Bronx 40.871667 -73.903513
	234483	108597 2092463	3605 SEDGWICK AVENUE 10463 Bronx 40.881542 -73.896483
	148906	121353 2112572	231 WEST 230 STREET 10463 Bronx 40.877587 -73.907309
	258135	89714 2083150	3024 KINGSBRIDGE AVENUE 10463 Bronx 40.878771 -73.907550
1	2752071	100042120040021	SOUL OVENDO AVENHELIADACSI DESCRILAD COMOEST TO CONTEST

3758

Took 15 sec. Last updated by nbuser at December 12 2021, 5:44:02 AM.

%pyspark

SPARK JOB (http://hc11.nyu.cluster:44088/jobs/job?id=5) FINISHED housing_maintenance_code_violations_sample = housing_maintenance_code_violations_df.sample(False, 0.1, 1121 housing_maintenance_code_violations_sample.count()

688358

Took 10 sec. Last updated by nbuser at December 12 2021, 5:22:58 AM.

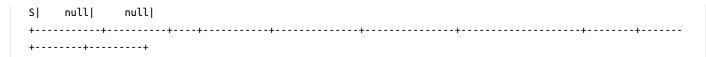
```
%pyspark
housing_maintenance_code_violations_sample.select('ViolationID', 'BuildingID', 'BIN', 'HouseNumber', 'LowHo
.filter(housing_maintenance_code_violations_sample['Postcode'] == 10463)\
.filter(housing_maintenance_code_violations_sample['Borough'] != 'BRONX')\
.show()
```

+						
1		ildingID BIN Hous	•	useNumber HighHo	useNumber StreetName Po	·
1	13780863	28645 1064543	110	110	116 TERRACE VIEW AVENUE	10463 MANHA
T	TAN 40.876444	-73.912829				
- 1	13827795	21653 1064651	26	26	28 FT CHARLES PLACE	10463 MANHA
T	TAN 40.875883	-73.910634				
- 1	10091660	25378 1064679	58	58	60 MARBLE HILL AVENUE	10463 MANHA
T	TAN 40.875720	-73.909702				
	11220171	28648 1064524	135	131	135 TERRACE VIEW AVENUE	10463 MANHA
T	TAN 40.876439	-73.912854				
	11604953	25387 1064584	99	99	101 MARBLE HILL AVENUE	10463 MANHA
T	TAN 40.877059	-73.909248				
	14040212	25387 1064584	99	99	101 MARBLE HILL AVENUE	10463 MANHA
	TANI 40 0770F01	72 2222421				

%pyspark SPARK JOB FINISHED housing_maintenance_code_violations_sample.filter(housing_maintenance_code_violations_sample['Postcode'].isNu |ViolationID|BuildingID| BIN|HouseNumber|LowHouseNumber|HighHouseNumber| StreetName|Postcode| Borough|La titude|Longitude| 14276447| 967507|null| 532| 532 532| EAST 142 STREET| null| BRONX| null| null| 10611493 940715|null| 201 201 201 | NORTH 11 STREET | null|BROOKLYN| null| null| 10612379 942882|null| 57-37 57-37 57-37 VAN DOREN STREET null| QUEENS| null| null 14299130 997630|null| 223 223 | MILFORD STREET | null|BROOKLYN| 221 null 10741476 956610|null| 91| 91| 91 | CRYSTAL STREET | null|BROOKLYN| null| null| 10971774| 90-47 90-47 null| QUEENS| 955005|null| 90-47 198 STREET Took 25 sec. Last updated by anonymous at December 11 2021, 9:18:58 PM.

```
%pyspark
housing_maintenance_code_violations_sample.filter(housing_maintenance_code_violations_df['StreetName'] == '
    'Latitude', 'Longitude').show()

+-----+
| ViolationID|BuildingID| BIN|HouseNumber|LowHouseNumber|HighHouseNumber| StreetName|Postcode|Boroug
h|Latitude|Longitude|
+-----+
| 10314587| 850855|null| 109-55| 109-55| 109-61|109-55 VAN WYCK E...| null| QUEEN
```



Took 22 sec. Last updated by nbuser at December 12 2021, 5:23:29 AM

%pyspark SPARK JOB (http://hc05.nyu.cluster:33536/jobs/job?id=19) FINISHED housing_maintenance_code_violations_sample.select('StreetName').distinct().count()

3295

Took 17 sec. Last updated by anonymous at December 11 2021, 9:19:36 PM.

%sh
/usr/bin/hadoop fs -rm -r housing_maintenance_code_violations_sample.csv

FINISHED

21/12/11 18:21:38 INFO fs.TrashPolicyDefault: Moved: 'hdfs://horton.hpc.nyu.edu:8020/user/ky572/housing_maint enance_code_violations_sample.csv' to trash at: hdfs://horton.hpc.nyu.edu:8020/user/ky572/.Trash/Current/user/ky572/housing_maintenance_code_violations_sample.csv

Took 2 sec. Last updated by anonymous at December 11 2021, 10:21:39 PM.

%pyspark SPARK JOB (http://hc05.nyu.cluster:33536/jobs/job?id=24) FINISHED housing_maintenance_code_violations_sample.write.option("header", True).csv('./housing_maintenance_code_violations_sample.write.option("header", True).csv('./housing_maintenance_code_violations_sample.write.option("header").

Took 38 sec. Last updated by anonymous at December 11 2021, 10:22:20 PM.

%sh
rm housing_maintenance_code_violations_sample.csv

FINISHED

Took 3 sec. Last updated by anonymous at December 11 2021, 10:22:46 PM. (outdated)

%sh /usr/bin/hadoop fs -getmerge housing_maintenance_code_violations_sample.csv housing_maintenance_code_violation

```
%python
                                                                                                      FINISHED
from openclean.pipeline import stream
from openclean.function.eval.base import Col
from openclean.function.eval.logic import And
from openclean.function.eval.null import IsNotEmpty, IsEmpty
from openclean.operator.map.violations import fd_violations
from openclean.cluster.key import key_collision
from openclean_geo.address.usstreet import USStreetNameKey
def lat_long_boro_fd(ds, latCol, longCol, boroCol):
    data = ds.select([boroCol, latCol, longCol]).where(And(IsNotEmpty(latCol), IsNotEmpty(longCol), IsNotEmpty(longCol), IsNotEmpty(longCol)
    df = data
        .select([boroCol, latCol, longCol])\
        .to df()
    groups = fd_violations(df, lhs=[latCol, longCol], rhs=boroCol)
    return groups
def cleanstreet(street,ds_full):
    streets = ds full.select(street).distinct()
```

Took 3 sec. Last updated by nbuser at December 12 2021, 5:17:28 AM.

%python
 FINISHED
housing_maintenance_code_violations_sample_ds = stream('./housing_maintenance_code_violations_sample.csv')

Took 0 sec. Last updated by nbuser at December 12 2021, 5:17:38 AM.

%python
housing_maintenance_code_violations_sample_ds.where(Col('StreetName') == '27TH STREET').select(['HouseNumber'

Empty DataFrame

Columns: [HouseNumber, StreetName, Postcode]

Index: []

Took 6 sec. Last updated by anonymous at December 11 2021, 10:28:36 PM.

%python
cleaned = cleanstreet('StreetName', housing_maintenance_code_violations_sample_ds)

Took 6 sec. Last updated by nbuser at December 12 2021, 5:17:48 AM.

%python
cleaned.where(IsEmpty('Postcode')).count()

378

Took 8 sec. Last updated by nbuser at December 12 2021, 5:18:01 AM.

%python
finished
cleaned.where(IsEmpty('Postcode')).select(['ViolationID', 'BIN', 'HouseNumber', 'StreetName', 'Postcode', 'Boundary')

	ViolationID E	BIN HouseNumber	Stre	eetName	Postcode	Borough
927	10859232	1754	ANTHONY	AVENUE		BRONX
1514	11323977	238-11	HILLSIDE	AVENUE		QUEENS
6238	13317952	881	ALBANY	AVENUE		BROOKLYN
6679	13417672	111-03	38	AVENUE		QUEENS
9041	13681082	235	SKILLMAN	STREET		BROOKLYN
11812	14276447	532	EAST 142	STREET		BRONX
11837	14276451	532	EAST 142	STREET		BRONX
12889	13767458	152	SACKMAN	STREET		BROOKLYN
12890	13767468	152	SACKMAN	STREET		BROOKLYN
12891	13767470	152	SACKMAN	STREET		BROOKLYN

Took 0 sec. Last updated by anonymous at December 12 2021, 12:21:48 AM.

%python FINISHED

Took 19 sec. Last updated by nbuser at December 12 2021, 5:18:34 AM.

```
      %python
      FINISHED

      for key in list(groups.keys())[:10]:
      print(groups.values(key=key, columns='Postcode'))

      print('\n')
      Counter({'11219': 2, '': 1})

      Counter({'19468': 2, '': 1})
      Counter({'11373': 1, '': 1})

      Counter({'': 1, '11226': 1})
      Counter({'': 1, '11226': 1})

      Counter({'': 1, '11221': 1})
      Counter({'': 1, '11221': 1})

      Took 0 sec. Last updated by anonymous at December 12 2021, 12:58:53 AM.
      AM.
```

```
%python
group_mapping = dict()
for key in groups.keys():
    values = groups.values(key=key, columns='Postcode')
    max=0
    maxValue=""
    for v in values:
        if v is None or len(v) == 0:
            continue
        if max<values[v]:
            max=values[v]
            maxValue=v
            group manning[kev] = maxValue</pre>
Took 0 sec. Last updated by nbuser at December 12 2021, 5:18:41 AM.
```

```
%python
for key in list(group_mapping.keys())[:10]:
    print('{} = {}'.format(key, group_mapping[key]))

('BROOKLYN', '204', 'LINDEN BOULEVARD') = 11226
('QUEENS', '84-47', 'KNEELAND AVENUE') = 11373
('BROOKLYN', '1064', '57 STREET') = 11219
('BROOKLYN', '41414', '47 STREET') = 11219
('BROOKLYN', '834', 'QUINCY STREET') = 11221
```

```
('BRONX', '2302', 'MORRIS AVENUE') = 10468
```

Took 0 sec. Last updated by nbuser at December 12 2021, 5:18:45 AM.

DOB Permit Issuance

%python from openclean.function.eval.domain import Lookup **FINISHED**

cleaned2 = cleaned.where(IsEmpty('Postcode')).update('Postcode', Lookup(columns=['Borough', 'HouseNumber',

Took 0 sec. Last updated by nbuser at December 12 2021, 5:18:58 AM.

```
%python
                                                                                                 FINISHED
cleaned2.where(IsEmpty('Postcode')).count()
```

372

Took 8 sec. Last updated by nbuser at December 12 2021, 5:19:10 AM.

```
%python
                                                                                                     FINISHED
cleaned2.where(IsEmpty('Postcode')).select(['ViolationID', 'BIN', 'HouseNumber', 'StreetName', 'Postcode',
       ViolationID BIN HouseNumber
                                                Borough Latitude Longitude
                              1754 ...
927
          10859232
                                                  BRONX
1514
          11323977
                            238-11 ...
                                                 QUEENS
          13317952
                                               BROOKLYN
6238
                               881 ...
6679
          13417672
                            111-03 ...
                                                 QUEENS
9041
          13681082
                               235 ...
                                               BROOKLYN
               ... ..
                                . . . . . . . .
                                                             . . .
205019
          11606836
                              2408 ...
                                                  BRONX
209011
          13843240
                               441
                                               BROOKLYN
209141
          14479940
                                590 ...
                                          STATEN ISLAND
218166
          14524400
                               114
                                               BROOKLYN
222810
          11807813
                               943 ...
                                                  BRONX
```

```
[100 rows x 8 columns]
```

Took 3 sec. Last updated by nbuser at December 12 2021, 5:20:15 AM.

```
FINISHED
cleaned.where(And(Col('HouseNumber') == '532', Col('StreetName') == 'EAST 142 STREET', IsNotEmpty('Postcode')
Empty DataFrame
Columns: [HouseNumber, StreetName, Postcode, Borough, BIN]
Index: []
Took 9 sec. Last updated by anonymous at December 12 2021, 1:03:10 AM.
```

```
%python
cleaned.where(Col('StreetName') == 'VAN WYCK EXPRESSWAY').select([ HouseNumber', 'StreetName', 'Postcode',
      HouseNumber
                            StreetName Postcode Borough
127656
            95-12 VAN WYCK EXPRESSWAY
                                          11419 QUEENS
252759
            95-12 VAN WYCK EXPRESSWAY
                                          11419 QUEENS
466509
            87-53 VAN WYCK EXPRESSWAY
                                          11435 QUEENS
```

12/12/21, 8:01 AM 10 of 33

D⊕B	Perm	M	SK.	Suanc	e 11419	QUEENS
535068				EXPRESSWAY		QUEENS
497254	95-18	VAN	${\tt WYCK}$	EXPRESSWAY	11419	QUEENS
494195	95-12	VAN	WYCK	EXPRESSWAY	11419	QUEENS
494194	95-12	VAN	WYCK	EXPRESSWAY	11419	QUEENS

Took 7 sec. Last updated by nbuser at December 12 2021, 5:19:27 AM.

```
%pyspark
                                                                ■ SPARK JOB FINISHED
housing_maintenance_code_violations_df.select('ViolationID', 'BuildingID', 'BIN', 'HouseNumber', 'LowHouseN
    (housing_maintenance_code_violations_df['BuildingID'] != 1).show(n=600)
+------
----+
|ViolationID|BuildingID| BIN|HouseNumber|LowHouseNumber|HighHouseNumber|
                                                              StreetName|Postcode|
orough|Latitude|Longitude|
+-----
-----+
                            273
                                       273|
                                                    273
  14265577
           904119|null|
                                                           WEST 122 STREET
                                                                          null|
                                                                                 MAN
HATTAN|
      null| null|
           958166|null|
                             42|
                                        42|
                                                           SWEETBROOK ROAD
                                                                          null|STATEN
  10521517
                                                    42|
ISLAND|
       null|
             null|
  10594623
             956415|null|
                          27-16
                                      27-16
                                                  27-16
                                                                18 STREET
                                                                          null|
QUEENS |
        null
                null|
  10623685
            803906|null|
                                                   2411|FREDERICK DOUGAL ...|
                                                                                 MAN
                           2411
                                       2411
                                                                          10027
HATTAN|
        null|
               null|
  10694011
            936460|null|
                           1431
                                       1431
                                                   1431
                                                            ZEREGA AVENUE
                                                                          null|
       null|
BRONX
               null
 10694022
             936460|null|
                           1431
                                       1431
                                                   1431
                                                            ZEREGA AVENUE
                                                                          null|
Took 8 sec. Last updated by anonymous at December 10 2021, 10:35:14 PM. (outdated)
```

+					+	
	+ iolationID Bui		seNumber LowH	ouseNumber HighH	ouseNumber StreetName Po	ostcode Borough Lati
tuc	B•Per	mit Issua	nce			
+					+	
			47.051	47.051	47 051 70 6705571	442771 005505140 72
1	10970791	890887 4536785	47-05	47-05	47-05 72 STREET	11377 QUEENS 40.73
188	51 -73.892238 11545234	890887 4536785	47-05	47-05	47-05 72 STREET	11377 QUEENS 40.73
 	11343234 51 -73.892238	•	47-05	47-03	47-05 72 SIREET	113// QUEENS 40./3
I	14531260	809275 3337150	190	190	198 72 STREET	11209 BROOKLYN 40.63
420	94 -74.030261	·	1901	190	190 72 31KLL1	11209 DROOKE N 40.05
	12136761	890887 4536785	47-05	47-05	47-05 72 STREET	11377 QUEENS 40.73
886	51 -73.892238	•	551	551	55 12 5111221	223.77
ı	12660175	169543 3146999	871	871	871 72 STREET	11228 BROOKLYN 40.62
775	56 -74.015102		·			
ı	12716984	890887 4536785	47-05	47-05	47-05 72 STREET	11377 QUEENS 40.73
200	(4) 73 0000001					
Too		ated by nbuser at Decemb	er 12 2021, 7:17:43	3 AM.		

+				+	+-	
+	+					
ViolationID Buil	dingID BIN Hous	eNumber LowHo	useNumber HighHo	ouseNumber	StreetName Po	ostcode Borou
gh Latitude Longi	•					
					+-	
10526625	1 null	401	401	403 ST (GEORGE'S CRESCENT	10016 MANHATT
AN null		4011	4011	405 51	LONGE 5 CHESCENT	10010 10010
11491071	1 null	401	401	403 ST C	GEORGE'S CRESCENT	10016 MANHATT
AN null	null	·	·	•	·	·
11491082	1 null	401	401	403 ST C	GEORGE'S CRESCENT	10016 MANHATT
AN null	null					
11491177	1 null	401	401	403 ST C	GEORGE'S CRESCENT	10016 MANHATT
AN null	null					
11491148	1 null	401	401	403 ST C	EORGE'S CRESCENT	10016 MANHATT
AN null	null					
11550294	1 null	401	401	403 ST C	EORGE'S CRESCENT	10016 MANHATT