Initializing System

```
1 !apt-get install openjdk-8-jdk-headless -qq > /dev/null
```

```
1 !tar xf "/content/drive/My Drive/BigDataAssignment3Files/spark-2.4.5-bin-hadoop2.7.tgz"
2 !pip install -q findspark
```

```
1 import os
2 os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"
3 os.environ["SPARK_HOME"] = "/content/spark-2.4.5-bin-hadoop2.7"
```

```
1 import findspark
2 findspark.init()
3 from pyspark.sql import SparkSession
4 from pyspark.context import SparkContext
5 spark = SparkSession.builder.master("local[*]").getOrCreate()
```

Reading Preprocessed Dataset

```
1 preprocessed_data = spark.read.json("hdfs://udit_gupta_1/processed_data")
```

Extract Data for EDA

```
1 #Get top 30 crime types in complaints
2 top30_crime_type = preprocessed_data.rdd \
3 .filter(lambda row : row['RECORD_TYPE'] == 'C') \
4 .map(lambda row : (row['OFNS_DESC'],1)) \
5 .reduceByKey(lambda key1, key2 : key1 + key2) \
6 .takeOrdered(30,lambda atuple: -atuple[1])
```

```
1 complaints_crime_list = [ele[0] for ele in top30_crime_type if ele[0] is not None]
```

```
1 #Get top 30 arrests crime types
2 top30_arrests_crime_type = preprocessed_data.rdd \
3 .filter(lambda row : row['RECORD_TYPE'] == 'A') \
4 .map(lambda row : (row['OFNS_DESC'],1)) \
5 .reduceByKey(lambda key1, key2 : key1 + key2) \
6 .takeOrdered(30,lambda atuple: -atuple[1])
```

```
1 arrests_crime_list = [ele[0] for ele in top30_arrests_crime_type if ele[0] is not None]
```

```
1 #Get top 30 location types for crime complaints
2 top30_crime_locations = preprocessed_data.rdd \
3 .filter(lambda row : row['RECORD_TYPE'] == 'C') \
4 .map(lambda row : (row['PREM_TYP_DESC'],1)) \
5 .reduceByKey(lambda key1, key2 : key1 + key2) \
6 .takeOrdered(30,lambda atuple: -atuple[1])
```

```
1 complaints_location_list = [ele[0] for ele in top30_crime_locations if ele[0] is not None]
```

Generate Pivot Tables for Different Attributes

```
1 #Pivot Table based on Location
2 location_groupby = preprocessed_data \
3 .filter(preprocessed_data['RECORD_TYPE'] == 'C') \
4 .filter(preprocessed_data.PREM_TYP_DESC.isin(complaints_location_list)) \
5 .groupby('PREM_TYP_DESC')
6 location_crime_pivot = location_groupby.pivot("OFNS_DESC", complaints_crime_list).agg({"*": "count"}).fillna(0).toPandas()
7 location_crime_pivot = location_crime_pivot.set_index('PREM_TYP_DESC')
```

```
1 #Pivot Table based on ethinicity
2 ethinicity_crime_pivot =  preprocessed_data \
3 .filter(preprocessed_data['RECORD_TYPE'] == 'C') \
4 .groupby('SUSP_RACE').pivot("OFNS_DESC", complaints_crime_list).agg({"*": "count"}).fillna(0).toPandas()
5
6 ethinicity_crime_pivot = ethinicity_crime_pivot.set_index('SUSP_RACE')
```

```
1 #Pivot Table for age-group
2 age_crime_pivot =  preprocessed_data \
3 .filter(preprocessed_data['RECORD_TYPE'] == 'C') \
4 .groupby('SUSP_AGE_GROUP').pivot("OFNS_DESC", complaints_crime_list).agg({"*": "count"}).fillna(0).toPandas()
5
6 age_crime_pivot = age_crime_pivot.set_index('SUSP_AGE_GROUP')
```

```
1 #Pivot Table for gender
2 gender_crime_pivot = preprocessed_data \
3 .filter(preprocessed_data['RECORD_TYPE'] == 'C') \
4 .groupby('SUSP_SEX').pivot("OFNS_DESC", complaints_crime_list).agg({"*": "count"}).fillna(0).toPandas()
5
6 gender_crime_pivot = gender_crime_pivot.set_index('SUSP_SEX')
```

```
1 #Pivot Table for precicnt codes
2 borough_crime_pivot =  preprocessed_data \
3 .filter(preprocessed_data['RECORD_TYPE'] == 'C') \
4 .groupby('BORO_NM').pivot("OFNS_DESC", complaints_crime_list).agg({"*": "count"}).fillna(0).toPandas()
5
6 borough_crime_pivot = borough_crime_pivot.set_index('BORO_NM')
```

Corrospondence Analysis

```
1 from numpy.linalg import svd
2 import numpy as np
3 from matplotlib.pyplot import figure
4 import numpy as np
5 import pandas as pd
6 import matplotlib.pyplot as plt
7
8 class CA(object):
```

```python
 9      """Simple corresondence analysis.
10
11      Inputs
12      ------
13      ct : array_like
14        Two-way contingency table. If `ct` is a pandas DataFrame object,
15        the index and column values are used for plotting.
16      Notes
17      -----
18      The implementation follows that presented in 'Correspondence
19      Analysis in R, with Two- and Three-dimensional Graphics: The ca
20      Package,' Journal of Statistical Software, May 2007, Volume 20,
21      Issue 3.
22      """
23
24      def __init__(self, ct):
25          self.rows = ct.index.values if hasattr(ct, 'index') else None
26          self.cols = ct.columns.values if hasattr(ct, 'columns') else None
27
28          # contingency table
29          N = np.matrix(ct, dtype=float)
30
31          # correspondence matrix from contingency table
32          P = N / N.sum()
33
34          # row and column marginal totals of P as vectors
35          r = P.sum(axis=1)
36          c = P.sum(axis=0).T
37
38          # diagonal matrices of row/column sums
39          D_r_rsq = np.diag(1. / np.sqrt(r.A1))
40          D_c_rsq = np.diag(1. / np.sqrt(c.A1))
41
42          # the matrix of standarized residuals
43          S = D_r_rsq * (P - r * c.T) * D_c_rsq
44
45          # compute the SVD
46          U, D_a, V = svd(S, full_matrices=False)
47          D_a = np.asmatrix(np.diag(D_a))
48          V = V.T
49
50          # principal coordinates of rows
51          F = D_r_rsq * U * D_a
52
53          # principal coordinates of columns
54          G = D_c_rsq * V * D_a
55
56          # standard coordinates of rows
57          X = D_r_rsq * U
58
59          # standard coordinates of columns
60          Y = D_c_rsq * V
61
62          # the total variance of the data matrix
63          inertia = sum([(P[i,i] - r[i,0] * c[i,0])**2 / (r[i,0] * c[i,0])
```
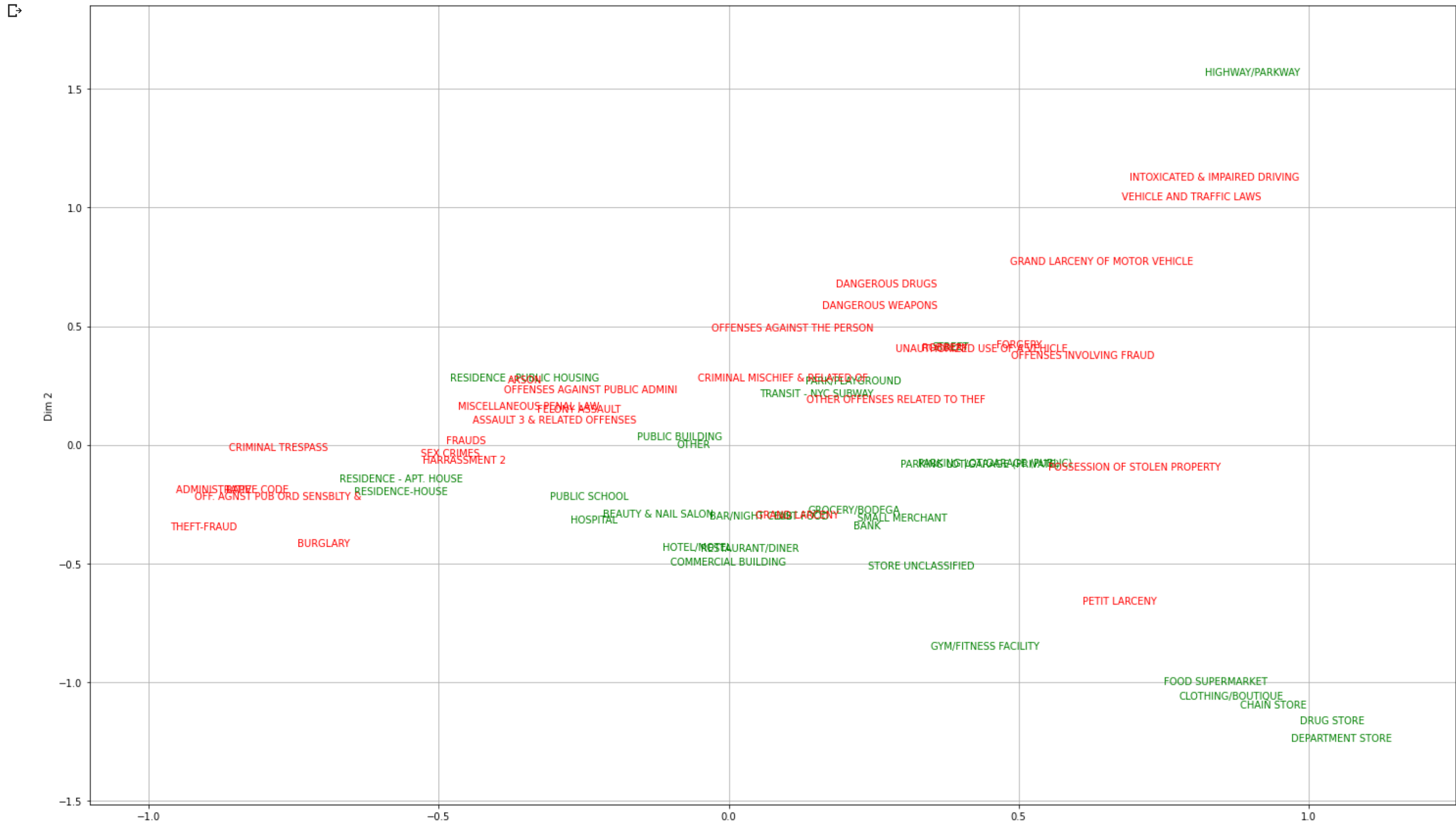
```python
64                          for i in range(N.shape[0])
65                          for j in range(N.shape[1])])
66
67          self.F = F.A
68          self.G = G.A
69          self.X = X.A
70          self.Y = Y.A
71          self.inertia = inertia
72          self.eigenvals = np.diag(D_a)**2
73
74      def plot(self):
75          """Plot the first and second dimensions."""
76          xmin, xmax = None, None
77          ymin, ymax = None, None
78          if self.rows is not None:
79              for i, t in enumerate(self.rows):
80                  x, y = self.F[i,0], self.F[i,1]
81                  plt.text(x, y, t, va='center', ha='center', color='green')
82                  xmin = min(x, xmin if xmin else x)
83                  xmax = max(x, xmax if xmax else x)
84                  ymin = min(y, ymin if ymin else y)
85                  ymax = max(y, ymax if ymax else y)
86          else:
87              plt.plot(self.F[:, 0], self.F[:, 1], 'ro')
88
89          if self.cols is not None:
90              for i, t in enumerate(self.cols):
91                  x, y = self.G[i,0], self.G[i,1]
92                  plt.text(x, y, t, va='center', ha='center', color='r')
93                  xmin = min(x, xmin if xmin else x)
94                  xmax = max(x, xmax if xmax else x)
95                  ymin = min(y, ymin if ymin else y)
96                  ymax = max(y, ymax if ymax else y)
97          else:
98              plt.plot(self.G[:, 0], self.G[:, 1], 'bs')
99
100         if xmin and xmax:
101             pad = (xmax - xmin) * 0.1
102             plt.xlim(xmin - pad, xmax + pad)
103         if ymin and ymax:
104             pad = (ymax - ymin) * 0.1
105             plt.ylim(ymin - pad, ymax + pad)
106
107         plt.grid()
108         plt.xlabel('Dim 1')
109         plt.ylabel('Dim 2')
110
111     def scree_diagram(self, perc=True, *args, **kwargs):
112         """Plot the scree diagram."""
113         eigenvals = self.eigenvals
114         xs = np.arange(1, eigenvals.size + 1, 1)
115         ys = 100. * eigenvals / eigenvals.sum() if perc else eigenvals
116         plt.plot(xs, ys, *args, **kwargs)
117         plt.xlabel('Dimension')
118         plt.ylabel('Eigenvalue' + (' [%]' if perc else ''))
```
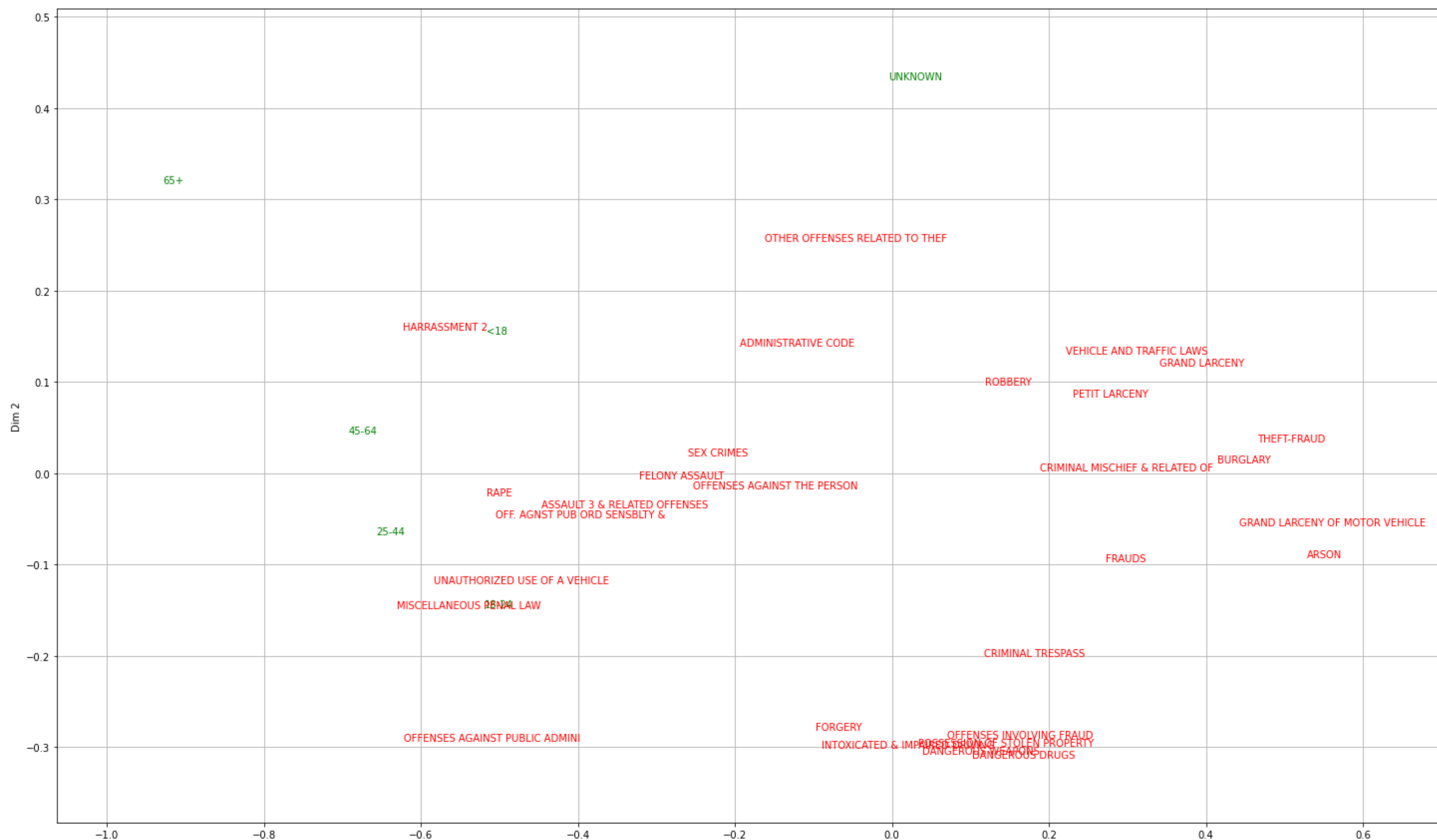
Generating Plots

```
1 ca = CA(location_crime_pivot)
2 plt.figure(figsize=(25,15))
3 ca.plot()
```
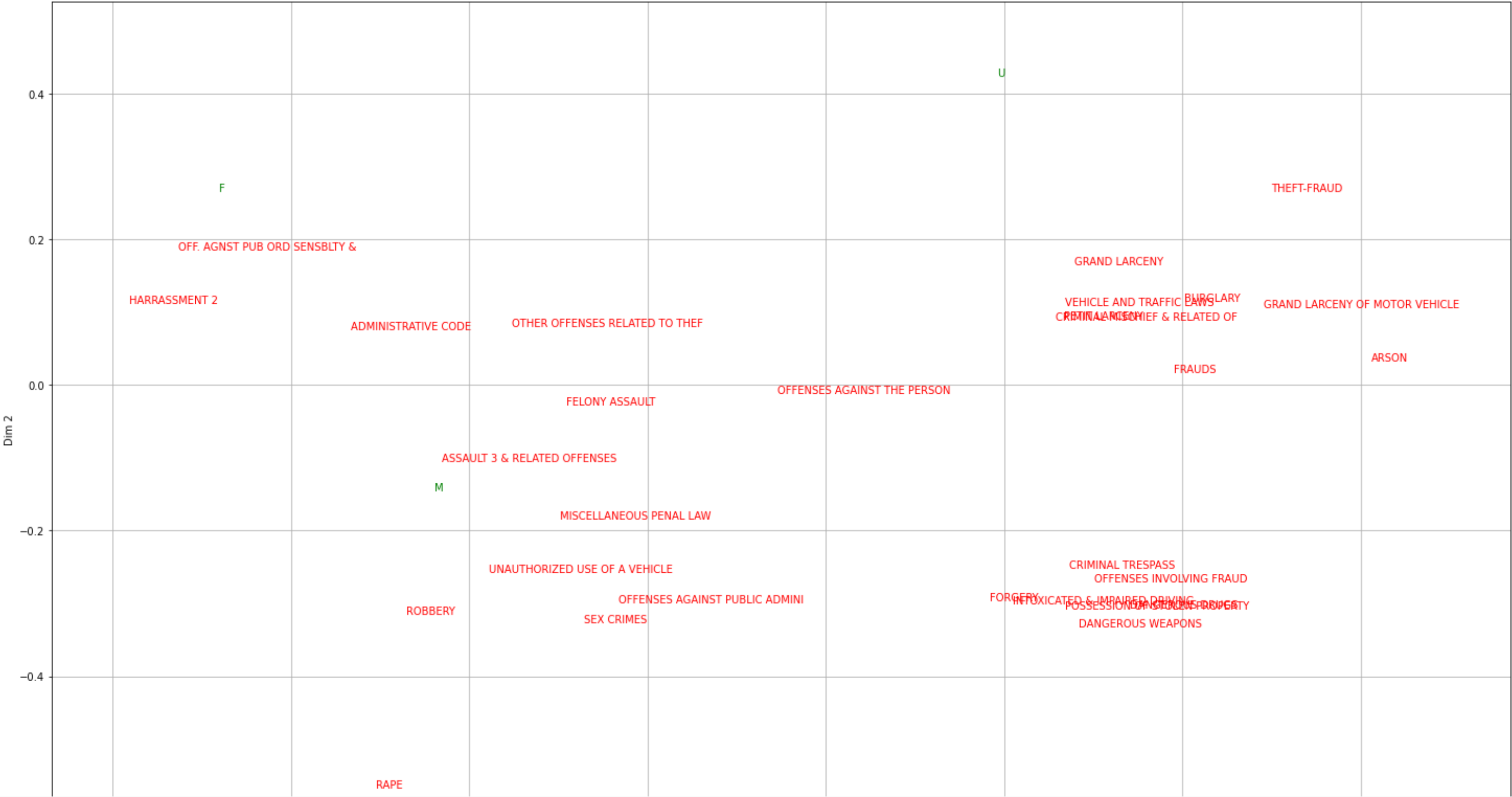


```
1 ca = CA(age_crime_pivot)
2 plt.figure(figsize=(25,15))
```

```
3 ca.plot()
```



```
1 ca = CA(gender_crime_pivot)
2 plt.figure(figsize=(25,15))
3 ca.plot()
```

```
1 ca = CA(borough_crime_pivot)
2 plt.figure(figsize=(25,15))
3 ca.plot()
```

```
1 ca = CA(ethinicity_crime_pivot)
2 plt.figure(figsize=(25,15))
3 ca.plot()
```