



EL OJO DE BUCAR

Carlos Javier Aceros Blanco

Inteligencia artificial I - Proyecto final
Universidad Industrial de Santander

PROBLEMA A RESOLVER

- La inseguridad es un problema que en los últimos años ha afectado a Bucaramanga, la percepción de inseguridad llegó hasta el 79,2%, es decir, 8 de cada 10 ciudadanos se sienten inseguros en Bucaramanga y su área metropolitana, según el periódico Vanguardia

SOLUCIÓN PROPUESTA

- Implementar una herramienta web basada en inteligencia artificial que permita clasificar las zonas más peligrosas de la ciudad para poder implementar políticas públicas basadas en datos y lograr que la tranquilidad vuelva a la capital de Santander

METODOLOGÍA

- A través de una búsqueda en los datasets disponibles en Datos abiertos de Colombia(<https://www.datos.gov.co/>), se encontró un historial que consolida los delitos reportados desde Enero de 2010 hasta diciembre de 2021
- Se notó que los datos se actualizaron hasta el 10 octubre de 2023
- A través de técnicas aprendidas y guiadas por el docente Gustavo Garzón de análisis de datos se implementaron acciones de limpieza de datos, visualización y algoritmos de clasificación para lograr el objetivo planteado previamente

ALGUNOS ANÁLISIS

columns and shape Index(['ORDEN', 'ARMAS_MEDIOS', 'BARRIOS_HECHO', 'LATITUD', 'LONGITUD', 'ZONA', 'NOM_COMUNA', 'ANO', 'MES', 'DIA', 'DIA_SEMANA', 'DESCRIPCION_CONDUCTA', 'CONDUCTA', 'CLASIFICACIONES_DELITO', 'EDAD', 'CURSO_DE_VIDA', 'ESTADO_CIVIL_PERSONA', 'GENERO', 'MOVIL_AGRESOR', 'MOVIL_VICTIMA'], dtype='object') (135076, 20)

ORDEN	ARMAS_MEDIOS	BARRIOS_HECHO	LATITUD	LONGITUD	ZONA	NOM_COMUNA	ANO	MES	DIA	DIA_SEMANA	DESCRIPCION_CONDUCTA	CONDUCTA	CLASIFICACIONES_DELITO	EDAD	CURSO_DE_VIDA	ESTADO_CIVIL_PERSONA	GENERO	MOVIL_AGRESOR	MOVIL_VICTIMA	
0	1	ARMA BLANCA / CORTOPUNZANTE	BUENOS AIRES	7.170557	-73.135108	URBANA	14. Morrorico	2010	01. Enero	1	05. Viernes	ARTÍCULO 111. LESIONES PERSONALES	LESIONES PERSONALES	Lesiones no fatales	30	05. Adultez	UNION LIBRE	MASCULINO	A PIE	A PIE
1	2	ARMA BLANCA / CORTOPUNZANTE	CAMPO HERMOSO	7.120645	-73.12605	URBANA	05. García Rovira	2010	01. Enero	1	05. Viernes	ARTÍCULO 111. LESIONES PERSONALES	LESIONES PERSONALES	Lesiones no fatales	21	04. Jovenes	SOLTERO	MASCULINO	A PIE	A PIE
2	3	ARMA BLANCA / CORTOPUNZANTE	CAMPO HERMOSO	7.120645	-73.12605	URBANA	05. García Rovira	2010	01. Enero	1	05. Viernes	ARTÍCULO 111. LESIONES PERSONALES	LESIONES PERSONALES	Lesiones no fatales	23	04. Jovenes	SOLTERO	MASCULINO	A PIE	A PIE
3	4	ARMA BLANCA / CORTOPUNZANTE	COMUNEROS	7.151359	-73.145705	URBANA	03. San Francisco	2010	01. Enero	1	05. Viernes	ARTÍCULO 111. LESIONES PERSONALES	LESIONES PERSONALES	Lesiones no fatales	36	05. Adultez	CASADO	MASCULINO	A PIE	A PIE
4	5	ARMA BLANCA / CORTOPUNZANTE	GIRARDOT	7.170557	-73.135108	URBANA	04. Occidental	2010	01. Enero	1	05. Viernes	ARTÍCULO 111. LESIONES PERSONALES	LESIONES PERSONALES	Lesiones no fatales	20	04. Jovenes	UNION LIBRE	MASCULINO	A PIE	A PIE
5	6	ARMA BLANCA / CORTOPUNZANTE	GIRARDOT	7.170557	-73.135108	URBANA	04. Occidental	2010	01. Enero	1	05. Viernes	ARTÍCULO 111. LESIONES PERSONALES	LESIONES PERSONALES	Lesiones no fatales	20	04. Jovenes	UNION LIBRE	MASCULINO	A PIE	A PIE
6	7	ARMA BLANCA /	LOS ANGELES	7.187455	-73.131727	URBANA	02. Nor	2010	01.	1	05. Viernes	ARTÍCULO 111. LESIONES	LESIONES	Lesiones no fatales	28	04. Jovenes	CASADO	MASCULINO	A PIE	A PIE

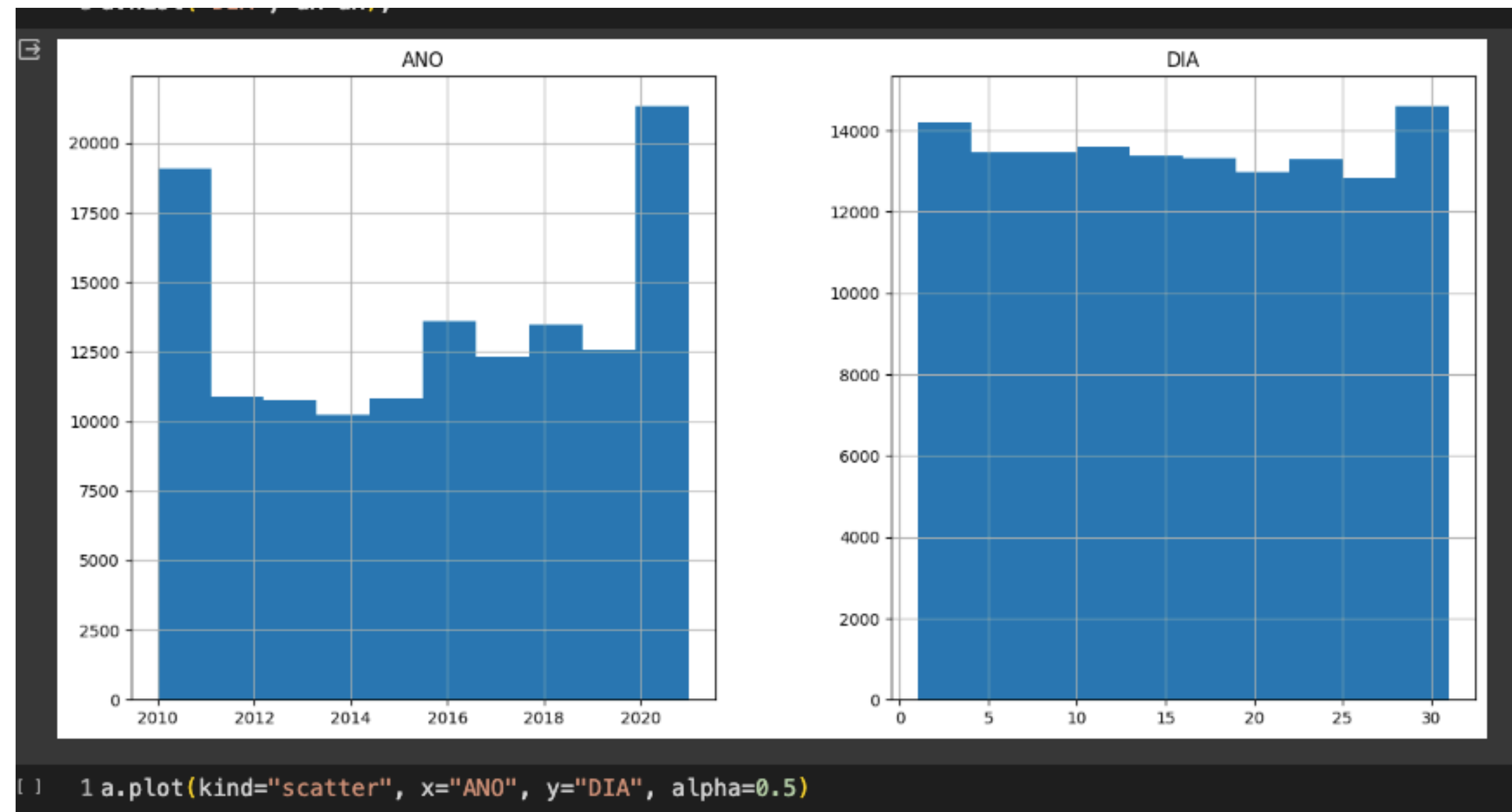
- Exploración del dataset: 20 columnas, 135.076 registros(delitos reportados)

ALGUNOS ANÁLISIS

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 135076 entries, 0 to 135075
Data columns (total 20 columns):
#   Column              Non-Null Count  Dtype
---  ---
0   ORDEN                135076 non-null  int64
1   ARMAS_MEDIOS         135076 non-null  object
2   BARRIOS_HECHO        135076 non-null  object
3   LATITUD              128713 non-null  object
4   LONGITUD             128713 non-null  object
5   ZONA                 135076 non-null  object
6   NOM_COMUNA           135076 non-null  object
7   AÑO                  135076 non-null  int64
8   MES                  135076 non-null  object
9   DIA                  135076 non-null  int64
10  DIA_SEMANA           135076 non-null  object
11  DESCRIPCION_CONDUCTA 135076 non-null  object
12  CONDUCTA             135076 non-null  object
13  CLASIFICACIONES_DELITO 135076 non-null  object
14  EDAD                 135076 non-null  object
15  CURSO_DE_VIDA        135076 non-null  object
16  ESTADO_CIVIL_PERSONA 135076 non-null  object
17  GENERO               135076 non-null  object
18  MOVIL_AGRESOR        135076 non-null  object
19  MOVIL_VICTIMA         135076 non-null  object
dtypes: int64(3), object(17)
memory usage: 20.6+ MB
```

- La mayoría de tipos de datos categóricos, implicaba etiquetado para transformar a numéricos y poder realizar un análisis interesante

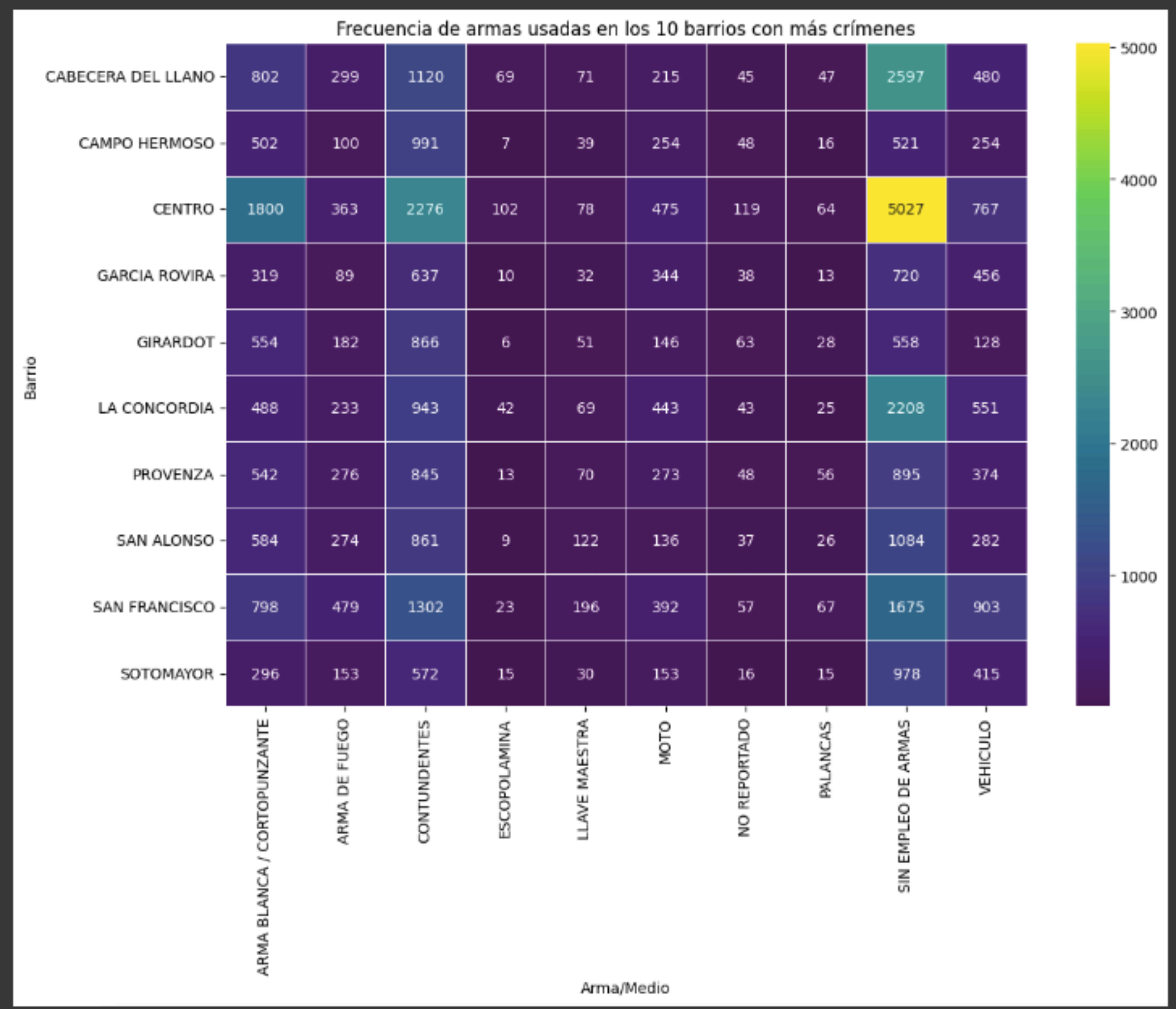
ALGUNOS ANÁLISIS



- Con los datos numéricos iniciales se lograron obtener algunas conclusiones interesantes: Los años en dónde más hubo delitos, los días del mes en donde más se cometen delitos(inicio y fin de mes) ¿Coincide con los días de pago?

ALGUNOS ANÁLISIS

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 # Vamos a encontrar los 10 barrios más comunes y las 10 armas más comunes.
6 top_barrios = a['BARRIOS_HECHO'].value_counts().head(10).index
7 top_armas = a['ARMAS_MEDIOS'].value_counts().head(10).index
8
9 # Filtramos el DataFrame para incluir solo los datos más comunes.
10 filtered_data = a[a['BARRIOS_HECHO'].isin(top_barrios) & a['ARMAS_MEDIOS'].isin(top_armas)]
11
12 # Calcular la frecuencia de cada combinación de arma y barrio.
13 conteo_armas_barrios = filtered_data.groupby(['BARRIOS_HECHO', 'ARMAS_MEDIOS']).size().reset_index(name='conteo')
14
15 # Pivotar el DataFrame para preparar los datos para el heatmap.
16 pivot_table = conteo_armas_barrios.pivot(index='BARRIOS_HECHO', columns='ARMAS_MEDIOS', values='conteo').fillna(0).astype(int)
17
18 # Crear un heatmap para visualizar la frecuencia de las armas usadas en cada barrio.
19 plt.figure(figsize=(12, 8)) # Ajusta el tamaño si es necesario
20 sns.heatmap(pivot_table, annot=True, fmt="d", linewidths=.5, cmap='viridis')
21 plt.title('Frecuencia de armas usadas en los 10 barrios con más crímenes')
22 plt.ylabel('Barrio')
23 plt.xlabel('Arma/Medio')
24 plt.show()
25
```



- Se determina con precisión el tipo de arma que más se usa en los delitos cometidos por zona

ALGUNOS ANÁLISIS

```
✓ DELITOS MAS COMUNES POR BARRIO

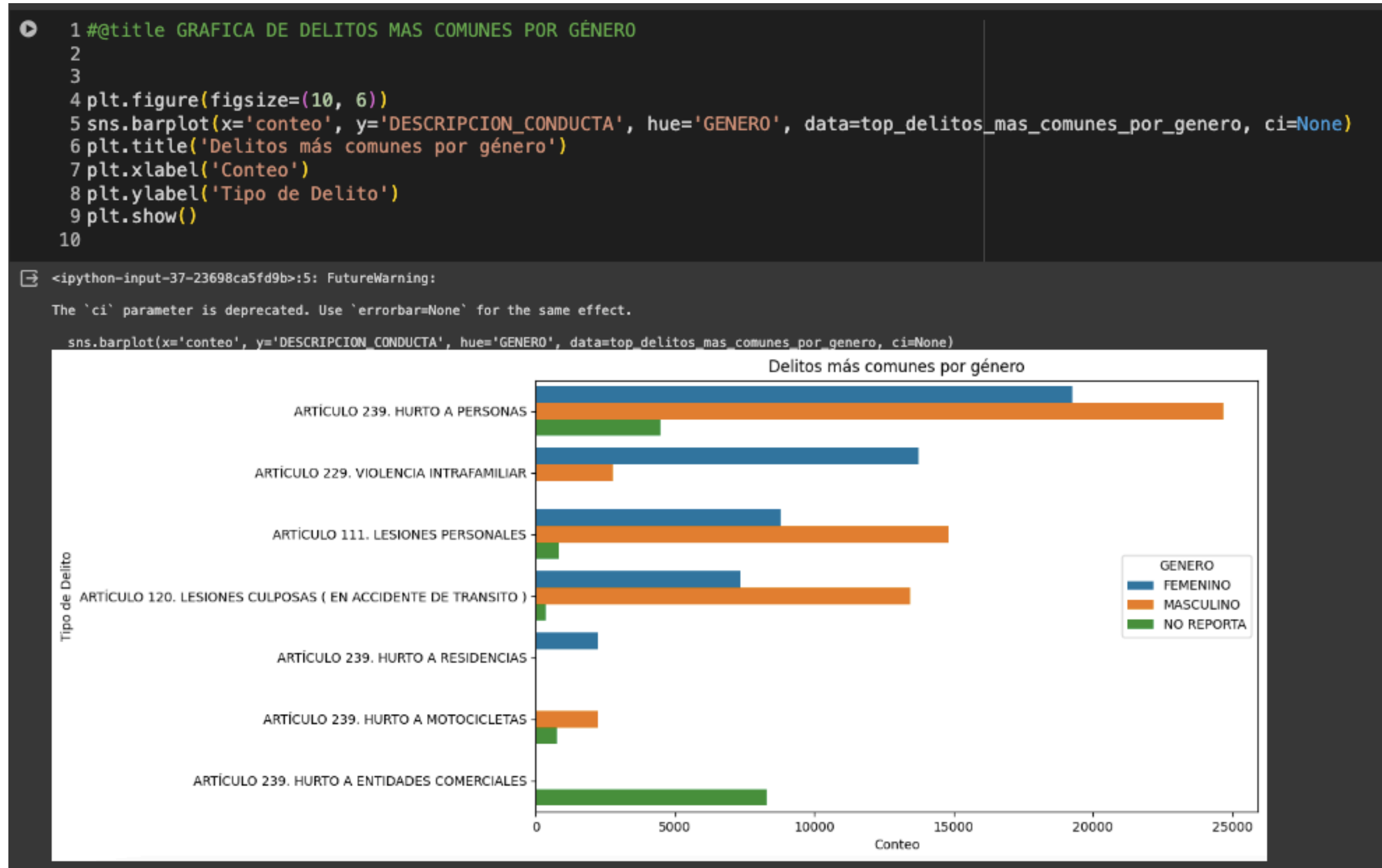
1 #@title DELITOS MAS COMUNES POR BARRIO
2
3 # Agrupamos por barrio y por conducta del delito y contamos las ocurrencias.
4 delitos_por_barrio = a.groupby(['BARRIOS_HECHO', 'CONDUCTA']).size().reset_index(name='conteo')
5
6 # Ordenamos los resultados para obtener los delitos más comunes en cada barrio.
7 delitos_mas_comunes_por_barrio = delitos_por_barrio.sort_values(by=['BARRIOS_HECHO', 'conteo'], ascending=[True, False])
8
9 # Ahora queremos tener solo el delito más común por barrio, por lo que agrupamos y tomamos el primero.
10 delito_mas_comun_por_barrio = delitos_mas_comunes_por_barrio.groupby('BARRIOS_HECHO').head(1)
11
12 # Finalmente, mostramos los resultados.
13 print(delito_mas_comun_por_barrio)
14
```

	BARRIOS_HECHO	CONDUCTA	conteo
13	12 DE OCTUBRE	VIOLENCIA INTRAFAMILIAR	85
18	13 DE JUNIO	HURTO A PERSONAS	13
28	20 DE JULIO	LESIONES PERSONALES	13
41	23 DE JUNIO	VIOLENCIA INTRAFAMILIAR	30
46	5 DE ENERO	LESIONES PERSONALES	7
...
4178	nuevo soto mayor	HURTO A PERSONAS	1
4179	pablo VI	VIOLENCIA INTRAFAMILIAR	2
4183	provenza	LESIONES PERSONALES	13
4185	san rafael	VIOLENCIA INTRAFAMILIAR	1
4187	sin informacion	EXTORSIÓN	22

[494 rows x 3 columns]

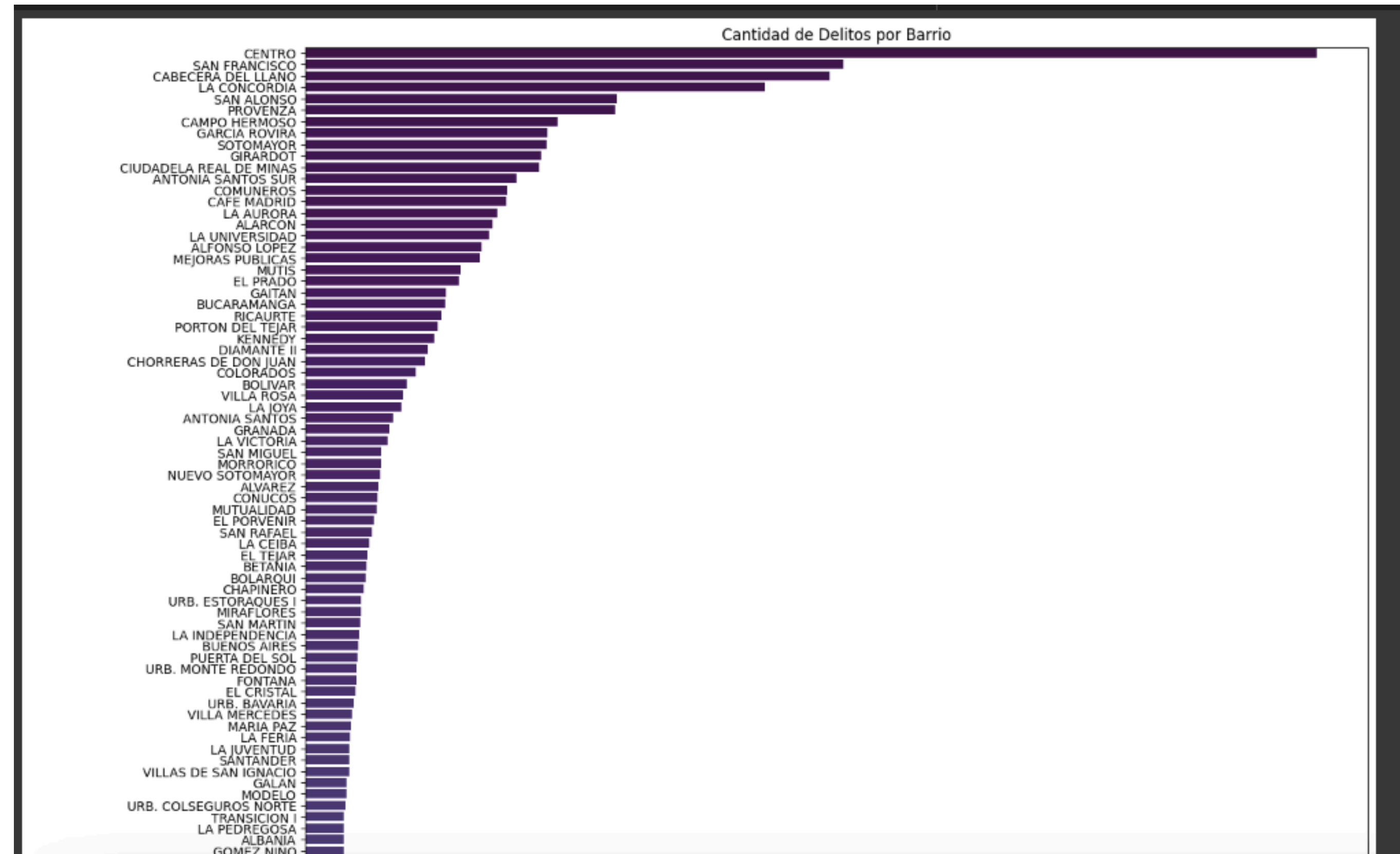
- Se determinan los delitos más comunes por barrio

ALGUNOS ANÁLISIS



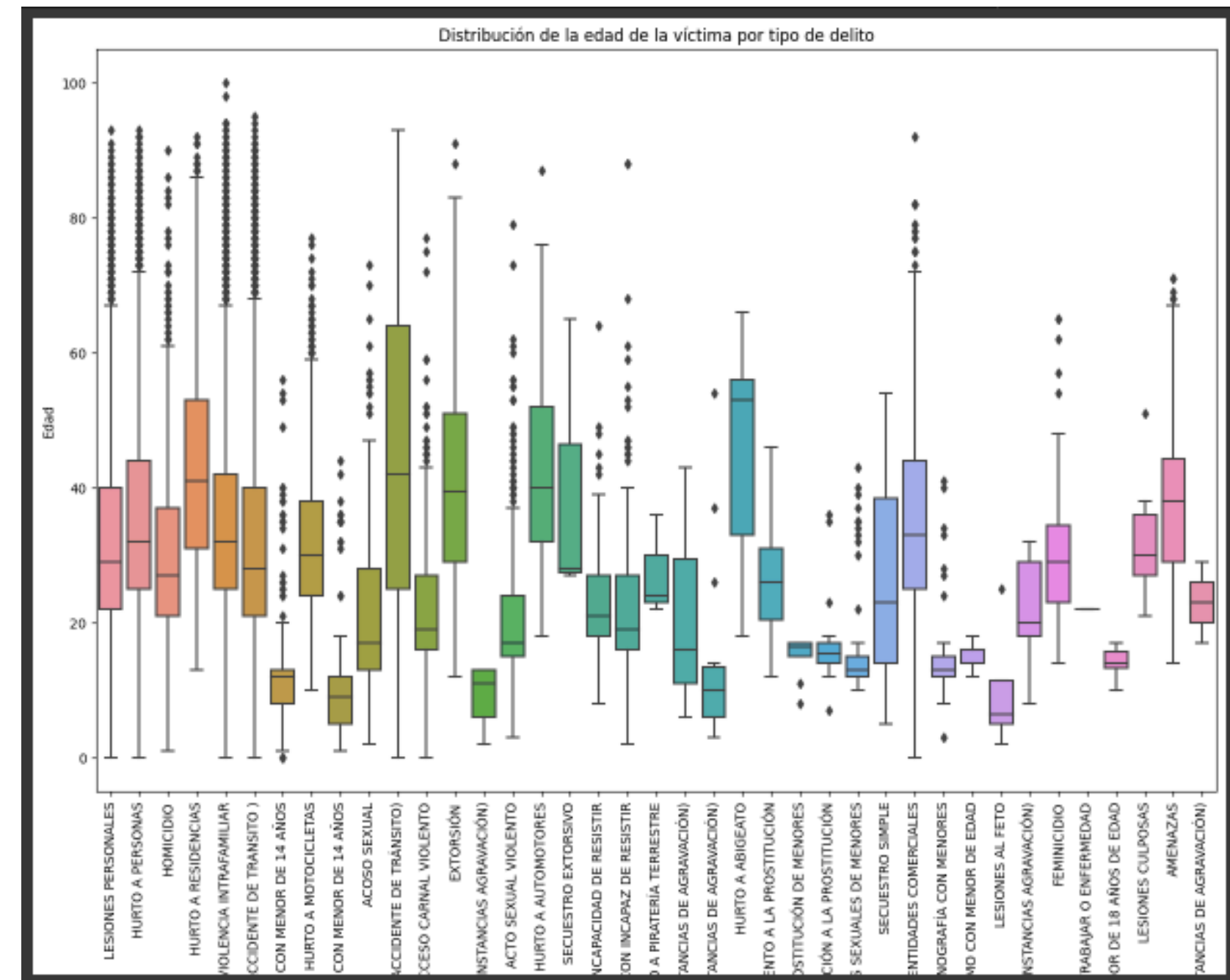
- Los delitos más comunes por género

ALGUNOS ANÁLISIS



- Se obtiene un primer acercamiento al objetivo deseado, determinar las zonas más peligrosas en Bucaramanga

ALGUNOS ANÁLISIS



- Se discriminan los delitos por edades de víctima usando boxplot

**SE APLICA UN PRIMER MODELO DE CLASIFICACIÓN
RANDOM FOREST**

APLICANDO MODELOS DE CLASIFICACIÓN

```
[ ] 1 #librerías
    2 from sklearn.model_selection import train_test_split
    3 from sklearn.ensemble import RandomForestClassifier
    4 from sklearn.metrics import classification_report
    5 from sklearn.model_selection import cross_val_score
    6 from sklearn.preprocessing import LabelEncoder

    1 # Preprocesamiento: Convertir categorías a números
    2 le_barrios = LabelEncoder()
    3 a['BARRIOS_HECHO'] = le_barrios.fit_transform(a['BARRIOS_HECHO'])
    4
    5 le_armas = LabelEncoder()
    6 a['ARMAS_MEDIOS'] = le_armas.fit_transform(a['ARMAS_MEDIOS'])
    7
    8 le_conducta = LabelEncoder()
    9 a['CONDUCTA'] = le_conducta.fit_transform(a['CONDUCTA'])
    10
    11 # División del conjunto de datos
    12 X = a[['BARRIOS_HECHO', 'ARMAS_MEDIOS']] # Características
    13 y = a['CONDUCTA'] # Variable objetivo
    14
    15 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=21)
    16
    17 # Construcción del modelo de clasificación
    18 clf = RandomForestClassifier(n_estimators=100, random_state=21)
    19 clf.fit(X_train, y_train)
    20
    21 # Predicciones
    22 y_pred = clf.predict(X_test)
    23
    24 # Evaluación del modelo
    25 print(classification_report(y_test, y_pred))
    26
    27 # Validación cruzada
    28 scores = cross_val_score(clf, X, y, cv=5)
    29 print("Precisión promedio de validación cruzada: %.2f (+/- %.2f)" % (scores.mean(), scores.std() * 2))
    30
```

■ Implementación de RF


```
warnings.warn(
precision    recall  f1-score   support

     0       0.15     0.02     0.04       189
     1       0.00     0.00     0.00         3
     2       0.07     0.02     0.03        45
     3       0.00     0.00     0.00        24
     4       0.11     0.01     0.02        87
     5       0.00     0.00     0.00         1
     6       0.04     0.02     0.03        43
     7       0.00     0.00     0.00        80
     8       0.00     0.00     0.00         1
     9       0.29     0.17     0.21       312
    10       0.00     0.00     0.00         2
    11       0.25     0.02     0.03        55
    12       0.00     0.00     0.00         2
    13       0.00     0.00     0.00         5
    14       0.00     0.00     0.00         1
    15       0.91     0.83     0.87        93
    16       0.00     0.00     0.00         7
    17       0.35     0.02     0.05       241
    18       0.00     0.00     0.00       128
    19       0.00     0.00     0.00         3
    20       0.00     0.00     0.00        33
    21       0.50     0.02     0.04       215
    22       0.86     0.61     0.71       641
    23       0.71     0.84     0.77     9559
    24       0.00     0.00     0.00         2
    25       0.69     0.12     0.21       929
    27       0.25     0.11     0.15         9
    28       0.00     0.00     0.00         1
    29       0.00     0.00     0.00         1
    30       0.97     0.98     0.97      4198
    31       0.48     0.61     0.54     4849
    32       0.00     0.00     0.00         1
    33       0.50     0.05     0.08        22
    37       0.00     0.00     0.00        15
    38       0.46     0.32     0.38     3218

 accuracy          0.67    25015
 macro avg         0.22     0.14    25015
 weighted avg      0.65     0.67     25015

Precisión promedio de validación cruzada: 0.66 (+/- 0.02)
```

```
38       0.46     0.32     0.38     3218

 accuracy          0.67    25015
 macro avg         0.22     0.14    25015
 weighted avg      0.65     0.67     25015

Precisión promedio de validación cruzada: 0.66 (+/- 0.02)
```

**SE APLICA UN SEGUNDO MODELO DE
CLASIFICACIÓN
SUPPORT VECTOR MACHINE(SVM)**


```
1 999 # División del conjunto de datos
2 from sklearn.svm import SVC
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.pipeline import make_pipeline
5 X = a[['BARRIOS_HECHO', 'ARMAS_MEDIOS']] # Características
6 y = a['CONDUCTA'] # Variable objetivo
7
8 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
9
10 # Para cada kernel, crea un pipeline con un StandardScaler y un SVC con el kernel correspondiente.
11 # Entrena y evalúa el modelo.
12 kernels = ['linear', 'poly', 'rbf']
13 for kernel in kernels:
14     clf = make_pipeline(StandardScaler(), SVC(kernel=kernel))
15     clf.fit(X_train, y_train)
16     y_pred = clf.predict(X_test)
17     print(f"Resultados para kernel {kernel}:")
18     print(classification_report(y_test, y_pred))
19
```

— LINEAR

Resultados para kernel linear:				
	precision	recall	f1-score	support
0	0.00	0.00	0.00	174
1	0.00	0.00	0.00	1
2	0.00	0.00	0.00	43
3	0.00	0.00	0.00	15
4	0.00	0.00	0.00	92
6	0.00	0.00	0.00	53
7	0.00	0.00	0.00	82
8	0.00	0.00	0.00	2
9	0.00	0.00	0.00	278
10	0.00	0.00	0.00	3
11	0.00	0.00	0.00	48
12	0.00	0.00	0.00	2
13	0.00	0.00	0.00	1
14	0.00	0.00	0.00	2
15	0.00	0.00	0.00	79
16	0.00	0.00	0.00	8
17	0.00	0.00	0.00	257
18	0.00	0.00	0.00	127
19	0.00	0.00	0.00	3
20	0.00	0.00	0.00	32
21	0.00	0.00	0.00	226
22	0.00	0.00	0.00	628
23	0.39	1.00	0.56	9701
25	0.00	0.00	0.00	926
27	0.00	0.00	0.00	5
28	0.00	0.00	0.00	2
29	0.00	0.00	0.00	1
30	0.00	0.00	0.00	4205
31	0.00	0.00	0.00	4736
32	0.00	0.00	0.00	1
33	0.00	0.00	0.00	12
34	0.00	0.00	0.00	1
35	0.00	0.00	0.00	1
37	0.00	0.00	0.00	7
38	0.00	0.00	0.00	3261
accuracy			0.39	25015
macro avg	0.01	0.03	0.02	25015
weighted avg	0.15	0.39	0.22	25015

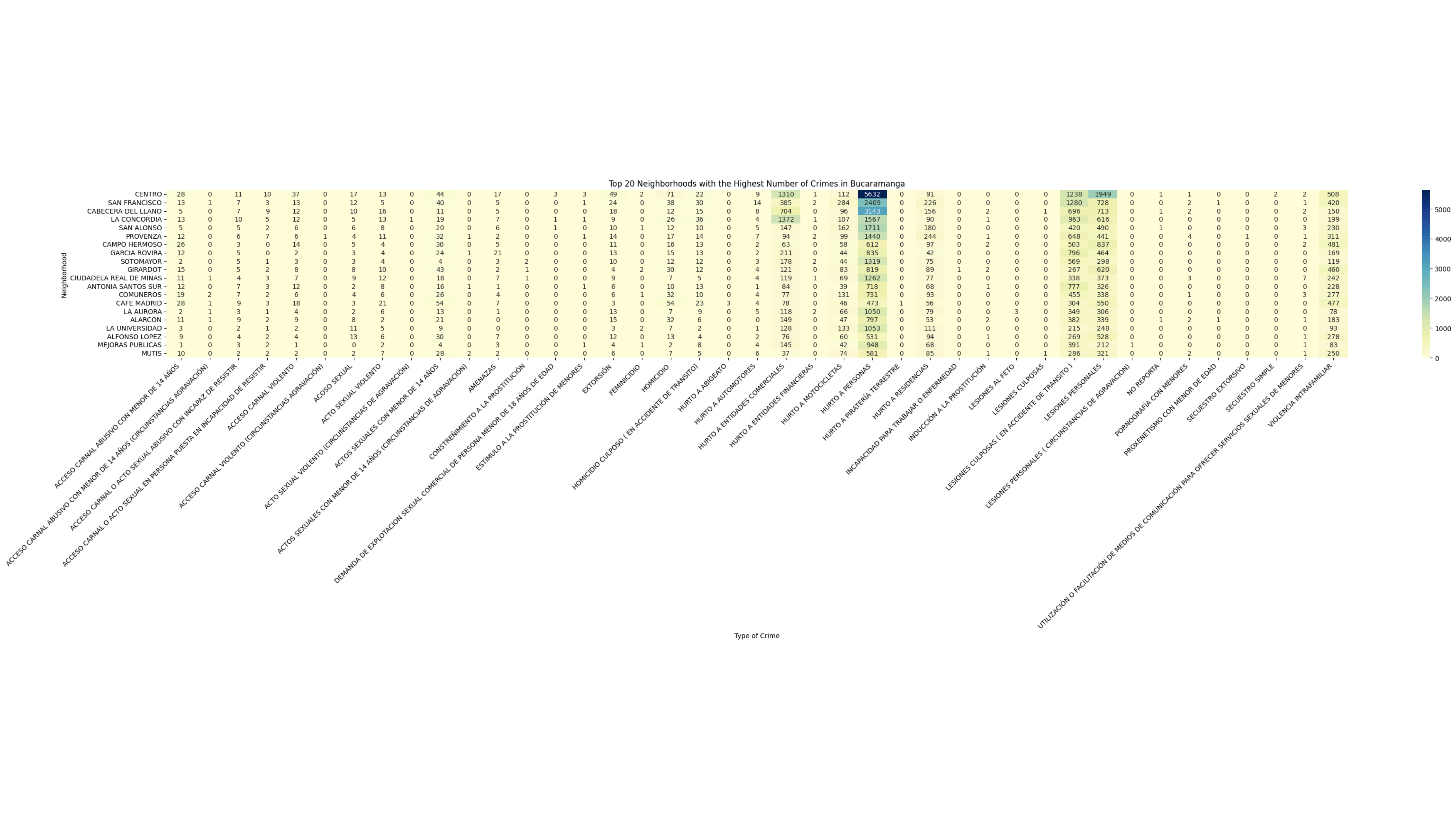
— POLY

Resultados para kernel poly:				
	precision	recall	f1-score	support
0	0.00	0.00	0.00	174
1	0.00	0.00	0.00	1
2	0.00	0.00	0.00	43
3	0.00	0.00	0.00	15
4	0.00	0.00	0.00	92
6	0.00	0.00	0.00	53
7	0.00	0.00	0.00	82
8	0.00	0.00	0.00	2
9	0.00	0.00	0.00	278
10	0.00	0.00	0.00	3
11	0.00	0.00	0.00	48
12	0.00	0.00	0.00	2
13	0.00	0.00	0.00	1
14	0.00	0.00	0.00	2
15	0.00	0.00	0.00	79
16	0.00	0.00	0.00	8
17	0.00	0.00	0.00	257
18	0.00	0.00	0.00	127
19	0.00	0.00	0.00	3
20	0.00	0.00	0.00	32
21	0.00	0.00	0.00	226
22	0.00	0.00	0.00	628
23	0.39	1.00	0.56	9701
25	0.00	0.00	0.00	926
27	0.00	0.00	0.00	5
28	0.00	0.00	0.00	2
29	0.00	0.00	0.00	1
30	0.00	0.00	0.00	4205
31	0.00	0.00	0.00	4736
32	0.00	0.00	0.00	1
33	0.00	0.00	0.00	12
34	0.00	0.00	0.00	1
35	0.00	0.00	0.00	1
37	0.00	0.00	0.00	7
38	0.00	0.00	0.00	3261
accuracy			0.39	25015
macro avg	0.01	0.03	0.02	25015
weighted avg	0.15	0.39	0.22	25015

— RBF

Resultados para kernel rbf:				
	precision	recall	f1-score	support
0	0.00	0.00	0.00	174
1	0.00	0.00	0.00	1
2	0.00	0.00	0.00	43
3	0.00	0.00	0.00	15
4	0.00	0.00	0.00	92
6	0.00	0.00	0.00	53
7	0.00	0.00	0.00	82
8	0.00	0.00	0.00	2
9	0.00	0.00	0.00	278
10	0.00	0.00	0.00	3
11	0.00	0.00	0.00	48
12	0.00	0.00	0.00	2
13	0.00	0.00	0.00	1
14	0.00	0.00	0.00	2
15	0.05	0.08	0.06	79
16	0.00	0.00	0.00	8
17	0.00	0.00	0.00	257
18	0.00	0.00	0.00	127
19	0.00	0.00	0.00	3
20	0.00	0.00	0.00	32
21	0.00	0.00	0.00	226
22	0.79	0.64	0.71	628
23	0.69	0.88	0.77	9701
25	0.58	0.13	0.21	926
27	0.00	0.00	0.00	5
28	0.00	0.00	0.00	2
29	0.00	0.00	0.00	1
30	0.89	0.98	0.93	4205
31	0.44	0.66	0.53	4736
32	0.00	0.00	0.00	1
33	0.00	0.00	0.00	12
34	0.00	0.00	0.00	1
35	0.00	0.00	0.00	1
37	0.00	0.00	0.00	7
38	0.00	0.00	0.00	3261
accuracy			0.65	25015
macro avg	0.10	0.10	0.09	25015
weighted avg	0.54	0.65	0.58	25015

**VISUALIZANDO LOS 20 BARRIOS MÁS PELIGROSOS
Y LOS DELITOS QUE MÁS SE COMETEN**



**PROBANDO MODELOS, COMPARÁNDOS Y
HACIENDO PREDICCIONES**

```

8 from sklearn.compose import ColumnTransformer
9
10 # Preparando los datos para la clasificación
11
12 # Selecciono las columnas que necesito y elimino filas vacías
13 classification_data = clean_data[['BARRIOS_HECHO', 'DIA_SEMANA', 'GENERO', 'CONDUCTA']].dropna()
14
15 # Transformando datos categóricos a numéricos
16 label_encoder = LabelEncoder()
17 onehot_encoder = ColumnTransformer(transformers=[('onehot', OneHotEncoder(), ['BARRIOS_HECHO', 'DIA_SEMANA', 'GENERO'])],
18                                     remainder='passthrough')
19
20 X = onehot_encoder.fit_transform(classification_data[['BARRIOS_HECHO', 'DIA_SEMANA', 'GENERO']])
21 y = label_encoder.fit_transform(classification_data['CONDUCTA'])
22
23 # Dividiendo datos de entrenamiento y prueba
24 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
25
26 # Se inicializan modelos
27 dt_classifier = DecisionTreeClassifier(random_state=42)
28 rf_classifier = RandomForestClassifier(random_state=42)
29 svm_classifier = SVC(random_state=42)
30 mlp_classifier = MLPClassifier(random_state=42)
31
32 # Entrenando modelos
33 dt_classifier.fit(X_train, y_train)
34 rf_classifier.fit(X_train, y_train)
35 svm_classifier.fit(X_train, y_train)
36 mlp_classifier.fit(X_train, y_train)
37
38 # Haciendo predicciones
39 dt_pred = dt_classifier.predict(X_test)
40 rf_pred = rf_classifier.predict(X_test)
41 svm_pred = svm_classifier.predict(X_test)
42 mlp_pred = mlp_classifier.predict(X_test)
43
44 # Evaluando los modelos
45 models = {'Decision Tree': dt_pred, 'Random Forest': rf_pred, 'SVM': svm_pred, 'Neural Network': mlp_pred}
46 performance = {model: accuracy_score(y_test, prediction) for model, prediction in models.items()}
47
48 performance
49
50

```



```

8 from sklearn.compose import ColumnTransformer
9
10 # Preparando los datos para la clasificación
11
12 # Selecciono las columnas que necesito y elimino filas vacías
13 classification_data = clean_data[['BARRIOS_HECHO', 'DIA_SEMANA', 'GENERO', 'CONDUCTA']].dropna()
14
15 # Transformando datos categóricos a numéricos
16 label_encoder = LabelEncoder()
17 onehot_encoder = ColumnTransformer(transformers=[('onehot', OneHotEncoder(), ['BARRIOS_HECHO', 'DIA_SEMANA', 'GENERO'])],
18                                     remainder='passthrough')
19
20 X = onehot_encoder.fit_transform(classification_data[['BARRIOS_HECHO', 'DIA_SEMANA', 'GENERO']])
21 y = label_encoder.fit_transform(classification_data['CONDUCTA'])
22
23 # Dividiendo datos de entrenamiento y prueba
24 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
25
26 # Se inicializan modelos
27 dt_classifier = DecisionTreeClassifier()
28 rf_classifier = RandomForestClassifier()
29 svm_classifier = SVC(random_state=42)
30 mlp_classifier = MLPClassifier()
31
32 # Entrenando modelos
33 dt_classifier.fit(X_train, y_train)
34 rf_classifier.fit(X_train, y_train)
35 svm_classifier.fit(X_train, y_train)
36 mlp_classifier.fit(X_train, y_train)
37
38 # Haciendo predicciones
39 dt_pred = dt_classifier.predict(X_test)
40 rf_pred = rf_classifier.predict(X_test)
41 svm_pred = svm_classifier.predict(X_test)
42 mlp_pred = mlp_classifier.predict(X_test)
43
44 # Evaluando los modelos
45 models = {'Decision Tree': dt_pred, 'Random Forest': rf_pred, 'SVM': svm_pred, 'Neural Network': mlp_pred}
46 performance = {model: accuracy_score(y_test, prediction) for model, prediction in models.items()}
47
48 performance
49
50

```

```

{'Decision Tree': 0.44160961931947357,
 'Random Forest': 0.4441357011140021,
 'SVM': 0.4528759441230707,
 'Neural Network': 0.4445146133831813}

```

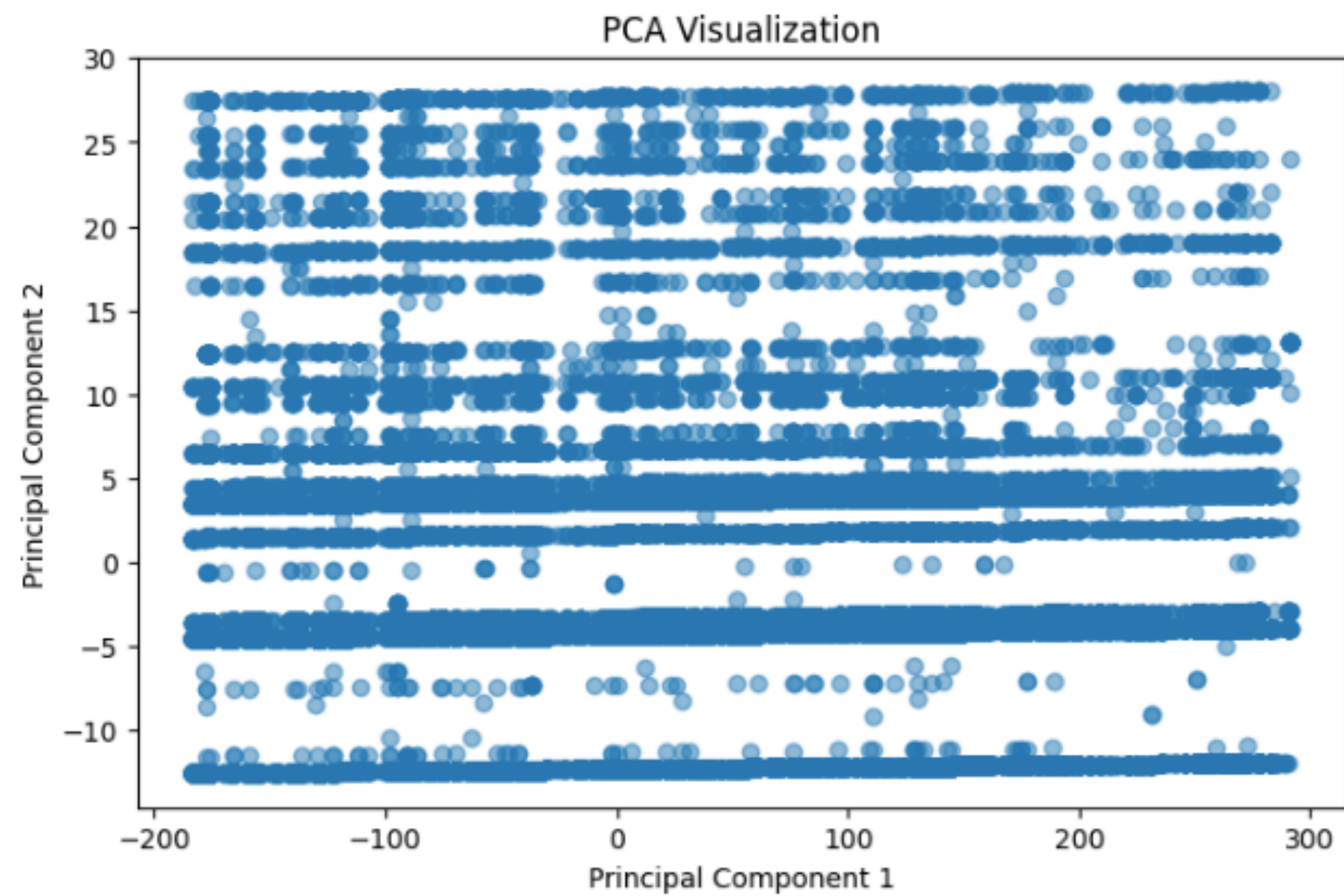
APLICANDO KMEANS, PCA Y T-SNE

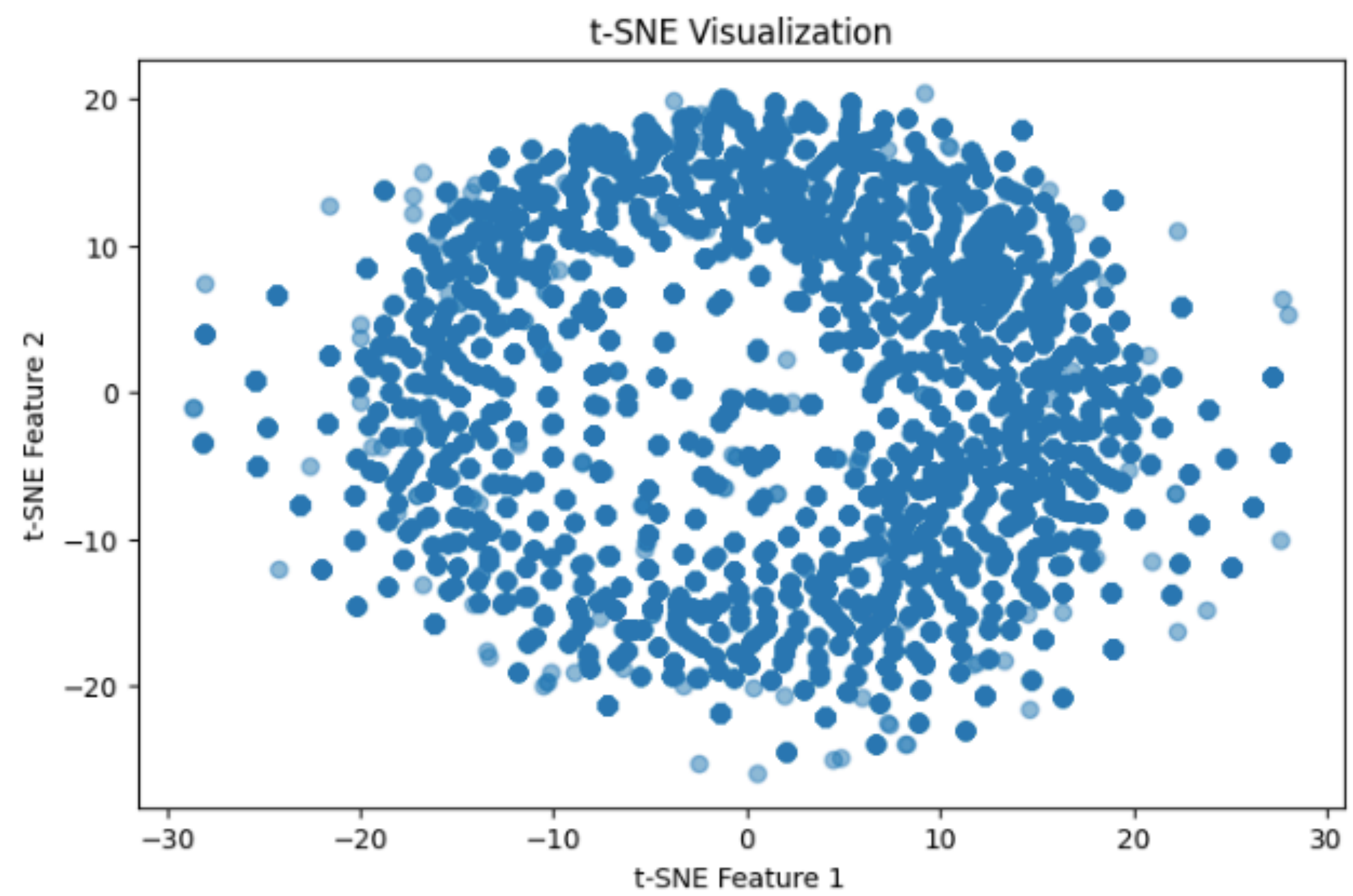
```

1 def get_similar_records(row, model, train_data, le_conducta, le_barrios):
2     # Predecir el cluster para el registro dado
3     conducta_encoded = le_conducta.transform([row['CONDUCTA']])[0]
4     barrios_hecho_encoded = le_barrios.transform([row['BARRIOS_HECHO']])[0]
5     cluster = model.predict([conducta_encoded, barrios_hecho_encoded])[0]
6
7     # Obtener registros del mismo cluster
8     same_cluster_data = train_data[train_data_kmeans.apply(lambda x: model.predict([x])[0], axis=1) == cluster]
9
10    # Seleccionar tres registros aleatorios del mismo cluster
11    if len(same_cluster_data) > 3:
12        return same_cluster_data.sample(3)
13    else:
14        return same_cluster_data
15
16    # Eliminar registros con valores desconocidos en 'BARRIOS_HECHO'
17    known_barrios = set(le_barrios.classes_)
18    test_data = test_data[test_data['BARRIOS_HECHO'].isin(known_barrios)]
19
20    # Aplicar LabelEncoder al conjunto de prueba corregido
21    test_data['BARRIOS_HECHO_encoded'] = le_barrios.transform(test_data['BARRIOS_HECHO'])
22
23    # Probar la función con un registro de prueba
24    example = test_data.iloc[0]
25    similar_records = get_similar_records(example, kmeans_with_categorical, train_data, le_conducta, le_barrios)
26

```

- Inicialmente apliqué k means para los valores numéricos solamente, pero no eran muy significativos para un análisis interesante(orden, año y dia)
- Decidí encontrar relaciones entre tipo de delito y ubicación del delito para eso tomé las columnas correspondientes y apliqué todo el procesamiento en estas columnas categóricas convirtiéndolas en numéricas con label encoding.
- Finalmente creo que se pudo encontrar casos similares o representativos dentro de un conjunto de datos grande. Esto podrá ser útil en aplicaciones como análisis criminal donde encontrar incidentes similares puede ayudar a entender mejor un caso específico.(de acuerdo a análisis previos coincide con los delitos más comunes y las ubicaciones más peligrosas)





CONCLUSIONES

- Se logró extraer información relevante para la toma de decisiones futura
- La visualización de datos aporta en gran medida a la comprensión del dataset y el contexto de la situación

- La limpieza es fundamental para aumentar la confiabilidad de los análisis y resultados
- Se logró identificar las zonas acorde al objetivo del proyecto. Sin embargo se desea implementar un mapa interactivo a nivel web a futuro.

GRACIAS



EL OJO DE BUCAR

Carlos Javier Aceros Blanco
Inteligencia artificial I - Proyecto final
Universidad Industrial de Santander