# MÁSTER EN DATA SCIENCE, BIG DATA & BUSINESS ANALYTICS 2023/2024

# TAREA DE BASE DE DATOS NoSQL

**INSTITUCIÓN**: UCM

ESTUDIANTE: Hao, Qi Xu

## Índice

INTRODUCCIÓN	1
IMPORTACIÓN DE DATOS	1
LOS CAMPOS DE CADA DOCUMENTO (DELITO)	2
MODIFICACIÓN Y CREACIÓN DE UNA NUEVA COLECCIÓN	6
FASE DE ANÁLISIS	8
Dimensión Temporal	8
Dimensión Geográfica	12
Tipología del delito	15
Perfil de las víctimas	18
CONCLUSIÓN	21

## INTRODUCCIÓN

El informe tiene como objetivo el análisis de delitos/crímenes que tomaron lugar y fueron denunciados en la ciudad de Los Ángeles, Estados Unidos, con el cual se quiere identificar patrones en estos delitos bajo varias dimensiones, como la temporal, la ubicación geográfica, la tipología de los delitos y el perfil de las víctimas. Dicho análisis puede ayudar a mejorar a comprender la dinámica de los delitos en esta ciudad e identificar áreas de preocupación y contribuir a la formulación de estrategias de prevención.

Para ello, se han recopilado datos sobre delitos denunciados que tomaron lugar desde el 2020 hasta el 1 de noviembre de 2023 (última actualización) proporcionados por el departamento de policía de la ciudad de Los Ángeles: <a href="https://data.lacity.org/Public-Safety/Crime-Data-from-2020-to-Present/2nrs-mtv8">https://data.lacity.org/Public-Safety/Crime-Data-from-2020-to-Present/2nrs-mtv8</a>

#### IMPORTACIÓN DE DATOS

La colección de datos cuenta con cerca de 830 mil registros sobre delitos denunciados de los últimos 4 años, desde el 2020 hasta la última actualización de datos: 1 de noviembre de 2023.

Como no es posible descargar el fichero en formato JSON directamente, se ha accedido al API<sup>2</sup> para solicitar los datos requeridos para el estudio. Para evitar posibles problemas de rendimiento relacionados con la solicitud de grandes conjuntos de datos, en vez de una sola consulta, se ha realizado la solicitud en 5 consultas SoQL<sup>3</sup> distintas, pidiendo 166,000 registros en cada una, que posteriormente fueron guardados en 5 ficheros en formato JSON.

Una vez importados los 5 ficheros JSON en MongoDBCompass, observamos que hubo en total 829,778 delitos denunciados que ocurrieron durante los últimos 4 años:

 $db.crimes.find(\{\}).count()$ 

1 829778

<sup>&</sup>lt;sup>1</sup> Es importante tener en consideración que el número real de delitos cometidos es mayor que los delitos denunciados dado que no siempre llegan a ser denunciados por diversos motivos.

<sup>&</sup>lt;sup>2</sup> The Socrata Open Data API (SODA). API Endpoint: https://data.lacity.org/resource/2nrs-mtv8.json

<sup>&</sup>lt;sup>3</sup> 1: https://data.lacity.org/resource/2nrs-mtv8.json?\$limit=166000

<sup>2:</sup> https://data.lacity.org/resource/2nrs-mtv8.json?\$limit=166000&\$offset=166000

<sup>3:</sup> https://data.lacity.org/resource/2nrs-mtv8.json?\$limit=166000&\$offset=332000

<sup>4:</sup> https://data.lacity.org/resource/2nrs-mtv8.json?\$limit=166000&\$offset=498000

<sup>5:</sup> https://data.lacity.org/resource/2nrs-mtv8.json?\$limit=166000&\$offset=664000

#### LOS CAMPOS DE CADA DOCUMENTO (DELITO)

Aquellos campos con fondo gris son los campos de interés que serán utilizados para el análisis posterior. No se va a incluir detalles como la cardinalidad y ejemplo para campos excluidos para análisis posterior.

- dr\_no: número de expediente oficial compuesto por un año de 2 dígitos, un identificador de área y 5 dígitos.
- **date rptd**: fecha de la denuncia del delito.
  - **\$** Ejemplo: '2019-04-19T00:00:00.000'
- **date occ**: fecha cuando ocurrió el delito.
  - **\\$** Ejemplo: '2019-04-19T00:00:00.000'
- time occ: la hora cuando ocurrió el delito, en formato militar.
  - ❖ Ejemplo: '1900'<sup>4</sup>
- area: código de las 21 Áreas Geográficas o Divisiones Policiales que hace referencia a un punto de referencia o a la comunidad circundante de la cual son responsables.
  - **❖** Ejemplo: '01'
- **area name**: nombre de cada una de las 21 Áreas Geográficas.
  - ❖ Cardinalidad: 21
  - ❖ Ejemplo: 'Central' (hace referencia al código 01 en "area")
- rpt\_dist\_no: Código de cuatro cifras que representa una subzona dentro de un Área Geográfica.
- part\_1\_2: toma como valores 1 o 2. Delitos de tipo 1 (part I) son aquellos delitos contra la persona o contra la propiedad (delitos de homicidio, violación, tráfico de humanos, robo, entre otros). Delitos tipo 2 se refieren a los delitos restantes (vandalismo, prostitución abuso de drogas, otros asaltos, entre otros). 

  •
- crm\_cd: código que indica el delito principal cometido.
- **crm cd desc**: descripción del código del delito principal cometido.
  - Cardinalidad: 138
  - ❖ Ejemplo: 'BATTERY SIMPLE ASSAULT'
- mocodes: código de 4 dígitos asociado al Modus Operandi; actividades asociadas con el sospechoso en la comisión del delito y los contextos en los que se produjo. Un registro puede tener varios mocodes.
  - **Ejemplo:** '0319 0416'
    - o 0319: Profanity Used
    - o 0416: Hit-Hit w/ weapon
    - 0379: Turns off lights/electricity

<sup>&</sup>lt;sup>4</sup> El equivalente en el formato estándar: 19:00

<sup>&</sup>lt;sup>5</sup> Ver Complete\_user\_manual.pdf para más información.

<sup>&</sup>lt;sup>6</sup> Las descripciones de los códigos están en MO\_CODES\_Numerical.pdf.

- vict age: edad de la víctima.
- vict sex: el sexo de la víctima.
  - Cardinalidad: 5
    - o 'F' (Female), 'M' (Male), 'X' (Unknown), 'H', 'N'
- vict descent: código de una letra que representa cada grupo étnico.
  - Cardinalidad: 20
    - A Other Asian ,B Black, C Chinese, D Cambodian, F Filipino G Guamanian, H Hispanic/Latin/Mexican, I American Indian/Alaskan Native, J Japanese, K Korean, L Laotian, O Other, P Pacific Islander, S Samoan, U Hawaiian, V Vietnamese, W White, X Unknown, Z Asian Indian
- premis\_cd: código de 3 dígitos que representa el tipo de estructura, vehículo o ubicación donde ocurrió el delito.
- **premis\_desc**: la descripción del tipo de estructura, vehículo o ubicación donde ocurrió el delito.
  - Cardinalidad: 306
  - ❖ Ejemplo: 'BANK'
- weapon\_used\_cd: código de 3 dígitos que representa el tipo de arma utilizada en el delito.
- weapon\_desc: descripción del tipo de arma utilizada en el delito.
  - ❖ Cardinalidad: 79
  - ❖ Ejemplo: 'RIFLE'
- **status**: abreviatura del estado en el que se encuentra el caso de delito.
- **status desc**: descripción del estado en el que se encuentra el delito.
  - Cardinalidad: 6
    - o 'Invest Cont': investigación en curso, recopilación pruebas, no hay arresto.
    - 'Adult Arrest', 'Juv Arrest': adulto/joven arrestado en relación con el caso.
    - o 'Adult Other', 'Juv Other': otros estados que no implican arresto.
    - 'UNK': estado desconocido / no especificado.
- **location**: dirección de la calle aproximado donde ocurrió el delito.
- lat: coordenada de latitud donde ocurrió el delito.
- **lon**: coordenada de longitud donde ocurrió el delito.

Se han elegido todos los campos con la descripción de los valores en vez de sus versiones codificados en dígitos.

Se escoge el campo "area\_name" como campo principal geográfico por descarte. Por un lado, "location" presenta una cardinalidad muy alta, llegando a tener 63,821 valores únicos, mientras que "area name" solo presenta 21 categorías. Por otro lado, los campos asociados a las coordenadas

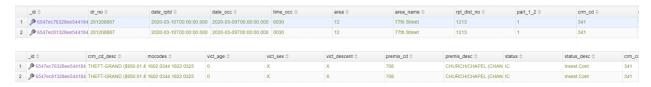
<sup>&</sup>lt;sup>7</sup> No se especifica en ningún lado de la documentación lo que H y N representan.

geográficas de los delitos no se van a poder aprovechar sin poder visualizarlas gráficamente en MongoDB.

Finalmente, se quiere comprobar si cada documento representa un delito distinto asociado a una sola víctima o un mismo documento (delito) puede estar asociado a varias víctimas; es decir, si la colección presenta duplicados del mismo delito con distintas víctimas. Se procede a comprobar con el conteo por cada "dr no", el número de expediente oficial de cada delito denunciado:

Al parecer, hay 3,466 registros duplicados por el campo "dr\_no".

db.crimes.find({dr no: "201208887"})



Después de haber revisado un par de aquellos documentos duplicados al azar, se observa que todos ellos contienen el mismo valor para todos los campos disponibles, incluso en las características de la víctima asociada a cada crimen. Esto sugiere que cada documento (delito) es único y está asociada a una sola víctima. Como los delitos con el mismo número de expediente no aportan ninguna información extra, se procede a eliminar dichos documentos duplicados.

<sup>&</sup>lt;sup>8</sup> La query para eliminar duplicados fue posible gracias en gran medida a la ayuda de CHATGPT.

```
]).forEach(function (doc) {
doc.uniqueIds.shift()
bulkOps.push({
  deleteOne: {
   filter: { _id: { $in: doc.uniqueIds } }
  }
})
})
if (bulkOps.length > 0) {
db.crimes.bulkWrite(bulkOps)
};
      "acknowledged" : true,
"deletedCount" : 3466,
      "insertedCount" : 0,
      "matchedCount" : 0,
      "upsertedCount" : 0,
      "insertedIds" : {
      "upsertedIds" : {
db.crimes_modificado.count()
1 826312
```

Una vez eliminado los registros duplicados, contamos con un total de 826,312 documentos.

#### MODIFICACIÓN Y CREACIÓN DE UNA NUEVA COLECCIÓN

Nótese que algunos campos podrían ser modelizados de manera más adecuada y que todos los campos son de tipo string:

 $db.crimes.find(\{\}).limit(1)$ 

```
"_id" : ObjectId("6547ec5e328ee54418463f00"),
    "dr no" : "010304468",
    "date_rptd" : "2020-01-08T00:00:00.000",
    "date_occ" : "2020-01-08T00:00:00.000".
    "time_occ" : "2230",
    "area" : "03",
    "area name" : "Southwest",
    "rpt_dist_no" : "0377",
    "part_1_2" : "2",
    "crm_cd" : "624",
    "crm_cd_desc" : "BATTERY - SIMPLE ASSAULT",
    "mocodes": "0444 0913",
     "vict_age" : "36",
    "vict_sex" : "F",
    "vict_descent" : "B",
    "premis_cd" : "501",
"premis_desc" : "SINGLE FAMILY DWELLING",
    "weapon_used_cd" : "400"
    "weapon desc" : "STRONG-ARM (HANDS, FIST, FEET OR BODILY FORCE)",
    "status" : "A0",
    "status desc" : "Adult Other".
    "crm_cd_1" : "624",
    "location" : "1100 W 39TH
                                                          PL",
    "lat" : "34.0141",
    "lon": "-118.2978"
}
```

Por ello, se realizarán los siguientes cambios en los campos de interés:

- Los campos relacionados a las características de las víctimas pasan a formar un nuevo campo "victim" con los campos de sus características embedidos. Además, se le cambia también el tipo del campo "vict age" previamente a integer.
- Se crea un nuevo campo "MOcodes" a partir del campo "mocodes" que contiene en su array los códigos de 4 dígitos que representan los métodos y las actividades que han utilizado el criminal en el delito, así como los contextos en los que se produjo.
- Se crea un nuevo campo "**time**" a partir de "time\_occ" en el cual se convierte la hora en formato militar al formato estándar; por ejemplo, de "2230" a "22:30".
- Una vez creada "time", se le concatena con el campo "date\_occ" para formar un nuevo campo "date\_occurrence", el cual contiene la fecha y la hora que tuvo lugar cada delito.
- Por último, creamos un nuevo campo "date\_reported" a partir de "date\_rptd" que tiene datetime como tipo de datos.

Cabe destacar que todas las modificaciones y nuevos campos creados se han realizado en un único pipeline de agregación. Además, se ha optado por crear una nueva colección llamada "crimes\_modificado" con las modificaciones mencionadas aplicadas y filtradas solamente por los campos de interés, listo para la fase de análisis:

```
db.crimes.aggregate([
          {$addFields: {
                  "victim.age": {$toInt: "$vict_age"},
                  "victim.sex": "$vict sex",
                  "victim.descent": "$vict descent",
                 "MOcodes": {$split: ["$mocodes", " "]},
                 "time": {\conormalsection{$\conormalsection{ \color=0,0,2]},":", {\substr: ["$time_occ",2,2]}]}, \color=0,0,2]}, \color=0,0,2], \color=0,0,2]}, \color=0,0,2], \color=0,0,2], \color=0,0,2],
                 "date_reported": {$toDate: "$date_rptd"},
                  "date_occ": {$toDate: "$date_occ"}
         }},
         {$addFields: {
                  "date occurrence": {$toDate: {$concat: [{$substr: ["$date occ", 0, 10]}, "T", "$time", ":00"]}}
         }},
         {$project: {
                 "_id": false,
                  "date occurrence": true,
                  "date occ": true,
                  "date_reported": true,
                  "area name": true,
                 "status desc": true,
                  "crm cd desc": true,
                 "part_1_2": true,
                 "MOcodes": true,
                  "weapon desc": true,
                 "premis desc": true,
                 "victim.age": true,
                 "victim.sex": true,
                  "victim.descent": true
         }},
         {$out: "crimes modificado"}
])
```

db.crimes modificado.find({}).limit(1)

```
_id" : ObjectId("6547f076d9eee3f6725f002f"),
    "date_occ" : ISODate("2020-01-08T01:00:00.000+01:00"),
    "area_name" : "Southwest",
    "part_1_2" : "2",
    "crm_cd_desc" : "BATTERY - SIMPLE ASSAULT",
    "premis_desc" : "SINGLE FAMILY DWELLING",
    "weapon_desc" : "STRONG-ARM (HANDS, FIST, FEET OR BODILY FORCE)",
    "status_desc" : "Adult Other",
     victim" : {
        "age": 36,
        "sex" : "F",
        "descent" : "B"
    "MOcodes" : [ "0444", "0913" ],
    "date_reported" : ISODate("2020-01-08T01:00:00.000+01:00"),
    "date_occurrence" : ISODate("2020-01-08T23:30:00.000+01:00")
}
```

#### FASE DE ANÁLISIS

El estudio de los crímenes en LA (Los Ángeles) se realizará por sus diferentes dimensiones mencionadas anteriormente:

#### Dimensión Temporal

• En primer lugar, queremos conocer cuántos delitos ocurrieron desde 2020 hasta la fecha de la última actualización, 1 de noviembre del 2023:

```
db.crimes_modificado.aggregate([
  {$group: {
    _id: {
      año: {$year: "$date occurrence"}
    count: {$sum: 1}
  }},
  {$sort: {_id: 1}}
])
                                   Value 🕥
▶ [ (1) { "año" : 2020 }
                                   { count : 195961 }
{ count : 209379 }
 { count : 234289 }
▶ [ (4) { "año" : 2023 }
                                   { count : 186683 }
```

Se contempla que la incidencia absoluta de delitos aumenta gradualmente cada año desde el 2020 hasta al 2022, incluidos. Hay cierta incertidumbre en si la incidencia de delitos en el año 2023 va a ser mayor o menor que el año anterior, ya que no se han registrado todavía los datos de los meses de noviembre y diciembre. No obstante, la ausencia de los datos del tamaño de la población de la ciudad de Los Ángeles (LA) en esos años para ajustar la incidencia de crímenes no nos permite

obtener conclusiones válidas. La incidencia criminalidad absoluta, a diferencia de la tasa de criminalidad per cápita, no tiene en cuenta la relación que existe entre el nivel de la población y la tasa de criminalidad. El aumento de la tasa de criminalidad en esos años podría ser atribuido a un aumento significativo de la población en esos períodos.

• ¿La incidencia de crímenes tiende a concentrarse en algunos meses del año?:

```
db.crimes modificado.aggregate([
   {$group: {
      id: {
        año: {$year: "$date occurrence"},
        mes: {$month: "$date occurrence"}
      },
      count: {$sum: 1}
   }},
   {$sort: { id: 1}}
])
▶ [ (1) { "año" : 2020, mes : 1 }
                                          { count : 15890 }
                                                             ▶ [ (25) { "año" : 2022, mes : 1 }
                                                                                                       { count : 18395 }
                                          { count : 16310 }
 ▶ [ (2) { "año" : 2020, mes : 2 }
                                                             ▶ [ (26) { "año" : 2022, mes : 2 }
                                                                                                       { count : 17704 }
 ▶ [ (3) { "año" : 2020, mes : 3 }
                                                             ▶ [ (27) { "año" : 2022, mes : 3 }
                                                                                                       { count : 19794 }
 ▶ [ (4) { "año" : 2020, mes : 4 }
                                                             ▶ [ (28) { "año" : 2022, mes : 4 }
                                                                                                       { count : 19765 }
 ▶ [ (5) { "año" : 2020, mes : 5 }
                                                             ▶ [till (29) { "año" : 2022, mes : 5 }
                                                                                                       { count : 20590 }
 ▶ [ (6) { "año" : 2020, mes : 6 }
                                           { count : 16666 }
                                                             ▶ [ (30) { "año" : 2022, mes : 6 }
                                                                                                       { count : 20255 }
 ▶ [ (7) { "año" : 2020, mes : 7 }
                                          { count : 16630 }
                                                             ▶ [ (31) { "año" : 2022, mes : 7 }
                                                                                                       { count : 19751 }
 ▶ [ (8) { "año" : 2020, mes : 8 }
                                          { count : 16336 }
                                                             ▶ [ (32) { "año" : 2022, mes : 8 }
                                                                                                       { count : 20423 }
 ▶ [ (9) { "año" : 2020, mes : 9 }
                                          { count : 15414 }
                                                             ▶ [ (33) { "año" : 2022, mes : 9 }
                                                                                                       { count : 19304 }
 ▶ [ (10) { "año" : 2020, mes : 10 }
                                          { count : 15876 }
                                                             ▶ [ (34) { "año" : 2022, mes : 10 }
                                                                                                       { count : 20047 }
 ▶ [ (11) { "año" : 2020, mes : 11 }
                                          { count : 15091 }
                                                             ▶ 🛅 (35) { "año" : 2022, mes : 11 }
                                                                                                       { count : 18739 }
 ▶ [ (12) { "año" : 2020, mes : 12 }
                                          { count : 15225 }
                                                             ▶ [ (36) { "año" : 2022, mes : 12 }
                                                                                                       { count : 20031 }
 ▶ [ (13) { "año" : 2021, mes : 1 }
                                          { count : 16113 }
                                                             ▶ [ (37) { "año" : 2023, mes : 1 }
                                                                                                       { count : 20032 }
 ▶ [ (14) { "año" : 2021, mes : 2 }
                                          { count : 15380 }
                                                             ▶ [ (38) { "año" : 2023, mes : 2 }
                                                                                                       { count : 18722 }
 ▶ [ (15) { "año" : 2021, mes : 3 }
                                                             ▶ [ (39) { "año" : 2023, mes : 3 }
                                                                                                       { count : 19338 }
 ▶ [ (16) { "año" : 2021, mes : 4 }
                                                             ▶ [ (40) { "año" : 2023, mes : 4 }
                                                                                                       { count : 19006 }
 ▶ [ (17) { "año" : 2021, mes : 5 }
                                          { count : 16779 }
                                                             ▶ [ (41) { "año" : 2023, mes : 5 }
                                                                                                       { count : 19409 }
 ▶ [ (18) { "año" : 2021, mes : 6 }
                                          { count : 17118 }
                                                             ▶ [ (42) { "año" : 2023, mes : 6 }
                                                                                                       { count : 18938 }
 ▶ [ (19) { "año" : 2021, mes : 7 }
                                          { count : 18740 }
                                                             ▶ [ (43) { "año" : 2023, mes : 7 }
                                                                                                       { count : 20051 }
                                          { count : 18285 }
 ▶ [ (20) { "año" : 2021, mes : 8 }
                                                             ▶ [ (44) { "año" : 2023, mes : 8 }
                                                                                                       { count : 20271 }
 ▶ (21) { "año" : 2021, mes : 9 }
                                          { count : 18256 }
                                                             ▶ [ (45) { "año" : 2023, mes : 9 }
                                                                                                       { count : 19385 }
 ▶ [ (22) { "año" : 2021, mes : 10 }
                                          { count : 18464 }
```

No parece ser el caso. No parece haber una diferencia significativa en la incidencia de delitos en algún mes en particular del año en los últimos 4 años.

• ¿Es más probable que ocurra un delito en una franja horaria específica a lo largo del día?:



- ❖ Mañana (6-12): 27.60%
- **A** Tarde (13-18): 31.59%
- ❖ Noche (19-0): 29.27%
- **❖** Madrugada (1-5): 11.46%

Aparte de la madrugada, período el cual hay menos interacción social, no hay una diferencia significativa a primera vista en la incidencia de delitos entre las distintas franjas horarias del día.

Dado la disponibilidad de la fecha de la denuncia de los delitos se quiere aprovechar para identificar el promedio de días que uno (víctima o testigo) tarda en denunciar el delito después de que dicho evento haya ocurrido:

<sup>&</sup>lt;sup>9</sup> Se utiliza "date\_occ" en vez de "date\_occurrence" en la resta porque tanto el campo "date\_rptd" como "date\_occ" no contienen la información de la hora.

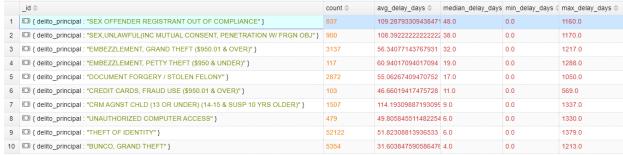
<sup>&</sup>lt;sup>10</sup> La resta de dos fechas en Mongo es devuelta en milisegundos y un día está compuesto de 8,6400,000 milisegundos.

```
}},
  {$group: {
    id: null,
    avg delay days: {$avg: "$rpt delay days"},
    median delay days: {$median: {
     input: "$rpt delay days",
     method: 'approximate'
    }},
    min_delay_days: {\$min: "\$rpt_delay_days"},
    max_delay_days: {$max: "$rpt_delay_days"}
  }}
])
      median_delay_days 4 min_delay_days 4
                                                                                       max_delay_days $
                              avg_delay_days $
                              10.14161963035754 1.0
                                                                   0.0
                                                                                       1379.0
  1 null
```

Se intuye la presencia de valores extremos en el retraso de la denuncia de los delitos, puesto que su promedio difiere significativamente de su mediana. Podemos confirmarlo al observar el retraso mínimo y máximo, de 0 a 1,379 días (casi cuatro años), respectivamente. Por ello, la medida de centralidad más adecuada para la variable "rpt\_delay\_days" es la mediana, ya que no es afectada por los valores extremos. Las víctimas o testigos de los delitos denuncian los crímenes a lo largo de un día de mediana, después de su incidencia.

Por último, se quiere identificar cuáles son los principales delitos cometidos con al menos
 100 observaciones que tuvieron de mediana mayores días de retraso en su denuncia:

```
db.crimes modificado.aggregate([
  {$addFields: {
    rpt_delay_days: {$divide: [{$subtract: ["$date_reported", "$date_occ"]}, 86400000]}
  }},
  {$group: {
    id: {
       delito principal: "$crm cd desc"
    },
    count: {$sum: 1},
    avg_delay_days: {$avg: "$rpt_delay_days"},
    median delay days: {$median: {
     input: "$rpt delay days",
     method: 'approximate'
    }},
    min_delay_days: {$min: "$rpt_delay_days"},
    max delay days: {$max: "$rpt delay days"}
```



La consulta señala que los crímenes que presentan mayores días de retraso en la denuncia fueron principalmente aquellos relacionados con abusos sexuales (1 y 2) y malversación de fondos (3 y 4), así como hurtos mayores de \$950. Los tres primeros delitos principales fueron denunciados de mediana más de un mes después de que ocurrieran. En el caso de los abusos sexuales puede deberse a la vergüenza y al estigma asociado con ser víctimas de tales crímenes, y/o la dependencia económica o emocional de las víctimas del agresor. En cuanto a las malversaciones y hurtos mayores, la tardía denuncia podría ser causado por la complejidad de detección y la recopilación de suficientes pruebas de dichos delitos, especialmente en el caso de las malversaciones puesto que suelen ser personas de altos cargos dentro de las instituciones.

#### Dimensión Geográfica

¿Cuáles son las áreas de la ciudad que concentraron mayor incidencia crímenes):

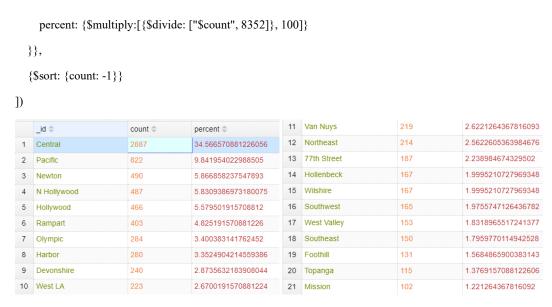
	_id	count \$	percent \$				
1	Central	55862	6.760400	12	West LA	37963	4.594269
2	77th Street	52111	6.306455				
3	Pacific	48258	5.840166	13	Northeast	35670	4.316771
4	Southwest	46358	5.610229	14	Van Nuys	34962	4.231089
5	Hollywood	43810	5.301871	15	West Valley	34573	4.1840128
6	Southeast	41948	5.076532	16	Harbor	34209	4.1399616
7	Olympic	41745	5.051965	17	Topanga	33594	4.065534
8	Newton	41338	5.002710	18	Devonshire	33413	4.0436300
9	N Hollywood	41065	4.969672	19	Mission	32863	3.977069 <sup>-</sup>
10	Wilshire	39361	4.763454	20	Hollenbeck	31014	3.7533038
11	Rampart	38698	4.683218	21	Foothill	27497	3.3276770

Se aprecia una cierta patrón en el resultado de la consulta. Las áreas de las divisiones policiales con más incidencias absolutas de delitos también son las zonas con mayor densidad de población en la ciudad de Los Ángeles: Central, 77th Street, Pacific, Southwest, Hollywood, entre otros. Similarmente, Áreas como Devonshire o Foothill, con menores incidencias de delitos, están bastante alejadas del centro de la ciudad y tienen menor densidad de población.

Especial atención recibe la zona Central de LA ya que en ella se encuentra "Skid Row", un barrio en el que reside la mayor concentración de gente sin techo del estado de California (y de todo Estados Unidos)<sup>11</sup>. Parte de los delitos en esa area tiene lugar entre la comunidad de gente sin techo en las calles de dicho barrio.

Se quiere comprobar la hipótesis de que una fracción significativa de los crímenes en LA Central tiene lugar entre la gente sin techo. Esta consulta es posible gracias a que la variable "MOcodes", la cual contiene una array de códigos que representan los contextos bajo los cuales cada delito tuvo lugar, contiene las siguientes detalles:

<sup>&</sup>lt;sup>11</sup> Skid Row: el infierno de la mayor concentración de indigentes de EE.UU. - BBC News Mundo



Hubo un total de 8,352 delitos en los que tanto las víctimas como los sospechosos fueron gente sin techo. Del total de delitos cometidos entre gente sin techo, más de un tercio de ellos, 34.57%, se concentran en LA Central, confirmando la relativa mayor influencia de los sin techos en los crímenes en dicha área. No obstante, estos delitos constituyen solamente el (2,887 / 55,862) x 100 = 5.17% del total de crímenes cometidos en LA Central en estos últimos 4 años.

• ¿En qué espacio/lugar en concreto tuvieron lugar gran parte de los crímenes?:

	_id	count \$	percent \$
1	STREET	208344	25.213720
2	SINGLE FAMILY DWELLIN	140253	16.973370
3	MULTI-UNIT DWELLING (#	101664	12.303343
4	PARKING LOT	57668	6.9789619
5	OTHER BUSINESS	38920	4.7100852
6	SIDEWALK	35703	4.3207650
7	VEHICLE, PASSENGER/TI	24582	2.9749053
8	GARAGE/CARPORT	16205	1.9611236
9	DRIVEWAY	13436	1.6260201
10	RESTAURANT/FAST FOOI	10579	1.2802670

Se observa que más de la mitad de los crímenes, 54.48%, tienen lugar en 3 lugares específicos: en la calle/vías públicas (25.21%), en viviendas unifamiliares (16.97%), y edificios de viviendas múltiples (12.3%).

#### Tipología del delito

Los delitos cometidos son en Estados Unidos son clasificados por dos categorías: delitos de Parte I (Part I Offenses) y delitos de Parte II (Part II Offenses). Los delitos de Parte I son considerados más graves: homicidio criminal, violación, asalto agravado, trata de personas, allanamiento de morada, robo de automóviles, entre otros. Por otra parte, los delitos de Parte II son el resto considerados menos graves que los de Parte I: asalto no agravado, vandalismo, prostitución, malversación, abuso de drogas, entre otros.

• ¿Cómo se distribuyen los delitos por esta tipología?:

```
db.crimes modificado.aggregate([
     {$group: {
        _id: "$part_1_2",
       count: {$sum: 1}
     }},
     {$project: {
        _id: true,
       count: true,
       percent: {$multiply:[{$divide: ["$count", 826312]}, 100]}
     }},
     {$sort: {percent: -1}}
  ])
                            count $
                                              percent $
1 1
                             483717
                                              58.5392684
2 2
                             342595
                                              41.4607315
```

Según el resultado, hubo más delitos de Parte I que de Parte II en los últimos 4 años, un 17.08% más en concreto.

• ¿Cuáles son los delitos cometidos con mayor incidencia en ambos tipos?:

	_id \$	count -
1	VEHICLE - STOLEN	87106
2	BURGLARY FROM VEHICLE	50916
3	BURGLARY	50506
4	ASSAULT WITH DEADLY WEAPON, AGGRAVATED ASSAULT	47500
5	THEFT PLAIN - PETTY (\$950 & UNDER)	42144
6	THEFT FROM MOTOR VEHICLE - PETTY (\$950 & UNDER)	31744
7	THEFT FROM MOTOR VEHICLE - GRAND (\$950.01 AND OVER)	29969
8	ROBBERY	28336
9	THEFT-GRAND (\$950.01 & OVER)EXCPT,GUNS,FOWL,LIVESTK,PROD	27026
10	SHOPLIFTING - PETTY THEFT (\$950 & UNDER)	19574

El resultado indica que la mayoría de los delitos de Parte I están relacionados con delitos de propiedad. Muchos de ellos están asociados a robos de vehículos en sí y sus motores.

	_id	count 🕶
1	BATTERY - SIMPLE ASSAULT	66121
2	THEFT OF IDENTITY	52122
3	VANDALISM - FELONY (\$400 & OVER, ALL CHURCH VANDALISMS)	50536
4	INTIMATE PARTNER - SIMPLE ASSAULT	41737
5	VANDALISM - MISDEAMEANOR (\$399 OR UNDER)	22148
6	CRIMINAL THREATS - NO WEAPON DISPLAYED	17173
7	TRESPASSING	11738
8	VIOLATION OF RESTRAINING ORDER	10527
9	LETTERS, LEWD - TELEPHONE CALLS, LEWD	6840
10	OTHER MISCELLANEOUS CRIME	5984

A diferencia de los de Parte I, los delitos de Parte II presentan más variedad. Los más comunes son las agresiones físicas, las amenazas de causar daño físico, robo de identidad y vandalismo.

• Finalmente, se quiere conocer las actividades y los contextos más comunes bajos los cuales ocurrieron los delitos durante los últimos 4 años:

	_id	count \$	percent \$
1	1822	284031	34.37333597962997
2	0344	243906	29.51742199072505
3	0913	134352	16.259233800307875
4	0329	109814	13.28965330286865
5	0416	107358	12.992429009865521
6	1300	81022	9.80525515785805
7	2000	66020	7.989718169408165
8	0400	65655	7.945545992312832
9	1402	51497	6.232149599666954
10	2004	45759	5.537738771795642

- ❖ 1822 (Stranger<sup>12</sup>): hubo una persona desconocida involucrada en el 34.37% de los delitos.
- 0344 (Removes vict property): hubo un robo/hurto de la propiedad de la víctima en el 29.52% de los casos.
- ❖ 0913 (Victim knew suspect): la víctima conocía a su agresor/sospechoso en el 16.26% de los casos.
- ❖ 0329 (Vandalized): hubo destrucción de propiedad en el 13.29% de los casos.
- 0416 (Hit-Hit w/weapon): la víctima fue agredida con el uso de cualquier arma en el 13% de los casos.
- ❖ 1300 (Vehicle involved): algún vehículo estuvo involucrado en el 9.8% de los casos.
- ❖ 2000 (Domestic violence): hubo violencia doméstica en el 8% de los casos.
- ❖ 0400 (Force used): la fuerza policial fue utilizado en el 7.94% de los casos.
- ❖ 1402 (Evidence booked): se han registrado pruebas/evidencias en relación con el delito en el 6.23% de los casos.
- ❖ 2004 (Suspect is homeless): el sospechoso era una persona sin hogar en el 5.54% de los casos.

#### Perfil de las víctimas

Por último, se quiere identificar el perfil de víctimas más comunes en los delitos cometidos en los últimos 4 años en LA:

• ¿Cuál es la edad media de las víctimas?:

```
db.crimes modificado.aggregate([
  {$group: {
     id: null,
    avg age: {\$avg: "\$victim.age"},
     median age: {$median: {
     input: "$victim.age",
     method: 'approximate'
     }},
     min age: {$min: "$victim.age"},
    max age: {$max: "$victim.age"}
  }}
])
     _id 🗘
                              avg_age $
                                                 median_age $
                                                                     min_age $
                                                                                       max_age $
                              29.862662045329124 31.0
```

Se percibe un resultado extraño en la edad mínima de las víctimas; -3. No se conoce si se trata de un problema relacionado con la calidad de datos o un código con un significado específico. Por

<sup>&</sup>lt;sup>12</sup> No se especifica lo que realmente representa. Intuyo que se refiere a que hubo personas desconocidas involucradas en el crimen.

ello, se utilizará la mediana para evitar la influencia de datos atípicos. Por consiguiente, las víctimas de los delitos tienen de mediana 31 años.

• ¿Cómo es la distribución de los delitos respecto al sexo de la víctima?:

	_id	count \$	percent \$
1	М	341789	41.363189691061
2	F	304850	36.892844349349886
3	null	107076	12.958301464822004
4	X	72506	8.774651705409095
5	Н	90	0.01089176969474
6	-	1	0.00012101966327488891

No se documenta lo que los valores "H" y "-" representan. Consecuentemente, estos valores no serán interpretados. Se observa que aproximadamente el 41.36% de las víctimas son hombres (M) y el 36.9%, mujeres (F). Por lo tanto, hubo un 4.46% o 36,939 más víctimas hombres que mujeres.

 Por último, se quiere conocer los grupos étnicos más vulnerables en relación con los crímenes:

#### {\$limit: 10}

])

	_id	count \$	percent \$
1	Н	254050	30.745045454985526
2	W	168723	20.41880064672908
3	В	117743	14.249218212975245
4	null	107084	12.959269622128202
5	X	80066	9.689560359767254
6	0	65645	7.944335795680082
7	Α	18143	2.195659750796309
8	К	4419	0.5347858920117341
9	F	3449	0.4173968186350918
10	С	3191	0.38617374551017053

### De los 20 grupos étnicos, los 4 más vulnerables son:

❖ H (Hispanic/Latin/Mexican): 30.75%

W(White): 20.42%B (Black): 14.25%O (Other): 7.44%

#### CONCLUSIÓN

- En los últimos 4 años, desde 2020 hasta el 1 de noviembre de 2023, tuvieron lugar 826,312 delitos denunciados en la ciudad de Los Ángeles, Estados Unidos.
- La incidencia absoluta de delitos ha aumentado gradualmente desde el 2020 hasta al 2022, quedando inconclusa para el 2023, puesto que faltan los datos para los meses de noviembre y diciembre.
  - ❖ No se han identificado patrones estacionales a lo largo del año ni durante el día.
- Los delitos han sido denunciados por las víctimas o testigos de promedio a lo largo de un solo día, posterior a su incidencia.
  - Los delitos principales relacionados con abusos sexuales y malversación de fondos son los que más se han tardado en denunciar, con unas medianas superiores a un mes.
- Las áreas de la ciudad más pobladas presentan mayor incidencia de crímenes que aquellas menos pobladas.
  - ❖ El mayor número de delitos se concentra en LA Central con un total de 55,862 delitos.
  - ❖ LA Central cuenta con la mayor incidencia de delitos entre las personas sintecho de toda la ciudad; el 34.57% de los casos. No obstante, ésta representa solamente el 5.17% del total de delitos cometidos en dicha zona.
- Cerca del 54.48% de los delitos tienen lugar en las calles y las viviendas familiares.
- Los delitos principales más graves están relacionados mayoritariamente con los robos de vehículos y sus motores, así como los hurtos, y los asaltos agravados.
- En cuanto al Modus Operandi y los contextos bajo los cuales los delitos tuvieron lugar:
  - Cerca del 30% de los delitos han supuesto un robo/hurto de la propiedad de la víctima.
  - ❖ La víctima conocía a su agresor/sospechoso en el 16.26% de los casos.
  - ❖ La víctima fue agredida con el uso de cualquier arma en el 13% de los casos.
- El perfil más común de las víctimas es:
  - ❖ Aquellas en un rango de edad cerca de los 31.
  - ♦ Hombres (4.46% más que las mujeres).
  - Grupo étnico latino, el cual constituye cerca del 31% de los casos totales.

<sup>&</sup>lt;sup>13</sup> Se debe tener en consideración la composición demográfica de la ciudad de Los Ángeles antes de sacar conclusiones respecto al perfil de las víctimas más vulnerables. Puede ser el caso de que dicha ciudad tenga de una población significativa de origen latino relativamente a los otros grupos étnicos.