

Procesamiento de Datos para el Análisis Social y Territorial de Nueva York

Primera Entrega – Informe de Avance



Pontificia Universidad
JAVERIANA
Colombia

Autores

Julián Camilo Ramos Granada
María Fernanda Rodríguez Ospina
Sebastián Andrés Rodríguez Pérez

Profesor

John Corredor Franco, PhD

Pontificia Universidad Javeriana Facultad de Ingeniería
Departamento de Ingeniería de Sistemas
Curso: Procesamiento de Alto Volumen de Datos
Octubre de 2025

Índice

1. Introducción	3
2. Entendimiento del negocio	3
2.1. Contexto general de Nueva York	3
2.2. Indicadores macroeconómicos de interés	4
2.3. Objetivos	5
3. Selección de los datos a utilizar	5
3.1. Datasets propuestos por el gobierno	6
3.2. Justificación de selección de datos	6
4. Colección y descripción de datos	7
4.1. Características del clúster Spark	7
4.2. Descripción técnica de los datasets	7
4.2.1. Datos de arrestos en Nueva York (NYPD Arrest Data)	8
4.2.2. Datos de pobreza en Nueva York (NYCgov Poverty Measure Data)	9
4.2.3. Datos de accidentes viales (Motor Vehicle Collisions – Vehicles)	9
4.2.4. Datos de educación y salud escolar (SAT NYC)	10
5. Exploración de los datos	11
5.1. Análisis estadístico descriptivo	11
5.2. Visualizaciones	13
5.3. NYPD Arrest Data	14
5.3.1. Distribución de arrestos por grupo de edad y sexo	14
5.3.2. Top 15 delitos más frecuentes	15
5.3.3. Número de arrestos por distrito (borough)	15
5.3.4. Distribución de arrestos por raza	16
5.3.5. Distribución por tipo legal del arresto	17
5.4. Motor Vehicle Collisions	17
5.4.1. Tipos de vehículos más involucrados	18
5.4.2. Principales factores contribuyentes	19
5.4.3. Puntos de impacto más comunes	20
5.4.4. Condición del vehículo antes del choque	21
5.4.5. Evolución de colisiones por año	22
5.5. Poverty Data	22
5.5.1. Cantidad de registros por distrito	23
5.5.2. Tasa de pobreza por distrito	23
5.5.3. Distribución de ingresos familiares estimados	24
5.5.4. Boxplot de ingreso familiar	25
5.5.5. Ingreso por nivel educativo	25
5.6. SAT NYC	26
5.6.1. Distribución de puntajes de lectura crítica	26

5.6.2. Distribución de puntajes de matemáticas	27
5.6.3. Distribución de puntajes de escritura	27
5.6.4. Comparación de puntajes por componente	28
5.6.5. Relación entre puntajes de lectura y matemáticas	29
5.7. Hallazgos preliminares	29
6. Reporte de calidad de datos	30
6.1. Análisis de valores faltantes	30
6.2. Detección de valores no numéricos en columnas numéricas	31
6.3. Propuesta de tratamiento	32
7. Planteamiento de preguntas sobre los datos	33
7.1. Preguntas principales	33
7.2. Justificación de su relevancia	33
8. Transformaciones, filtrado y limpieza inicial	34
8.1. Transformaciones preliminares	34
8.2. Filtrados aplicados	35
8.3. Limpiezas realizadas	35
9. Web scraping de datos poblacionales	36
10.Consulta climática con OpenWeatherMap	37
11.Conclusiones y recomendaciones	38
Referencias	39

1. Introducción

El presente proyecto tiene como propósito aplicar técnicas de procesamiento de alto volumen de datos para analizar diversas problemáticas sociales y urbanas en la ciudad de Nueva York. A partir de fuentes oficiales de datos abiertos, se busca comprender fenómenos asociados a la seguridad, la movilidad y las condiciones socioeconómicas de la población, con el fin de apoyar la toma de decisiones basada en evidencia. Este trabajo forma parte de un ejercicio consultor en el cual se emplean herramientas de Big Data, usando Apache Spark para procesar, limpiar y explorar grandes volúmenes de información pública, permitiendo identificar patrones relevantes y relaciones entre indicadores territoriales.

La metodología utilizada sigue el enfoque CRISP-DM (Cross Industry Standard Process for Data Mining), que estructura el desarrollo del proyecto en seis fases principales: entendimiento del negocio, entendimiento de los datos, preparación de los datos, modelado, evaluación y despliegue. En esta primera entrega, el énfasis se encuentra en las dos primeras etapas, centradas en conocer el contexto de la ciudad, seleccionar las fuentes de información más pertinentes, realizar la carga y descripción técnica de los datos, y efectuar un análisis exploratorio inicial.

Con ello, se busca establecer una base sólida para las etapas posteriores del proyecto, donde se profundizará en el modelamiento y análisis predictivo. En última instancia, el objetivo es generar hallazgos de valor que contribuyan al diseño de estrategias orientadas a mejorar la seguridad ciudadana y reducir la incidencia de accidentes viales, promoviendo así un desarrollo urbano más seguro, equitativo y sostenible para Nueva York.

2. Entendimiento del negocio

El entendimiento del negocio constituye la primera fase del proyecto y tiene como propósito contextualizar el caso de estudio en el entorno real de la ciudad de Nueva York. En esta sección se analiza la situación actual de la ciudad desde una perspectiva económica, social y territorial, identificando los principales retos que enfrenta en materia de seguridad, movilidad y equidad. Comprender este contexto es importante para orientar el análisis de datos hacia la generación de hallazgos útiles para la toma de decisiones públicas. Así, se parte del conocimiento del territorio y de sus indicadores clave para formular estrategias basadas en evidencia que contribuyan a mejorar los indicadores priorizados por el gobierno: la cantidad de arrestos y la frecuencia de accidentes viales.

2.1. Contexto general de Nueva York

Nueva York es una ciudad emblemática y una de las más influyentes a nivel mundial. Está situada en la costa noreste de Estados Unidos, en la desembocadura del río Hudson, frente al océano Atlántico. Su ubicación estratégica la convierte en un punto clave tanto para la economía estadounidense como para los negocios internacionales. La ciudad se destaca como el principal centro financiero del mundo, albergando la Bolsa de Valores de Nueva York y numerosas instituciones bancarias y financieras globales. Además, posee una riqueza cultural y social extraordinaria gracias a la concentración de comunidades de todo el planeta, lo que ha

moldeado su identidad diversa y multicultural. Esta mezcla cultural se refleja en sus museos, teatros, festivales y en la vibrante vida urbana que atrae a millones de turistas cada año. Nueva York también cumple un papel vital en la política global al ser sede de la Organización de las Naciones Unidas, reafirmando su importancia como un epicentro internacional para la cooperación y la diplomacia mundial.

La relevancia de Nueva York trasciende su economía, ya que representa un símbolo global de diversidad y resiliencia. La ciudad alberga sectores que abarcan desde las finanzas, la tecnología y el comercio, hasta las artes y los medios de comunicación, conformando un ecosistema urbano multifacético y dinámico. Su influencia se extiende a la cultura popular, la moda, la educación superior y la innovación tecnológica, lo que le otorga un lugar privilegiado en el escenario mundial. Conocida como “la ciudad que nunca duerme”, Nueva York marca tendencias y movimientos sociales a nivel global, impactando las políticas económicas, las migraciones y la cultura contemporánea.

A pesar de su estatus y prosperidad, Nueva York enfrenta desafíos sociales significativos que afectan a diversos sectores de su población. La inseguridad, aunque ha disminuido considerablemente en las últimas décadas gracias a políticas públicas y estrategias de vigilancia, sigue siendo un reto en ciertas áreas urbanas donde los índices de criminalidad y violencia son más elevados, impactando especialmente a las comunidades más vulnerables. Así mismo, los accidentes viales constituyen un problema social y de salud pública de gran importancia. A pesar de las medidas pioneras en seguridad vial y del liderazgo de la ciudad en la protección de peatones y ciclistas, las cifras de amenaza continúan generando preocupación, lo que exige esfuerzos constantes para mejorar la movilidad y la seguridad urbana.

Otro problema importante es la pobreza, visible en la marcada desigualdad económica entre distintos sectores de la ciudad. Mientras algunas zonas muestran altos niveles de vida y desarrollo, otras enfrentan limitaciones, falta de vivienda adecuada e inseguridad. Esta brecha socioeconómica repercute en el acceso y la calidad de la educación: a pesar de contar con una amplia oferta educativa y prestigiosas instituciones, las desigualdades afectan principalmente a los barrios más desfavorecidos.

Por lo tanto, Nueva York es una ciudad de contrastes donde conviven la innovación, la riqueza cultural y económica con desafíos sociales profundos. Su posición como nodo económico global y centro cultural de referencia es indiscutible, pero también lo es la importancia de abordar las problemáticas que afectan a su población para garantizar un desarrollo más equitativo y sostenible en el futuro.

2.2. Indicadores macroeconómicos de interés

Nueva York, con una población estimada en 8,48 millones de habitantes a julio de 2024, es una de las ciudades más densamente pobladas de Estados Unidos, lo que implica una gran concentración de recursos y desafíos urbanos. La población creció en aproximadamente 87.000 personas entre julio de 2023 y julio de 2024, un aumento que refleja una recuperación tras las pérdidas generadas por la pandemia, destacándose especialmente el crecimiento en los distritos de Manhattan y Brooklyn. Esta dinámica demográfica se ve influenciada tanto por la migración internacional como por movimientos internos dentro del país.

En materia de empleo, la tasa de desempleo en Nueva York para agosto de 2025 se sitúa en torno al 4,7 %, ligeramente por encima del promedio nacional (4,3 %). Esto refleja

una recuperación laboral gradual, aunque persisten brechas que afectan a ciertos sectores y grupos poblacionales, lo que demanda políticas activas de empleo y capacitación. La tasa de participación laboral, que mide la proporción de personas activas respecto a la población en edad de trabajar, se ubica en torno al 62,3 %.

En el ámbito educativo, la ciudad presenta un alto nivel de formación académica en su población. Para el año 2023, alrededor de 3,08 millones de residentes cuentan con un título universitario, mientras que 2,88 millones poseen educación secundaria completa. Asimismo, cerca de 1,77 millones han alcanzado un título de máster, 245.000 cuentan con un doctorado, y aproximadamente 373.000 personas no tienen educación formal. Estos datos reflejan una ciudad con una amplia base educativa, aunque las desigualdades sociales y económicas aún influyen en el acceso y la calidad de la formación en distintos barrios.

En cuanto a la pobreza, las Directrices Federales de Pobreza de Estados Unidos establecieron en 2023 un umbral de 14.580 dólares anuales para un hogar unipersonal y de 30.000 dólares para uno conformado por cuatro personas. Sin embargo, estos valores no se ajustan regionalmente, por lo que en una ciudad con un costo de vida elevado como Nueva York, las cifras oficiales pueden subestimar la magnitud real de la pobreza. Según estimaciones del Censo (ACS 2017–2021), aproximadamente el 17 % de la población —es decir, unos 1,5 millones de personas— vive por debajo del umbral federal de pobreza. De esta población, cerca del 29 % son menores de 18 años, y las mayores concentraciones de bajos ingresos se encuentran en el sur del Bronx, el Alto Manhattan y diversas zonas de Brooklyn.

En conjunto, estos indicadores evidencian una ciudad con un dinamismo económico y educativo destacado, pero también con desafíos persistentes en materia de desigualdad y acceso equitativo a oportunidades.

2.3. Objetivos

El propósito principal de este proyecto es apoyar al gobierno de la ciudad de Nueva York en la comprensión y mejora de dos indicadores prioritarios: la cantidad de arrestos y la frecuencia de accidentes viales. A partir del análisis de grandes volúmenes de datos públicos, se busca identificar patrones, factores asociados y posibles relaciones entre variables sociales, económicas y territoriales.

El objetivo es utilizar herramientas de procesamiento de datos a gran escala, específicamente Apache Spark, para generar hallazgos que sirvan como base para la toma de decisiones. Con esto se podrán diseñar soluciones concretas y eficientes que contribuyan a reducir los niveles de inseguridad y mejorar la movilidad en la ciudad.

Para esta primera entrega, el alcance del trabajo se centra en el entendimiento del negocio y de los datos disponibles. Esto incluye seleccionar las fuentes de información más relevantes, describir sus características, realizar una exploración inicial y formular preguntas de análisis que orienten las etapas posteriores del proyecto.

3. Selección de los datos a utilizar

Para el desarrollo del proyecto, se seleccionaron diversas fuentes oficiales de datos públicas del portal *Open Data NYC*, priorizando aquellas que permiten analizar los indicadores de

interés definidos por el gobierno de Nueva York: la cantidad de arrestos y los accidentes viales. Adicionalmente, se incorporaron conjuntos de datos socioeconómicos y educativos que facilitan comprender las condiciones del entorno y su posible relación con los fenómenos en estudio.

3.1. Datasets propuestos por el gobierno

Los conjuntos de datos seleccionados y empleados durante esta etapa son los siguientes:

- **NYPD Arrest Data (Year-to-Date)**

Fuente: <https://data.cityofnewyork.us/Public-Safety/NYPD-Arrest-Data-Year-to-Date-/uip8-fykc>

Contiene información sobre los arrestos realizados por el Departamento de Policía de Nueva York (NYPD), incluyendo tipo de delito, edad, sexo, raza, fecha y ubicación del suceso.

- **NYCgov Poverty Measure Data (2018)**

Fuente: <https://data.cityofnewyork.us/City-Government/NYCGov-Poverty-Measure-Data-2018-cts7-vksw>

Proporciona indicadores de pobreza, ingresos y distribución de la población por zonas geográficas.

- **Motor Vehicle Collisions – Vehicles**

Fuente: <https://data.cityofnewyork.us/Public-Safety/Motor-Vehicle-Collisions-Vehicles/bm4k-52h4>

Reúne información sobre vehículos involucrados en accidentes de tránsito, como tipo de vehículo, dirección del accidente, daños e información del conductor.

- **Health and Education Report 2016–2017**

Fuente: <https://data.cityofnewyork.us/Education/2016-2017-Health-Education-Report/2dzy-e7cu>

Incluye indicadores de bienestar y educación escolar en los distritos de Nueva York.

3.2. Justificación de selección de datos

Cada conjunto de datos aporta información específica para la resolución de los objetivos planteados. Los registros de arrestos permiten identificar patrones delictivos y zonas con mayor incidencia de detenciones, lo que permite crear estrategias de seguridad pública efectivas. Los datos de accidentes viales cuentan con evidencia cuantitativa sobre la siniestralidad en el tránsito urbano y sus causas facilitando la propuesta de medidas preventivas y de control. Por otro lado, la información sobre pobreza brinda el contexto económico y social necesario para analizar posibles correlaciones entre condiciones de vulnerabilidad y los indicadores de seguridad o accidentes. Por último, los datos educativos y de salud escolar brindan una visión sobre los factores que pueden influir en la reducción de comportamientos delictivos y en la mejora de las condiciones sociales a largo plazo.

En general, la unión de estos datos permiten tener una visión integral del territorio y establecer relaciones entre variables sociales, económicas y de seguridad . Todo esto con el propósito de contar con información suficiente para formular soluciones de calidad , orientadas a reducir los niveles de inseguridad y la incidencia de accidentes en la ciudad de Nueva York.

4. Colección y descripción de datos

En esta etapa se llevó a cabo la recopilación, carga y revisión inicial de los conjuntos de datos seleccionados. Todos los procesos de lectura y análisis se realizaron sobre un clúster de procesamiento al cual se le realizó una configuración previa, utilizando Apache Spark como herramienta principal.

El objetivo de esta fase es garantizar que los datos estén correctamente integrados en el entorno de trabajo distribuido y que puedan ser procesados de manera eficiente. Para ello, se verificaron los formatos de los archivos, los tipos de datos y la estructura general de cada conjunto.

A continuación, se presentan las especificaciones del clúster utilizado y la descripción técnica de cada uno de los datasets empleados en el proyecto.

4.1. Características del clúster Spark

Para el procesamiento de los datos se implementó un clúster propio sobre máquinas virtuales con sistema operativo Rocky Linux. La infraestructura se configuró manualmente con el fin de simular un entorno distribuido de cómputo y evaluar el rendimiento del procesamiento paralelo en tareas de análisis de datos a gran escala.

El clúster está compuesto por tres nodos: un nodo *master*, que también actúa como *worker*, y dos nodos *worker* adicionales encargados de ejecutar las tareas distribuidas. La comunicación entre los nodos se estableció a través de red local privada, garantizando estabilidad en la transmisión de datos y control sobre los recursos asignados.

Desde el entorno de trabajo en Jupyter Notebook se estableció la sesión de Spark utilizando la biblioteca **PySpark**. La conexión se realizó mediante la inicialización del entorno con `findspark`, la importación de librerías principales y la creación de una sesión a través de `SparkConf()` y `SparkSession`. En esta configuración se definió el modo de planificación FAIR, el cual distribuye las tareas de forma equitativa entre los nodos disponibles, y se asignó el nombre de aplicación `Proyecto_BigData_NY_2025`.

De esta manera, se garantiza una ejecución distribuida eficiente de los procesos de lectura, transformación y análisis de los datos.

4.2. Descripción técnica de los datasets

Una vez configurado el entorno de trabajo, se procedió a la carga y exploración inicial de los conjuntos de datos seleccionados. Cada dataset fue analizado a partir de su estructura, tipo de variables y volumen de registros con el fin de comprender su composición y evaluar la calidad de la información disponible.

A continuación, se describen las principales características técnicas de cada base de datos utilizada en el proyecto.

4.2.1. Datos de arrestos en Nueva York (NYPD Arrest Data)

El primer conjunto de datos corresponde al archivo `nypd.arrests.csv`, el cual contiene información detallada sobre los arrestos efectuados por el Departamento de Policía de Nueva York (NYPD). Los datos fueron obtenidos desde el portal oficial de datos abiertos de la ciudad de Nueva York y cargados en el clúster configurado para el proyecto.

El archivo contiene un total de **143.494 registros** distribuidos en **19 columnas**. Cada registro representa un arresto individual con información demográfica del implicado, el tipo de delito, la localización, entre otros. A continuación, se presenta la descripción general de sus variables principales:

- **ARREST_KEY**: identificador único del arresto.
- **ARREST_DATE**: fecha en que se efectuó el arresto.
- **PD_CD / KY_CD**: códigos de clasificación del delito según el NYPD.
- **PD_DESC / OFNS_DESC**: descripciones internas y estandarizadas del delito.
- **LAW_CODE / LAW_CAT_CD**: código legal violado y categoría del delito (*F*: felonía, *M*: delito menor, *V*: violación).
- **ARREST_BORO / ARREST_PRECINCT**: condado y precinto donde se realizó el arresto.
- **AGE_GROUP, PERP_SEX, PERP_RACE**: variables demográficas del arrestado.
- **Latitude, Longitude**: coordenadas geográficas del arresto.

Durante la revisión inicial se verificó que la mayoría de las columnas no presentan valores nulos. Las únicas excepciones fueron:

- **KY_CD**: 10 valores nulos.
- **LAW_CAT_CD**: 687 valores nulos (equivalentes al 0.48 % del total de registros).

Por otro lado, no se encontraron registros duplicados en el conjunto.

El dataset ofrece un alto nivel de detalle sobre los arrestos realizados en la ciudad, abarcando distintos tipos de delitos y características poblacionales. Mediante este conjunto de datos se busca analizar la distribución territorial de los incidentes y los perfiles más frecuentes de arrestos dentro de los diferentes distritos de Nueva York.

4.2.2. Datos de pobreza en Nueva York (NYCgov Poverty Measure Data)

El segundo conjunto de datos utilizado corresponde al archivo `nycgov_poverty_data.csv`, el cual contiene indicadores oficiales del nivel de pobreza en la ciudad de Nueva York. La información proviene del *NYC Center for Economic Opportunity* y se encuentra publicada en el portal de datos abiertos de la ciudad.

Cada registro representa un año de medición e incluye datos por grupo poblacional, tipo de hogar y zona geográfica. Los indicadores permiten analizar la evolución de la pobreza, la distribución del ingreso y el impacto de los costos de vida en los diferentes distritos de la ciudad.

Entre las variables más relevantes se incluyen:

- **Year:** año de referencia de la medición.
- **Poverty_Rate:** porcentaje de la población que vive bajo la línea oficial de pobreza.
- **Near_Poverty_Rate:** proporción de personas con ingresos apenas por encima del umbral de pobreza.
- **Deep_Poverty_Rate:** porcentaje de la población con ingresos extremadamente bajos.
- **Median_Income:** ingreso familiar medio ajustado por tamaño del hogar.
- **Borough:** distrito o zona geográfica correspondiente (Bronx, Brooklyn, Manhattan, Queens o Staten Island).
- **Population:** población total considerada para el cálculo.

El conjunto de datos cuenta con un total de 55 registros y 7 columnas. Este conjunto permite observar diferencias marcadas entre distritos. El Bronx, por ejemplo, presenta históricamente los índices más altos de pobreza, mientras que Manhattan y Staten Island registran los valores más bajos. Con esto, se busca abordar los temas de seguridad y movilidad desde una perspectiva socioeconómica, facilitando el análisis de posibles correlaciones entre pobreza, vulnerabilidad y niveles de criminalidad o accidentalidad en la ciudad.

4.2.3. Datos de accidentes viales (Motor Vehicle Collisions – Vehicles)

El tercer conjunto de datos corresponde al archivo `motor_vehicle_collisions_vehicles.csv`, el cual contiene información sobre los vehículos involucrados en accidentes de tránsito registrados en la ciudad. La fuente original es el portal de datos abiertos del NYPD, bajo la categoría de seguridad pública.

Cada registro representa un vehículo asociado a un evento de colisión, identificado por un código único. El conjunto de datos cuenta con más de 4 millones de registros y 25 columnas, siendo este uno de los más extensos y relevantes para el proyecto. La información se considera estructurada y cuenta con un nivel de detalle alto, abarcando variables técnicas, direccionales y descriptivas sobre el incidente.

Entre las variables principales se incluyen:

- `COLLISION_ID`: identificador único del evento de colisión.
- `CRASH_DATE` y `CRASH_TIME`: fecha y hora del accidente.
- `VEHICLE_TYPE`, `VEHICLE_MAKE`, `VEHICLE_MODEL` y `VEHICLE_YEAR`: características del vehículo involucrado.
- `DRIVER_SEX`, `DRIVER_LICENSE_STATUS` y `TRAVEL_DIRECTION`: información sobre el conductor y su comportamiento previo al accidente.
- `POINT_OF_IMPACT`, `VEHICLE_DAMAGE` y `PUBLIC_PROPERTY_DAMAGE`: detalles del impacto y nivel de daño.
- `CONTRIBUTING_FACTOR_1` a `CONTRIBUTING_FACTOR_5`: factores que contribuyeron al accidente según los reportes del NYPD.

Durante la revisión inicial, se observó que el dataset incluye registros con valores nulos en algunas columnas descriptivas, principalmente en los factores contribuyentes (`CONTRIBUTING_FACTOR_4` y `CONTRIBUTING_FACTOR_5`), e indica que no todos los reportes contienen información completa. No se identificaron duplicados en el campo `COLLISION_ID`.

Con este conjunto se busca identificar patrones de siniestralidad vial, tipos de vehículos más involucrados y condiciones asociadas a los accidentes.

4.2.4. Datos de educación y salud escolar (SAT NYC)

El último conjunto de datos corresponde al archivo `SAT_Results_NYC.csv`, el cual contiene los resultados promedio de las pruebas SAT aplicadas a estudiantes de educación secundaria en la ciudad de Nueva York durante el año 2012. La información fue descargada desde el portal *NYC Open Data*, cuya última actualización se registró el 26 de noviembre de 2024.

El archivo fue cargado exitosamente en el entorno de Spark, con un total de **478 registros** y **6 columnas**. A diferencia de los datasets anteriores, este conjunto tiene un tamaño menor, pero ofrece una visión agregada del desempeño académico de los estudiantes en tres áreas clave: lectura crítica, matemáticas y escritura. Esta información puede servir como referencia para explorar correlaciones entre el nivel educativo y las tasas de criminalidad o accidentalidad por distrito.

Las variables incluidas en el conjunto son las siguientes:

- `DBN`: identificador único de la escuela (District Borough Number).
- `SCHOOL_NAME`: nombre oficial de la institución educativa.
- `Num_of_SAT_Test_Takers`: número de estudiantes que presentaron el examen SAT.
- `SAT_Critical_Reading_Avg_Score`: puntaje promedio de lectura crítica.
- `SAT_Math_Avg_Score`: puntaje promedio de matemáticas.
- `SAT_Writing_Avg_Score`: puntaje promedio de escritura.

Durante la carga, se observó que las columnas correspondientes a los puntajes y al número de estudiantes fueron interpretadas como tipo `string`, lo cual indica la presencia de valores no numéricos. Este comportamiento se debe a la existencia de registros con caracteres como `'s'` o celdas vacías, que impiden la conversión automática a tipo numérico.

El análisis de calidad evidenció que no existen valores nulos explícitos (`null`) ni registros duplicados. Sin embargo, se identificaron **57 registros con valores no numéricos** en las columnas relacionadas con los puntajes y el número de participantes, específicamente en:

- `Num_of_SAT_Test_Takers`
- `SAT_Critical_Reading_Avg_Score`
- `SAT_Math_Avg_Score`
- `SAT_Writing_Avg_Score`

Estos valores serán tratados como datos faltantes en la fase de limpieza, mediante su conversión a `null`.

5. Exploración de los datos

La fase de exploración de los datos tiene como objetivo comprender en profundidad la estructura, el comportamiento y las características generales de cada conjunto de datos utilizado en el proyecto. A través de un análisis estadístico descriptivo y la generación de diversas visualizaciones, se buscó identificar patrones, tendencias y posibles anomalías que sirvan como base para el análisis posterior. Esta etapa permitió reconocer diferencias relevantes entre distritos, categorías demográficas y niveles socioeconómicos, así como evidenciar relaciones preliminares entre las variables de interés.

5.1. Análisis estadístico descriptivo

Con el propósito de comprender mejor la estructura y el comportamiento general de los datos, se realizó un análisis estadístico descriptivo sobre los cuatro conjuntos trabajados. Esta etapa busca identificar patrones relevantes, diferencias entre grupos y valores predominantes en las variables clave, sirviendo además como una verificación final de la calidad de los datos antes de avanzar hacia etapas más complejas como visualizaciones o modelado.

Arrestos en Nueva York (NYPD Arrest Data).

El análisis reveló que el grupo de edad con mayor incidencia de arrestos fue el de personas entre 25 y 44 años, seguido por los adultos de 45 a 64. En contraste, los menores de edad y los adultos mayores de 65 años presentaron proporciones significativamente menores. Al examinar la variable de sexo, se observó que la diferencia es bastante evidente: más del 80 % de los arrestos fueron a hombres, mientras que las mujeres representaron poco menos del 20 %. Por su parte, la variable de raza mostró que las personas identificadas como negras y como hispanas (tanto blancas como negras) concentraban la mayoría de los registros. En términos territoriales, Brooklyn, Manhattan y el Bronx fueron los distritos que agruparon el

mayor volumen de arrestos, siendo los precintos 14, 40 y 75 algunos de los más recurrentes. Finalmente, al observar la variable de tipo de delito, se encontró que las agresiones menores, el hurto y los delitos relacionados con drogas fueron los más comunes. Cerca del 60 % de los delitos correspondían a la categoría legal de “delito menor” (M), mientras que el 40 % restante fueron clasificados como felonías (F).

Resultados académicos (SAT NYC).

Este conjunto contiene los puntajes promedio del examen SAT para 421 instituciones educativas de la ciudad. Una vez procesado el conjunto y corregidos los errores de tipado, se realizó un análisis de los puntajes por área. Los resultados generales se resumen en la siguiente tabla:

Área	Media	Mediana	Máximo
Lectura crítica	400.9	391	679
Matemáticas	413.4	395	735
Escritura	394.0	381	682

Como se puede evidenciar, las tres áreas mostraron un rendimiento similar, aunque matemáticas tuvo un desempeño ligeramente más alto. Las distribuciones fueron relativamente simétricas, con medianas cercanas a las medias. Sin embargo, también se identificaron escuelas con puntajes muy por debajo o muy por encima del promedio, lo que refleja diferencias significativas entre instituciones.

El número de estudiantes que presentó el examen varió ampliamente. Si bien el promedio fue de 110 estudiantes por colegio, la mediana fue de solo 62, lo que indica una distribución asimétrica. Hubo instituciones con menos de 10 participantes y otras con más de 1.200, lo que puede estar asociado al tamaño de la institución.

Accidentes viales (Motor Vehicle Collisions).

Este conjunto es el más extenso, con más de 4 millones de registros. Luego de su limpieza, se trabajó sobre 17 variables seleccionadas. Entre los tipos de vehículos más involucrados se encuentran los sedanes, station wagons, SUVs y taxis. Cabe destacar que una parte considerable de los registros originalmente no especificaba el tipo de vehículo, con etiquetas como *unknown* o *unspecified*, las cuales fueron unificadas bajo la categoría **NO_INFO** para facilitar su análisis y contabilización. Esta consolidación permitió visualizar con mayor claridad la proporción de datos faltantes en variables clave sin dispersar los resultados en múltiples formas de “sin información”.

En cuanto a las causas de los accidentes, más del 60 % de los registros indicaban que el “factor contribuyente” era “no especificado”. Estos también fueron agrupados bajo la misma etiqueta **NO_INFO**. Aun así, dentro de los factores que sí fueron reportados, los más frecuentes fueron la distracción del conductor, el no ceder el paso, y el exceso de velocidad. Estos resultados sugieren que una buena parte de los accidentes podría estar relacionada con comportamientos evitables.

También se analizaron los años de fabricación de los vehículos. El promedio fue cercano a 2014, con una mediana similar. Sin embargo, se detectaron valores extremos como 1000 o 20063, lo cual claramente corresponde a errores de ingreso. En cuanto a la ubicación del impacto, los golpes frontales fueron los más comunes, seguidos por impactos laterales y traseros. De forma similar a lo anterior, los valores no informados en esta categoría fueron

homogeneizados como **NO_INFO**. Solo un pequeño porcentaje de los registros confirmó daños a propiedad pública, mientras que la gran mayoría marcó esta variable como “no” o “desconocido”, también agrupados bajo la misma etiqueta para mantener consistencia en el tratamiento de los datos faltantes.

Condiciones de pobreza (Microdatos NYCgov).

El conjunto de pobreza utilizado corresponde a una base de microdatos individuales. En total se analizaron 68.273 registros, distribuidos principalmente entre los distritos de Brooklyn, Queens y el Bronx. Al usar la variable `NYCgov_Pov_Stat`, que identifica si una persona se encuentra bajo la línea de pobreza, se estimó una tasa de pobreza general del 17.7 %. La distribución por distrito se resume en la siguiente tabla:

Distrito	Código	Tasa de pobreza (%)
Bronx	1	25.8
Brooklyn	2	18.3
Manhattan	3	13.8
Queens	4	15.8
Staten Island	5	13.8

El Bronx se consolida como el distrito con mayor proporción de personas en situación de pobreza, mientras que Manhattan y Staten Island presentan las tasas más bajas.

En cuanto al ingreso ajustado por la metodología del gobierno de la ciudad (`NYCgov_Income`), el promedio fue de aproximadamente \$79.120, pero la mediana fue notablemente más baja (\$60.122), lo que indica una distribución desigual. Se observaron valores extremos, desde ingresos negativos hasta cifras que superaban los \$899.000. Estos casos podrían corresponder a familias con múltiples fuentes de ingreso.

En conjunto, este análisis permitió conocer mejor la forma en que están estructurados los datos y los patrones generales presentes en cada conjunto. Las diferencias encontradas entre distritos, categorías demográficas y niveles de ingreso generan una base para continuar con la fase de visualización e interpretación de correlaciones entre variables sociales, educativas y territoriales.

5.2. Visualizaciones

A partir de los distintos conjuntos de datos utilizados, se generaron gráficos que complementan los análisis estadísticos previos y facilitan la comprensión de fenómenos sociales, territoriales y demográficos en la ciudad.

En cada caso se seleccionaron variables clave, tanto cuantitativas como categóricas, y se representaron mediante histogramas, gráficos de barras, mapas de calor o boxplots, dependiendo del tipo de dato y del objetivo de análisis.

Cada subsección presenta una serie de cinco visualizaciones por conjunto de datos (arrestos, accidentes, SAT, pobreza), acompañadas de una breve interpretación que resalta los principales hallazgos. Estas gráficas permiten comparar distritos, identificar rangos etarios más afectados, explorar distribuciones económicas, evaluar diferencias raciales, y detectar comportamientos atípicos o categorías dominantes.

5.3. NYPD Arrest Data

Las visualizaciones seleccionadas para este conjunto buscan resaltar los aspectos más representativos de los arrestos en la ciudad durante 2025, complementando el análisis estadístico previo. Se eligieron variables clave como edad, sexo, tipo de delito, distrito y raza del arrestado, ya que permiten identificar patrones demográficos y geográficos relevantes. Además, se incluyó una gráfica que diferencia entre delitos mayores y menores, para tener una visión más clara del tipo de conductas sancionadas.

5.3.1. Distribución de arrestos por grupo de edad y sexo

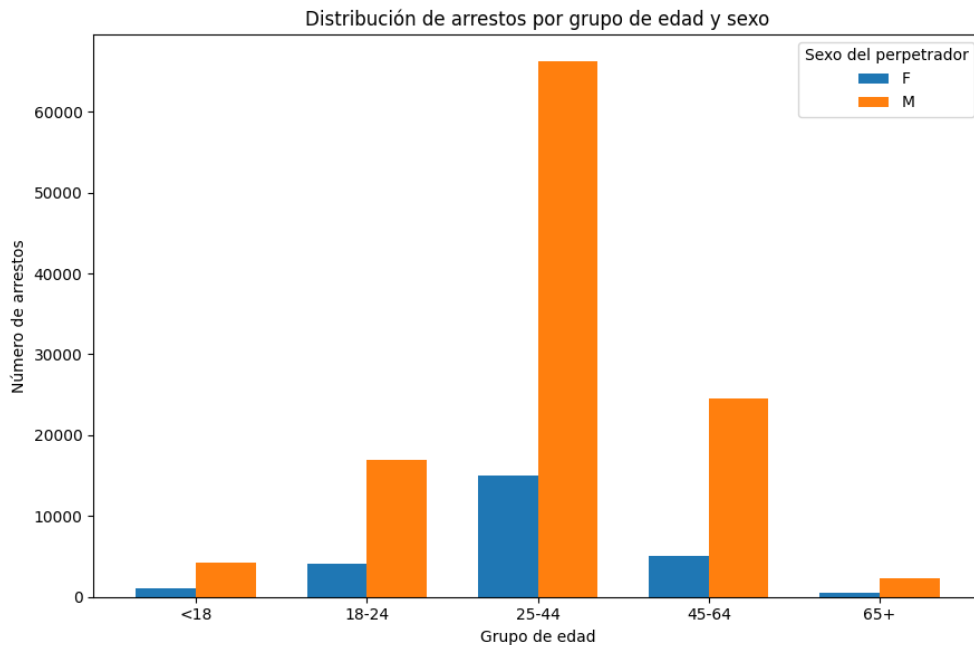


Figura 1: Distribución de arrestos por grupo de edad y sexo.

Este gráfico muestra la cantidad de arrestos según distintos grupos etarios, diferenciando por sexo. Se observa que la mayoría de los arrestos ocurre entre personas de 25 a 44 años, con una notable predominancia de hombres en todos los rangos de edad.

5.3.2. Top 15 delitos más frecuentes

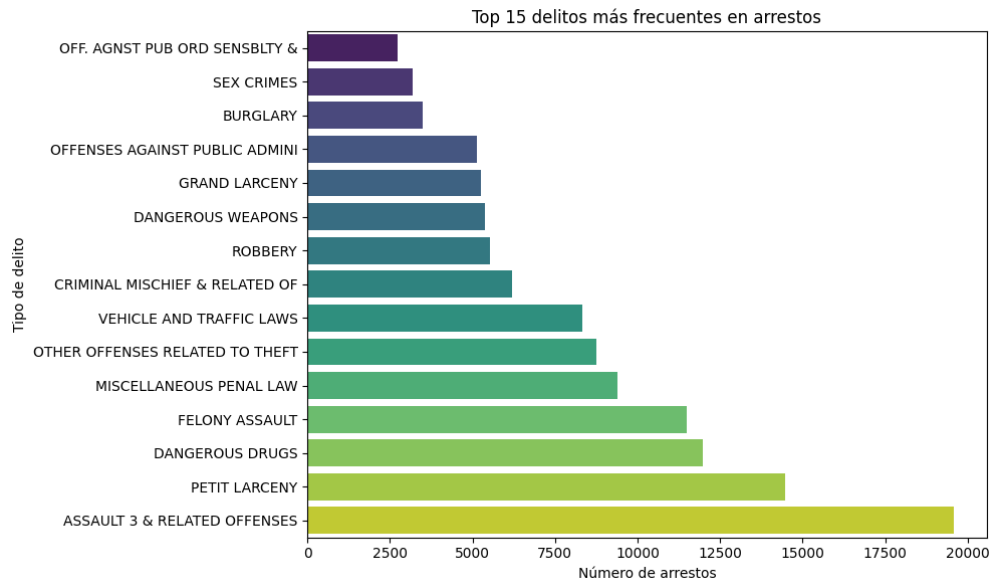


Figura 2: Quince delitos más frecuentes en los arrestos.

Se representa la frecuencia de arrestos según el tipo de delito. Encabezan la lista el asalto en tercer grado, el hurto menor y los delitos relacionados con drogas.

5.3.3. Número de arrestos por distrito (borough)

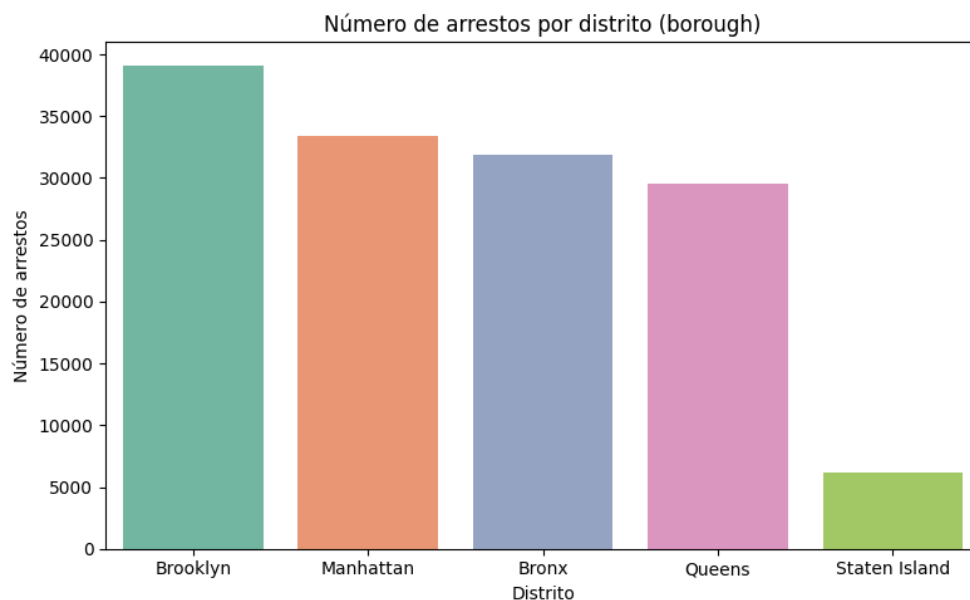


Figura 3: Total de arrestos por borough.

Brooklyn es el distrito con mayor cantidad de arrestos, seguido por Manhattan, Bronx y Queens. Esta distribución territorial puede estar asociada a factores como densidad poblacional o características socioeconómicas.

5.3.4. Distribución de arrestos por raza

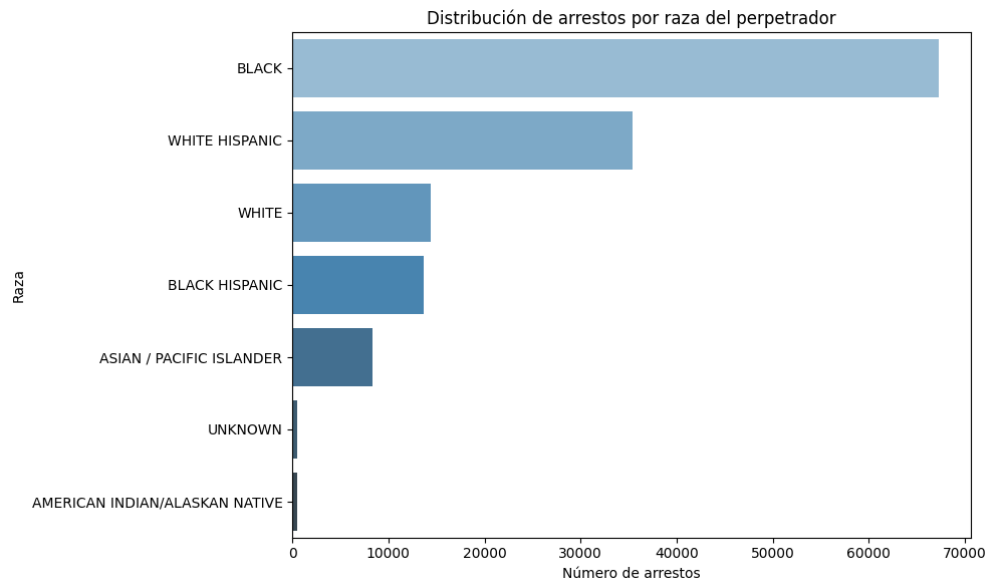


Figura 4: Distribución de arrestos según la raza del perpetrador.

En esta visualización se aprecia que las personas identificadas como Black y White Hispanic son las más representadas en los registros de arresto.

5.3.5. Distribución por tipo legal del arresto

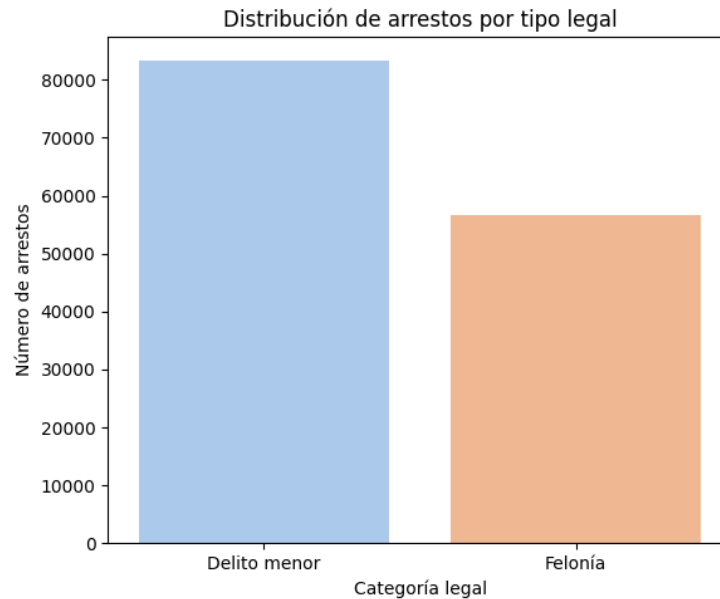


Figura 5: Distribución por tipo legal: felonía o delito menor.

La mayoría de los arrestos registrados corresponden a delitos menores. Esta gráfica ayuda a dimensionar qué tan grave es, en promedio, la conducta sancionada, y puede ser útil para reflexionar sobre prioridades en políticas de justicia y seguridad.

5.4. Motor Vehicle Collisions

En el caso del conjunto de accidentes vehiculares registrados en la ciudad de Nueva York, las visualizaciones permiten observar con mayor detalle las características más relevantes de los siniestros y las tendencias que se desprenden de ellos. Estas gráficas complementan los resultados estadísticos al representar de forma visual los tipos de vehículos más involucrados, las causas más comunes, las zonas de impacto, la condición de los vehículos al momento del choque y la evolución temporal de los accidentes.

Cabe señalar que, durante el proceso de limpieza, se identificó una gran cantidad de registros con valores no informados o ambiguos en distintas columnas. Por coherencia y para facilitar la interpretación, todas las categorías de este tipo fueron unificadas bajo la etiqueta `NO_INFO`. Esto permitió manejar de manera más clara los datos faltantes y distinguir entre la información efectivamente registrada y aquella que no fue reportada.

5.4.1. Tipos de vehículos más involucrados

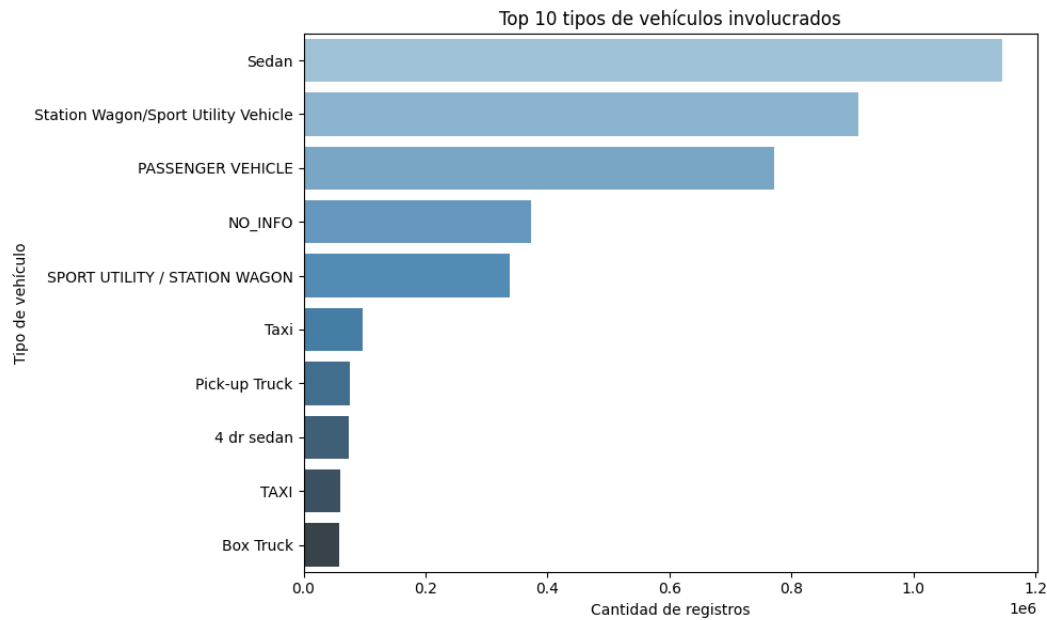


Figura 6: Top 10 tipos de vehículos más involucrados en accidentes.

El primer gráfico muestra los tipos de vehículos más frecuentemente implicados en colisiones. Los automóviles tipo *sedán* encabezan la lista, seguidos por los *station wagons* o SUVs y los vehículos clasificados genéricamente como *passenger vehicle*. La presencia de un número considerable de casos marcados como *NO_INFO* refleja limitaciones en la calidad del registro y la falta de precisión en algunos reportes, aunque no altera las tendencias generales. Este resultado indica que los vehículos de uso cotidiano son los que más se ven expuestos al riesgo de colisión dentro del entorno urbano.

5.4.2. Principales factores contribuyentes

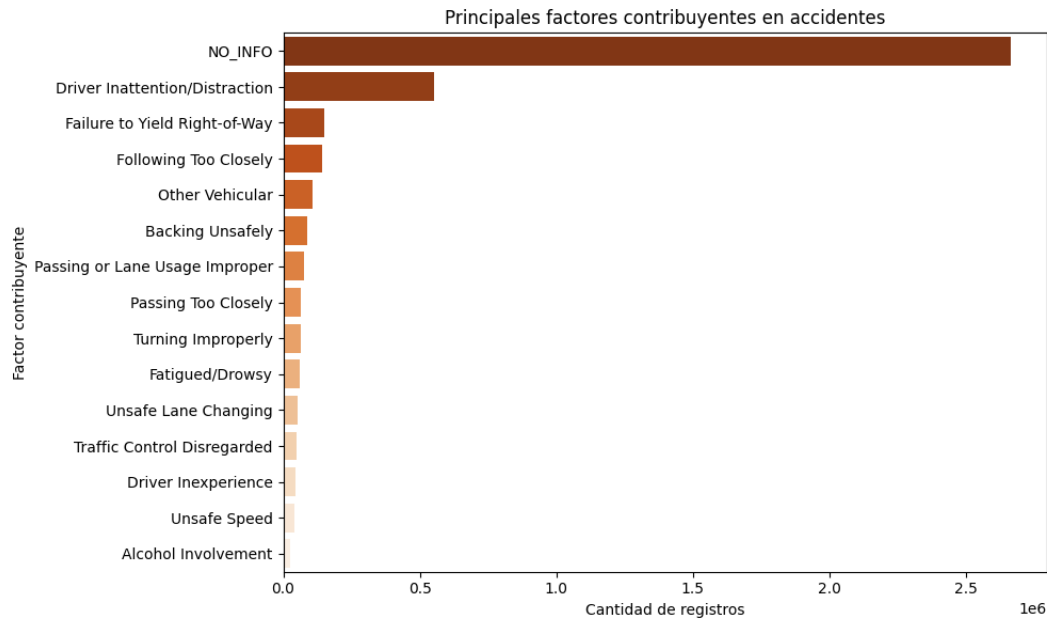


Figura 7: Factores contribuyentes más frecuentes en los accidentes.

El segundo gráfico ilustra los factores atribuidos como causa principal de los accidentes. Aunque la categoría NO_INFO concentra una proporción significativa de registros, entre los casos reportados sobresalen la distracción del conductor, el no ceder el paso y la conducción a corta distancia respecto a otros vehículos. Estas causas evidencian que la mayoría de los accidentes tienen un origen humano, asociado a fallas de atención, imprudencia o desconocimiento de las normas básicas de seguridad vial.

5.4.3. Puntos de impacto más comunes

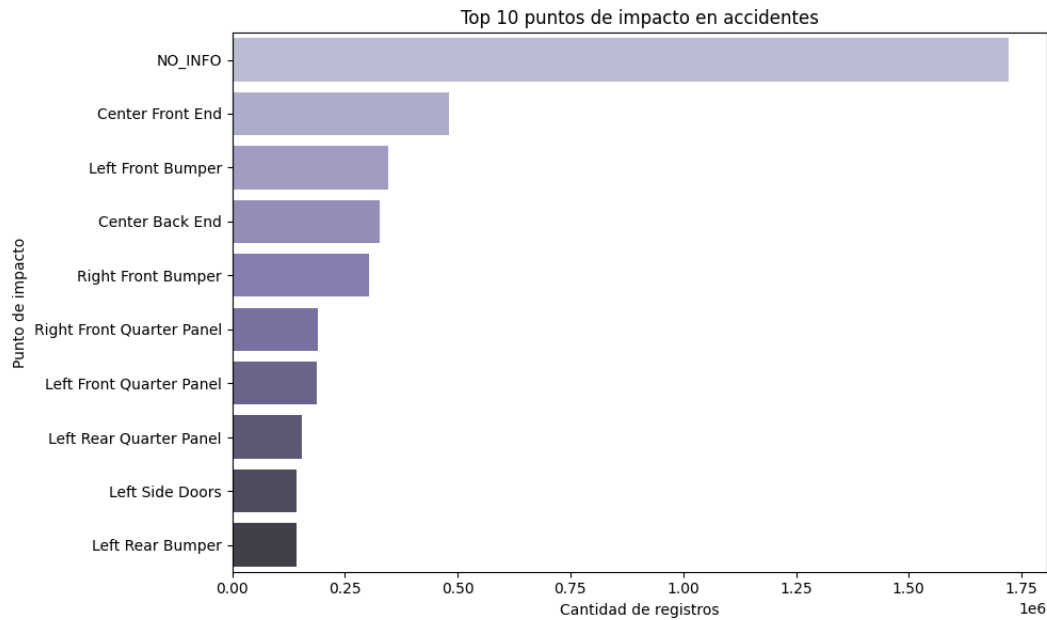


Figura 8: Distribución de puntos de impacto en los vehículos involucrados.

En este gráfico se observan las zonas del vehículo más afectadas por los choques. Los impactos frontales son los más frecuentes, seguidos por los golpes en las partes traseras y laterales. Este patrón coincide con la naturaleza de los siniestros más comunes, como colisiones por alcance o choques frontales en intersecciones. La abundancia de valores NO_INFO también pone en evidencia la dificultad de registrar de manera precisa ciertos detalles técnicos en los informes policiales.

5.4.4. Condición del vehículo antes del choque

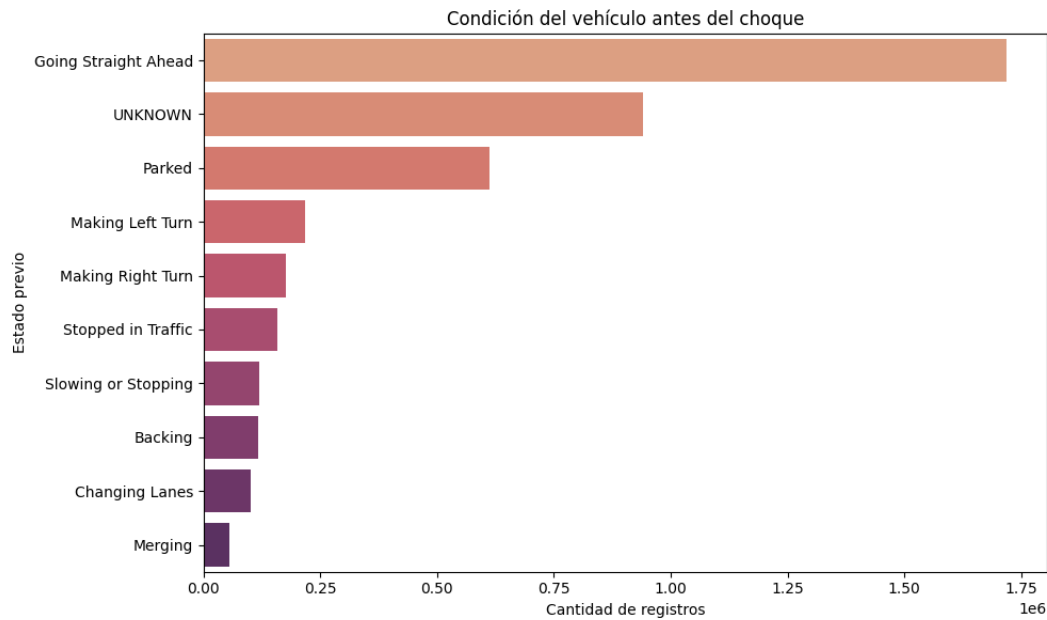


Figura 9: Condición del vehículo antes del choque.

La mayoría de los vehículos involucrados en accidentes estaban circulando en línea recta al momento del impacto, seguidos por aquellos que se encontraban estacionados o realizando maniobras de giro. Esto indica que muchos incidentes se producen en condiciones normales de desplazamiento, posiblemente por distracciones o fallas en la atención del conductor.

5.4.5. Evolución de colisiones por año

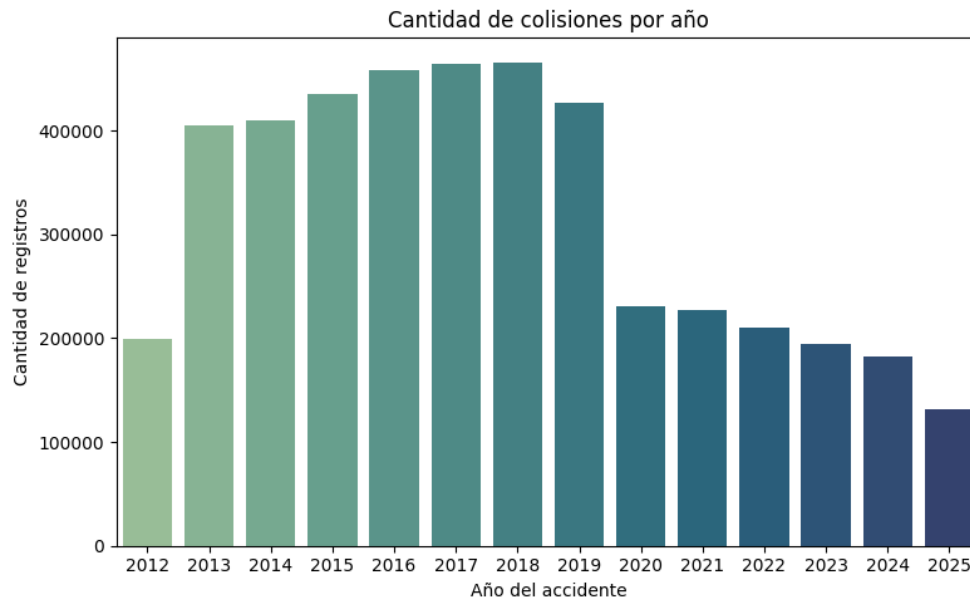


Figura 10: Evolución del número de colisiones registradas por año.

Finalmente, la última gráfica muestra la evolución anual del número de colisiones registradas. Se observa un incremento sostenido entre 2013 y 2019, seguido por una disminución marcada a partir de 2020. Este descenso podría estar asociado a la pandemia de COVID-19 y las restricciones de movilidad que redujeron la circulación de vehículos,. En conjunto, esta visualización proporciona una perspectiva temporal que ayuda a contextualizar las variaciones en la frecuencia de accidentes dentro del periodo analizado.

5.5. Poverty Data

A partir de esta base de datos, se construyeron diversas visualizaciones que permiten observar con mayor claridad la distribución de la pobreza en la ciudad, así como su relación con variables como el ingreso familiar, el distrito de residencia y el nivel educativo.

5.5.1. Cantidad de registros por distrito

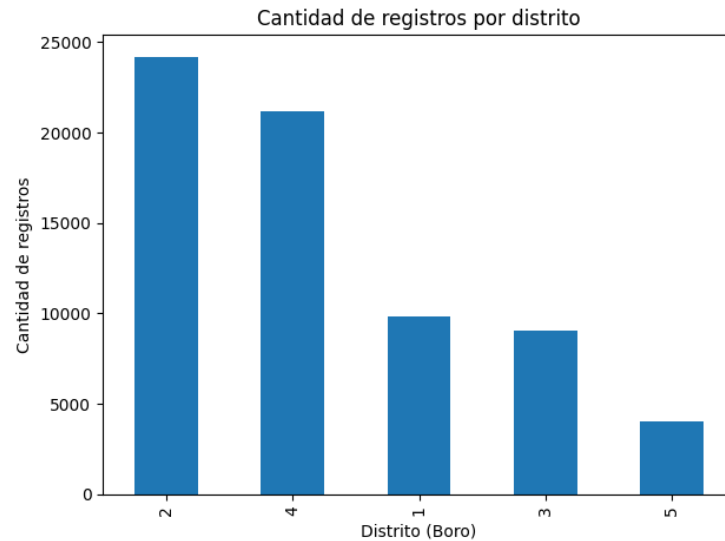


Figura 11: Cantidad de registros por distrito.

Este gráfico muestra el volumen de registros en cada uno de los cinco distritos (boroughs) de Nueva York. Aunque no representa una variable analítica directamente, es útil para tener contexto sobre la cobertura y densidad de información proveniente de cada zona. Los distritos 2 y 4 concentran el mayor número de registros, seguidos por el 1, 3 y 5.

5.5.2. Tasa de pobreza por distrito

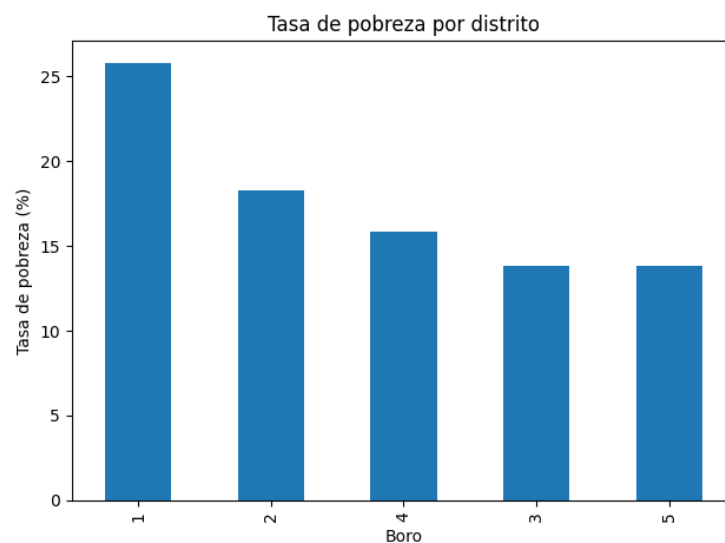


Figura 12: Tasa de pobreza por distrito.

Este gráfico presenta el porcentaje de personas clasificadas como en condición de pobreza, según su distrito de residencia. Se evidencia una desigualdad notable, siendo el distrito 1 el que presenta la mayor tasa de pobreza (por encima del 25 %), mientras que los distritos 3 y 5 muestran los niveles más bajos. Esta información es clave para identificar zonas con mayores necesidades económicas y orientar estrategias de intervención.

5.5.3. Distribución de ingresos familiares estimados

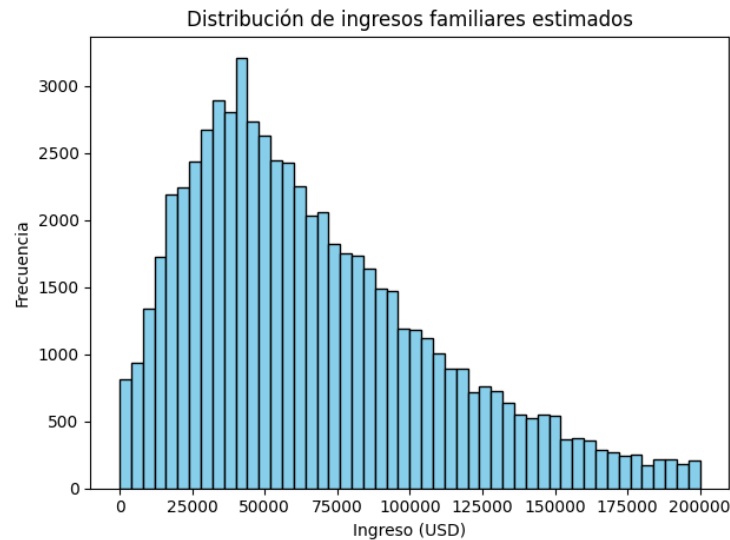


Figura 13: Distribución estimada de ingresos familiares.

El histograma de ingresos familiares muestra que la mayoría de los hogares se concentran en rangos bajos o medios, particularmente entre los \$20.000 y \$60.000 anuales. A medida que el ingreso aumenta, la frecuencia disminuye, reflejando una distribución desigual donde los ingresos elevados son mucho menos comunes.

5.5.4. Boxplot de ingreso familiar

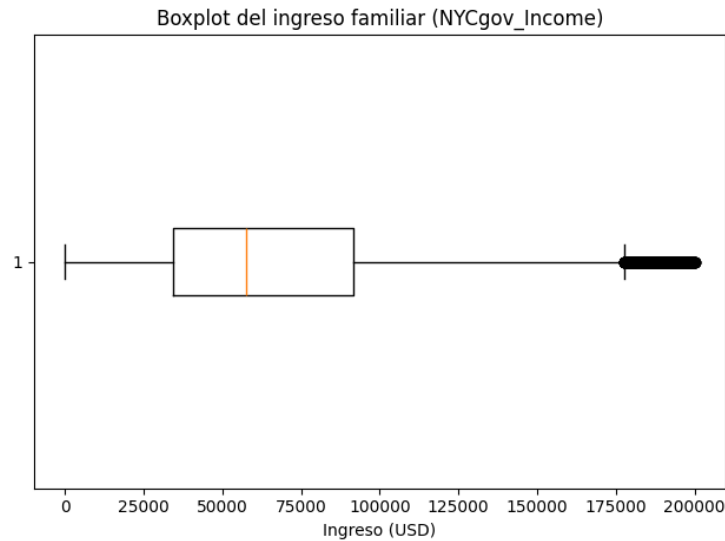


Figura 14: Boxplot del ingreso familiar estimado.

El boxplot refuerza lo observado anteriormente, mostrando una mediana cercana a los \$55.000 y una amplia dispersión de los datos. También se evidencian múltiples valores atípicos hacia el extremo superior, lo cual sugiere la presencia de hogares con ingresos significativamente más altos que el promedio general.

5.5.5. Ingreso por nivel educativo

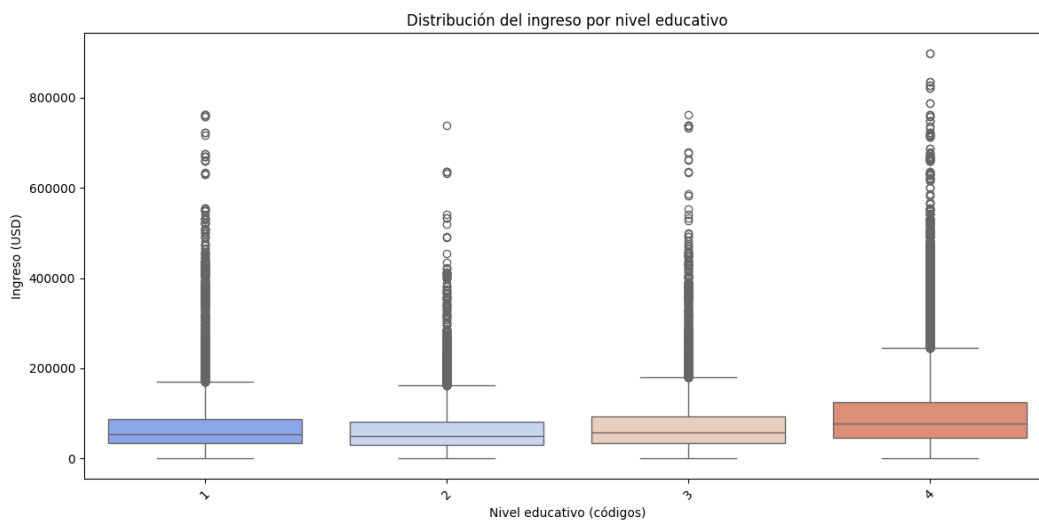


Figura 15: Distribución del ingreso familiar según nivel educativo.

Este gráfico compara el ingreso estimado entre los distintos niveles educativos registrados. A mayor nivel educativo, mayor es la mediana de ingreso y mayor es también el rango intercuartílico. Esto sugiere una clara relación entre el nivel educativo y el bienestar económico.

5.6. SAT NYC

En el caso del conjunto de datos de resultados del SAT para escuelas públicas de Nueva York, se elaboraron visualizaciones que permiten examinar con mayor detalle el rendimiento promedio de los estudiantes en las tres áreas evaluadas: lectura crítica, matemáticas y escritura. Estas gráficas complementan el análisis estadístico previo, facilitando la identificación de patrones y posibles correlaciones entre los puntajes.

Se incluyen histogramas para cada componente del examen, un boxplot comparativo y un gráfico de dispersión que permite observar la relación entre lectura y matemáticas.

5.6.1. Distribución de puntajes de lectura crítica

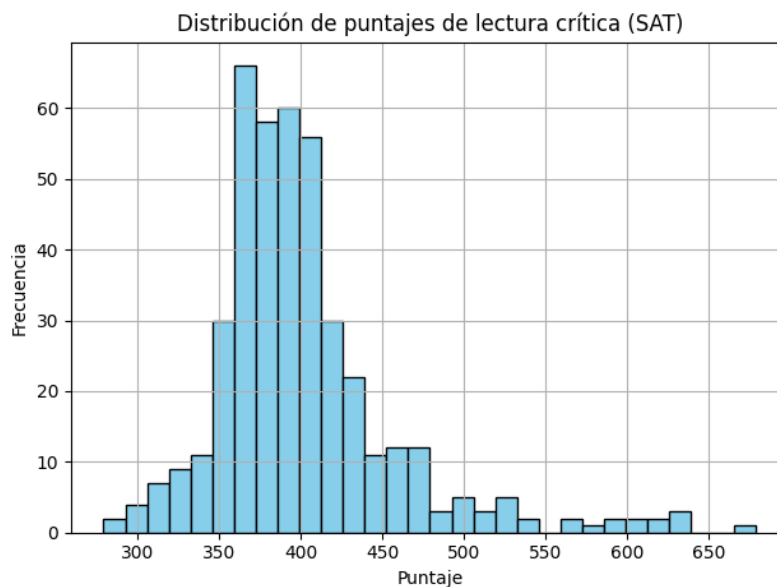


Figura 16: Distribución de puntajes de lectura crítica (SAT).

El histograma refleja que la mayoría de las escuelas presentan puntajes promedio entre 350 y 400 puntos en lectura crítica. Se trata de una distribución ligeramente sesgada hacia la derecha, con pocas instituciones alcanzando valores cercanos o superiores a 600 puntos.

5.6.2. Distribución de puntajes de matemáticas

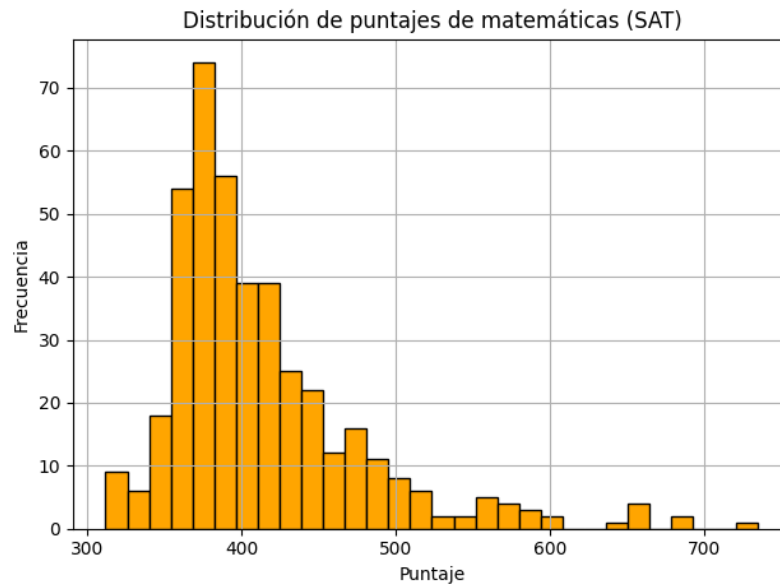


Figura 17: Distribución de puntajes de matemáticas (SAT).

La distribución en matemáticas muestra una mayor dispersión, con una concentración similar a la de lectura en el rango medio, pero con más casos que alcanzan puntajes elevados, superando incluso los 700 puntos. Esto indica un mejor desempeño en esta sección.

5.6.3. Distribución de puntajes de escritura

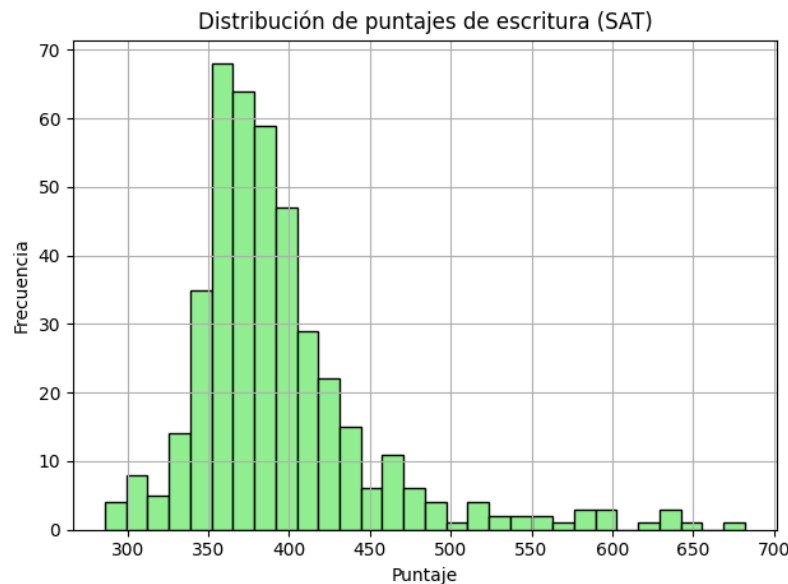


Figura 18: Distribución de puntajes de escritura (SAT).

Los resultados de escritura presentan una forma similar a la de lectura crítica, con una fuerte concentración entre los 350 y 400 puntos. La menor dispersión sugiere una menor variabilidad en el desempeño de las escuelas en esta área.

5.6.4. Comparación de puntajes por componente

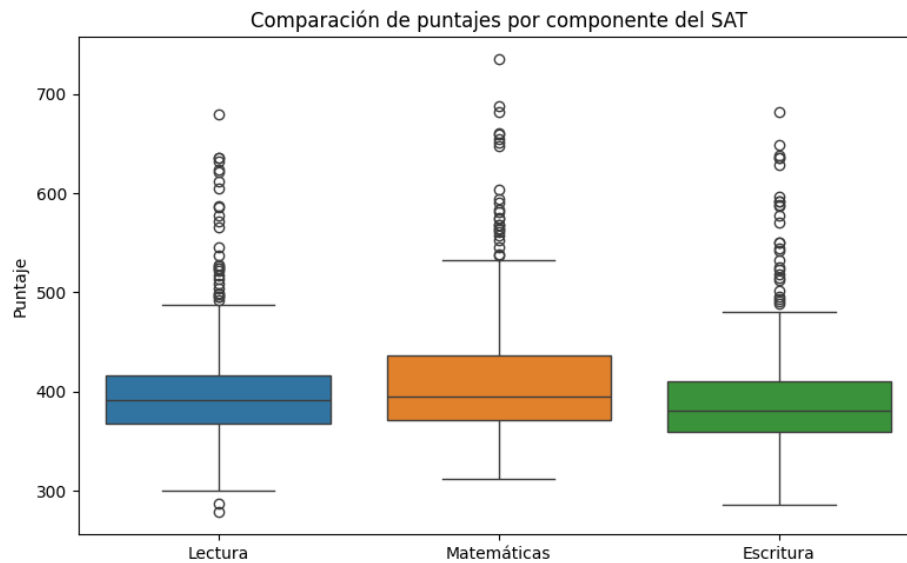


Figura 19: Comparación de puntajes por componente del SAT.

Este boxplot permite comparar visualmente la distribución de puntajes entre las tres secciones del examen. Se observa que matemáticas presenta una mediana ligeramente superior, así como una mayor presencia de valores atípicos.

5.6.5. Relación entre puntajes de lectura y matemáticas

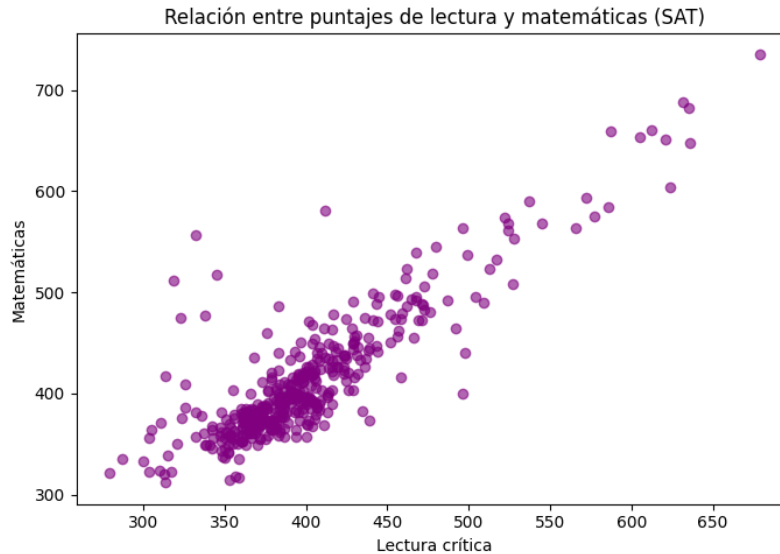


Figura 20: Relación entre puntajes de lectura crítica y matemáticas.

El gráfico de dispersión revela una correlación positiva entre los puntajes en lectura y matemáticas. Las escuelas con buen desempeño en una de estas áreas tienden también a obtener buenos resultados en la otra, lo que sugiere la posible existencia de factores comunes que afectan ambas dimensiones del rendimiento académico.

5.7. Hallazgos preliminares

Durante la fase de exploración de los datos, se identificaron una serie de patrones que pueden aportar al análisis territorial y social de la ciudad. Si bien aún no se ha realizado un cruce formal entre los conjuntos, algunos hallazgos que destacan son:

- **Concentración etaria en los arrestos:** La mayoría de los arrestos se concentran en personas entre los 25 y 44 años, con una proporción mucho mayor de hombres que de mujeres.
- **Alta presencia de registros no informados en accidentes:** En el conjunto de accidentes vehiculares, varias columnas presentan valores desconocidos o sin especificar. Aun así, se evidencian causas frecuentes como distracción del conductor o maniobras inseguras. También se detectó una alta proporción de impactos frontales y vehículos sedán como los más comúnmente involucrados.
- **Diferencias distritales en pobreza:** En el análisis del conjunto de pobreza, se observan diferencias importantes entre los distritos. Por ejemplo, Manhattan (distrito 1) presenta la tasa de pobreza más alta, mientras que Staten Island (distrito 5) tiene la más baja. Además, se aprecia que las personas en condición de pobreza tienen ingresos medios sustancialmente menores, como era de esperarse.

- **Relaciones entre desempeño académico:** El conjunto SAT muestra una correlación clara entre los puntajes en lectura crítica y matemáticas. Las escuelas que tienen mejores resultados en una sección tienden a tenerlos también en las otras. Esto puede reflejar tanto la calidad educativa como factores socioeconómicos más amplios.

Estos hallazgos servirán como punto de partida para análisis más complejos, especialmente en etapas posteriores donde se crucen variables entre conjuntos para detectar correlaciones o diferencias significativas entre zonas o poblaciones.

6. Reporte de calidad de datos

Como se mencionó anteriormente, antes de avanzar con el análisis exploratorio y la fase de transformación, se realizó una evaluación de la calidad de los datos disponibles. Esta revisión tuvo como propósito identificar problemas de integridad, consistencia, completitud y tipificación incorrecta que pudieran interferir con los procesos analíticos posteriores.

La validación de calidad incluyó tanto técnicas cuantitativas (conteo de nulos, duplicados, detección de valores no válidos) como observaciones cualitativas basadas en los tipos de datos de cada conjunto. En términos generales, el reporte se organizó en dos fases principales: análisis de valores faltantes y propuesta de tratamiento para las columnas afectadas.

6.1. Análisis de valores faltantes

Para cada uno de los conjuntos de datos seleccionados se realizó un diagnóstico inicial de calidad basado en el conteo de valores nulos, vacíos y duplicados. Esta revisión permitió identificar columnas potencialmente problemáticas, con el fin de definir estrategias de tratamiento en las etapas posteriores.

NYPD Arrest Data.

Para este conjunto sólo se detectaron valores nulos en dos columnas clave: LAW_CAT_CD (categoría legal del delito), con 687 registros nulos, lo que representa aproximadamente el 0,48 %, y KY_CD (código estandarizado del delito), con 10 registros nulos, equivalentes al 0,007. El resto de las columnas no contiene valores faltantes ni registros duplicados, por lo que se considera un conjunto de alta calidad para el análisis.

NYCgov Poverty Measure Data.

Este conjunto presenta un mayor grado de incompletitud. Se identificaron valores nulos en 61 columnas, principalmente relacionadas con variables laborales y de transporte. Por ejemplo, ENG (nivel de inglés en el hogar) presenta un 55,01 % de nulos, JWTR (medio de transporte al trabajo) un 51,20 %, y WKW (semanas trabajadas durante el año) un 45,98 %. Sin embargo, las variables ajustadas por ingreso, como SEMP_adj (ingreso por trabajo por cuenta propia), SSIP_adj (ingreso por seguridad suplementaria) y WAGP_adj (ingreso por salario), no contienen valores faltantes. Tampoco se detectaron registros duplicados. El tratamiento posterior deberá enfocarse en filtrar o imputar columnas con alta proporción de nulos, priorizando aquellas con mayor relevancia para el modelado.

Motor Vehicle Collisions – Vehicles.

Este conjunto, el cual contiene más de 4,4 millones de registros, presenta un nivel considerable de datos faltantes. Entre las columnas más afectadas se encuentran `PUBLIC_PROPERTY_DAMAGE_TYPE` (tipo de daño a propiedad pública), con un 99,3 % de valores faltantes, `VEHICLE_MODEL` (modelo del vehículo) con un 98,8 %, y `VEHICLE_DAMAGE_3` (zona secundaria de daño) con un 76,7 %. Otras columnas como `DRIVER_LICENSE_STATUS` (estado de la licencia del conductor) y `DRIVER_SEX` (sexo del conductor) presentan vacíos por encima del 50 %. A pesar de ello, columnas fundamentales como `COLLISION_ID` (identificador del accidente), `CRASH_DATE` (fecha del accidente) y `CRASH_TIME` (hora del accidente) no contienen nulos ni duplicados. Debido al tamaño del conjunto, será necesario aplicar filtros y reducir dimensionalidad para facilitar su análisis

SAT NYC.

Si bien para este conjunto no se encontraron valores nulos explícitos, se detectó que las columnas `Num_of_SAT_Test_Takers` (número de estudiantes que presentaron el examen), `SAT_Critical_Reading_Avg_Score` (puntaje promedio en lectura crítica), `SAT_Math_Avg_Score` (puntaje promedio en matemáticas) y `SAT_Writing_Avg_Score` (puntaje promedio en escritura) fueron interpretadas como texto. Esto indica la presencia de valores no numéricos, lo que será abordado más adelante como parte del análisis semántico.

En resumen, únicamente el conjunto NYPD Arrest Data presenta una estructura completa y lista para el análisis inmediato. Los conjuntos NYCgov Poverty Measure Data y Motor Vehicle Collisions – Vehicles requieren un tratamiento más profundo de los valores faltantes, dada la alta proporción de vacíos en múltiples columnas. Por su parte, el conjunto SAT NYC requiere un proceso de depuración específico para normalizar las columnas que fueron interpretadas como texto y reemplazar los valores no numéricos identificados.

6.2. Detección de valores no numéricos en columnas numéricas

Además de los valores nulos explícitos, se realizó una inspección de aquellas columnas que deberían contener únicamente datos numéricos. Este paso fue necesario debido a que algunos conjuntos presentaban columnas con tipo de dato texto (`string`) a pesar de representar variables cuantitativas. El objetivo fue identificar valores no numéricos o codificaciones erróneas que impidieran una conversión correcta.

En los conjuntos NYPD Arrest Data, NYCgov Poverty Measure Data y Motor Vehicle Collisions – Vehicles se evaluaron todas las columnas que representan variables numéricas, con el fin de identificar posibles valores no válidos o inconsistencias de tipo. En los tres casos se confirmó que las variables numéricas están correctamente tipadas y no presentan valores alfabéticos, símbolos extraños ni codificaciones inválidas. Si bien algunas variables incluyen valores negativos o ceros, estos son coherentes con el significado de cada atributo y no requieren tratamiento adicional en esta etapa.

SAT NYC. En este conjunto se identificaron valores no válidos del tipo texto en todas las columnas que deberían ser numéricas: `Num_of_SAT_Test_Takers` (número de estudiantes que presentaron el SAT), `SAT_Critical_Reading_Avg_Score` (puntaje promedio en lectura crítica),

SAT_Math_Avg_Score (puntaje promedio en matemáticas) y SAT_Writing_Avg_Score (puntaje promedio en escritura). En cada una de estas columnas se detectó la presencia del carácter 's' como valor no numérico. Esto impide la conversión automática a tipo entero y representa un caso de falsos nulos, por lo que requerirá un proceso de limpieza que reemplace dichos valores por `null` antes de realizar cualquier análisis cuantitativo.

Con esto, se confirma que solo uno de los cuatro conjuntos de datos (SAT NYC), requiere una corrección directa de tipo para lograr su análisis numérico. Los demás conjuntos presentan una estructura numérica válida, sin presencia de datos atípicos o no numéricos en sus variables clave.

6.3. Propuesta de tratamiento

Con base en lo encontrado durante el análisis de calidad de datos, se plantean las siguientes estrategias de tratamiento para preparar los conjuntos antes de su exploración y modelado.

NYPD Arrest Data. Dado que este conjunto presenta únicamente una proporción menor de valores nulos en las columnas LAW_CAT_CD (categoría legal del delito) y KY_CD (código estandarizado del delito), se propone realizar una eliminación directa de los registros incompletos en estas variables, ya que representan menos del 1 % del total. No se requieren imputaciones ni transformaciones adicionales.

NYCgov Poverty Measure Data. Este conjunto presenta múltiples columnas con valores faltantes, especialmente en variables relacionadas con condiciones laborales y transporte. Se propone:

- Eliminar columnas con más del 50 % de nulos.
- Aplicar imputación por moda o mediana en variables numéricas o por la categoría UNKNOWN en variables categóricas.
- Mantener las variables económicas que no presentan nulos.

Motor Vehicle Collisions – Vehicles. Debido al alto volumen de datos y al gran número de columnas con vacíos, se plantea:

- Filtrar columnas con más del 50 % de valores nulos.
- Reemplazar valores nulos en campos categóricos con la etiqueta UNKNOWN..
- Mantener las columnas completas.

SAT NYC Este conjunto requiere una transformación específica en sus variables cuantitativas, que actualmente están representadas como cadenas de texto. Se propone el siguiente tratamiento:

- Reemplazar los valores s por null en las columnas de puntajes y número de estudiantes, ya que impiden el análisis numérico. - Convertir dichas columnas a tipo numérico, con el propósito de poder realizar operaciones estadísticas. - Una vez realizado el casteo, se repetirá el proceso de conteo de valores nulos y duplicados con el nuevo esquema de datos.

7. Planteamiento de preguntas sobre los datos

En esta sección se formulan una serie de preguntas clave orientadas al análisis territorial y social de la ciudad de Nueva York, a partir de los conjuntos de datos previamente explorados. Estas preguntas no solo reflejan patrones identificados en arrestos, accidentes, condiciones socioeconómicas y desempeño académico, sino que buscan generar hallazgos que puedan traducirse en acciones concretas por parte del equipo de gobierno. Todas las preguntas fueron diseñadas con un enfoque en el impacto social, la desigualdad territorial y la formulación de estrategias de prevención y mejora de calidad de vida.

7.1. Preguntas principales

1. ¿Qué tipos de delitos son más frecuentes en zonas de bajos ingresos, y qué diferencias hay respecto a zonas de mayor ingreso?
2. ¿Hay relación entre los puntajes promedio del SAT por distrito y la tasa de arrestos?
3. ¿El rendimiento académico está asociado a la frecuencia de accidentes viales en las zonas donde se ubican las escuelas?
4. ¿En qué horarios o temporadas se concentran más accidentes o arrestos?
5. ¿Qué relación existe entre el nivel educativo promedio y la condición de pobreza en los distintos distritos?
6. ¿Las zonas con mayores niveles de pobreza presentan también mayores tasas de arrestos? ¿Cómo varía esta relación entre distritos?
7. ¿Qué tipos de vehículos están más involucrados en colisiones con daño a propiedad pública?
8. ¿Las zonas con mayor cantidad de arrestos por delitos violentos presentan también una mayor frecuencia de accidentes viales?

7.2. Justificación de su relevancia

Cada una de las preguntas anteriores tiene el potencial de generar valor estratégico para el equipo de gobierno, al permitir una mejor comprensión del territorio y la orientación de políticas públicas más efectivas. A continuación, se presenta la justificación de cada una:

- **Pregunta 1:** Permite identificar patrones delictivos asociados a condiciones económicas, lo cual es fundamental para diseñar estrategias de seguridad diferenciadas según el contexto socioeconómico.
- **Pregunta 2:** Analizar la relación entre desempeño académico y arrestos podría revelar dinámicas de exclusión educativa y social, y abrir el camino a intervenciones escolares con enfoque preventivo.

- **Pregunta 3:** Permite explorar si el entorno escolar influye en la seguridad vial, lo que puede justificar programas educativos o infraestructura segura en zonas críticas.
- **Pregunta 4:** Identificar horarios o temporadas con mayor concentración de incidentes ayuda a planificar mejor la asignación de recursos policiales y de tránsito.
- **Pregunta 5:** Explorar la relación entre educación y pobreza ofrece insumos para políticas públicas de inclusión educativa y reducción de desigualdad.
- **Pregunta 6:** Busca comprender cómo el nivel educativo se relaciona con las condiciones de pobreza en los distintos distritos, lo que puede ayudar a identificar brechas de oportunidad y orientar programas educativos o sociales hacia las comunidades más vulnerables.
- **Pregunta 7:** Facilita la identificación de los tipos de vehículos más problemáticos desde una perspectiva de seguridad urbana y daño a infraestructura pública, lo que puede derivar en regulaciones más estrictas o campañas de prevención.
- **Pregunta 8:** Relacionar violencia y accidentalidad permite identificar zonas de alta vulnerabilidad urbana donde se concentran múltiples riesgos.

8. Transformaciones, filtrado y limpieza inicial

Una vez finalizada la etapa de evaluación de calidad de los datos, se procedió con el desarrollo de un conjunto de transformaciones, filtros y limpiezas iniciales orientadas a garantizar la consistencia, completitud y utilidad de las variables involucradas. Esta fase tuvo como propósito preparar cada conjunto de datos para su análisis posterior, eliminando inconsistencias, depurando registros incompletos y normalizando formatos. Las acciones realizadas se guiaron por la propuesta de tratamiento definida previamente y se según las características propias de cada conjunto. Se buscó intervenir lo mínimo necesario para conservar la mayor cantidad posible de información relevante.

8.1. Transformaciones preliminares

Durante esta etapa se realizaron transformaciones básicas sobre los conjuntos de datos, con el objetivo de dejar las variables clave listas para su análisis posterior. Estas transformaciones incluyeron cambios de tipo, creación de nuevas columnas y ajustes en el formato de ciertos valores que venían mal representados.

En el caso del conjunto SAT NYC, algunas columnas que deberían contener números estaban registradas como texto debido a la presencia del carácter **s**, que indicaba falta de reporte. Se reemplazaron estos valores por nulos, se convirtieron las columnas a tipo numérico, y se creó una nueva columna que suma los puntajes por área. Luego de esta transformación, se eliminaron 57 registros que no contenían ningún dato útil, quedando un total de 421 observaciones válidas.

En el conjunto de pobreza se eliminaron columnas con más del 50 % de nulos y se imputaron los valores faltantes en las variables restantes. Para las columnas numéricas se usó la

mediana, y para las categóricas una etiqueta genérica unknown. El resultado fue una tabla sin datos faltantes y con 59 columnas útiles.

El conjunto de colisiones vehiculares fue tratado de forma similar. Se eliminaron las columnas más incompletas y se imputaron los valores faltantes en las variables restantes. Además, se unificaron los distintos valores que representaban ausencia de información (como “Unspecified”, “N/A”, “Unknown”, entre otros) bajo una sola categoría estándar: unknown. En total se conservaron 17 columnas, y no fue necesario eliminar filas, ya que la limpieza se centró en estandarizar y completar los valores existentes.

Por último, en el conjunto de arrestos se encontraron nulos en dos columnas clave que definen el tipo de delito. Como estas variables son esenciales y los valores faltantes eran pocos (menos del 0.5 %), se decidió eliminar directamente esas filas. El conjunto final quedó con 142.797 registros completos.

8.2. Filtrados aplicados

En esta etapa se aplicaron filtros sobre los conjuntos, con el objetivo de enfocar el análisis en registros relevantes y reducir información poco útil o poco representativa, sin comprometer la cobertura general de los datos. Se mantuvo el principio de intervenir lo menos posible, filtrando únicamente cuando era necesario y justificado.

En el caso del conjunto de arrestos, todos los registros correspondían al año 2025, por lo que no fue necesario filtrar por fechas. Sin embargo, se consideró importante restringir el análisis a delitos de mayor gravedad. Por esta razón, se conservaron únicamente los registros cuya categoría legal (LAW_CAT_CD) correspondía a felonías (F) o delitos menores (M), excluyendo las violaciones menores (V). Este filtro redujo el tamaño de la base de 142.797 a 140.086 registros, lo cual representa una pérdida menor al 2 % del total, pero permite centrar el análisis en eventos de mayor impacto social y territorial.

Para los demás conjuntos de datos —pobreza, colisiones vehiculares y SAT NYC— no se aplicaron filtros adicionales en esta etapa, ya que presentaban una cobertura amplia, datos recientes o bien definidos, y variedad suficiente. Se consideró que aplicar filtros adicionales podría reducir innecesariamente el volumen de información disponible sin aportar un valor claro al análisis.

8.3. Limpiezas realizadas

En términos generales, los conjuntos de datos fueron sometidos a un proceso de limpieza enfocada en mejorar su calidad y garantizar su utilidad analítica. Este proceso incluyó la eliminación de columnas con altos porcentajes de valores nulos, algunos filtros específicos orientados a mejorar el enfoque del análisis, y transformaciones necesarias para corregir formatos, normalizar tipos de datos y crear variables clave. Además, en ciertos conjuntos se realizó una estandarización de valores categóricos, unificando distintas expresiones de datos no informados bajo una misma categoría. Se procuró intervenir lo menos posible, aplicando solo las acciones indispensables para asegurar la consistencia, completitud y coherencia de la información disponible.

9. Web scraping de datos poblacionales

Para cumplir con el requerimiento de realizar un proceso de web scraping sobre la población de Nueva York, se intentó inicialmente acceder al enlace oficial proporcionado por el enunciado del proyecto:

`https://www.health.ny.gov/statistics/cancer/registry/appendix/neighborhoodpop.htm`

Sin embargo, este enlace actualmente se encuentra fuera de servicio (error 404), por lo que se optó por otra alternativa: el uso de la API oficial del U.S. Census Bureau, específicamente la del *American Community Survey (ACS) 5-Year 2020*.

Mediante esta API, se consultó la variable B01003.001E, correspondiente al total de población por condado. En el caso de la ciudad de Nueva York, los condados consultados corresponden a sus cinco distritos: Bronx, Brooklyn (Kings), Manhattan (New York), Queens y Staten Island (Richmond). Se usaron sus respectivos códigos FIPS para realizar la consulta, y posteriormente se construyó un *DataFrame* con la información.

El resultado de la consulta fue el siguiente:

- **Brooklyn (Kings):** 2,576,771 habitantes
- **Queens:** 2,270,976 habitantes
- **Manhattan (New York):** 1,629,153 habitantes
- **Bronx:** 1,427,056 habitantes
- **Staten Island (Richmond):** 475,596 habitantes

A partir de estos datos se generó la siguiente visualización:

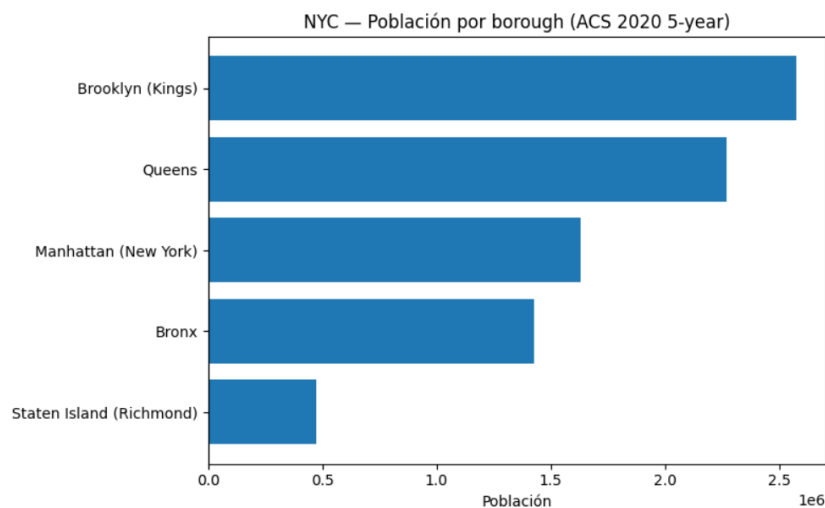


Figura 21: Distribución poblacional por distrito en la ciudad de Nueva York (ACS 2020 5-Year)

Este gráfico permite visualizar cómo se distribuye la población en los diferentes distritos de la ciudad. Brooklyn y Queens destacan como las zonas más densamente pobladas, mientras que Staten Island presenta una concentración significativamente menor.

Esta información será fundamental en etapas posteriores del análisis, ya que permitirá normalizar indicadores como arrestos o accidentes por cada 100,000 habitantes y hacer comparaciones más equitativas entre zonas con diferentes tamaños poblacionales.

10. Consulta climática con OpenWeatherMap

Como parte del segundo bono opcional de la entrega, se realizó una consulta a la API de OpenWeatherMap utilizando el endpoint de pronóstico a 5 días con cortes cada 3 horas. La ciudad consultada fue *New York, US*, y se extrajeron variables como temperatura, humedad, precipitación (lluvia y nieve), velocidad del viento y el estado general del clima.

La respuesta fue exitosa (status 200), y los datos fueron procesados para construir un *DataFrame* ordenado cronológicamente. En las primeras observaciones ya se evidencian patrones relevantes: cielos despejados durante el 21 de octubre, seguidos de episodios de lluvia ligera durante la madrugada del día 22, con temperaturas que oscilan entre los 19°C y los 11°C.

Se generaron dos visualizaciones a partir de esta información. La primera muestra la evolución de la temperatura durante los próximos días junto con las precipitaciones registradas en cada intervalo de 3 horas. Esta visualización permite identificar los ciclos térmicos diarios, así como los momentos de mayor precipitación. Por ejemplo, el 22 de octubre se concentran los valores más altos de lluvia, con varios picos cercanos a los 2 mm por bloque de 3 horas.

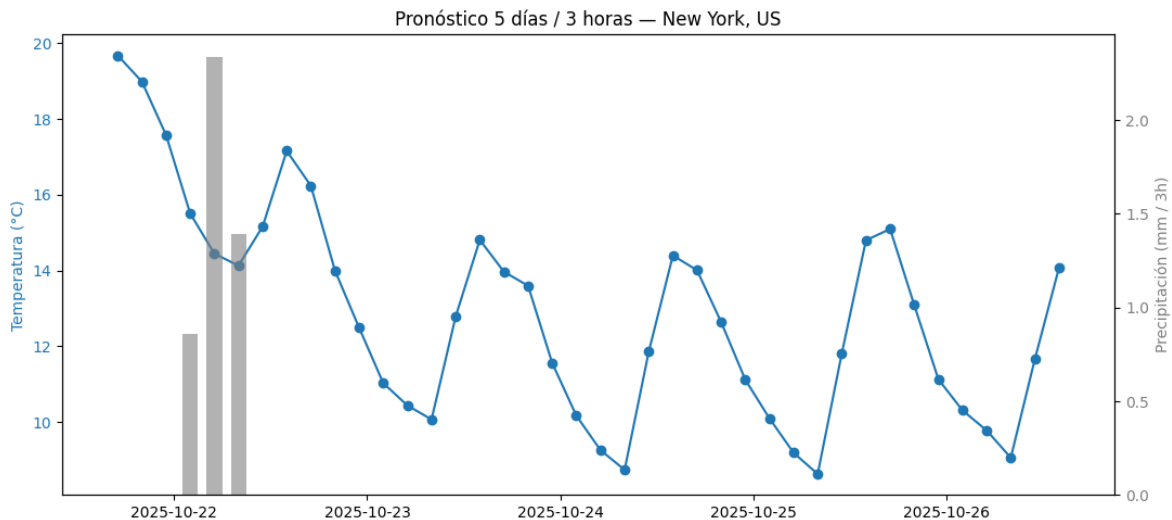


Figura 22: Pronóstico de temperatura y precipitación por intervalos de 3 horas — New York, US

La segunda visualización agrega la información a nivel diario, mostrando la temperatura máxima y mínima por jornada, así como la lluvia total diaria acumulada. En esta figura

se observa un descenso progresivo de la temperatura a lo largo del periodo, comenzando cerca a los 20°C y descendiendo hasta valores cercanos a 14°C. La única jornada con lluvia significativa fue el 22 de octubre, con una acumulación superior a los 4.5 mm, mientras que el resto de los días se mantuvieron secos según el pronóstico.

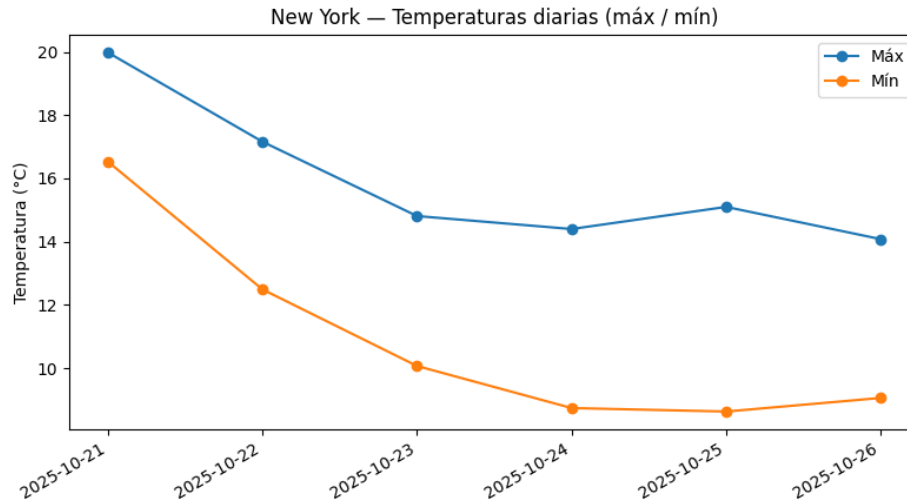


Figura 23: Temperaturas máximas y mínimas diarias — New York, US

Estas visualizaciones pueden servir más adelante para analizar posibles relaciones entre condiciones climáticas y fenómenos como la movilidad o la accidentalidad vial en la ciudad.

11. Conclusiones y recomendaciones

Durante esta primera entrega se logró consolidar el entendimiento del negocio, la recopilación y descripción técnica de los conjuntos de datos, así como su exploración inicial mediante análisis estadísticos y visualizaciones. Estos avances permitieron obtener una comprensión integral del contexto urbano y social de la ciudad de Nueva York, identificando los principales retos asociados a la seguridad, la movilidad y las condiciones socioeconómicas de su población.

El procesamiento distribuido en el clúster de Apache Spark permitió validar la infraestructura propuesta y garantizar la capacidad para manejar grandes volúmenes de información de manera eficiente. A su vez, la revisión de calidad de datos permitió detectar problemas comunes como valores nulos, registros no numéricos y categorías inconsistentes, los cuales fueron tratados mediante procesos de limpieza y transformación que dejaron las bases listas para su análisis posterior.

En términos de resultados, se evidenciaron patrones relevantes: la concentración de arres-tos en hombres adultos jóvenes, la mayor siniestralidad vial asociada a la distracción del conductor, la persistente desigualdad económica entre distritos y la correlación positiva entre los puntajes académicos del SAT. Estos hallazgos preliminares sientan las bases para los análisis cruzados y modelamientos que se desarrollarán en la segunda entrega.

Referencias

- Datos oficiales: <https://data.cityofnewyork.us>
- CRISP-DM: https://www.ibm.com/docs/en/spss-modeler/SaaS?topic=SS3RA7_Cloud/com.ibm.spss.modeler.help/idh_crispdm_main.htm
- OpenWeatherMap API: <https://openweathermap.org/api>
- Wikipedia - Nueva York: https://es.wikipedia.org/wiki/Nueva_York
- Blog SparklyMaid NYC: <https://www.sparklymaidnyc.com/blog/que-lugar-ocupa-nueva-york->
- Información general: <https://www.nuevayork.net/informacion-general>
- Seguridad vial en NY: <https://www.semana.com/mundo/noticias-estados-unidos/articulo/nueva-york-lidera-la-seguridad-vial-en-estados-unidos-logro-convertirse-en-202517/>
- Crisis de personas sin techo: <https://www.france24.com/es/programas/enlace/20240906-desolacion-y-desamparo-la-crisis-de-los-sintecho-en-nueva-york>
- Desigualdad en NY: https://www.bbc.com/mundo/noticias/2016/04/160418_nueva_york_ricos_pobres_primarias_ps
- Riesgos sociales en NY: <https://nychazardmitigation.com/es/documentation/nyc-hazard-environmental-social/>
- Tasa de desempleo (EE.UU.): <https://es.tradingeconomics.com/united-states/unemployment-rate>
- Estadísticas de NY: <https://datacommons.org/place/geoId/36?hl=es>
- Estimaciones de población NYC 2025: https://www.nyc.gov/assets/planning/downloads/pdf/our-work/reports/new-york-city-population-estimates-and-trends_may-2025.pdf