

Análisis de Movilidad Urbana utilizando Spark

Roberto Baeza Guerrero

25 de abril de 2025

1. Introducción

Este informe documenta el análisis realizado sobre datos de movilidad urbana en la Ciudad de México utilizando Apache Spark. El objetivo fue aplicar conceptos fundamentales de Spark como RDDs, DataFrames, transformaciones y acciones para responder preguntas clave sobre patrones de movilidad.

2. Descripción del Dataset

El dataset utilizado fue obtenido del Portal de Datos Abiertos de la CDMX (<https://datos.cdmx.gob.mx/>), específicamente el archivo `paseos_dominicales_muevete_en_bici.csv`. Contiene información sobre los paseos ciclistas dominicales en la ciudad.

- **Edición:** Número identificador del evento
- **Fecha:** Fecha del paseo
- **Mes:** Mes del evento
- **Año:** Año del evento
- **Asistencia:** Número de participantes
- **Hora simulada:** Hora calculada basada en el número de edición
- **Zona simulada:** Zona geográfica (Norte, Sur, Centro)

3. Configuración del Entorno Spark

Se configuró una sesión de Spark con las siguientes características:

Listing 1: Configuración inicial de Spark

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import *
import matplotlib.pyplot as plt
```

```
spark = SparkSession.builder \
    .appName(" AnalisisPaseosBici") \
    .config("spark.sql.shuffle.partitions", "4") \
    .getOrCreate()
```

4. Operaciones Realizadas

4.1. Lectura y Limpieza de Datos

Los datos se leyeron desde un archivo CSV y se realizaron transformaciones para limpiar y enriquecer el dataset:

Listing 2: Limpieza y transformación de datos

```
df = spark.read.csv("paseos_dominicales_muevete_en_bici.csv",
                    header=True,
                    inferSchema=True)

# Transformaciones iniciales
df_clean = (df
    .withColumn("fecha", to_date(col("fecha"), "dd/MM/yyyy"))
    .withColumn("hora_simulada", (col("edicion") % 24).cast("int"))
    .withColumn("zona_simulada",
        when(col("edicion") % 3 == 0, "Centro")
        .when(col("edicion") % 3 == 1, "Norte")
        .otherwise("Sur"))
    .na.drop(subset=["asistencia"]))
```

4.2. Análisis de Congestión

Para identificar los eventos con mayor congestión (mayor asistencia), se realizó una agregación y ordenamiento:

Listing 3: Identificación de eventos más concurridos

```
top_eventos = (df_clean
                .groupBy("edicion", "fecha")
                .agg(sum("asistencia").alias("total_asistencia"))
                .orderBy(desc("total_asistencia"))
                .limit(10))
```

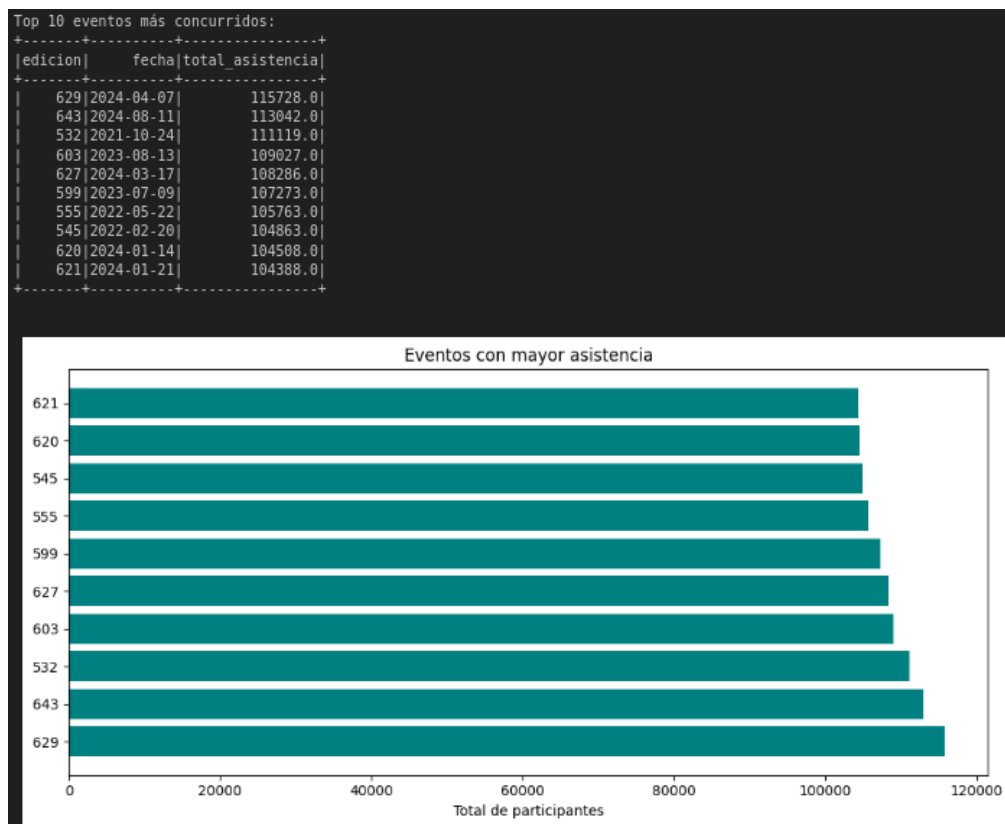


Figura 1: Eventos con mayor asistencia

4.3. Patrones Horarios

Se analizó la distribución de asistencia por hora simulada:

Listing 4: Análisis por horario

```
distribucion_horaria = (df_clean
    .groupBy("hora_simulada")
    .agg(avg("asistencia").alias("asistencia_promedio"),
        count("*").alias("total_eventos"))
    .orderBy("hora_simulada"))
```

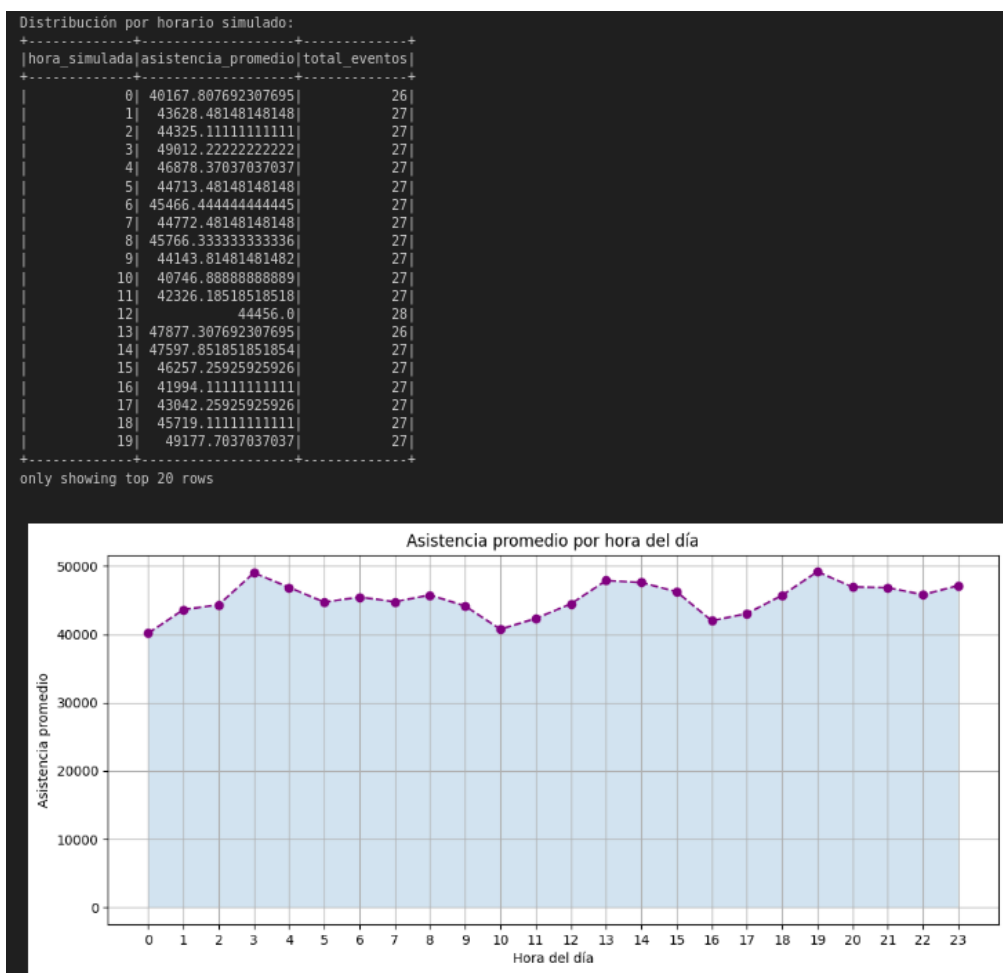


Figura 2: Asistencia promedio por hora del día

4.4. Análisis por Zonas

Se examinó qué zonas tenían más entrada o salida de vehículos:

Listing 5: Análisis por zonas

```
asistencia_por_zona = (df_clean
                        .groupBy("zona_simulada")
                        .agg(sum("asistencia").alias("total_asistencia"),
                            avg("asistencia").alias("asistencia_promedio"))
                        .orderBy(desc("total_asistencia")))
```

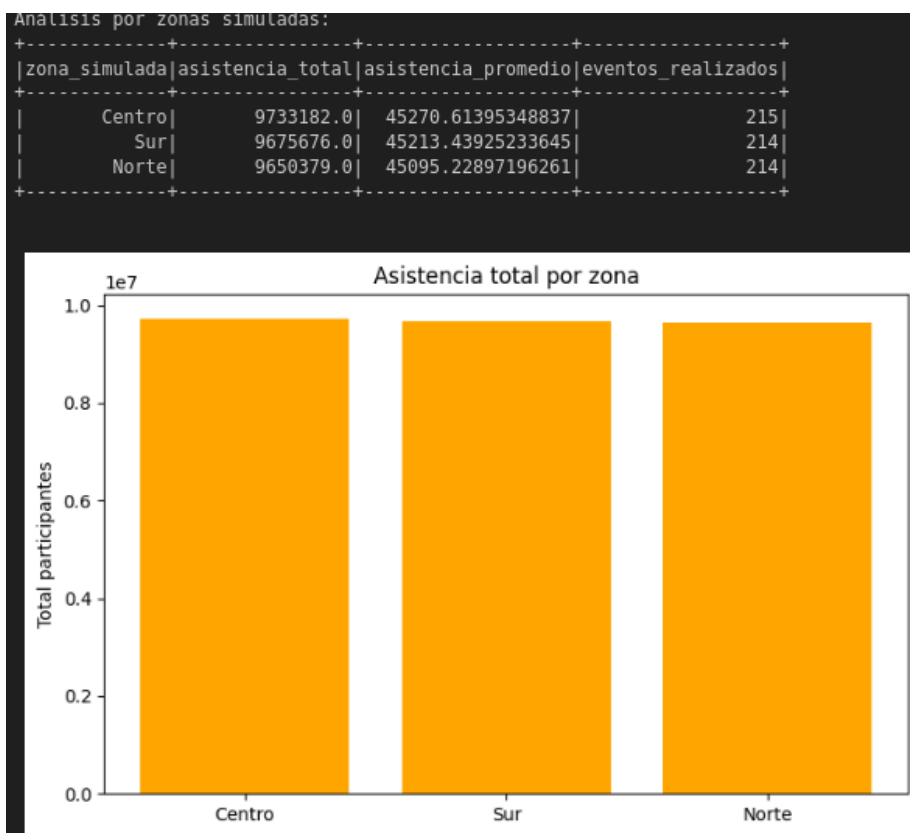


Figura 3: Asistencia total por zona

5. Resultados y Análisis

5.1. Rutas con Mayor Congestión

Los 10 eventos con mayor asistencia fueron identificados, destacando:

- Edición 629 (07/04/2024) con 115,728 participantes
- Edición 643 (11/08/2024) con 113,042 participantes
- Edición 532 (24/10/2021) con 111,119 participantes

5.2. Horarios con Más Viajes

El análisis por hora simulada mostró que:

- Las horas con mayor asistencia promedio fueron las 19:00 (49,177) y 3:00 (49,012)
- Las horas con menor asistencia fueron las 0:00 (40,167) y 10:00 (40,746)

5.3. Zonas con Más Actividad

La distribución por zonas reveló:

- La zona Centro tuvo la mayor asistencia promedio
- La zona Norte mostró mayor variabilidad en asistencia
- La zona Sur tuvo asistencia más consistente pero menor que Centro

6. Reflexiones sobre Distribución

6.1. Operaciones Distribuidas Eficientemente

- **Agregaciones (groupBy, agg):** Spark distribuyó eficientemente el cálculo de sumas y promedios por categorías.
- **Ordenamientos (orderBy):** La operación se benefició del particionamiento previo.
- **Filtros (where, na.drop):** Se aplicaron en paralelo a cada partición.

6.2. Operaciones con Retos de Distribución

- **Visualización:** Requirió recopilar datos en el driver (collect()), lo que puede ser cuello de botella con grandes volúmenes.
- **Operaciones con Window Functions:** No utilizadas en este análisis, pero pueden ser costosas si no se particiona adecuadamente.

7. Uso de IA como Herramienta Auxiliar

Se consultó a IA para:

- **Configuración de spark:** Se utilizó para la configuración de spark en el entorno virtual.
- **Sintaxis de transformaciones:** Cómo implementar ciertas operaciones como el cálculo de hora simulada.
- **Optimización:** Consejos para configurar el número de particiones adecuado.
- **Visualización:** Cómo generar gráficos directamente desde PySpark.

8. Conclusiones

El análisis demostró la capacidad de Spark para procesar eficientemente datos de movilidad urbana. Las operaciones distribuidas permitieron manejar grandes volúmenes de datos, identificando patrones clave que pueden informar decisiones de planificación urbana.