



Caracterización de datos de trayectorias individuales

Presentado por:
Jorge Rafael Martínez Buenrostro

Asesora: Dra. Elizabeth Pérez Cortés

México, CDMX, a 31 de julio de 2025

Resumen

Una trayectoria se define como la secuencia de desplazamientos realizada por un individuo en movimiento, compuesta por tres elementos principales: **puntos de recorrido**, **tiempos de pausa** y **longitudes de vuelo**. Los puntos de recorrido corresponden a las ubicaciones o pasos específicos por los que transita el individuo. Los tiempos de pausa representan los intervalos durante los cuales el individuo permanece detenido en un mismo punto de recorrido. Por último, la longitud de vuelo se refiere a la distancia recorrida entre dos puntos de recorrido consecutivos. En el presente trabajo se describe el proceso de caracterización de datos de movilidad, es decir, el proceso de limpieza y depuración de la información, mediante el cual elimina aquellos campos y registros que no le aportan valor a la trayectoria individual. El objetivo es identificar la mayor cantidad de trayectorias individuales, para así poder crear un modelo que permita simular el movimiento de individuos.

Contenido

Lista de Códigos	IV
1. Introducción del Proyecto	1
1.1. Descripción general del proyecto	1
1.2. Objetivos y propósitos	1
1.3. Alcance del sistema	1
2. Requisitos del sistema	2
2.1. Requisitos	2
2.1.1. Instrucciones de instalación	2
3. Caracterización de datos de trayectorias individuales	3
3.1. Caracterización de datos de trayectorias individuales	3
3.1.1. Exploración inicial del conjunto de datos	4
3.1.2. Dimensiones del conjunto de datos	5
3.1.3. Depuración de columnas	5
3.1.4. Depuración de filas	6
A. Uso de Docker y Docker Compose	9
A.1. ¿Qué son Docker y Docker Compose?	9
A.2. Instalación en Linux (Ubuntu/Debian)	9
A.3. Instalación en Windows	10
A.4. Descripción del archivo <code>docker-compose.yml</code>	11
A.5. Scripts de control del contenedor	13
A.5.1. <code>start_container.sh</code>	13
A.5.2. <code>restart_container.sh</code>	13
A.5.3. <code>stop_container.sh</code>	14
A.6. Proceso de uso y desarrollo del contenedor	14
B. Scripts para la caracterización de datos	16

Lista de Figuras

3.1.	Frecuencia de aparición de los valores de 'device_horizontal_accuracy'.	7
3.2.	Frecuencia de aparición de los identificadores únicos.	7
3.3.	Comparación de histogramas por rangos de repeticiones.	8

Lista de Códigos

2.1. Iniciar contenedor del proyecto.	2
2.2. Iniciar contenedor del proyecto.	2
2.3. Reiniciar contenedor del proyecto.	2
A.1. Actualizar el sistema.	9
A.2. Instalar Docker.	10
A.3. Verificar instalación de Docker.	10
A.4. Instalar Docker Compose.	10
A.5. Verificar instalación de Docker Compose.	10
A.6. Verificar instalación de Docker y Docker Compose.	10
A.7. Archivo docker-compose.yml	11
A.8. Script para iniciar el contenedor.	13
A.9. Script para reiniciar el contenedor.	13
A.10. Script para detener y eliminar el contenedor y sus volúmenes.	14
A.11. Dar permisos de ejecución a los scripts.	14
A.12. Iniciar contenedor y ejecutar el proyecto.	14
A.13. Reiniciar el contenedor completamente.	14
A.14. Eliminar el contenedor y limpiar el entorno.	14
B.1. csv_glance.py, exploración inicial del conjunto de datos.	16
B.2. csv_count_registers.py, conteo de registros en el conjunto de datos.	17
B.3. remove_columns.py, eliminación de campos innecesarios en el conjunto de datos.	18
B.4. unique_values.py, obtención de valores únicos de la columna 'device_horizontal_accuracy'.	19
B.5. accuracy_histogram.py, creación de un histograma de frecuencias de la columna 'device_horizontal_accuracy'.	21
B.6. identifier_histogram.py, creación de un histograma de frecuencias de la columna 'identifier'.	24
B.7. identifier_histogram_detailed.py, análisis de frecuencias de la columna 'identifier'.	27
B.8. csv_deduplicate.py, eliminación de duplicados en el conjunto de datos.	31

Capítulo 1

Introducción del Proyecto

1.1 Descripción general del proyecto

La simulación de una red de comunicaciones con dispositivos personales requiere modelos que representen fielmente los patrones de movimiento de las personas. De lo contrario, las conclusiones derivadas de dicha simulación pueden ser poco útiles. Para avanzar hacia la definición de un modelo de trayectorias individuales, se propone caracterizar los datos de una base existente que permita modelar trayectorias de forma eficaz.

1.2 Objetivos y propósitos

El objetivo principal del proyecto es obtener una caracterización estadística de las trayectorias individuales.

Los propósitos específicos son:

- Caracterizar la base de datos para extraer las trayectorias contenidas.
- Aplicar un modelo de inteligencia artificial para identificar y analizar dichas trayectorias.

1.3 Alcance del sistema

El sistema se enfoca en la identificación de trayectorias peatonales individuales y su análisis mediante herramientas de IA. El alcance incluye:

- Caracterización de la base de datos existente.
- Identificación de trayectorias individuales.
- Generación de reportes y visualizaciones de los resultados.

No se incluye la creación de modelos de IA desde cero; se emplearán herramientas y modelos ya existentes.

Capítulo 2

Requisitos del sistema

2.1 Requisitos

Docker: version 28.2.2, build e6534b4

Docker Compose: version 1.29.2, build unknown

2.1.1 Instrucciones de instalación

1. Clonar el repositorio del proyecto desde el siguiente enlace el proyecto se encuentra dentro de la carpeta **Implementación**.
2. Para levantar y acceder al contenedor, ejecutar el siguiente script:¹

```
1 ./start_container.sh      # Linux/Mac
2 .\start_container.bat     # Windows
```

Código 2.1: Iniciar contenedor del proyecto.

3. Para cerrar el contenedor del proyecto, ejecutar el siguiente script:

```
1 ./stop_container.sh      # Linux/Mac
2 .\stop_container.bat     # Windows
```

Código 2.2: Iniciar contenedor del proyecto.

4. Para ver reflejados los cambios realizados en el código, ejecuta el siguiente script:

```
1 ./restart_container.sh   # Linux/Mac
2 .\restart_container.bat  # Windows
```

Código 2.3: Reiniciar contenedor del proyecto.

¹Para más detalles sobre el uso de Docker y Docker Compose, consulte el Apéndice A.

Capítulo 3

Caracterización de datos de trayectorias individuales

3.1 Caracterización de datos de trayectorias individuales

El análisis de datos comienza con una etapa fundamental: la caracterización del conjunto de datos. Esta fase tiene como objetivo examinar y comprender la estructura, el contenido y las principales propiedades de los datos antes de aplicar técnicas analíticas más complejas. En el caso de los datos de trayectorias individuales, la caracterización permite identificar posibles inconsistencias, redundancias y elementos irrelevantes que puedan afectar la calidad del análisis. Las tareas principales llevadas a cabo en esta etapa son las siguientes:

- Explorar las primeras filas del conjunto de datos para obtener una visión general de su estructura.
- Verificar la cantidad total de registros y columnas disponibles.
- Identificar y eliminar columnas que no aportan información relevante para el análisis o inconsistentes.
- Identificar y eliminar las filas que no aportan información relevante para el análisis o inconsistentes.

A continuación, se describen en detalle las acciones específicas realizadas durante el proceso de caracterización.

3.1.1 Exploración inicial del conjunto de datos

Como primer paso en la caracterización, se realizó una exploración preliminar del conjunto de datos con el fin de comprender su estructura general. Para ello, se inspeccionaron las primeras dos filas, lo cual permitió identificar las columnas presentes y observar ejemplos representativos de sus valores. El código utilizado para realizar esta exploración se encuentra en el Apéndice B.1. A continuación, se presenta un resumen de las columnas detectadas junto con una muestra de sus respectivos valores:

1. `id`: Identificador numérico único por registro
[`'34284565'`, `'34284566'`]
2. `identifier`: UUID del dispositivo
[`'f2640430-7e39-41b7-80bb-3fddaa44779c'`]
3. `identifier_type`: Tipo de ID (ej. `'gaid'` para Android)
[`'gaid'`, `'gaid'`]
4. `timestamp`: Fecha-hora del registro
[`'2022-11-07 02:04:21'`]
5. `device_lat/device_lon`: Coordenadas GPS
[`'21.843149'`], [`'-102.196838'`]
6. `country_short/province_short`: Códigos de ubicación
[`'MX'`], [`'MX.01'`]
7. `ip_address`: Dirección IPv6
[`'2806:103e:16::'`]
8. `device_horizontal_accuracy`: Precisión GPS en metros
[`'8.0'`]
9. `source_id`: Hash de la fuente de datos
[`'449d086d...344'`]
10. `record_id`: Hash único por registro
[`'77d795df...'`]
11. `home_country_code`: País de residencia
[`'MX'`]
12. `home_geog_point/work_geog_point`: Coordenadas en WKT
[`'POINT(-102.37038 22.20753)'`]
13. `home_hex_id/work_hex_id`: ID hexagonal (H3)
[`'85498853fffffff'`]
14. `data_execute`: Fecha de procesamiento
[`'2023-05-30'`]
15. `time_zone_name`: Zona horaria
[`'America/Mexico.City'`]

3.1.2 Dimensiones del conjunto de datos

Para verificar las dimensiones del conjunto de datos, se utilizó la biblioteca Dask, que permite trabajar con grandes volúmenes de datos de manera eficiente. Junto con Python se usó el código en el Apéndice B.2. Como resultado ahora sabemos que el conjunto de datos contiene un total de **69,980,000** registros y **19** campos. Esto indica que hay una cantidad significativa de datos disponibles para el análisis.

3.1.3 Depuración de columnas

Dado que el conjunto de datos original contiene 19 campos, es fundamental identificar y eliminar aquellas columnas que no aportan valor al análisis. Para ello, se realizó una revisión de los valores únicos presentes en cada campo, con el objetivo de detectar información redundante o irrelevante. A partir de este análisis, se identificaron las siguientes columnas como innecesarias para los fines del estudio:

- `id`
- `identifier_type`
- `country_short`
- `province_short`
- `ip_address`
- `source_id`
- `home_country_code`
- `home_geog_point`
- `work_geog_point`
- `home_hex_id`
- `work_hex_id`
- `data_execute`

En lugar de eliminar columnas explícitamente, se optó por seleccionar únicamente aquellas que se desean conservar. El código utilizado para esta tarea se encuentra incluido en el Apéndice B.3. Dicho script emplea la biblioteca `dask` para cargar y guardar una nueva versión del conjunto de datos que contiene exclusivamente las siguientes columnas relevantes:

- `identifier`
- `timestamp`
- `device_lat`
- `device_lon`
- `device_horizontal_accuracy`
- `record_id`
- `time_zone_name`

Como resultado, se genera un nuevo archivo **CSV** que conserva únicamente la información útil para el análisis posterior, optimizando así el tamaño y la calidad del conjunto de datos.

3.1.4 Depuración de filas

Una vez obtenida una versión más ligera del conjunto de datos, el siguiente paso consiste en identificar y eliminar aquellas filas que no aportan valor al análisis. Para ello, se generaron representaciones gráficas que permiten observar la distribución de los datos y facilitar la toma de decisiones. Las columnas seleccionadas para este proceso fueron:

- **identifier**: Identificador único del dispositivo.
- **device_horizontal_accuracy**: Precisión del GPS en metros. A menor valor, mayor precisión.

La primera columna a analizar será **device_horizontal_accuracy**, que refleja la precisión del GPS en metros. Este valor depende tanto del sistema de medición como de la fuente de datos, y suele clasificarse según la siguiente escala:

- GPS puro (satelital): 1–20 metros.
- A-GPS (asistido por red): 5–50 metros.
- Triangulación por WiFi o redes móviles: 20–500 metros.
- Geolocalización por IP: 1000–5000 metros.

Con base en esta escala, primero hay que identificar el rango de valores presentes en la columna. Para ello se utilizó el código mostrado en el Apéndice B.4, el cual extrae los valores únicos de **device_horizontal_accuracy** y los guarda en un archivo de texto. El resultado indicó que los valores oscilan entre 0.916 y 199.9, lo que permitió construir un histograma (Apéndice B.5) para analizar la frecuencia de cada valor y así evaluar su relevancia para el análisis. El resultado se muestra en la siguiente figura:

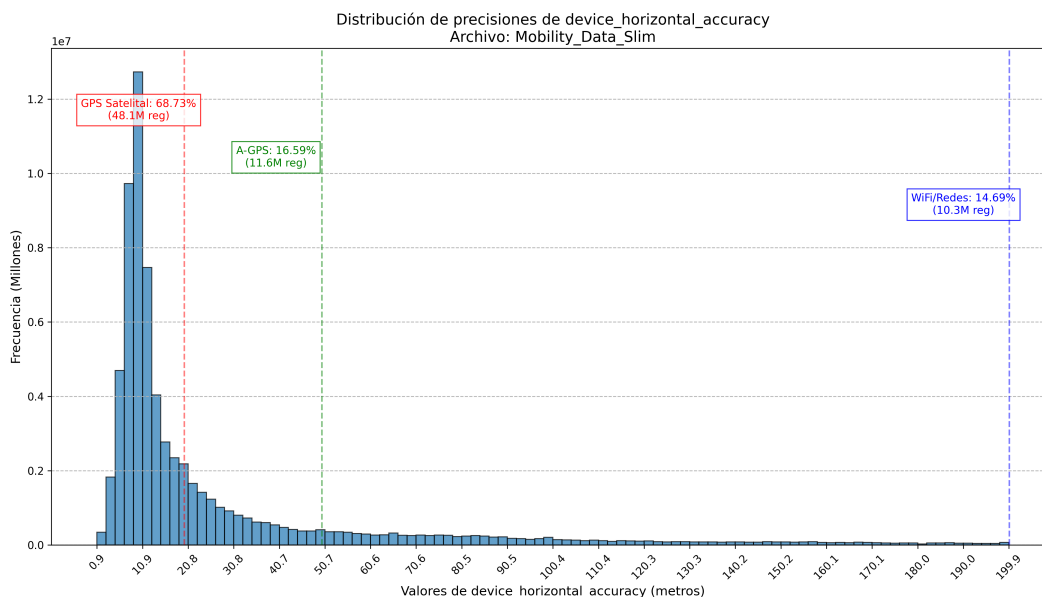


Figura 3.1: Frecuencia de aparición de los valores de 'device_horizontal_accuracy'.

Para el objetivo de este proyecto, se busca que la configuración del GPS sea lo más precisa posible, por lo que aquellos que estén dentro del rango del GPS puro (1-20 metros) son los más relevantes. Como se puede ver en la Figura 3.1, el **68.73%** de los valores se encuentran dentro de este rango. Sin embargo, el **31.27%** de registros con están por encima de este rango, precisión A-GPS (5-50 metros) y triangulación por WiFi/red móvil (20-500 metros).

La siguiente columna a evaluar es **identifier**, corresponde al identificador único de cada dispositivo. Para analizar la frecuencia de aparición de estos valores se empleó un script que agrupa las repeticiones por rangos y grafica la cantidad de valores únicos usando escala logarítmica (ver Apéndice B.6).

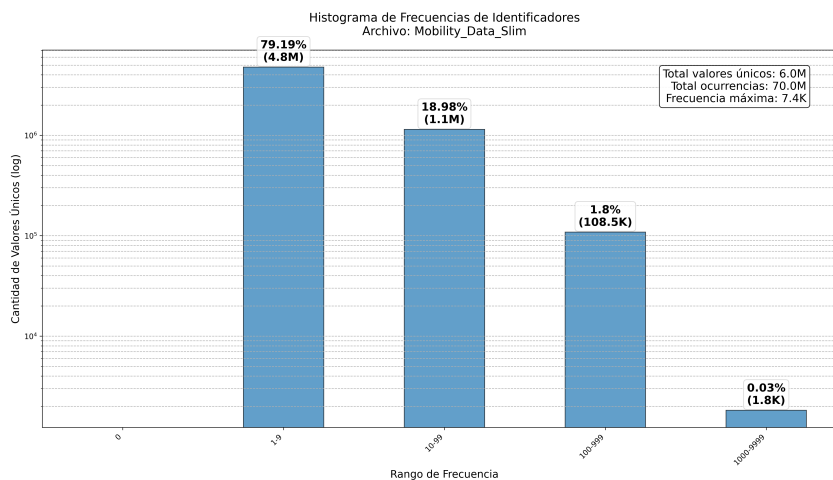
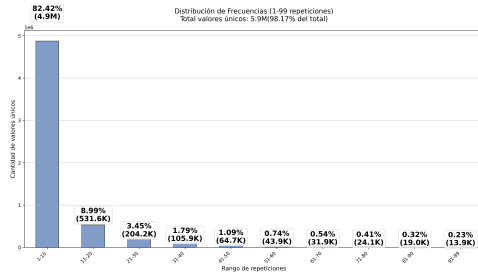
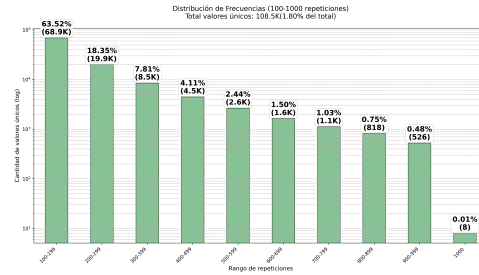


Figura 3.2: Frecuencia de aparición de los identificadores únicos.

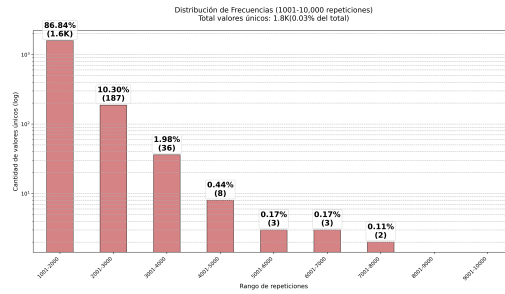
De este script sabemos que el total de individuos es de **6,022,772** de los cuales el **79.19 %** tienen una frecuencia de aparición de una a nueve veces, esto es **4,769,317** de individuos. Así mismo de la Figura 3.2 podemos observar que hay poco más de un **20 %** de individuos con más de 99 repeticiones. Por lo que se necesita hacer un análisis más detallado, para ello se ejecuta el código del Apéndice B.7, el cual segmenta los datos en tres rangos: 1-99, 100-1000 y 1001-10000 repeticiones.



(a) Histograma 1-99 repeticiones



(b) Histograma 100-1000 repeticiones



(c) Histograma 1001-10000 repeticiones

Figura 3.3: Comparación de histogramas por rangos de repeticiones.

Con la información obtenida de los histogramas de la figura anterior, se puede observar que el **98.17 %** de los identificadores únicos tienen entre 1 y 99 repeticiones, lo que equivale a **5,912,437** individuos. Por otro lado, el **1.83 %** restante tiene entre 100 y 10,000 repeticiones, lo que equivale a **110,335** individuos. Con base en esta información aún no se puede determinar que registros eliminar. Por lo que ahora se eliminarán aquellos registros que sean duplicados, es decir, aquellos que tengan el mismo `timestamp`, `device_lat` y `device_lon`. Para ello se utilizó el código del Apéndice B.8, que elimina los duplicados y genera un nuevo archivo CSV con los registros únicos.

Apéndice A

Uso de Docker y Docker Compose

En la sección 2.1 se establece como requisito el uso de Docker y Docker Compose para la ejecución del proyecto. A continuación, se detallan las instrucciones necesarias para su instalación, ya que ambas herramientas son fundamentales para la implementación. Además, se describe el archivo `docker-compose.yml`, el cual permite crear un contenedor que incluye todas las dependencias requeridas para el correcto funcionamiento del sistema.

A.1 ¿Qué son Docker y Docker Compose?

Docker es una plataforma de virtualización ligera que permite desarrollar, empaquetar y ejecutar aplicaciones en contenedores aislados. Un contenedor incluye el código, las dependencias y configuraciones necesarias para que la aplicación se ejecute de manera consistente en cualquier entorno. Esto facilita la portabilidad, escalabilidad y despliegue de software.

Docker Compose es una herramienta que permite definir y ejecutar aplicaciones multicontenedor mediante archivos de configuración YAML. A través de un solo archivo `docker-compose.yml`, es posible especificar los servicios, redes y volúmenes que componen una aplicación, simplificando así su orquestación.

Estas herramientas son fundamentales en este proyecto para garantizar que el entorno de ejecución sea replicable y controlado, independientemente del sistema operativo o configuración local del usuario.

A.2 Instalación en Linux (Ubuntu/Debian)

Para instalar Docker y Docker Compose en un sistema Linux basado en Debian o Ubuntu, siga los siguientes pasos:

1. Actualizar los paquetes del sistema:

```
1 sudo apt update
2 sudo apt upgrade
```

Código A.1: Actualizar el sistema.

2. Instalar Docker:

```
1 sudo apt install docker.io
2 sudo systemctl enable docker
3 sudo systemctl start docker
```

Código A.2: Instalar Docker.

3. Verificar que Docker está instalado correctamente:

```
1 docker --version
```

Código A.3: Verificar instalación de Docker.

4. Instalar Docker Compose:

```
1 sudo apt install docker-compose
```

Código A.4: Instalar Docker Compose.

5. Verificar la instalación:

```
1 docker-compose --version
```

Código A.5: Verificar instalación de Docker Compose.

A.3 Instalación en Windows

Para instalar Docker y Docker Compose en Windows, se recomienda utilizar Docker Desktop, que incluye ambas herramientas de forma integrada.

1. Acceder al sitio oficial: <https://www.docker.com/products/docker-desktop/>
2. Descargar el instalador correspondiente para Windows.
3. Ejecutar el instalador y seguir el asistente de instalación.
4. Reiniciar el sistema si es necesario.
5. Verificar que Docker y Docker Compose estén correctamente instalados desde la terminal de Windows (PowerShell o CMD):

```
1 docker --version
2 docker-compose --version
```

Código A.6: Verificar instalación de Docker y Docker Compose.

Nota: Docker Desktop requiere que la virtualización esté habilitada en la BIOS del sistema. También es necesario contar con Windows 10 o superior.

A.4 Descripción del archivo `docker-compose.yml`

El archivo `docker-compose.yml` permite definir y configurar el entorno de ejecución del proyecto utilizando un contenedor de Docker. A continuación, se presenta su contenido y una explicación de cada uno de sus elementos:

```
1 version: "3.8"
2
3 services:
4   data-analysis:
5     image: python:3.13-bookworm
6     container_name: data-analysis
7     runtime: nvidia
8     tty: true
9     stdin_open: true
10    volumes:
11      - ./:/app
12      - python-packages:/usr/local/lib/python3.13/site-packages
13    command: sh -c "pip install -r requirements.txt && pip install -e . && python3 /app/src/main.py"
14    working_dir: /app
15    environment:
16      - PYTHONPATH=/app
17
18 volumes:
19   python-packages:
```

Código A.7: Archivo `docker-compose.yml`

A continuación se explica el propósito de cada sección:

- **version: "3.8"**
Define la versión del esquema de Docker Compose utilizado. La versión 3.8 es compatible con la mayoría de las características modernas de Docker.
- **services → data-analysis**
Se define un servicio llamado `data-analysis`, que representa el contenedor principal del proyecto.
- **image: python:3.13-bookworm**
Utiliza una imagen oficial de Python 3.13 basada en Debian Bookworm como entorno base.
- **container_name: data-analysis**
Asigna un nombre personalizado al contenedor para facilitar su identificación.
- **runtime: nvidia**
Indica que el contenedor utilizará el runtime de NVIDIA para permitir acceso a la GPU. Requiere tener instalado `nvidia-docker`.
- **tty: true y stdin_open: true**
Habilitan la interacción con el terminal del contenedor, lo que es útil para

ejecutar comandos manuales si es necesario.

- **volumes**

- `./:/app`: Monta el directorio actual del proyecto como `/app` dentro del contenedor.
- `python-packages:/usr/local/lib/python3.13/site-packages`: Crea un volumen persistente para las bibliotecas de Python instaladas.

- **command**

Ejecuta una serie de comandos cuando el contenedor inicia: instala las dependencias del archivo `requirements.txt`, instala el proyecto en modo editable (`pip install -e .`) y ejecuta el archivo `main.py`.

- **working_dir**: `/app`

Establece el directorio de trabajo dentro del contenedor como `/app`.

- **environment**

Define la variable de entorno `PYTHONPATH` para que Python pueda encontrar correctamente los módulos dentro del proyecto.

- **volumes** → **python-packages**

Declara un volumen persistente llamado `python-packages`, que se utiliza para almacenar los paquetes instalados sin perderlos entre reinicios del contenedor.

Este archivo permite que el entorno de desarrollo sea fácilmente replicable y ejecutable, sin necesidad de instalar manualmente dependencias o configurar rutas en el sistema anfitrión.

A.5 Scripts de control del contenedor

Para facilitar el manejo del contenedor durante el desarrollo del proyecto, se han creado tres scripts auxiliares en Bash que automatizan las operaciones más comunes: iniciar, reiniciar y detener el contenedor.

A.5.1 start_container.sh

Este script verifica si el contenedor `data-analysis` ya se encuentra en ejecución. En caso de que no esté activo, lo inicia utilizando `docker-compose up -d`. Posteriormente, ejecuta el archivo `main.py` dentro del contenedor.

```
1 #!/bin/bash
2
3 if ! docker ps --filter "name=~/data-analysis$" --filter "
   status=running" | grep -q data-analysis; then
4     echo "Contenedor no está corriendo. Levantando con docker-
       compose..."
5     docker-compose up -d
6     echo "Esperando que se instalen las dependencias..."
7     while ! docker exec data-analysis pip show colorama &> /dev
       /null; do
8         sleep 2
9     done
10    echo "Dependencias instaladas correctamente."
11 else
12    echo "Contenedor ya está corriendo. Usando instancia
       existente."
13 fi
14
15 echo "Ejecutando script..."
16 docker exec -it data-analysis python3 /app/src/main.py
```

Código A.8: Script para iniciar el contenedor.

A.5.2 restart_container.sh

Este script reinicia completamente el contenedor (equivalente a detenerlo y volverlo a levantar), lo cual resulta útil cuando se han modificado archivos como `requirements.txt` o `setup.py`. Tras reiniciar, vuelve a ejecutar el archivo principal del proyecto.

```
1 #!/bin/bash
2 docker restart data-analysis
3 sleep 2
4 docker exec -it data-analysis python3 /app/src/main.py
```

Código A.9: Script para reiniciar el contenedor.

A.5.3 stop_container.sh

Este script detiene y elimina el contenedor junto con los volúmenes asociados. Debe utilizarse con precaución, ya que elimina todas las dependencias instaladas en el entorno del contenedor. Solo es necesario en casos donde se requiere limpiar completamente el entorno.

```
1 #!/bin/bash
2 docker-compose down --volumes
```

Código A.10: Script para detener y eliminar el contenedor y sus volúmenes.

A.6 Proceso de uso y desarrollo del contenedor

A continuación se describe el flujo recomendado para desarrollar y ejecutar el sistema dentro del contenedor de Docker:

1. Verifique que Docker y Docker Compose están instalados (Apéndice A.6).

Nota para usuarios de Windows: Si se utiliza Windows como sistema operativo, se deben usar los archivos con extensión `.bat` en lugar de `.sh`, y deben ser ejecutados desde la terminal de Windows (por ejemplo, CMD o PowerShell).

2. Asigne permisos de ejecución a los scripts:

```
1 chmod +x start_container.sh restart_container.sh
   stop_container.sh
```

Código A.11: Dar permisos de ejecución a los scripts.

3. Para iniciar el contenedor y ejecutar el proyecto con los cambios más recientes del código fuente:

```
1 ./start_container.sh
```

Código A.12: Iniciar contenedor y ejecutar el proyecto.

4. Si se realizan cambios en las dependencias o archivos de configuración del entorno (como `requirements.txt`), utilice:

```
1 ./restart_container.sh
```

Código A.13: Reiniciar el contenedor completamente.

5. Para detener el contenedor y eliminar todos los volúmenes asociados:

```
1 ./stop_container.sh
```

Código A.14: Eliminar el contenedor y limpiar el entorno.

Este conjunto de scripts permite un desarrollo ágil dentro del contenedor, ya que los cambios realizados en el código fuente local se reflejan de inmediato gracias al uso de **volumes**. Además, se reduce la necesidad de ejecutar manualmente comandos repetitivos, facilitando el trabajo del usuario final y asegurando la correcta ejecución del proyecto.

Apéndice B

Scripts para la caracterización de datos

En la sección 3.1 se describen los pasos del proceso de caracterización de datos de trayectorias individuales. En este anexo se presentan los scripts utilizados para llevar a cabo dicho proceso.

```
1  import dask.dataframe as dd
2  import sys
3
4  print("Exploracion_inicial_de_datos_con_Dask\n")
5
6  if len(sys.argv) < 2:
7      print("Error:_Debe_especificar_un_archivo_CSV")
8      sys.exit(1)
9
10 ruta_archivo = sys.argv[1]
11
12 ddf = dd.read_csv(
13     ruta_archivo,
14     encoding="utf-8",
15     sep="," ,
16     dtype="object",
17 )
18
19 columnas = ddf.columns.tolist()
20
21 print("Columnas_y_2_ejemplos_por_cada_una:\n")
22 for col in columnas:
23     ejemplos = ddf[col].head(2).values.tolist()
24     print(f"-_{col}:_{ejemplos}")
25
26 input("Presiona_Enter_para_continuar...")
```

Código B.1: csv_glance.py, exploración inicial del conjunto de datos.

```
1  import dask.dataframe as dd
2  import sys
3  import os
4
5  def contar_registros(ruta_archivo):
6
7      columnas_usar = ["record_id"]
8      try:
9          print(f"\nCargando archivo {ruta_archivo}...")
10         ddf = dd.read_csv(
11             ruta_archivo,
12             usecols=columnas_usar,
13             sep=";",
14             dtype={"record_id": "str"},
15             blocksize="256MB",
16         )
17
18         print("Contando registros (paciencia para archivos grandes)...")
19         total_registros = ddf.shape[0].compute()
20
21         print(f"\nAnálisis completado:")
22         print(f"Archivo analizado: {ruta_archivo}")
23         print(f"Total de registros: {total_registros:,}")
24
25     except Exception as e:
26         print(f"\nOcurrió un error inesperado: {str(e)}")
27
28 if __name__ == "__main__":
29     print("=== Contador de registros en archivos CSV grandes ===")
30
31     if len(sys.argv) < 2:
32         print("Uso: python csv_count_registers.py <nombre_del_archivo.csv>")
33         sys.exit(1)
34
35     archivo = sys.argv[1]
36     contar_registros(archivo)
```

Código B.2: csv_count_registers.py, conteo de registros en el conjunto de datos.

```
1  import dask.dataframe as dd
2
3  columnas_deseadas = [
4      'identifier',
5      'timestamp',
6      'device_lat',
7      'device_lon',
8      'device_horizontal_accuracy',
9      'record_id',
10     'time_zone_name'
11 ]
12
13 df = dd.read_csv('Mobility_Data.csv', usecols=
14     columnas_deseadas)
15
16 df.to_csv('Mobility_Data_Slim.csv', index=False, single_file=
17     True, encoding='utf-8-sig')
```

Código B.3: remove_columns.py, eliminación de campos innecesarios en el conjunto de datos.

```
1  import pandas as pd
2  from tqdm import tqdm
3  import os
4  import sys
5  from src.menus.menu import MainMenu
6  def main():
7      print("\n" + "="*50)
8      print("└─EXTRACTOR└─DE└─VALORES└─UNICOS└─DE└─COLUMNAS└─CSV")
9      print("="*50 + "\n")
10
11     if len(sys.argv) < 2:
12         print("Uso:└─python└─extract_unique.py└─<archivo.csv>")
13         sys.exit(1)
14
15     csv_file = sys.argv[1]
16
17     if not os.path.exists(csv_file):
18         print(f"Error:└─El└─archivo└─'{csv_file}'└─no└─existe.")
19         sys.exit(1)
20
21     chunk_size = 1_000_000
22
23     try:
24         available_columns = pd.read_csv(csv_file, nrows=0).
25             columns.tolist()
26     except Exception as e:
27         print(f"Error└─leyendo└─el└─archivo:└─{e}")
28         sys.exit(1)
29
30     try:
31         selected_index = MainMenu.display_available_columns(
32             available_columns)
33         target_column = available_columns[selected_index]
34     except (ValueError, IndexError):
35         print("Selección└─inválida.")
36         sys.exit(1)
37     except Exception as e:
38         print(f"Error└─inesperado└─al└─seleccionar└─columna:└─{e}"
39             )
40         sys.exit(1)
41
42     safe_column_name = target_column.replace("└─", "└─").
43         replace("/", "└─")
44     output_file = f"valores_unicos_{safe_column_name}.txt"
45
46     unique_values = set()
47     print(f"\nProcesando└─columna:└─{target_column}\n")
48
49     try:
50         for chunk in tqdm(pd.read_csv(csv_file, usecols=[
```



```

47         target_column], chunksize=chunk_size)):
48             unique_values.update(chunk[target_column].dropna
49                                   ().astype(str))
50     except Exception as e:
51         print(f"Error durante el procesamiento: {e}")
52         sys.exit(1)
53
54     try:
55         numeric_values = sorted([float(v) for v in
56                                 unique_values])
57         is_numeric = True
58     except ValueError:
59         is_numeric = False
60
61     try:
62         with open(output_file, "w", encoding="utf-8") as f:
63             if is_numeric:
64                 min_val = numeric_values[0]
65                 max_val = numeric_values[-1]
66                 f.write(f"# Rango de valores: {min_val} - {
67                         max_val}\n")
68                 f.write("\n".join(str(v) for v in
69                                   numeric_values))
70             else:
71                 sorted_values = sorted(unique_values)
72                 f.write(f"# Rango de valores: No numerico\n")
73                 f.write("\n".join(sorted_values))
74     except Exception as e:
75         print(f"Error guardando los resultados: {e}")
76         sys.exit(1)
77
78     print(f"\nSe encontraron {len(unique_values):,} valores
79         unicos.")
80     print(f"Resultados guardados en: {output_file}")
81
82     print("\nMuestra de valores unicos (primeros 10):")
83     print("\n".join(sorted(unique_values)[:10]))
84
85 if __name__ == "__main__":
86     main()

```

Código B.4: unique_values.py, obtención de valores únicos de la columna 'device_horizontal_accuracy'.

```
1 import os
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import numpy as np
5 import sys
6 from tqdm import tqdm
7
8 def classify_tech(valor):
9     if 1 <= valor <= 20:
10         return 'GPS_Satelital'
11     elif 5 <= valor <= 50:
12         return 'A-GPS_(Asistido_por_red)'
13     elif 20 <= valor <= 500:
14         return 'Triangulacion_WiFi/Redes_Moviles'
15     else:
16         return 'Fuera_de_rango'
17
18 def format_count(count):
19     if count >= 1_000_000:
20         return f"{count/1_000_000:.1f}M"
21     elif count >= 1_000:
22         return f"{count/1_000:.1f}K"
23     return str(count)
24
25 def main():
26     if len(sys.argv) < 2:
27         print("Error: Debe especificar un archivo CSV como argumento")
28         sys.exit(1)
29
30     csv_file = sys.argv[1]
31     filename = os.path.splitext(os.path.basename(csv_file))
32     [0]
33     column = "device_horizontal_accuracy"
34     bins = 100
35
36     print(f"\nIniciando procesamiento del archivo: {csv_file}")
37     print(f"Columna analizada: {column}")
38
39     os.makedirs("img", exist_ok=True)
40     print("Directorio 'img' verificado/creado")
41
42     print("\nProcesando datos y clasificando tecnologías...")
43
44     frequency = pd.Series(dtype=float)
45     tech_counts = {
46         'GPS_Satelital': 0,
47         'A-GPS_(Asistido_por_red)': 0,
48         'Triangulacion_WiFi/Redes_Moviles': 0,
```

```

47         'Fuera_de_rango': 0
48     }
49
50     total_row = sum(1 for _ in pd.read_csv(csv_file,
51                                     usecols=[column], chunksize=1_000_000))
52
53     with tqdm(total=total_row, unit='M_rows') as pbar:
54         for chunk in pd.read_csv(csv_file, usecols=[column],
55                                 chunksize=1_000_000):
56             chunk_clean = chunk[column].dropna()
57
58             for valor in chunk_clean:
59                 tech = classify_tech(valor)
60                 tech_counts[tech] += 1
61
62             counts = chunk_clean.value_counts()
63             if not frequency.empty or not counts.empty:
64                 frequency = pd.concat([frequency, counts],
65                                     axis=0).groupby(level=0).sum()
66             pbar.update(1)
67
68     total = sum(tech_counts.values())
69     percentage = {k: (v/total)*100 for k, v in tech_counts.items()}
70
71     print("\nGenerando histograma con estadísticas...")
72     counts, edges = np.histogram(frequency.index, bins=bins,
73                                 weights=frequency.values)
74
75     plt.figure(figsize=(14, 8))
76     plt.bar(edges[:-1], counts, width=np.diff(edges), align='edge',
77            edgecolor='black', alpha=0.7)
78
79     plt.axvline(x=20, color='r', linestyle='--', alpha=0.5)
80     plt.axvline(x=50, color='g', linestyle='--', alpha=0.5)
81     plt.axvline(x=200, color='b', linestyle='--', alpha=0.5)
82
83     gps_str = f"GPS_Satelital: {percentage['GPS_Satelital']:.2f}%\n({format_count(tech_counts['GPS_Satelital'])})reg)"
84
85     agps_str = f"A-GPS: {percentage['A-GPS (Asistido por red)']:.2f}%\n({format_count(tech_counts['A-GPS (Asistido por red)'])})reg)"
86
87     wifi_str = f"WiFi/Redes: {percentage['Triangulacion WiFi/Redes Moviles']:.2f}%\n({format_count(tech_counts['Triangulacion WiFi/Redes Moviles'])})reg)"
88
89     plt.text(10, max(counts)*0.9, gps_str, ha='center', color='r', fontsize=10,

```

```

83         bbox=dict(facecolor='white', alpha=0.8,
84                   edgecolor='r'))
85         plt.text(40, max(counts)*0.8, agps_str, ha='center',
86                 color='g', fontsize=10,
87                 bbox=dict(facecolor='white', alpha=0.8,
88                           edgecolor='g'))
89         plt.text(190, max(counts)*0.7, wifi_str, ha='center',
90                 color='b', fontsize=10,
91                 bbox=dict(facecolor='white', alpha=0.8,
92                           edgecolor='b'))
93
94     plt.title(f"Distribución de precisiones de {column}\n"
95              f"Archivo: {filename}", fontsize=14)
96     plt.xlabel(f"Valores de {column} (metros)", fontsize
97               =12)
98     plt.ylabel("Frecuencia (Millones)", fontsize=12)
99     plt.xticks(edges[:5], rotation=45)
100    plt.grid(axis='y', linestyle='--')
101    plt.tight_layout()
102
103    output_path = os.path.join("img", f"histograma_{column}"
104                               f"_{filename}.png")
105    plt.savefig(output_path, dpi=300, bbox_inches='tight')
106    plt.close()
107
108    print("\n=== DISTRIBUCION DE TECNOLOGIAS DE GEOLOCALIZACION ===")
109    for tech, count in tech_counts.items():
110        print(f"{tech}: {count}, registros ({percentage[
111            tech]:.2f}%)")
112
113    print(f"\nHistograma generado exitosamente")
114    print(f"Archivo guardado en: {output_path}")
115    print(f"Total registros analizados: {total:,}\n")
116
117    if __name__ == "__main__":
118        main()

```

Código B.5: accuracy_histogram.py, creación de un histograma de frecuencias de la columna 'device_horizontal_accuracy'.

```
1  import os
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  import numpy as np
5  from collections import Counter
6  import sys
7  from tqdm import tqdm
8  import math
9
10 def format_count(count):
11     if count >= 1_000_000:
12         return f"{count/1_000_000:.1f}M"
13     elif count >= 1_000:
14         return f"{count/1_000:.1f}K"
15     return str(count)
16
17 def main():
18     if len(sys.argv) < 2:
19         print("Error: Debe especificar un archivo CSV como argumento")
20         sys.exit(1)
21
22     csv_file = sys.argv[1]
23     filename = os.path.splitext(os.path.basename(csv_file))
24         [0]
25     column = "identifier"
26     chunksize = 1_000_000
27
28     print(f"\nIniciando procesamiento del archivo: {csv_file}")
29     print(f"Columna analizada: {column}")
30     os.makedirs("img", exist_ok=True)
31     print("Directorio 'img' verificado/creado")
32
33     print("\nProcesando datos y contando frecuencias...")
34     counter = Counter()
35
36     total_chunks = sum(1 for _ in pd.read_csv(csv_file,
37         usecols=[column], chunksize=chunksize))
38
39     with tqdm(total=total_chunks, unit='chunk') as pbar:
40         for chunk in pd.read_csv(csv_file, usecols=[column],
41             chunksize=chunksize):
42             counter.update(chunk[column].dropna().astype(str))
43             pbar.update(1)
44
45     frequency = pd.Series(counter)
46     total_unique_values = len(frequency)
47     max_freq = frequency.max()
```

```
45
46     print(f"Datos▯procesados▯correctamente")
47     print(f"Total▯de▯valores▯Unicos:▯{total_unique_values
48           : ,}")
49     print(f"Frecuencia▯maxima:▯{max_freq: ,}")
50
51     bins = [0] + [10**i for i in range(0, int(np.log10(
52           max_freq)) + 2)]
53     group_freq = pd.cut(frecuency, bins=bins, right=False).
54           value_counts().sort_index()
55
56     total_ocurrence = frecuency.sum()
57     percentage_per_range = (group_freq /
58           total_unique_values * 100).round(2)
59
60     print("\nGenerando▯histograma▯con▯estadísticas...")
61     plt.figure(figsize=(16, 9))
62     ax = group_freq.plot(kind='bar', logy=True, alpha=0.7,
63           edgecolor='black')
64
65     formatted_labels = []
66     for interval in group_freq.index.categories:
67         left = int(interval.left)
68         right = int(interval.right - 1)
69         formatted_labels.append(f"{left}-{right}" if left
70               != right else f"{left}")
71
72     plt.xticks(range(len(formatted_labels)),
73           formatted_labels, rotation=45, ha='right')
74
75     plt.title(f"Histograma▯de▯Frecuencias▯de▯
76           Identificadores\nArchivo:▯{filename}", fontsize=16,
77           pad=20)
78     plt.xlabel("Rango▯de▯Frecuencia", fontsize=14)
79     plt.ylabel("Cantidad▯de▯Valores▯Unicos▯(log)", fontsize
80           =14)
81     plt.grid(True, which="both", ls="--", axis='y')
82
83     stats_text = (
84         f"Total▯valores▯unicos:▯{format_count(
85               total_unique_values)}\n"
86         f"Total▯ocurrencias:▯{format_count(total_ocurrence)
87               }\n"
88         f"Frecuencia▯maxima:▯{format_count(max_freq)}"
89     )
90     plt.annotate(stats_text,
91           xy=(0.95, 0.95),
92           xycoords='axes▯fraction',
93           fontsize=15,
94           ha='right',
95           va='top',
```

```

84         bbox=dict(boxstyle='round', facecolor='
            white', alpha=0.9))
85
86     max_val = group_freq.max()
87     min_y = 0.9
88
89     for i, (count, percent) in enumerate(zip(group_freq.
        values, percentage_per_range.values)):
90         if count > 0:
91             y_pos = count * 1.1 if count * 1.1 > min_y else
                min_y * 1.2
92
93             text = f"{percent}%\n({format_count(count)})"
94
95             ax.text(
96                 i, y_pos, text,
97                 ha='center', va='bottom',
98                 fontsize=15,
99                 fontweight='bold',
100                 bbox=dict(
101                     facecolor='white',
102                     alpha=0.85,
103                     edgecolor='lightgray',
104                     boxstyle='round,pad=0.3'
105                 )
106             )
107
108     output_path = os.path.join("img", f"histograma_{column}
        _{filename}.png")
109     plt.tight_layout()
110     plt.savefig(output_path, dpi=300, bbox_inches='tight')
111     plt.close()
112
113     print("\n=== DISTRIBUCION DE FRECUENCIAS ===")
114     for i, (intervalo, count) in enumerate(group_freq.items
        ()):
115         print(f"Rango {formatted_labels[i]}: {count:,}")
116
117     print(f"\nHistograma generado exitosamente")
118     print(f"Archivo guardado en: {output_path}")
119     print(f"Total ocurrencias analizadas: {total_occurrence
        :,}\n")
120
121     if __name__ == "__main__":
122         main()

```

Código B.6: `identifier_histogram.py`, creación de un histograma de frecuencias de la columna 'identifier'.

```
1  import os
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  import numpy as np
5  from collections import Counter
6  import sys
7  from tqdm import tqdm
8
9  def format_count(count):
10     if count >= 1_000_000:
11         return f"{count/1_000_000:.1f}M"
12     elif count >= 1_000:
13         return f"{count/1_000:.1f}K"
14     return str(count)
15
16  def create_histogram(data, bins, title, filename, color='
skyblue', log_scale=False):
17     grouped = pd.cut(data, bins=bins, right=False).
value_counts().sort_index()
18     total_values = len(data)
19     max_count = grouped.max()
20
21     plt.figure(figsize=(14, 8))
22     ax = grouped.plot(kind='bar', color=color, edgecolor='
black', alpha=0.7, logy=log_scale)
23
24     bin_labels = []
25     for interval in grouped.index.categories:
26         left = int(interval.left)
27         right = int(interval.right)
28         bin_labels.append(f"{left}-{right-1}" if right-left
> 1 else str(left))
29
30     plt.xticks(range(len(bin_labels)), bin_labels, rotation
=45, ha='right')
31     plt.title(f"{title}\nTotal valores únicos: {
format_count(total_values)}", fontsize=14, pad=20)
32     plt.xlabel("Rango de repeticiones", fontsize=12)
33     plt.ylabel("Cantidad de valores únicos" + (" (log)" if
log_scale else ""), fontsize=12)
34     plt.grid(True, which="both", ls="--", axis='y')
35
36     min_y = 0.9
37     for i, (count, interval) in enumerate(zip(grouped.
values, grouped.index)):
38         if count > 0:
39             percentage = (count / total_values) * 100
40             y_pos = count * 1.1 if count * 1.1 > min_y else
min_y * 1.2
41             text = f"{percentage:.2f}%\n({format_count(
```



```
count)}})"
42
43     ax.text(i, y_pos, text,
44             ha='center', va='bottom',
45             fontsize=15, fontweight='bold',
46             bbox=dict(facecolor='white', alpha=0.8,
47                       edgecolor='lightgray', boxstyle='round',
48                       pad=0.2'))
49
50     output_path = os.path.join("img", filename)
51     plt.tight_layout()
52     plt.savefig(output_path, dpi=300, bbox_inches='tight')
53     plt.close()
54     return output_path
55
56 def main():
57     if len(sys.argv) < 2:
58         print("Error: Debe especificar un archivo CSV")
59         sys.exit(1)
60
61     csv_file = sys.argv[1]
62     filename_base = os.path.splitext(os.path.basename(
63         csv_file))[0]
64     column = "identifier"
65     chunksize = 1_000_000
66     os.makedirs("img", exist_ok=True)
67
68     print(f"\nIniciando análisis de {csv_file}")
69     print(f"Columna analizada: {column}")
70
71     print("\nContando frecuencias...")
72     counter = Counter()
73
74     with tqdm(desc="Contando filas totales", unit='filas') as pbar:
75         total_rows = 0
76         for chunk in pd.read_csv(csv_file, usecols=[column],
77                                 chunksize=chunksize):
78             total_rows += len(chunk)
79             pbar.update(len(chunk))
80
81     with tqdm(total=total_rows, desc="Procesando datos",
82             unit='filas') as pbar:
83         for chunk in pd.read_csv(csv_file, usecols=[column],
84                                 chunksize=chunksize):
85             counter.update(chunk[column].dropna().astype(
86                 str))
87             pbar.update(len(chunk))
88
89     frequencies = pd.Series(counter)
90     total_unique = len(frequencies)
```

```

84         print(f"\nDatos procesados - Total valores unicos: {
            format_count(total_unique)}")
85
86     print("\nClasificando frecuencias...")
87     with tqdm(total=4, desc="Progreso") as pbar:
88         low_freq = frequencies[(frequencies >= 1) & (
            frequencies <= 99)]
89         pbar.update(1)
90         mid_freq = frequencies[(frequencies >= 100) & (
            frequencies <= 1000)]
91         pbar.update(1)
92         high_freq = frequencies[(frequencies >= 1001) & (
            frequencies <= 10000)]
93         pbar.update(1)
94
95     low_bin = list(range(1, 100, 10)) + [100]
96     mid_bin = list(range(100, 1001, 100)) + [1001]
97     high_bin = list(range(1001, 10001, 1000)) + [10001]
98
99     print("\n===Resumen de frecuencias===")
100    print(f"\nRango 1-99 repeticiones:")
101    print(f"Valores unicos: {format_count(len(low_freq))} ({len(low_freq)/total_unique:.1%}")
102
103    print(f"\nRango 100-1000 repeticiones:")
104    print(f"Valores unicos: {format_count(len(mid_freq))} ({len(mid_freq)/total_unique:.1%}")
105
106    print(f"\nRango 1001-10000 repeticiones:")
107    print(f"Valores unicos: {format_count(len(
        high_freq))} ({len(high_freq)/total_unique:.1%}")
108
109
110    print("\nGenerando graficos...")
111    with tqdm(total=3, desc="Progreso") as pbar:
112        low_path = create_histogram(
113            low_freq,
114            bins=low_bin,
115            title="Distribucion de Frecuencias (1-99 repeticiones)",
116            filename=f"histograma_1-99_{column}_{
                filename_base}.png",
117            color='#4C72B0'
118        )
119        pbar.update(1)
120
121        mid_path = create_histogram(
122            mid_freq,
123            bins=mid_bin,
124            title="Distribucion de Frecuencias (100-1000 repeticiones)",

```

```
125         filename=f"histograma_100-1k_{column}_{
126             filename_base}.png",
127         color='#55A868',
128         log_scale=True
129     )
130     pbar.update(1)
131     high_path = create_histogram(
132         high_freq,
133         bins=high_bin,
134         title="Distribucion de Frecuencias (1001-10,000
135             repeticiones)",
136         filename=f"histograma_1k-10k_{column}_{
137             filename_base}.png",
138         color='#C44E52',
139         log_scale=True
140     )
141     pbar.update(1)
142     print("\nGraficos generados exitosamente:")
143     print(f"{low_path}")
144     print(f"{mid_path}")
145     print(f"{high_path}")
146     if __name__ == "__main__":
147         main()
```

Código B.7: `identfier_histogram_detailed.py`, análisis de frecuencias de la columna 'identfier'.

```
1  import dask.dataframe as dd
2  import sys
3  import os
4
5  def delete_duplicates(input_file, output_file):
6
7      ddf = dd.read_csv(input_file)
8
9      print(f"\nProcesando archivo: {input_file}")
10     print(f"Numero inicial de registros: {len(ddf):,}")
11
12     ddf_deduplicate = ddf.drop_duplicates(
13         subset=['timestamp', 'device_lon', 'device_lat'],
14         keep='first'
15     )
16
17     print(f"Numero de registros despues de eliminar
18         duplicados: {len(ddf_deduplicate):,}")
19
20     ddf_deduplicate.to_csv(
21         output_file,
22         index=False,
23         single_file=True
24     )
25
26     print(f"\nArchivo sin duplicados guardado en: {
27         output_file}")
28
29     if __name__ == "__main__":
30         if len(sys.argv) < 2:
31             print("Error: Debe especificar un archivo CSV como
32                 argumento")
33             sys.exit(1)
34
35         input_csv = sys.argv[1]
36         base_name = os.path.splitext(input_csv)[0]
37         output_csv = f"{base_name}_DeDuplicate.csv"
38
39         delete_duplicates(input_csv, output_csv)
```

Código B.8: csv_deduplicate.py, eliminación de duplicados en el conjunto de datos.

Referencias

- [1] Autor referencia 1
- [2] Autor referencia 2
- [3] Autor referencia 3