

Universidad Nacional de San Agustín de Arequipa

AÑO DE LA RECUPERACIÓN Y CONSOLIDACIÓN DE LA  
ECONOMÍA PERUANA



ESCUELA PROFESIONAL DE CIENCIA DE LA COMPUTACIÓN  
Análisis Exploratorio de Datos -  
Data Wrangling

---

## INFORME DATA WRANGLING

A-2025

---

**Docente: Dra. Ana Maria CuadrosValdivia**

Javier Wilber Quispe Rojas

# Índice

Índice	1
<b>1 Contexto de mi DataSet</b>	<b>2</b>
1.1 Boston_feature_df: . . . . .	2
1.2 BostonMobility2021: . . . . .	2
<b>2 Descripción del data set</b>	<b>2</b>
2.1 Cuantos Registros tiene cada DataSet . . . . .	2
2.2 Descripcion DataSet Boston_feature_df.csv . . . . .	2
2.3 Descripción de las columnas del dataset BostonMobility2021.csv . . . . .	8
2.4 ¿Que contiene cada Registro . . . . .	11
2.4.1 Descripción de un Registro del dataset Boston_feature_df . . . . .	11

# 1. Contexto de mi DataSet

## 1.1. Boston\_feature\_df:

Este dataset contiene información sobre diferentes áreas de Boston, identificadas por códigos geográficos. Para cada área, muestra la proporción de diferentes grupos raciales como blancos, negros, asiáticos e hispanos, tanto de la población residente como de quienes llegan a esas zonas. También incluye datos sobre los ingresos de las personas en cada área, divididos en rangos económicos, y cómo estos ingresos se distribuyen entre quienes se trasladan hacia esas zonas. Además, registra la intensidad con la que la gente visita distintos tipos de lugares, como tiendas, restaurantes, hospitales o sitios religiosos. Cada área está geolocalizada con latitud y longitud, y se especifica la población total que vive ahí.

## 1.2. BostonMobility2021:

Este archivo muestra los movimientos de personas dentro de Boston durante 2021, indicando desde dónde salen y hacia dónde se dirigen. Cada fila representa un flujo de personas entre dos áreas geográficas, con las coordenadas de origen y destino, además del número de visitantes que se trasladaron en un período específico, generalmente semanal. También incluye estimaciones de la población en movimiento y, en algunos casos, datos sobre la cantidad de dispositivos móviles activos en el área durante el día. Este dataset es mucho más grande y permite analizar con detalle los patrones de movilidad dentro de la ciudad.

# 2. Descripción del data set

## 2.1. Cuantos Registros tiene cada DataSet

```
Boston_feature_df.csv:
Un registro representa un área pequeña (CBG) con datos demográficos, socioeconómicos y puntos de interés.
Número de registros: 462

BostonMobility2021.csv:
Un registro representa un flujo de movilidad humana entre dos áreas (origen y destino) en un rango temporal.
Número de registros: 236530
```

Figura 1: Cantidad de Registros por Data set

### Boston\_feature\_df.csv:

Este dataset tiene 462 registros, lo que es bastante manejable para casi cualquier computadora actual. No es un volumen de datos grande, por lo que se puede procesar fácilmente en memoria sin necesidad de herramientas especiales o grandes recursos de CPU y RAM. Ideal para análisis rápidos y exploratorios.

### BostonMobility2021.csv:

Con 236,530 registros, este dataset es considerablemente más grande, pero sigue siendo posible de manejar en la mayoría de las computadoras personales modernas, especialmente si cuentas con al menos 8 GB de RAM. Dependiendo del tipo de análisis, podrías experimentar un poco más de uso de CPU y memoria, pero en general no es un volumen masivo que requiera infraestructura de big data.

## 2.2. Descripción DataSet Boston\_feature\_df.csv

### ¿Cuál es el objeto u entidad de estudio?

El dataset estudia diferentes áreas o sectores geográficos (probablemente census tracts o zonas censales) dentro de la ciudad de Boston. Cada fila representa una unidad espacial que contiene información demográfica, socioeconómica y de actividad social de esa zona.

Cuadro 1: Identificación y composición racial

Atributo	Significado / Descripción	Tipo de dato	Rango aproximado / Valores
GEOID	Identificador único del área geográfica (sector o tract)	Categorico (texto)	Código único por zona
white	Proporción de población blanca en el área	Numérico decimal	0 a 1
black	Proporción de población negra en el área	Numérico decimal	0 a 1
asian	Proporción de población asiática en el área	Numérico decimal	0 a 1
hispanic	Proporción de población hispana en el área	Numérico decimal	0 a 1
white_inflow	Proporción de población blanca que ingresa a esa área	Numérico decimal	0 a 1
black_inflow	Proporción de población negra que ingresa a esa área	Numérico decimal	0 a 1
hispanic_inflow	Proporción de población hispana que ingresa a esa área	Numérico decimal	0 a 1
asian_inflow	Proporción de población asiática que ingresa a esa área	Numérico decimal	0 a 1

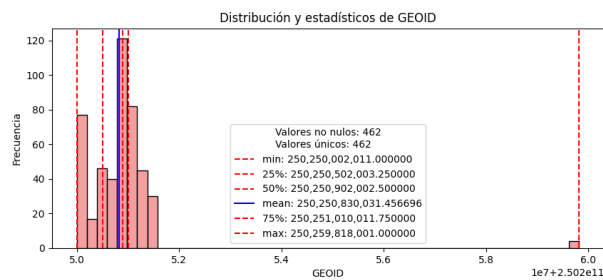


Figura 2: GEOID

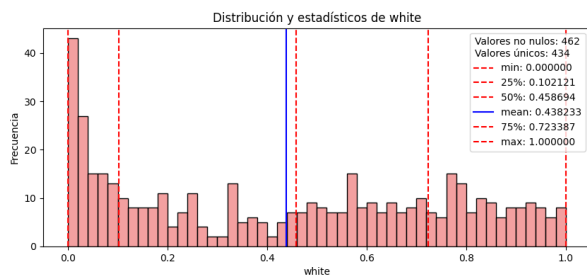


Figura 3: white

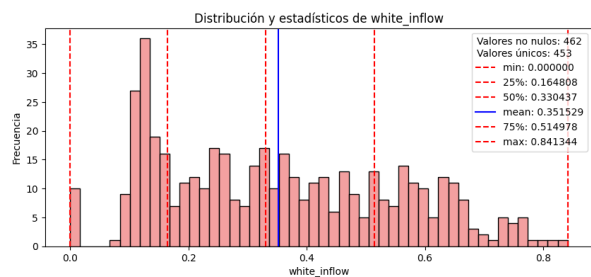


Figura 4: white\_inflow

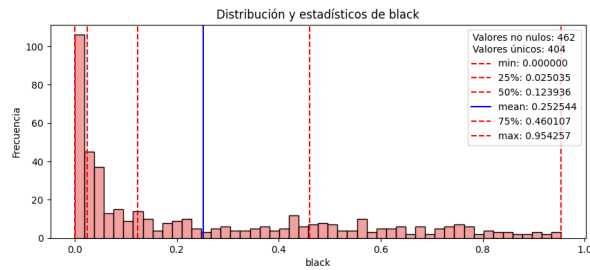


Figura 5: black

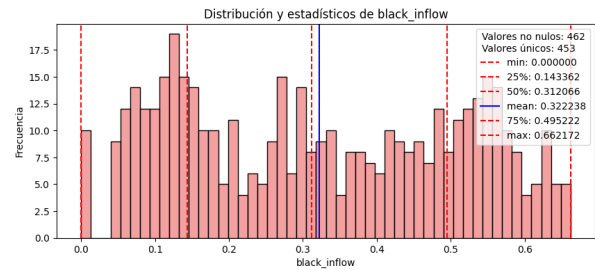


Figura 6: black\_inflow

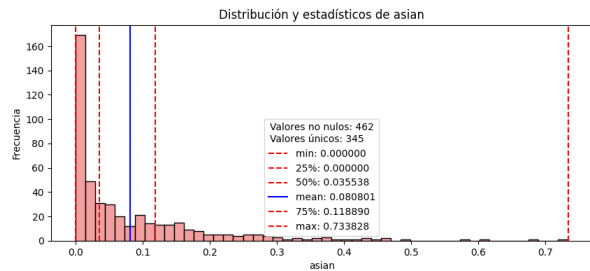


Figura 7: asian

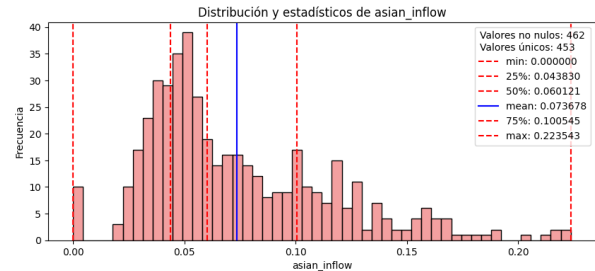


Figura 8: asian\_inflow

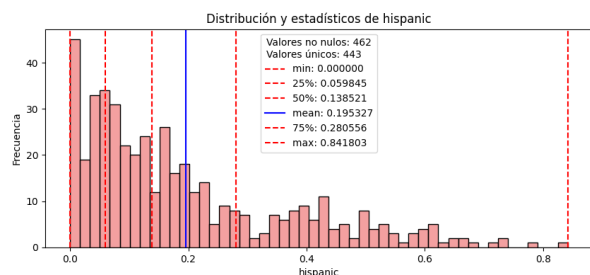


Figura 9: hispanic

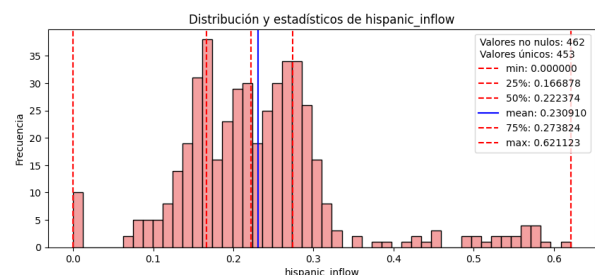


Figura 10: hispanic\_inflow

Cuadro 2: Ingreso y flujo de ingreso

Atributo	Significado / Descripción	Tipo de dato	Rango aproximado / Valores
Under \$50K	Proporción de población con ingresos menores a \$50,000	Númérico decimal	0 a 1
\$50K - \$100K	Proporción con ingresos entre \$50,000 y \$100,000	Númérico decimal	0 a 1
\$100K - \$200K	Proporción con ingresos entre \$100,000 y \$200,000	Númérico decimal	0 a 1
Over \$200K	Proporción con ingresos mayores a \$200,000	Númérico decimal	0 a 1
Under \$50K_inflow	Proporción con ingresos menores a \$50,000 que ingresan a esa área	Númérico decimal	0 a 1
\$50K - \$100K_inflow	Proporción con ingresos entre \$50,000 y \$100,000 que ingresan a esa área	Númérico decimal	0 a 1
\$100K - \$200K_inflow	Proporción con ingresos entre \$100,000 y \$200,000 que ingresan a esa área	Númérico decimal	0 a 1
Over \$200K_inflow	Proporción con ingresos mayores a \$200,000 que ingresan a esa área	Númérico decimal	0 a 1

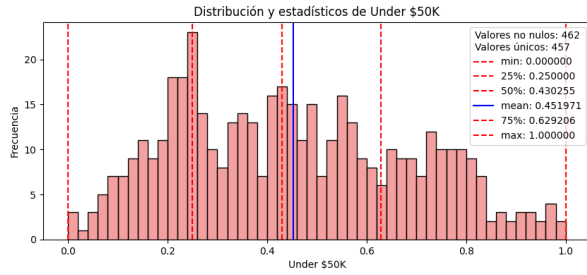


Figura 11: Under \$50K

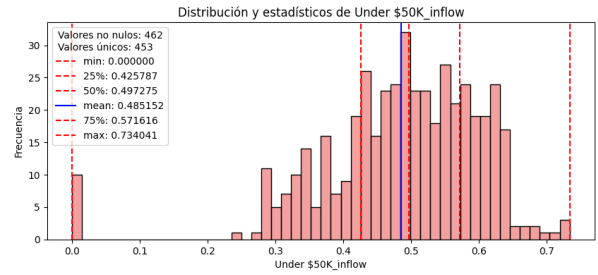


Figura 12: Under \$50K\_inflow

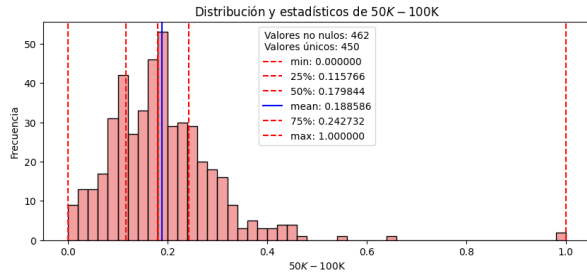


Figura 13: \$50K - \$100K

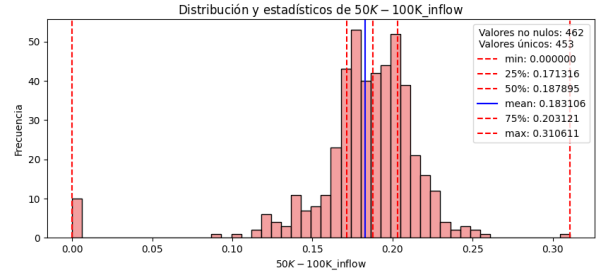


Figura 14: \$50K - \$100K\_inflow

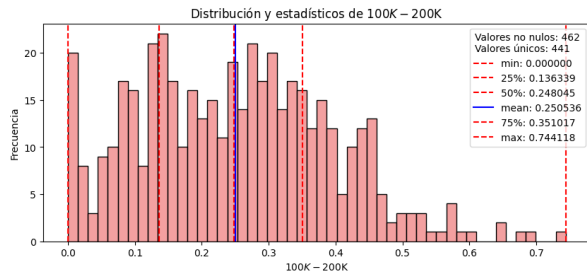


Figura 15: \$100K - \$200K

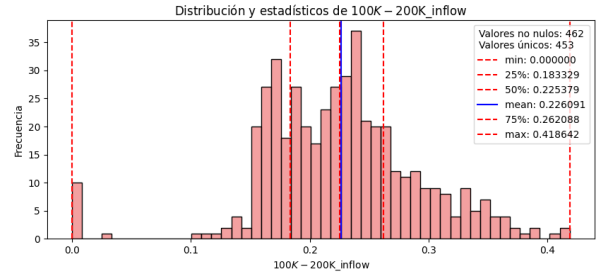


Figura 16: \$100K - \$200K\_inflow

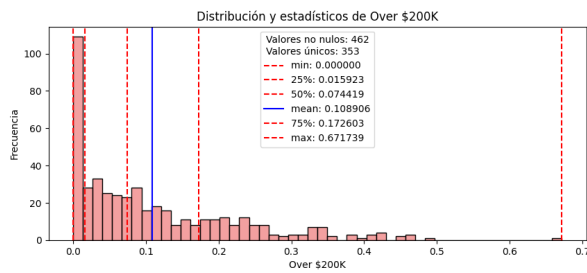


Figura 17: Over \$200K

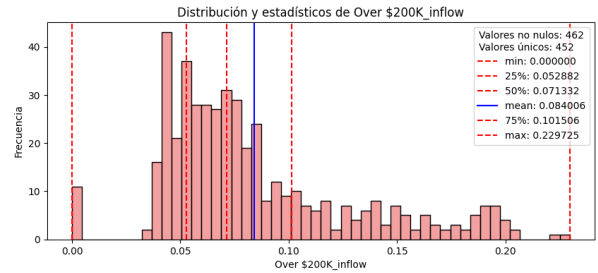


Figura 18: Over \$200K\_inflow

Cuadro 3: Actividades y geolocalización

Atributo	Significado / Descripción	Tipo de dato	Rango aproximado / Valores
Food	Nivel o cantidad de actividad relacionada con comida en el área	Numérico decimal	Valores positivos (escala no definida)
Shopping	Nivel o cantidad de actividad de compras en el área	Numérico decimal	Valores positivos (escala no definida)
Work	Nivel o cantidad de actividad laboral en el área	Numérico decimal	Valores positivos (escala no definida)
Health	Nivel o cantidad de actividad relacionada con salud en el área	Numérico decimal	Valores positivos (escala no definida)
Religious	Nivel o cantidad de actividad religiosa en el área	Numérico decimal	Valores positivos (escala no definida)
Service	Nivel o cantidad de actividad de servicios en el área	Numérico decimal	Valores positivos (escala no definida)
Entertainment	Nivel o cantidad de actividad de entretenimiento en el área	Numérico decimal	Valores positivos (escala no definida)
Grocery	Nivel o cantidad de actividad de supermercados/grocery en el área	Numérico decimal	Valores positivos (escala no definida)
Education	Nivel o cantidad de actividad educativa en el área	Numérico decimal	Valores positivos (escala no definida)
Arts/Museum	Nivel o cantidad de actividad cultural y museos en el área	Numérico decimal	Valores positivos (escala no definida)
Transportation	Nivel o cantidad de actividad de transporte en el área	Numérico decimal	Valores positivos (escala no definida)
Sports	Nivel o cantidad de actividad deportiva en el área	Numérico decimal	Valores positivos (escala no definida)
LATITUDE	Latitud geográfica del centro del área	Numérico decimal	Aproximadamente 42.x (Boston)
LONGITUDE	Longitud geográfica del centro del área	Numérico decimal	Aproximadamente -71.x (Boston)
total_population	Número total estimado de habitantes en el área	Numérico entero	Desde decenas hasta varios miles

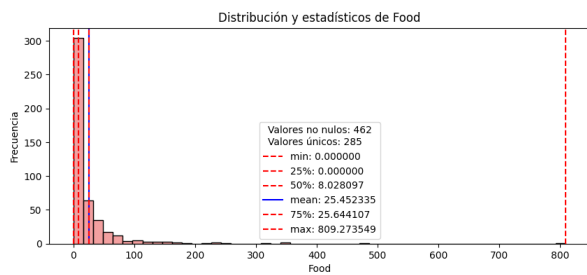


Figura 19: Food



Figura 20: Shopping

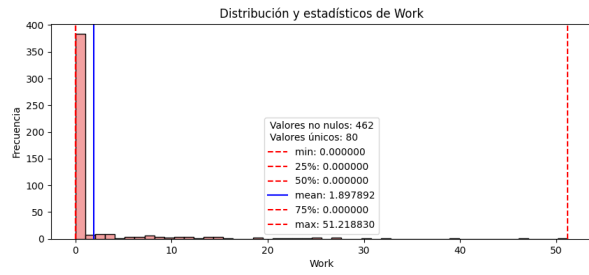


Figura 21: Work

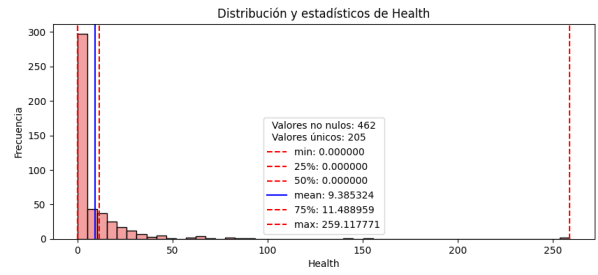


Figura 22: Health

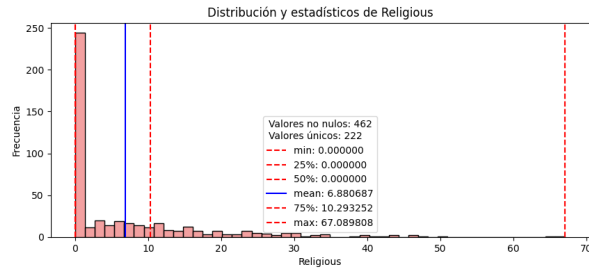


Figura 23: Religious

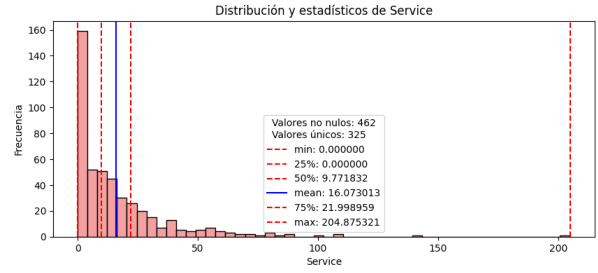


Figura 24: Service

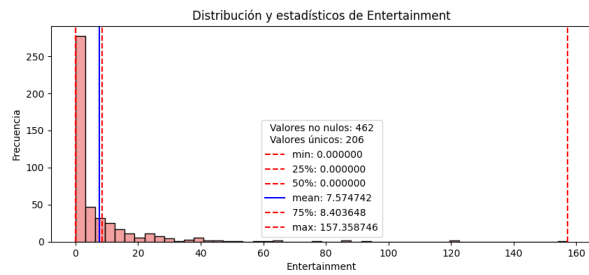


Figura 25: Entertainment

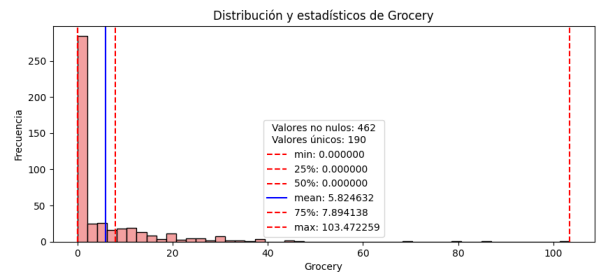


Figura 26: Grocery

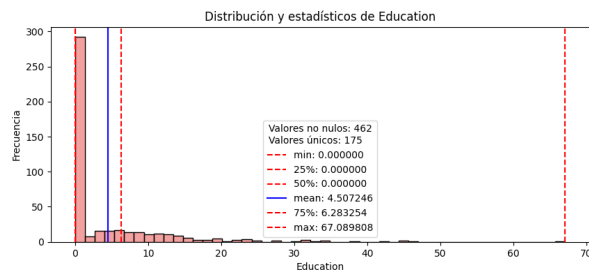


Figura 27: Education

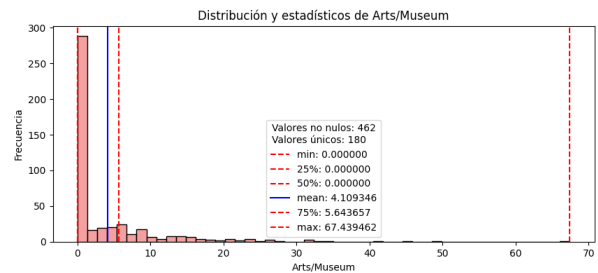


Figura 28: Arts\_Museum



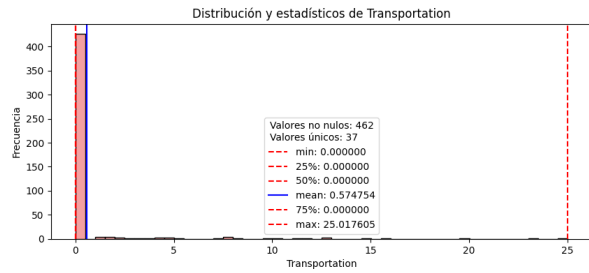


Figura 29: Transportation

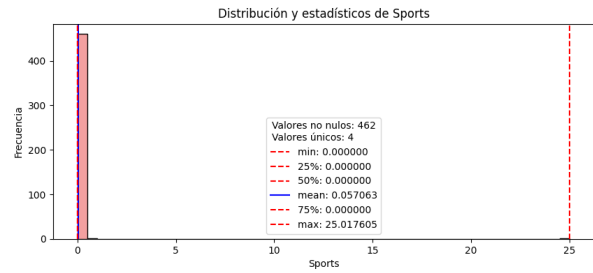


Figura 30: Sports

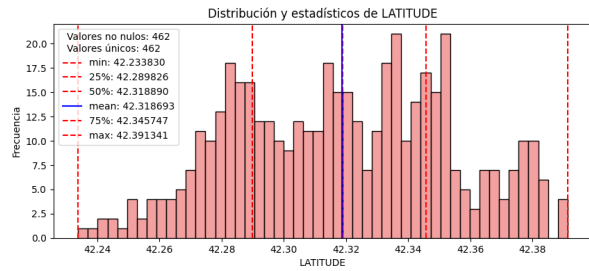


Figura 31: LATITUDE

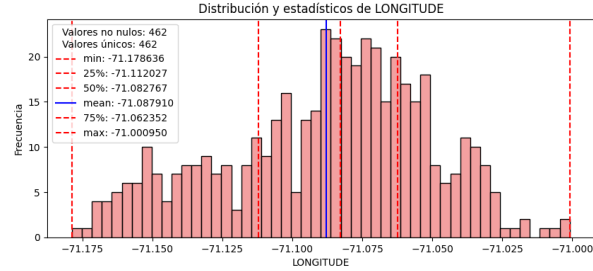


Figura 32: LONGITUDE

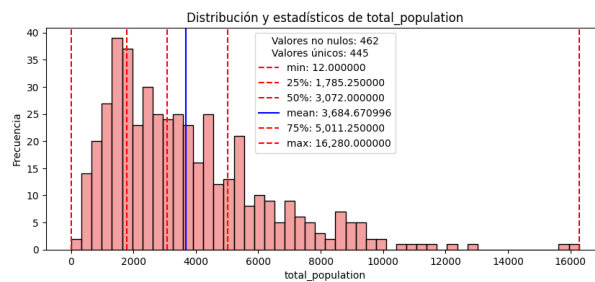


Figura 33: total\_population

## 2.3. Descripción de las columnas del dataset BostonMobility2021.csv

### ¿Cuál es el objeto u entidad de estudio?

El objeto de estudio son los flujos de movilidad de personas entre diferentes áreas geográficas dentro de la ciudad de Boston. Es decir, cada registro representa la cantidad de personas que se desplazan desde una zona de origen hacia una zona de destino en un rango de fechas determinado.

Cuadro 4: Resumen del objeto de estudio y atributos del dataset BostonMobility2021

Atributo	Significado / Descripción	Tipo de dato	Rango aproximado / Valores
geoid_o	Identificador del área de origen del flujo	Catégorico (texto)	Código único por zona
geoid_d	Identificador del área de destino del flujo	Catégorico (texto)	Código único por zona
lng_o	Longitud del área de origen	Numérico decimal	Aproximadamente -71.x (Boston)
lat_o	Latitud del área de origen	Numérico decimal	Aproximadamente 42.x (Boston)
lng_d	Longitud del área de destino	Numérico decimal	Aproximadamente -71.x (Boston)
lat_d	Latitud del área de destino	Numérico decimal	Aproximadamente 42.x (Boston)
date_range	Rango de fechas del registro de movilidad	Catégorico (texto)	Ejemplo: "01/04/21 - 01/10/21"
visitor_flows	Número de visitantes entre áreas en el rango de fechas	Numérico entero	Desde 0 hasta valores variables
pop_flows	Población total estimada asociada al flujo o área	Numérico decimal	Valores positivos variados
census_block	Agrupación censal (datos mayormente ausentes)	Desconocido / vacío	Mayormente NaN
number_devices	Número estimado de dispositivos activos en área origen durante el día	Numérico entero o NaN	Valores numéricos o vacíos

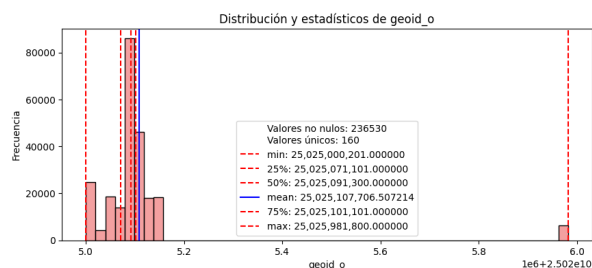


Figura 34: geoid origen

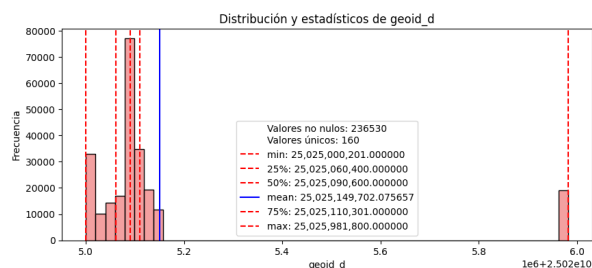


Figura 35: geoid destino

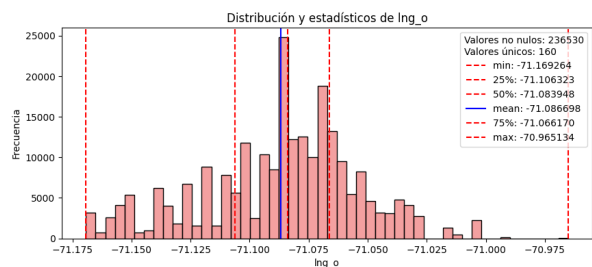


Figura 36: LONGITUDE Origen

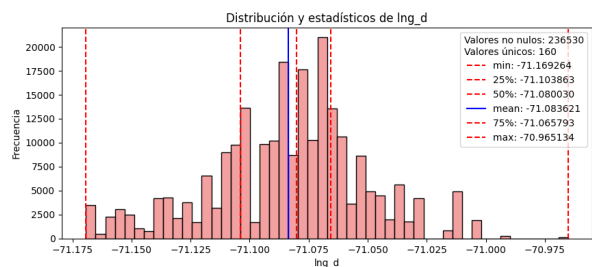


Figura 37: LONGITUDE Destino

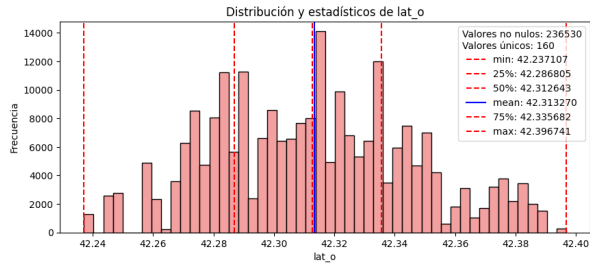


Figura 38: LATITUDE Origen

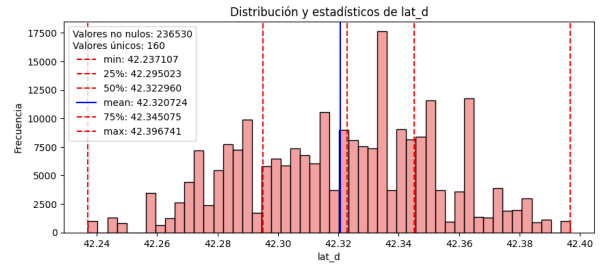


Figura 39: LATITUDE Destino

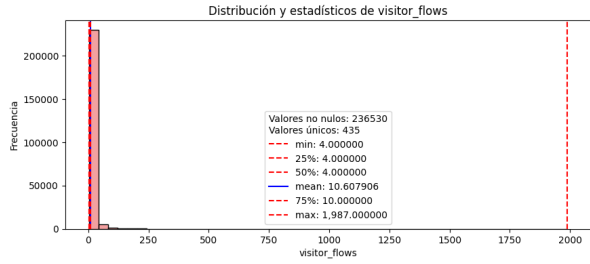


Figura 40: Flujo de persona

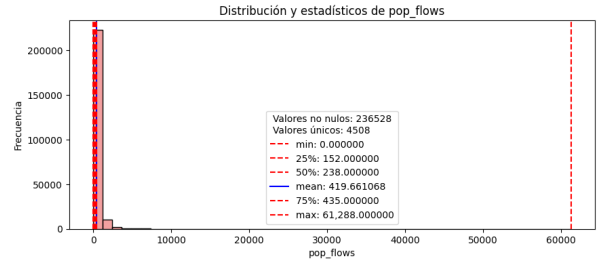


Figura 41: Población total

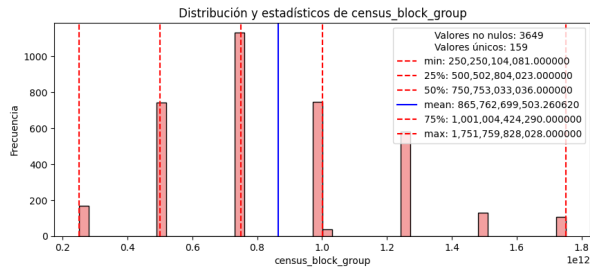


Figura 42: Agrupación censal

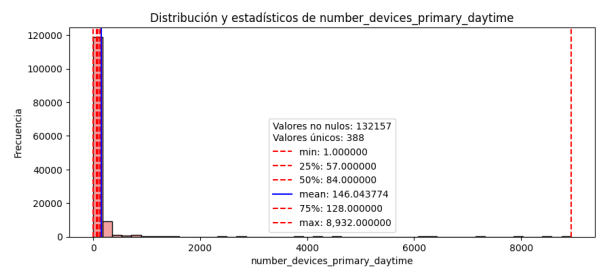


Figura 43: dispositivos activos

## 2.4. ¿Que contiene cada Registro

### 2.4.1. Descripción de un Registro del dataset `Boston.feature_df`

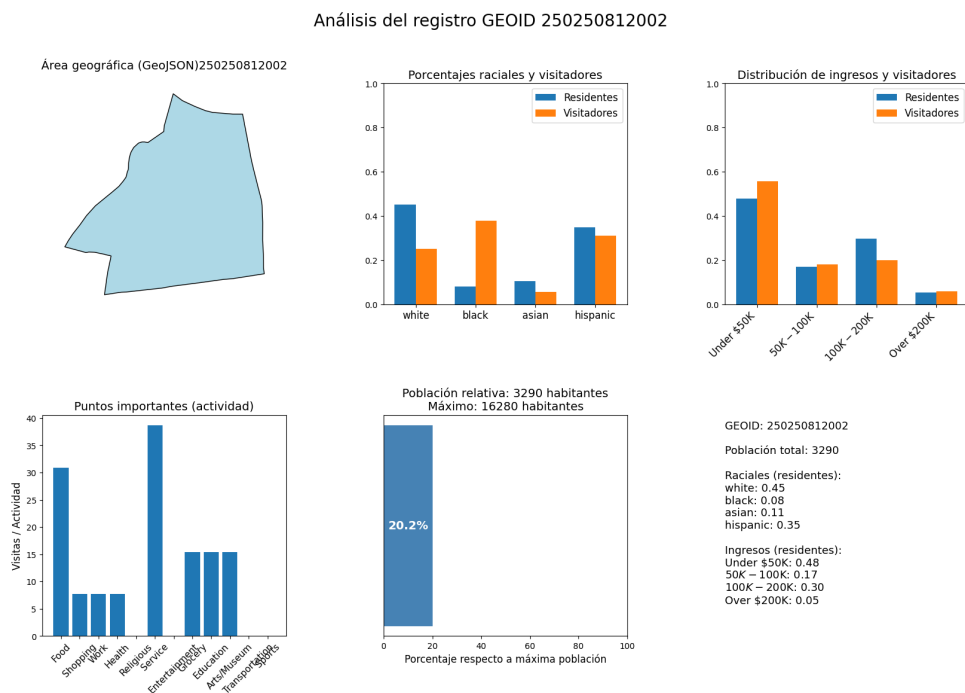


Figura 44: Registro Geoid:250250812002

Cada fila en el dataset representa un área geográfica pequeña dentro de Boston, identificada por un código único llamado **GEOID**. Este código corresponde a un bloque censal o grupo de bloques con características propias.

Tomando como ejemplo el registro con **GEOID = 250250812002**, podemos describir la fila de la siguiente manera:

- **Ubicación geográfica:** La fila corresponde a un polígono delimitado en el mapa que define el área específica de ese bloque censal.
- **Composición racial y visitantes:** Se presentan las proporciones de la población residente dividida en grupos raciales: *white*, *black*, *asian* y *hispanic*, así como la composición de los visitantes que llegan a esa área. Esto refleja la diversidad y dinámica demográfica de la zona.
- **Distribución de ingresos y visitantes:** La información incluye la distribución de ingresos de los residentes en cuatro rangos económicos y la composición similar de los visitantes, brindando un panorama socioeconómico detallado.
- **Puntos importantes y actividad:** Se muestran las intensidades de visitas o actividades en diferentes tipos de lugares, como comida, tiendas, trabajo, salud, servicios y entretenimiento, que reflejan el uso y atractivo de la zona.
- **Población relativa:** Se presenta como la población total de esta área.