

Análisis Exploratorio de Datos - Data Wrangling

Tópicos en Ciencia de Datos

Integrantes: Kemely Francis Castillo Caccire

Docente: Ana Maria Cuadros Valdivia

Fecha de entrega: 30 de mayo del 2025

Arequipa, Perú

1. Análisis del comportamiento de los datos

Imagina que una ciudad quiere contarle al mundo cómo funciona su sistema de transporte público: por dónde pasan los buses, a qué hora llegan, dónde paran y qué rutas existen. Para eso existe GTFS, un formato de datos que organiza toda esa información de forma clara, ordenada y fácil de usar por computadoras.

Un archivo GTFS es como una colección de hojas de Excel: una con las paradas, otra con las rutas, otra con los horarios, y así sucesivamente. Cada hoja tiene datos que se relacionan entre sí, como si fueran piezas de un rompecabezas.

En este trabajo nos centramos en GTFS Static(“instantánea” periódica sin datos en tiempo real). Este enfoque sigue lo propuesto por G2Viz [1] como una herramienta para analizar visualmente datos GTFS estáticos.

1.1. Descripción del los datos

Dataset: [muni_gtfs-current.zip](#)

El dataset utilizado contiene datos GTFS actualizados de la **San Francisco Municipal Transportation Agency (SFMTA)**[3]. Este conjunto incluye archivos como `routes.txt`, `stops.txt`, `trips.txt`, `stop_times.txt`, entre otros. Cada archivo representa una entidad o relación específica dentro del sistema de transporte, y su estructura permite consultas cruzadas para responder preguntas sobre operación y cobertura del servicio.

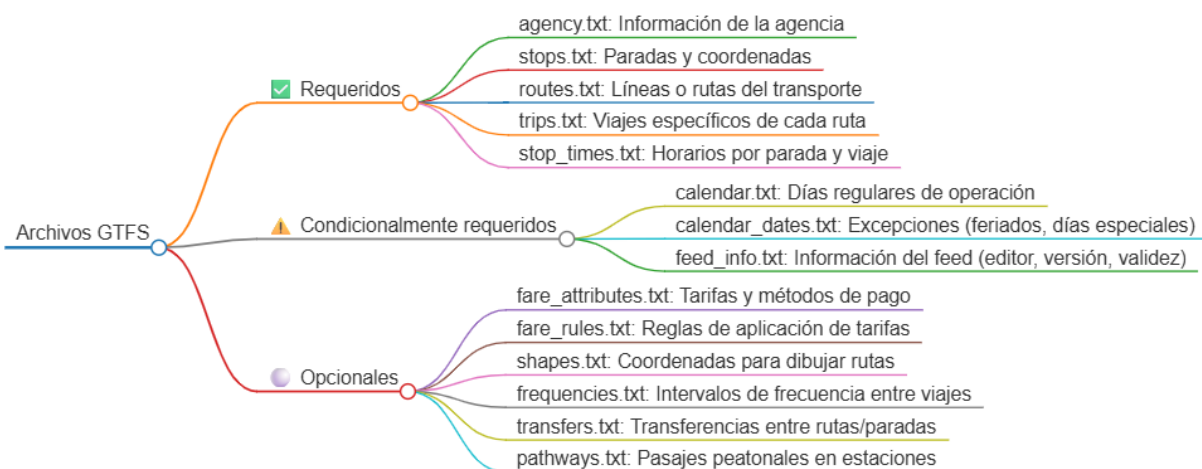
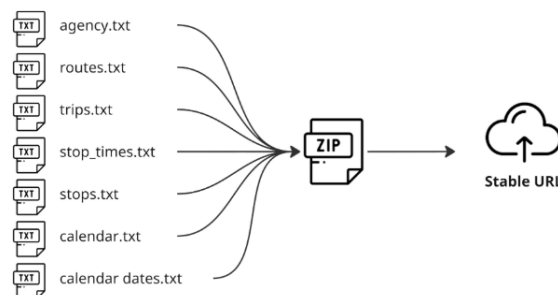


Figura 1: Archivos más comunes incluidos en un feed GTFS.

Este mapa ayuda a entender la estructura GTFS, inspirada en el artículo G2Viz.

1.2. GTFS SFMTA

Archivos: Codificación, Tamaño y Granularidad

Archivo	Codificación	Tamaño (MB)	Granularidad
agency.txt	ascii	0.000000	1 filas x 8 columnas
stops.txt	ascii	0.272000	3278 filas x 8 columnas
routes.txt	ascii	0.006000	69 filas x 10 columnas
trips.txt	ascii	1.885000	35017 filas x 9 columnas
stop_times.txt	ascii	54.836000	1313881 filas x 9 columnas
calendar.txt	ascii	0.000000	3 filas x 10 columnas
calendar_dates.txt	ascii	0.000000	26 filas x 3 columnas
fare_attributes.txt	ascii	0.000000	2 filas x 7 columnas
fare_rules.txt	ascii	0.000000	69 filas x 2 columnas
shapes.txt	ascii	1.798000	46249 filas x 5 columnas
directions.txt	ascii	0.002000	138 filas x 3 columnas
timepoints.txt	ascii	4.527000	249904 filas x 2 columnas

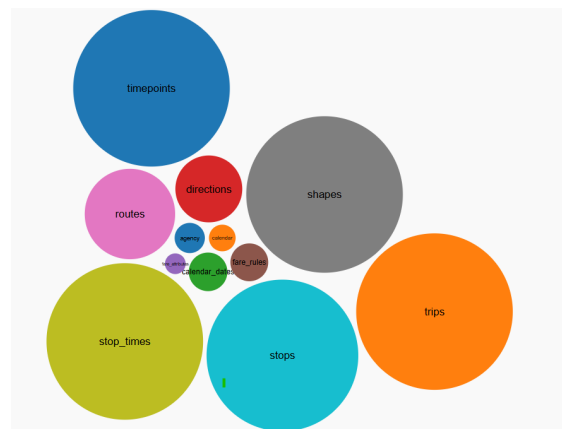


Figura 2: Visualización de archivos GTFS

- El conjunto de datos GTFS tiene 12 archivos con diferentes tamaños, codificación y granularidad, reflejando la complejidad del transporte público.
- Los archivos `stop_times.txt` y `trips.txt` son los más grandes y con mayor granularidad, mostrando detalles finos de horarios y viajes.
- Archivos como `calendar.txt` y `fare_rules.txt` son pequeños y cumplen funciones auxiliares.
- La mayoría de los archivos usan codificación ASCII, facilitando su lectura y manipulación.
- La granularidad varía desde pocos registros con muchas columnas hasta muchos registros con pocas columnas.

Modelo Entidad-Relación

Para facilitar la comprensión de cómo estos archivos se interconectan, se ha elaborado un diagrama Entidad-Relación (ER) que muestra las entidades principales, sus atributos clave y las relaciones entre ellas.

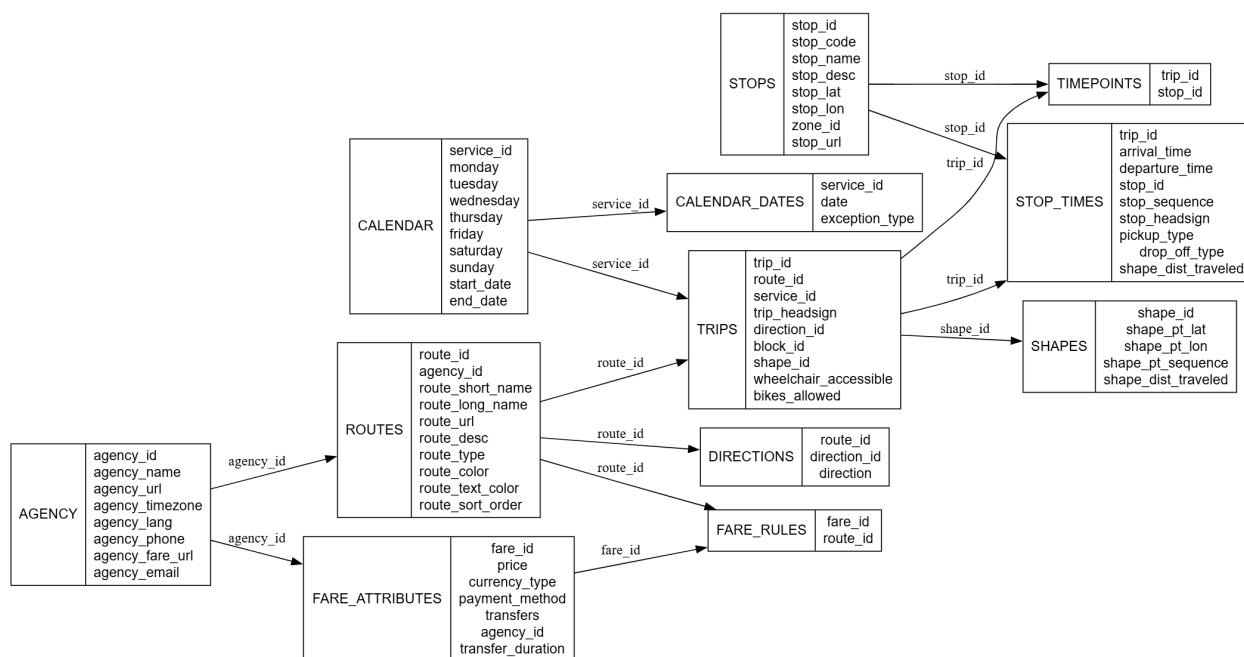


Figura 3: Modelo relacional de archivos de texto incluidos en el GTFS SFMTA

Analizamos las tablas en un `:colab` y sus contenidos en orden a su importancia.

Requeridos

▪ agency

	agency_id	agency_name	agency_url	agency_timezone	agency_lang	agency_phone	agency_fare_url	agency_email
0	SFMTA	San Francisco Municipal Transportation Agency	http://www.sfmta.com	America/Los_Angeles	en	311	https://SFMTA.com/Fares	munifeedback@sfmta.com

Figura 4: 3 primeros registros de la tabla agency

- **Una ciudad puede tener una o varias agencias.**
El archivo puede listar múltiples agencias. Por ejemplo, San Francisco usa solo SFMTA, pero ciudades como Arequipa en nuestro caso pueden tener varias hasta 10 [2].
- **La zona horaria es clave.**
El campo `agency_timezone` asegura que los horarios se muestren correctamente en apps. Si está mal configurado, los horarios se desajustan [5].
- **El número de teléfono no siempre es internacional.**
Por ejemplo, San Francisco usa "311", un número local. No todos los feeds usan formato internacional, lo cual puede generar confusión [6].
- **El idioma afecta cómo se ve todo.**
El campo `agency_lang` define el idioma predeterminado del feed, lo que influye en cómo se muestra la información en las apps [7].

- **No todos los campos son obligatorios.**

Campos como `agency_email` o `agency_fare_url` son opcionales, pero útiles para el contacto o información adicional [8].

- **stops**

--- stops.txt: primeros 3 registros ---

	stop_id	stop_code	stop_name	stop_desc	stop_lat	stop_lon	zone_id	stop_url
0	390	10390	19th Avenue & Holloway St		37.721190	-122.475153		https://SFMTA.com/10390
1	913	10913	Dublin St & La Grande Ave		37.719192	-122.425802		https://SFMTA.com/10913
2	3016	13016	3rd St & 4th St		37.772618	-122.389786		https://SFMTA.com/13016
3	3018	13018	Bacon St & San Bruno Ave		37.727859	-122.402994		https://SFMTA.com/13018
4	3019	13019	Bacon St & San Bruno Ave		37.727645	-122.403269		https://SFMTA.com/13019
5	3020	13020	Bacon St & Somerset St		37.726670	-122.407460		https://SFMTA.com/13020
6	3021	13021	Bacon St & Somerset St		37.726540	-122.407660		https://SFMTA.com/13021
7	3023	13023	Baker Beach Parking Lot SW		37.790249	-122.482263		https://SFMTA.com/13023
8	3024	13024	Baker St & Greenwich St		37.797606	-122.445693		https://SFMTA.com/13024
9	3031	13031	Bay St & Midway St		37.806094	-122.409161		https://SFMTA.com/13031

Figura 5: 3 primeros registros de la tabla stops

- Un mismo nombre de parada puede tener varios `stop_id` diferentes. Por ejemplo, *Bacon St & San Bruno Ave* aparece con varios `stop_id` distintos (3018 y 3019). Esto ocurre porque cada lado de la calle o sentido puede tener una parada distinta, aunque el nombre sea igual.
- `stop_code` puede ser distinto de `stop_id` y funciona como código para usuarios. `stop_code` suele ser un número más corto o amigable que los pasajeros pueden ver en señales o apps, mientras que `stop_id` es un identificador único para el sistema.
- Algunas paradas no tienen descripción (`stop_desc`) pero sí URL con más información. Aunque el campo `stop_desc` esté vacío, el campo `stop_url` ofrece un enlace donde se puede obtener más datos o mapas, lo que ayuda a usuarios que quieren detalles extras.
- Las coordenadas geográficas (`stop_lat` y `stop_lon`) no solo marcan la ubicación, sino que permiten calcular distancias y rutas automáticamente. Esto facilita la generación automática de mapas, cálculo de distancias a pie o la navegación GPS para apps.
- La ausencia de `zone_id` indica que el sistema no usa zonas tarifarias en algunas paradas. Si bien `zone_id` sirve para definir tarifas basadas en zonas geográficas, en este caso muchas paradas no tienen zona, probablemente porque la tarifa es fija o la ciudad no usa zonas para precios.
- Las paradas pueden ser puntos muy cercanos, incluso a pocos metros, para diferentes direcciones o rutas. Esto pasa mucho en cruces grandes, donde cada dirección de viaje tiene su parada propia, a veces justo en lados opuestos de la calle.
- El sistema permite representar desde paradas simples en la calle hasta lugares complejos como estacionamientos o terminales. Por ejemplo, *Baker Beach Parking Lot SW* no es

una parada normal sino un área específica para transporte, lo que muestra la flexibilidad del formato para distintos tipos de puntos de parada.

■ routes

--- routes.txt: primeros 3 registros ---

	route_id	agency_id	route_short_name	route_long_name	route_url	route_desc	route_type	route_color	route_text_color	route_sort_order
0	714	SFMTA	714	BART EARLY BIRD	http://www.sfmta.com/714	Weekdays 4am-5am	3	005B95	FFFFFF	NaN
1	8	SFMTA	8	BAYSHORE	http://www.sfmta.com/8	5am-12 midnight daily	3	005B95	FFFFFF	NaN
2	8AX	SFMTA	8AX	BAYSHORE A EXPRESS	http://www.sfmta.com/8AX	Weekdays 6:30-9:30am (northbound) 3:30-6:50pm (southbound)	3	005B95	FFFFFF	NaN
3	8BX	SFMTA	8BX	BAYSHORE B EXPRESS	http://www.sfmta.com/8BX	Weekdays 6:30-9am (northbound) 3:30-6:30pm (southbound)	3	005B95	FFFFFF	NaN
4	9	SFMTA	9	SAN BRUNO	http://www.sfmta.com/9	5am-12 midnight daily	3	005B95	FFFFFF	NaN
5	90	SFMTA	90	SAN BRUNO OWL	http://www.sfmta.com/90	12 midnight-5am daily	3	666666	FFFFFF	NaN
6	91	SFMTA	91	3RD-19TH AVE OWL	http://www.sfmta.com/91	12 midnight-5am daily	3	666666	FFFFFF	NaN
7	9R	SFMTA	9R	SAN BRUNO RAPID	http://www.sfmta.com/9R	Weekdays 7am-6pm	3	BF2B45	FFFFFF	NaN
8	CA	SFMTA	CA	CALIFORNIA STREET CABLE CAR	http://www.sfmta.com/CA	7 am-8:30 pm daily	5	B49A36	000000	NaN
9	F	SFMTA	F	MARKET & WHARVES	http://www.sfmta.com/F	7am-10pm daily	0	B49A36	000000	NaN

Figura 6: 3 primeros registros de la tabla routes

- Los IDs como 8, 8AX y 8BX son rutas relacionadas pero diferentes. Aunque comparten el número base (8), representan variaciones con diferentes horarios o paradas. El sufijo (AX, BX, R) indica rutas expresas o rápidas. Esto no siempre está estandarizado, así que las aplicaciones deben interpretar bien estos nombres.
- Hay rutas que sólo operan en horarios extremos. Por ejemplo, 714 BART EARLY BIRD opera sólo de 4 am a 5 am, y 90 SAN BRUNO OWL de 12 am a 5 am. Estas rutas “temprano” o “nocturnas” existen para cubrir la ciudad cuando el sistema principal (como BART) está cerrado.
- El mismo número puede aparecer en una versión diurna y otra nocturna. La ruta 9 es "SAN BRUNO", y 90 es "SAN BRUNO OWL", su equivalente nocturno. No siempre es obvio para el usuario casual que están relacionadas.
- Los colores ayudan a distinguir tipos de servicio. Rutas normales usan azul (005B95), rápidas usan rojo (BF2B45), y nocturnas usan gris (666666). Aunque el GTFS no impone reglas sobre colores, algunas agencias lo usan para comunicar visualmente el tipo de servicio.
- El tipo de transporte es más diverso de lo que parece. El campo `route_type = 3` indica autobuses, pero 5 es para tranvías (como el Cable Car), y 0 indica tranvía histórico (como la ruta F de Market & Wharves). Esta clasificación permite representar muchos modos en una sola estructura.
- Algunas rutas tienen nombres largos y promocionales. Ejemplo: CALIFORNIA STREET CABLE CAR no es un bus regular, sino una atracción turística reconocida, y su nombre lo deja claro. Esto es útil para turistas pero puede confundir al integrarse con rutas regulares.
- El campo `route_sort_order` está vacío, pero podría usarse para ordenar visualmente. Si se completa, permite que apps muestren las rutas en orden lógico, como primero locales, luego expresas, luego nocturnas. Que esté en blanco significa que el orden depende de cada app o queda al azar.

■ trips

--- trips.txt: primeros 3 registros ---

	route_id	service_id	trip_id	trip_headsign	direction_id	block_id	shape_id	wheelchair_accessible	bikes_allowed
0	1	M12	11708585_M12	Geary + 33rd Avenue	0	109_M12	102	1	1
1	1	M12	11708586_M12	Presidio Avenue	0	105_M12	101	1	1
2	1	M12	11708587_M12	Presidio Avenue	0	101_M12	101	1	1
3	1	M12	11708588_M12	Presidio Avenue	0	104_M12	101	1	1
4	1	M12	11708589_M12	Geary + 33rd Avenue	0	102_M12	102	1	1
5	1	M12	11708590_M12	Geary + 33rd Avenue	0	107_M12	102	1	1
6	1	M12	11708591_M12	Geary + 33rd Avenue	0	110_M12	102	1	1
7	1	M12	11708592_M12	Geary + 33rd Avenue	0	103_M12	102	1	1
8	1	M12	11708593_M12	Geary + 33rd Avenue	0	101_M12	102	1	1
9	1	M12	11708594_M12	Geary + 33rd Avenue	0	104_M12	102	1	1

Figura 7: 3 primeros registros de la tabla trips

- Un mismo **route_id** puede tener muchos **trip_id** distintos. Esto significa que una ruta puede tener varios viajes diferentes durante el día, cada uno con su propio identificador para organizar horarios y recorridos.
- El **service_id** define cuándo opera un viaje. Por ejemplo, puede indicar si el viaje es solo en días laborables, fines de semana o días festivos, permitiendo controlar el calendario del servicio.
- El **trip_headsign** es el texto que aparece en el bus o app para mostrar el destino. Aunque la ruta sea la misma, diferentes viajes pueden tener destinos distintos para informar a los pasajeros.
- El **direction_id** usualmente es 0 o 1 y marca ida o vuelta. Esto ayuda a diferenciar viajes en ambas direcciones dentro de una misma ruta.
- El **block_id** agrupa viajes que se hacen con el mismo bus consecutivamente. Así, un vehículo puede cubrir varios viajes sin regresar a la base, optimizando recursos.
- El **shape_id** indica la ruta exacta o forma que sigue el viaje en el mapa. Esto permite representar con precisión el recorrido en sistemas de mapas o apps.
- El campo **wheelchair_accessible** muestra si el viaje es accesible para personas con discapacidad motriz. Valor 1 indica accesible, 0 no accesible o desconocido.
- El campo **bikes_allowed** indica si está permitido llevar bicicletas en ese viaje. Esto es útil para pasajeros que usan bicicleta y el transporte público.
- Aunque los viajes compartan ruta y servicio, pueden tener destinos, accesibilidad o recorridos diferentes. Esto refleja la flexibilidad del sistema para manejar variantes dentro de una misma ruta.

■ stop_times

Cada registro en **stop_times** representa un paso de un viaje (**trip**) por una parada específica, con hora de llegada y salida, y el orden en la secuencia de paradas.

--- stop_times.txt: primeros 3 registros ---

	trip_id	arrival_time	departure_time	stop_id	stop_sequence	stop_headsign	pickup_type	drop_off_type	shape_dist_traveled
0	11708204_M21	04:32:00	04:32:00	3892	1	NaN	NaN	NaN	NaN
1	11708204_M21	04:32:41	04:32:41	3875	2	NaN	NaN	NaN	NaN
2	11708204_M21	04:33:28	04:33:28	3896	3	NaN	NaN	NaN	NaN
3	11708204_M21	04:34:12	04:34:12	3852	4	NaN	NaN	NaN	NaN
4	11708204_M21	04:34:45	04:34:45	3845	5	NaN	NaN	NaN	NaN
5	11708204_M21	04:35:25	04:35:25	3822	6	NaN	NaN	NaN	NaN
6	11708204_M21	04:36:00	04:36:00	3824	7	NaN	NaN	NaN	NaN
7	11708204_M21	04:36:36	04:36:36	7160	8	NaN	NaN	NaN	NaN
8	11708204_M21	04:37:15	04:37:15	3828	9	NaN	NaN	NaN	NaN
9	11708204_M21	04:37:43	04:37:43	3831	10	NaN	NaN	NaN	NaN

Figura 8: 3 primeros registros de la tabla stop_times

- El campo **trip_id** conecta cada tiempo con un viaje específico, permitiendo saber qué parada pertenece a qué recorrido.
- **arrival_time** y **departure_time** pueden ser iguales, lo que indica que el vehículo no espera en esa parada o la parada es muy rápida.
- Los tiempos pueden exceder las 24 horas (por ejemplo, 25:15:00), lo que indica viajes que pasan después de la medianoche y ayudan a manejar horarios nocturnos sin confusión.
- **stop_sequence** indica el orden exacto en que el vehículo visita las paradas dentro de un viaje, lo que es clave para reconstruir el recorrido.
- **stop_headsign** puede estar vacío, pero si se usa, indica un destino o información específica sobre esa parada en ese viaje particular.
- **pickup_type** y **drop_off_type** definen reglas especiales para subir o bajar pasajeros en esa parada (por ejemplo, si está permitido sólo bajar, o es parada a demanda).
- **shape_dist_traveled** puede contener la distancia recorrida desde el inicio del viaje hasta esa parada, útil para cálculos de posición o velocidad.
- Esta tabla es fundamental para calcular horarios exactos, tiempos de viaje entre paradas y generar mapas animados de los buses en tiempo real.

Condicionalmente requeridos

- calendar


```
--- calendar.txt: primeros 3 registros ---
```

	service_id	monday	tuesday	wednesday	thursday	friday	saturday	sunday	start_date	end_date
0	M21	1	1	1	1	1	0	0	20250516	20250620
1	M12	0	0	0	0	0	1	0	20250516	20250620
2	M13	0	0	0	0	0	0	1	20250516	20250620

Figura 9: 3 primeros registros de la tabla calendar

- **Servicios diferenciados por día:** Los servicios están claramente divididos entre días laborables y fines de semana. Por ejemplo, algunos `service_id` como M21 operan de lunes a viernes, mientras que otros como M12 o M13 lo hacen solo en sábado o domingo respectivamente.
- **Servicios con validez limitada:** Las fechas `start_date` y `end_date` indican que los servicios son válidos solo por un rango de tiempo específico (por ejemplo, del 16 de mayo al 20 de junio de 2025), lo cual puede representar una programación especial o temporal.
- **Servicios únicos por día:** Existen identificadores de servicio que operan únicamente en un día de la semana, lo que podría indicar rutas especiales o de baja demanda.
- **Posible planificación para contingencias o eventos:** La variedad en los patrones de servicio sugiere que SFMTA podría tener diferentes configuraciones de calendario para responder a eventos, obras o temporadas específicas.
- **Estructura binaria fácil de analizar:** Las columnas correspondientes a los días de la semana tienen valores binarios (1 o 0), lo que facilita el análisis de patrones operativos y la automatización de filtros por día.
- **Sin traslape completo entre servicios:** En los registros analizados, cada `service_id` tiene un patrón único de operación semanal, lo que refleja una planificación con bajo nivel de solapamiento total.

■ calendar_dates

```
--- calendar_dates.txt: primeros 3 registros ---
```

	service_id	date	exception_type
0	M21	20250526	2
1	M12	20250526	1
2	M21	20250605	2
3	M41	20250605	1
4	M21	20250606	2
5	M41	20250606	1
6	M21	20250609	2
7	M41	20250609	1
8	M21	20250610	2
9	M41	20250610	1

Figura 10: Primeros registros de la tabla calendar_dates

- **Modificaciones al servicio regular:** El archivo contiene fechas específicas en las que se realizan excepciones al calendario estándar definido en `calendar.txt`. Estas excepciones están identificadas por el campo `exception_type`.
- **Tipos de excepción:** El campo `exception_type` tiene valores binarios:
 - 1: El servicio es añadido en esa fecha.
 - 2: El servicio es eliminado en esa fecha.
- **Uso frecuente de eliminaciones:** Se observa que el `service_id` M21 tiene múltiples eliminaciones (`exception_type` = 2) en fechas consecutivas como el 5, 6, 9 y 10 de junio de 2025, lo que sugiere cancelaciones programadas de servicios.
- **Adiciones planificadas:** El `service_id` M12 y M41 aparecen con `exception_type` = 1, indicando que fueron añadidos especialmente para el 26 de mayo y fechas posteriores, posiblemente por ser días feriados u ocasiones especiales.
- **Relación directa con calendario base:** Estas excepciones ajustan el calendario base definido en `calendar.txt`, y son esenciales para una correcta planificación del servicio real de transporte.
- **Posibilidad de coincidencia con eventos locales:** Las fechas podrían coincidir con feriados o eventos en San Francisco, lo que justificaría tanto la cancelación como la adición de servicios.
- **Planificación flexible del sistema:** El uso activo de `calendar_dates.txt` demuestra que SFMTA implementa una planificación dinámica capaz de responder a cambios temporales en la demanda.

Opcionales

▪ fare_rules

--- fare_rules.txt: primeros 3 registros ---

	fare_id	route_id
0	1	1
1	1	12
2	1	14

Figura 11: Primeros registros de la tabla fare_rules

- **Relación entre tarifas y rutas:** El archivo `fare_rules.txt` establece qué tarifas (`fare_id`) se aplican a qué rutas específicas (`route_id`). Es clave para entender la estructura de precios del sistema.
- **Tarifa uniforme:** Todos los registros muestran el mismo `fare_id` = 1, lo que sugiere una tarifa estándar para múltiples rutas, independientemente de su número o sufijo (por ejemplo, rutas 1, 1X, 14R).
- **Cobertura amplia:** Esta tarifa uniforme aplica a una amplia gama de rutas con identificadores numéricos simples y extendidos (como 14R y 1X), lo cual puede simplificar el modelo tarifario para los usuarios.

- **Simplificación para el usuario:** La asignación de una única tarifa a múltiples rutas puede facilitar la comprensión del sistema tarifario por parte de los pasajeros, lo cual es importante para la accesibilidad y eficiencia del sistema.
- **Posible ausencia de zonificación:** No hay indicios de múltiples `fare_id`, lo cual sugiere que el sistema de SFMTA no está dividido en zonas tarifarias complejas, al menos en esta muestra.
- **Importancia para el modelado:** Este archivo es fundamental para cualquier análisis que requiera calcular costos de viaje o simular escenarios de impacto tarifario en rutas específicas.

■ fare_attributes

--- fare_attributes.txt: primeros 3 registros ---

	fare_id	price	currency_type	payment_method	transfers	agency_id	transfer_duration
0	1	3	USD	0	NaN	NaN	5400
1	2	8	USD	0	0.0	NaN	0

Figura 12: Primeros registros de la tabla fare_attributes

- **Definición de tarifas:** El archivo `fare_attributes.txt` proporciona detalles sobre cada tarifa identificada por `fare_id`, incluyendo el precio, el tipo de moneda y las reglas de transferencia.
- **Precios diferenciados:** Se observan al menos dos tarifas: una de \$3.00 (USD) y otra de \$8.00, lo que indica que existen diferentes tipos de tarifas, posiblemente para servicios especiales o pasajeros con condiciones específicas.
- **Método de pago:** El campo `payment_method` con valor 0 indica que el pago debe hacerse antes de abordar el transporte (según la especificación GTFS).
- **Transferencias:** El campo `transfers` tiene valores faltantes o NaN, lo cual puede indicar que no se especifican reglas claras para transferencias o que no están permitidas por defecto.
- **Duración de la transferencia:** Para la tarifa básica, se permite una transferencia durante un periodo de 5400 segundos (1.5 horas), lo cual puede ser útil para usuarios que hacen conexiones entre rutas.
- **Incompletitud de los datos:** Algunos campos como `agency_id` están vacíos o incompletos, lo que sugiere que la información no está completamente normalizada o no aplica para una sola agencia (como SFMTA).
- **Importancia en la planificación:** Esta tabla es esencial para entender el modelo tarifario, realizar simulaciones de ingresos o estudiar el impacto de precios sobre el uso del transporte.

■ shapes

--- shapes.txt: primeros 3 registros ---

	shape_id	shape_pt_lat	shape_pt_lon	shape_pt_sequence	shape_dist_traveled
0	23	37.803674	-122.443434	1	0.000000
1	23	37.803736	-122.443388	3	0.004962
2	23	37.803921	-122.441952	4	0.084341
3	23	37.803749	-122.441918	5	0.096362
4	23	37.803076	-122.441794	6	0.143322
5	23	37.802990	-122.441767	7	0.149441
6	23	37.802781	-122.443296	8	0.234099
7	23	37.802779	-122.443426	9	0.241193
8	23	37.801898	-122.443262	10	0.302680
9	23	37.801843	-122.443236	11	0.306734

Figura 13: Primeros registros de la tabla shapes

- **Trayectorias geográficas:** El archivo `shapes.txt` describe la geometría de las rutas a través de puntos geográficos secuenciales identificados por `shape_id`.
- **Puntos ordenados:** Cada punto incluye coordenadas de latitud (`shape_pt_lat`) y longitud (`shape_pt_lon`), y se ordena mediante `shape_pt_sequence` para trazar correctamente el recorrido.
- **Distancia acumulada:** La columna `shape_dist_traveled` indica la distancia acumulada recorrida desde el inicio de la forma, medida en unidades consistentes (generalmente kilómetros o millas), útil para estimar la longitud de cada trayecto.
- **Visualización de rutas:** Esta información permite reconstruir y visualizar rutas sobre un mapa, apoyando análisis espaciales de cobertura y trayectorias reales de los vehículos.
- **Múltiples trayectorias:** Un mismo `shape_id` puede representar la forma compartida por varios viajes, optimizando el almacenamiento y representando rutas con geometría común.
- **Importancia cartográfica:** Este archivo es esencial para herramientas de visualización como G2Viz o mapas interactivos, ya que permite representar la red de transporte con precisión espacial.

■ directions

--- directions.txt: primeros 3 registros:

	route_id	direction_id	direction
0	1	0	Outbound
1	1	1	Inbound
2	2	0	Outbound
3	2	1	Inbound
4	5	0	Outbound
5	5	1	Inbound
6	6	0	Outbound
7	6	1	Inbound
8	7	0	Outbound
9	7	1	Inbound

Figura 14: Primeros registros de la tabla directions

- **Direcciones de ruta:** El archivo `directions.txt` proporciona una descripción textual de las direcciones en que opera cada ruta, indicando si se trata de un servicio **Outbound** (salida) o **Inbound** (entrada).
- **Identificación de rutas:** Cada fila está asociada a una `route_id`, que permite vincular esta información con otras tablas como `trips.txt` o `routes.txt`.
- **Dirección numérica:** La columna `direction_id` usa valores binarios (0 o 1) para representar las dos posibles direcciones, lo cual es estándar en el formato GTFS.
- **Claridad semántica:** La columna `direction` entrega una etiqueta legible (como **Inbound** o **Outbound**) que facilita la comprensión del sentido de los trayectos tanto para usuarios como para desarrolladores.
- **Visualización de sentidos:** Esta información es útil en visualizaciones para diferenciar los sentidos de circulación de una misma ruta sobre el mapa o al analizar frecuencia por dirección.
- **Complemento a trips.txt:** Aunque GTFS ya define el campo `direction_id` en `trips.txt`, este archivo complementario mejora la legibilidad al asignar nombres más descriptivos a las direcciones.

■ timepoints

--- timepoints.txt: primeros 3 registros -

	trip_id	stop_id
0	11710030_M41	7914
1	11710030_M41	5689
2	11710030_M41	5688
3	11710030_M41	7563
4	11710030_M41	5404
5	11710030_M41	5391
6	11710030_M41	4229
7	11710030_M41	4732
8	11710030_M41	4748
9	11710030_M41	3927

Figura 15: 3 primeros registros de la tabla timepoints

- **Relación de puntos clave por viaje:** El archivo `timepoints.txt` contiene una lista de paradas clave (`stop_id`) para cada viaje específico (`trip_id`). Estos puntos suelen utilizarse como referencias principales para el control del horario.
- **Identificación única de viaje:** La columna `trip_id` permite vincular estas paradas con un viaje específico definido en `trips.txt`, facilitando el análisis del recorrido planeado.
- **Importancia operativa:** Los `timepoints` suelen representar paradas estratégicas donde se asegura el cumplimiento de horarios, por ejemplo, en intersecciones importantes o puntos de transferencia.
- **Sin información horaria directa:** Aunque este archivo no contiene tiempos explícitos, puede combinarse con `stop_times.txt` para extraer o verificar la precisión temporal en los puntos definidos.
- **Análisis de puntualidad y control de calidad:** Su uso es fundamental en análisis más detallados del cumplimiento de horarios y en la generación de gráficos de rendimiento si se dispone de datos en tiempo real (aunque este proyecto se centra en datos estáticos).
- **Selección parcial de paradas:** No todas las paradas de un viaje están presentes; sólo las consideradas como puntos de control, lo cual lo diferencia de `stop_times.txt`.

Valores nulos

Para evaluar la calidad de los datos del conjunto GTFS de la SFMTA, se realizó un análisis detallado de los valores nulos. Se identificaron campos específicos con presencia de datos faltantes, los cuales se resumen a continuación:

Tabla	Estado de Calidad	Observación
agency	Completa (sin nulos)	Datos listos para análisis
stops	Todas filas con nulos, requiere limpieza	Clave para ubicación, limpieza crítica
routes	Mayoría con nulos, revisar columnas	Nulos en algunas filas, evaluar columnas
trips	Completa (sin nulos)	Datos completos, buen nivel de calidad
stop_times	Todas filas con nulos, requiere limpieza	Tabla muy grande, limpieza imprescindible
calendar	Completa (sin nulos)	Datos completos, pocas filas
calendar_dates	Completa (sin nulos)	Datos completos, pocas filas
fare_attributes	Mayoría con nulos, revisar columnas	Pocas filas completas, evaluar columnas
fare_rules	Completa (sin nulos)	Datos completos, buen nivel de calidad
shapes	Completa (sin nulos)	Datos completos, buen nivel de calidad
directions	Completa (sin nulos)	Datos completos, buen nivel de calidad
timepoints	Completa (sin nulos)	Datos completos, buen nivel de calidad

Tabla 1: Resumen de columnas con valores nulos detectadas en archivos GTFS

Archivo	Columnas afectadas	Descripción del problema
stops.txt	stop_desc, zone_id	Muchas filas sin descripción o sin zona tarifaria asignada.
stop_times.txt	stop_headsign	Este campo suele estar vacío si no hay información adicional.
routes.txt	route_sort_order	Campo vacío que podría usarse para ordenar rutas en visualizaciones.
fare_attributes.txt	transfers, agency_id	Faltan valores clave para modelar transferencias o no aplican al caso de una sola agencia.

Comentarios adicionales por archivo:

- **stops.txt:** Algunas paradas no tienen descripción (**stop_desc**), pero sí cuentan con una URL (**stop_url**) que puede complementar la información. La ausencia de **zone_id** sugiere que no se utilizan zonas tarifarias diferenciadas.
- **stop_times.txt:** El campo **stop_headsign** está vacío en muchos casos, lo cual es aceptable si no hay información adicional por parada.
- **routes.txt:** El campo **route_sort_order** no se utiliza, pero podría emplearse para establecer un orden lógico de presentación en aplicaciones visuales.

- **fare_attributes.txt:** El campo **transfers** presenta valores NaN, lo que podría indicar que no existen reglas claras de transferencia. El campo **agency_id** también está incompleto o no aplica al tratarse de una única agencia (SFMTA).

La mayoría de las tablas están completas y listas para análisis. Los valores nulos observados no comprometen la estructura ni el uso funcional del GTFS, pero sí sugieren oportunidades de mejora en la documentación, la estandarización y la preparación de datos para visualizaciones más completas o simulaciones tarifarias avanzadas.

Limpieza y Transformación de Datos (Basado en G2Viz)

Como parte del preprocesamiento del conjunto de datos GTFS de la SFMTA, se aplicaron técnicas de limpieza inspiradas en el enfoque propuesto por **G2Viz**. Estas técnicas se enfocaron en los siguientes aspectos:

Eliminación de columnas irrelevantes

Se eliminaron atributos opcionales que no aportan valor directo al análisis visual ni a las métricas de rendimiento. Estas columnas suelen estar vacías o ser redundantes.

Tabla 2: Columnas eliminadas por irrelevancia

Tabla	Columnas eliminadas
stops	zone_id, stop_url, location_type, parent_station, stop_timezone, wheelchair_boarding, level_id, platform_code
trips	block_id, direction_id, wheelchair_accessible, bikes_allowed

Resultado de la transformación

Las siguientes tablas fueron exportadas con una estructura limpia y lista para análisis:

- **stops.csv:** estructura depurada de paradas con solo los campos esenciales.
- **trips.csv:** viajes con estructura reducida para análisis más eficiente.

Esta etapa optimiza el rendimiento del sistema y prepara los datos para visualizaciones dinámicas y cálculos de métricas como frecuencia, velocidad y *headway*, conforme a los principios definidos por la herramienta G2Viz.

	Tabla	Filas antes	Columnas antes	Filas después	Columnas después	Columnas eliminadas
0	stops	3278	8	3278	6	zone_id, stop_url
1	trips	35017	9	35017	5	wheelchair_accessible, bikes_allowed, block_id, direction_id

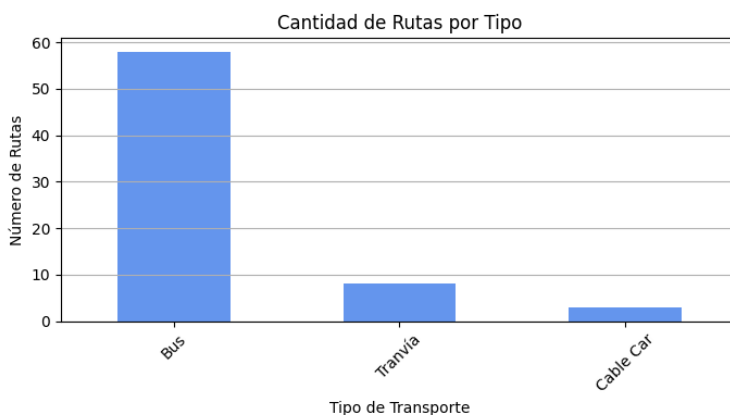
Figura 16: Comparación de Tablas Originales vs Limpias

1.3. Visualización

Para complementar el análisis exploratorio de los datos GTFS de SFMTA, se generaron visualizaciones simples que permiten observar patrones en la distribución de rutas, paradas y viajes programados. Estas representaciones ayudan a responder preguntas clave del sistema, facilitar la validación de los datos procesados y detectar posibles inconsistencias.

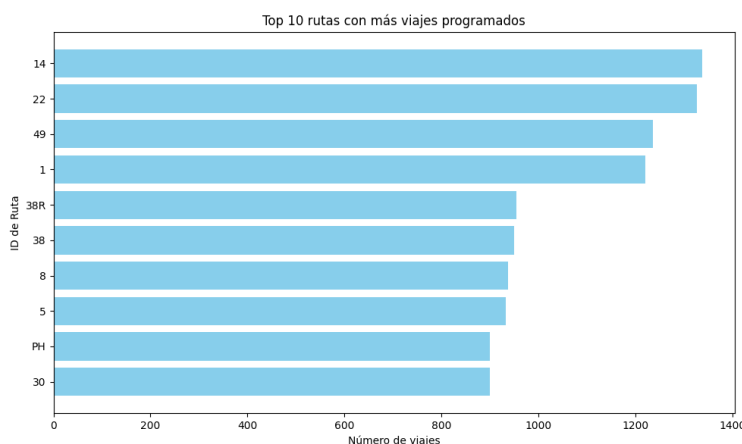
Las principales variables exploradas fueron:

- **Tipos de rutas:** Ayuda a entender la variedad de medios de transporte que opera la SFMTA. Se construyó un gráfico de barras con la cantidad de rutas por tipo (bus, tranvía, cable car), usando el campo `route_type` de `routes.txt`.



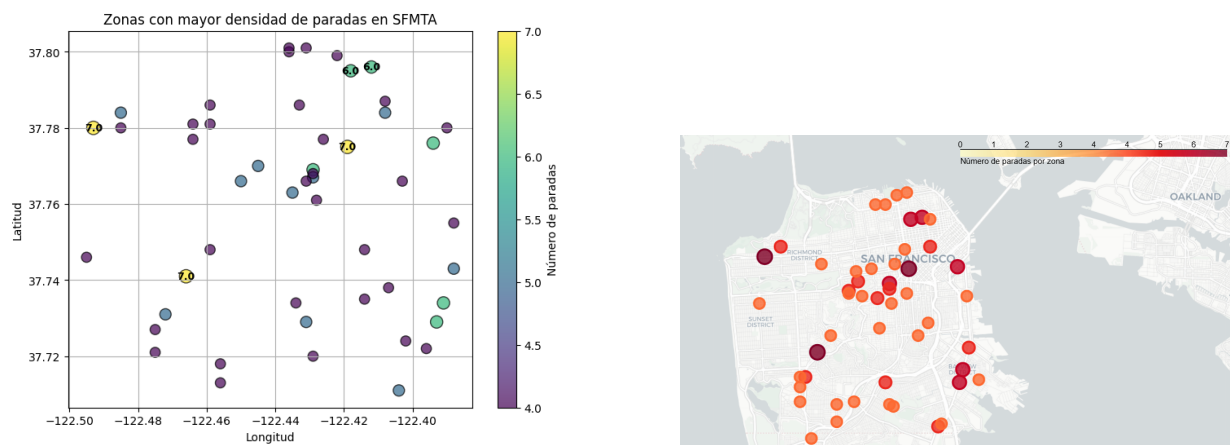
El gráfico muestra que la mayoría de las rutas en el sistema SFMTA corresponden al servicio de bus, mientras que los tranvías ligeros y los cable cars representan una proporción mucho menor. Esto indica que el sistema depende principalmente del transporte terrestre para cubrir la mayor parte de su red.

- **Rutas con mayor número de viajes:** Permite identificar qué rutas tienen mayor frecuencia o demanda prevista. A partir de `trips.txt`, se identificaron las rutas más frecuentes mediante un conteo de `trip_id` por `route_id`.



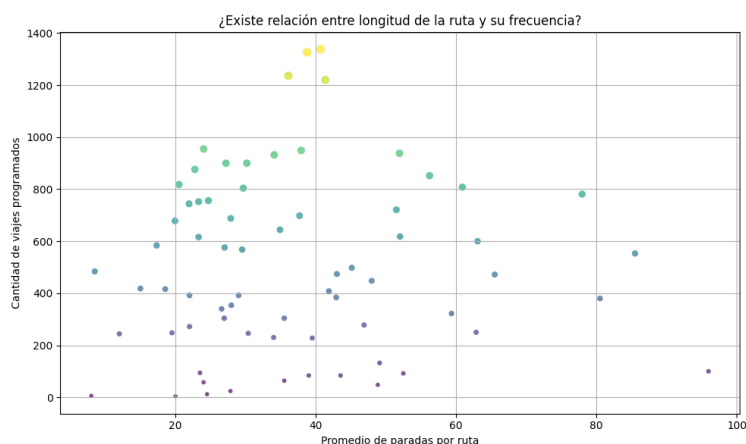
El gráfico muestra que algunas rutas tienen muchos más viajes programados que otras. Esto sugiere que esas rutas son más usadas o importantes dentro del sistema, probablemente porque pasan por zonas con más personas o más movimiento.

- **Densidad de paradas:** A partir de las coordenadas geográficas en `stops.txt` (`stop_lat` y `stop_lon`), se agruparon las paradas por zonas aproximadas mediante redondeo de coordenadas. Luego, se visualizaron sobre un mapa base para identificar áreas con mayor concentración de paradas, lo cual permite analizar la accesibilidad y cobertura del sistema.



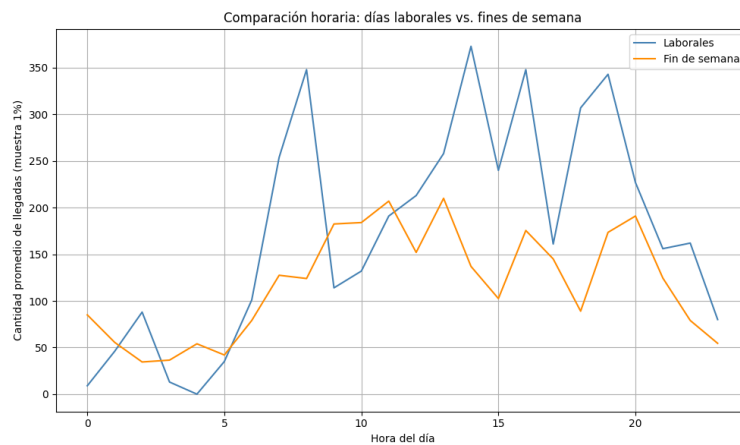
El gráfico muestra que las zonas con mayor densidad de paradas se encuentran principalmente en el centro de San Francisco. Estas áreas concentran muchas paradas cercanas entre sí, lo que sugiere una mayor accesibilidad al transporte público en zonas comerciales o con alta demanda de movilidad.

- **Relación longitud-frecuencia:** Se cruzaron datos de `stop_times.txt` y `trips.txt` para analizar si las rutas con más paradas también presentan mayor frecuencia de viajes.



El gráfico muestra que no existe una relación directa entre la longitud de una ruta (número de paradas) y su frecuencia (número de viajes). Algunas rutas largas tienen pocos viajes, mientras que rutas más cortas son más frecuentes, lo que sugiere que otros factores influyen en la programación del servicio.

- **Cobertura semanal:** A partir de los archivos `calendar.txt`, `trips.txt` y `stop_times.txt`, se analizó la frecuencia de llegadas por hora en cada día de la semana.



El gráfico muestra que los días laborales presentan una mayor frecuencia de llegadas, con picos marcados en horas punta (mañana y tarde), mientras que los fines de semana tienen un servicio más reducido y distribuido de forma más uniforme durante el día. Esto refleja una programación ajustada a patrones de movilidad distintos entre semana y fin de semana.

Estas visualizaciones, aunque básicas, fueron útiles para validar el preprocesamiento y apoyar la formulación de hipótesis. En etapas posteriores, se integrarán a un dashboard interactivo desarrollado con D3.js.

Referencias

- [1] Para, S., Wirotasathon, T., Jundee, T., Demissie, M. G., Sekimoto, Y., Biljecki, F., & Phithak-
kitnukoon, S. (2024b). G2Viz: an online tool for visualizing and analyzing a public transit
system from GTFS data. Public Transport.
<https://link.springer.com/article/10.1007/s12469-024-00362-x>
- [2] trufi-association. (s. f.). trufi-arequipa/GTFS-Peru-Arequipa/out/gtfs/agency.txt at main ·
trufi-association/trufi-arequipa
[https://github.com/trufi-association/trufi-arequipa/blob/main/GTFS-Peru-Arequipa/
out/gtfs/agency.txt](https://github.com/trufi-association/trufi-arequipa/blob/main/GTFS-Peru-Arequipa/out/gtfs/agency.txt)
- [3] GTFS transit data. (2024, 9 octubre). SFMTA.
<https://www.sfmta.com/reports/gtfs-transit-data>
- [4] Create - General Transit Feed Specification. (s. f.).
<https://gtfs.org/getting-started/create/>
- [5] Google. GTFS Reference - agency.txt.
<https://developers.google.com/transit/gtfs/reference#agencytxt>
- [6] 311 SF City Services.
<https://sf311.org/>
- [7] GTFS.org. GTFS Best Practices.
<https://gtfs.org/documentation/schedule/reference/#agencytxt>
- [8] Transitland - Feed Analysis.
<https://transit.land/feed-registry/>
- [9] Google Developers. GTFS Static Reference - stops.txt.
<https://gtfs.org/reference/static/>
- [10] MobilityData. Understanding GTFS: stops.txt.
<https://mobilitydata.org/understanding-gtfs-stops/>
- [11] OpenMobilityData (TransitFeeds). Ejemplos de feeds GTFS.
<https://transitfeeds.com/>
- [12] Transit Developers Group. GTFS Best Practices.
<https://github.com/MobilityData/gtfs-best-practices>
- [13] Google Colab. (s. f.).
[https://colab.research.google.com/drive/1fEFLHNArU7imlqQ-vPyBFht61JZJp1WG?
usp=sharing](https://colab.research.google.com/drive/1fEFLHNArU7imlqQ-vPyBFht61JZJp1WG?usp=sharing)