

# Documentación

## Resumen

En este proyecto desarrollé un pipeline de datos automatizado para ingestar, limpiar, transformar y almacenar datos de registros de vehículos eléctricos. A continuación, el resumen de mi enfoque:

1. **Automatización de descarga:** para ello usé BeautifulSoup en Python para hacer web scraping del CSV público con datos de Electric Vehicle Population Data.
2. **Limpieza de datos:** con Pandas leí el CSV, corregí tipos de datos (fechas, numéricos), eliminé registros nulos o inconsistentes, y normalicé formatos de texto (por ejemplo, nombres de condados).  
Resultado: un DataFrame limpio listo para análisis.
3. **Generación de archivo Parquet:** Exporté el DataFrame limpio a formato Parquet. Esto es porque Airflow no me permite pasar datos tan pesados por xcom.
4. **Modelado y creación de tablas:** Diseñé un esquema en estrella (star schema) con:
  - **Dimensiones:** dim\_vehicle, dim\_date, dim\_location, dim\_electric\_type, dim\_policy
  - **Tabla de hechos:** fact\_registration
5. **Carga de datos en PostgreSQL:** Mediante SQLAlchemy y COPY en Postgres, cargué los datos del Parquet en sus tablas correspondientes. Las dimensiones se llenan evitando duplicados (usando slow change dimension), y el hecho enlaza a los IDs.

## DAGS y Scripts

- **dags/electric\_vehicles\_etl.py:** define un DAG con 5 tareas: descargar CSV, limpiar, convertir a Parquet, crear esquema en Postgres, cargar datos.
- **Lógica en etl/:**
  - data\_reader.py: descarga y lee CSV.
  - data\_cleaner.py: aplica transformaciones Pandas.
  - schema\_creator.py: ejecuta sentencias DDL para tablas.
  - load\_data.py: realiza la carga eficiente del Parquet con SQLAlchemy y COPY.
  - db\_connection.py: configuración de la conexión a Postgres.

## Decisiones de diseño

- **Parquet vs CSV:** Parquet ofrece compresión y tipos nativos, acelerando cargas.

- **Star Schema:** mejora rendimiento de consultas analíticas y facilita visualizaciones.
- **Airflow:** permite programar re-ejecuciones, recuperaciones en fallos y monitorización de logs.

## **Transformaciones de datos**

1. **Fechas:** básicamente, se creó desde 0 la tabla dim\_date.
2. **Normalización de texto:** todas las columnas poseen nombres en snake\_case
3. **Extracción de coordenadas:** se obtuvo tanto la latitud como la longitud a partir de la columna POINT (...)
4. **Rellenado de valores:** algunos valores faltantes (como valores numéricos) fueron rellenados con 0

## Análisis de datos

Se configuró Docker para que ya traiga una base de datos integrada, hosteada en el puerto 5433 por lo que solo se tuvo que conectar a dicho host y luego seleccionar las tablas creadas en el proyecto.

### Base de datos PostgreSQL

Servidor

localhost:5433

Base de datos

promptior\_etl

Modo Conectividad de datos ⓘ

☒ Importar

☐ DirectQuery

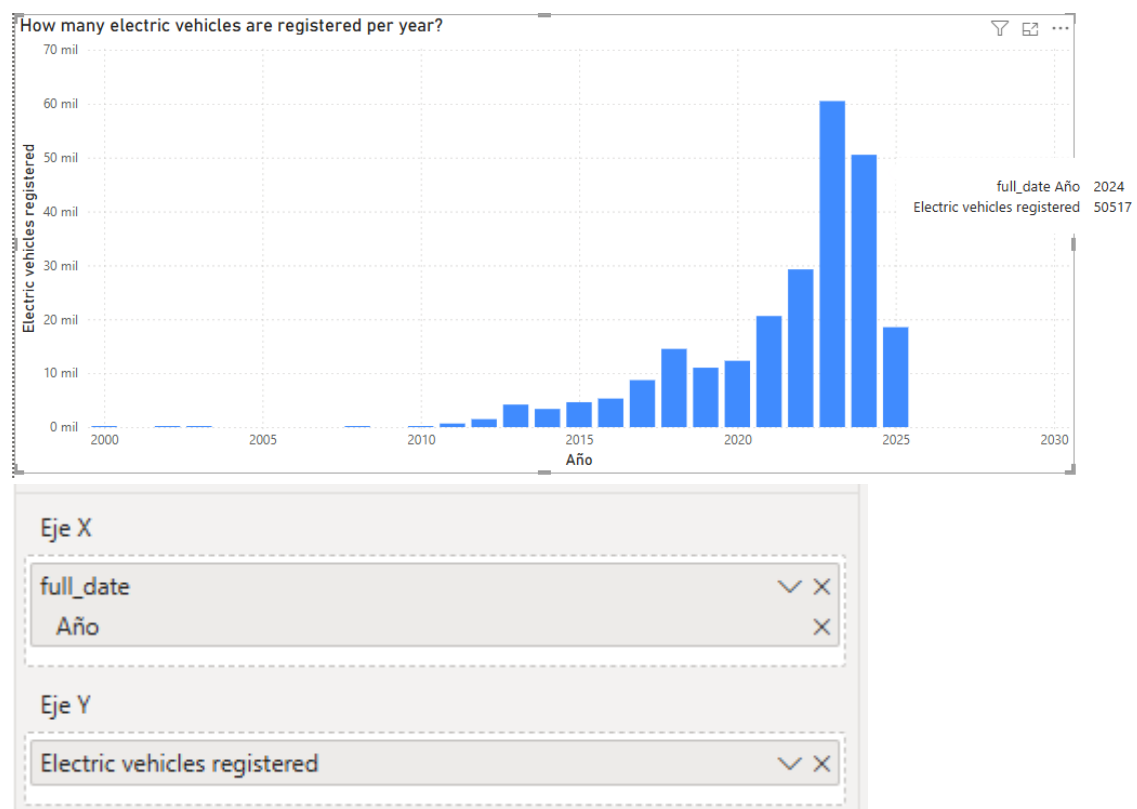
► Opciones avanzadas

Aceptar

Cancelar

### - How many electric vehicles are registered per year?

Para obtener este gráfico se creó la medida Total Registros que lo único que hace es contar las filas de la tabla de hechos

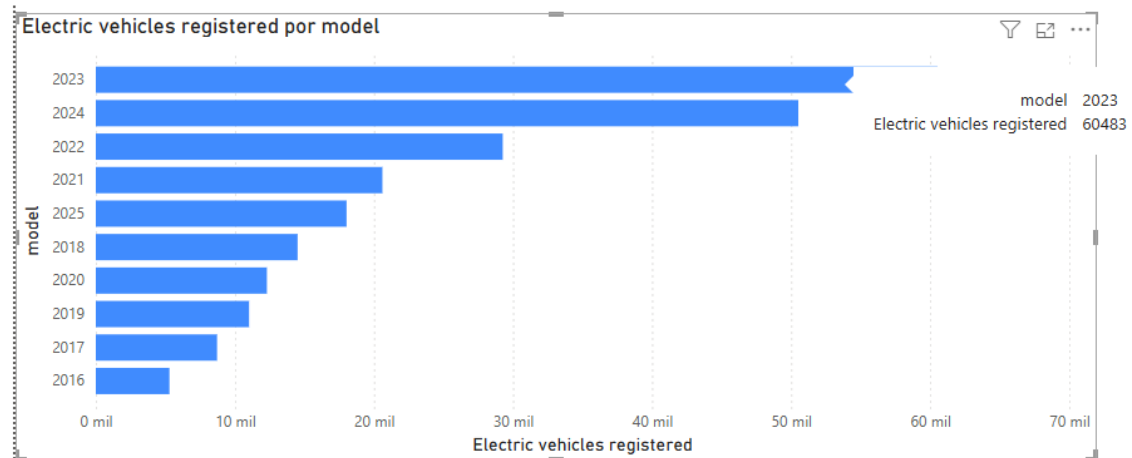


```
1 Total Registros =  
2 COUNTROWS('public fact_registration')
```

- What are the top 10 electric vehicle models by registration count?

Para obtener este gráfico se creó la medida Registros por Modelo que quita todos los filtros de la tabla excepto los que se aplican al modelo.

Además, se filtró por el top 10 de modelos según Registros por Modelo



**Filtros**

Buscar

Filtros de este objeto visual

Electric vehicles regist...  
es (todos)

**model**  
10 principales por Re...

Tipo de filtro ⓘ  
Top N

Mostrar artículos  
Superior 10

Por valor  
Registros por Modelo

Aplicar filtro

Agregar campos de datos ...

Filtros de esta página

Agregar campos de datos ...

Filtros de todas las páginas

Agregar campos de datos ...

**Visualizaciones**

Compilar visual

Eje Y  
model

Eje X  
Electric vehicles registered

Leyenda  
Agregar campos de datos aquí

Múltiplos pequeños  
Agregar campos de datos aquí

Información sobre herramientas  
Agregar campos de datos aquí

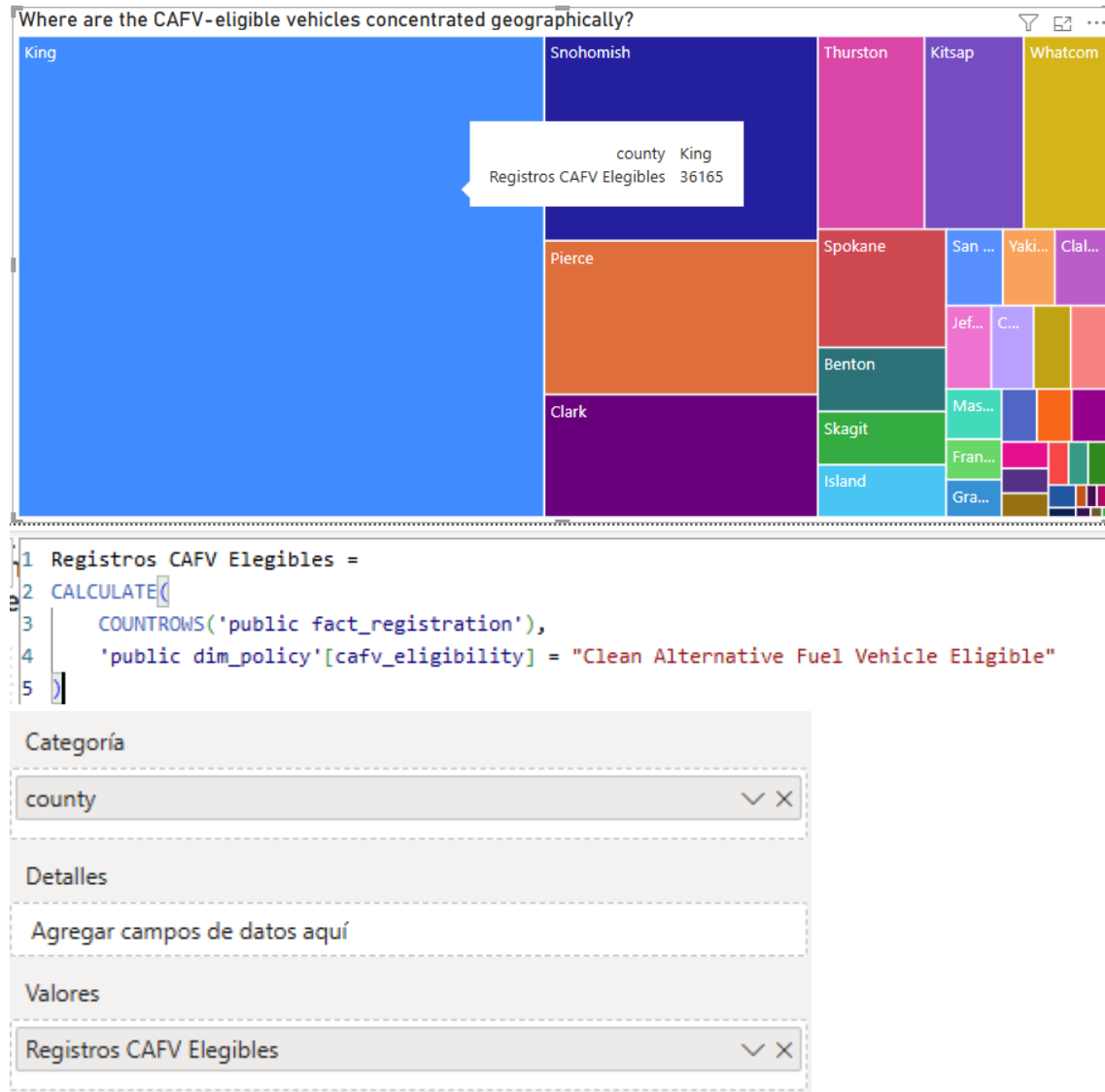
Obtener detalles  
Entre varios informes  
Mantener todos los filtros  
Agregue los campos de obtención de detalles aquí.

```
1 Registros por Modelo =  
2 CALCULATE(  
3     COUNTROWS('public fact_registration'),  
4     ALLEXCEPT('public dim_vehicle', 'public dim_vehicle'[model])  
5 )
```

- Where are the CAFV-eligible vehicles concentrated geographically?

Para obtener este gráfico se creó la medida Registros CAFV Elegibles cuenta todas las columnas de dim\_policy que contengan un cafv\_eligibility = 'Clean Alternative Fuel Vehicle Eligible'

Y se mostró por esos valores por condado

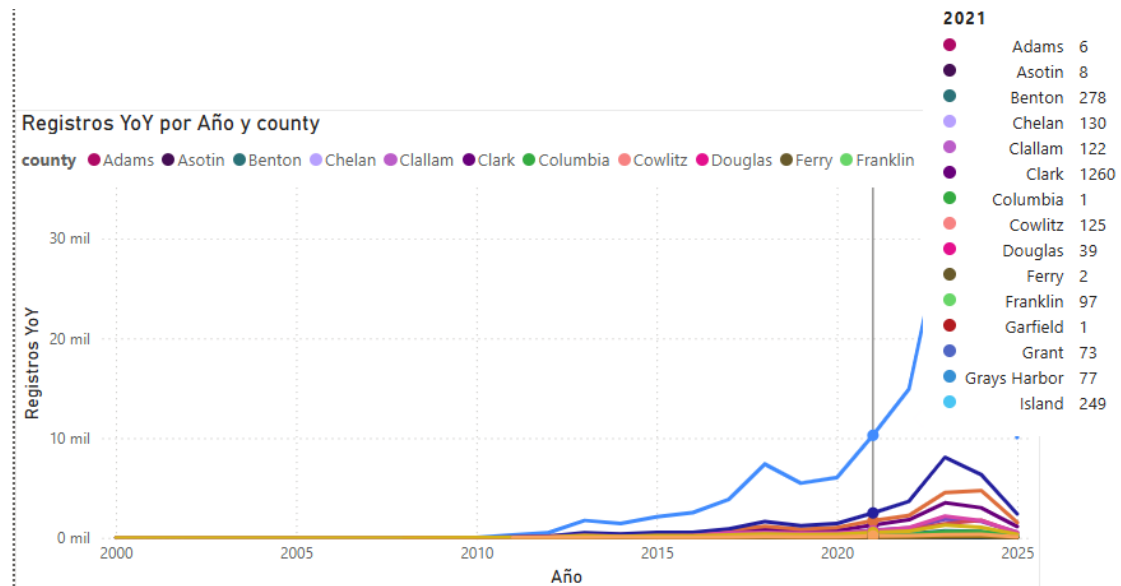


- What is the year-over-year change in electric vehicle registrations by county?

Para obtener este grafico se creó la medida Registros YoY que consta de 3 partes:

- 1) TotalAnioActual: cuenta todas filas de la tabla de hechos
- 2) TotalAnioAnterior: cuenta todas las filas de la tabla hechos pero de un año anterior
- 3) DiferenciaYoY: resta TotalAnioActual menos TotalAnioAnterior dando como resultado, para un año en especifico, la diferencia de registros que había entre ese año y el anterior.

Se usa como leyenda a los condados para verlos como líneas y se divide por años



```

Registros YoY =
VAR TotalAnioActual = COUNTROWS('public fact_registration')
VAR TotalAnioAnterior =
    CALCULATE(
        TotalAnioActual,
        SAMEPERIODLASTYEAR('public dim_date'[full_date])
    )
VAR DiferenciaYoY =
    TotalAnioActual - TotalAnioAnterior
RETURN
    TotalAnioActual
    
```