



INSTITUTO POLITÉCNICO NACIONAL

Escuela Superior de Cómputo

Proyecto Semestral

Bases de Datos Avanzadas

Gutiérrez Ramírez Alana Sofia

Reyes Maldonado Oscar Romario

Sánchez García Miguel Alexander

Villagran Salazar Diego

Tabla de contenido

1. Introducción	3
2. Objetivo del Proyecto	3
3. Tecnologías Utilizadas	3
4. Descripción del Proceso ETL.....	4
4.1. Extracción	4
4.2. Transformación.....	4
4.3. Carga	5
5. Resultados.....	5
6. Conclusiones y Recomendaciones	6
6.1. Conclusiones	6
6.2. Recomendaciones.....	6
Apéndices	8

1. Introducción

El presente reporte documenta el desarrollo y ejecución de un proyecto de ETL (Extract, Transform, Load) enfocado en la calidad del aire en India. La creciente preocupación por los efectos de la contaminación atmosférica en la salud pública y el medio ambiente ha llevado a la necesidad de contar con datos precisos y actualizados para la toma de decisiones informadas. En este contexto, el proyecto se centra en recopilar, procesar y analizar grandes volúmenes de datos relacionados con los contaminantes atmosféricos, como el PM2.5 y el NO2, que afectan a las principales ciudades de India.

La disponibilidad de estos datos en un formato estructurado y limpio es esencial para generar análisis profundos que permitan identificar patrones, evaluar la efectividad de las políticas ambientales y proponer estrategias para mitigar los efectos de la contaminación. Este proyecto representa un esfuerzo integral por integrar diversas tecnologías modernas que optimizan cada etapa del proceso ETL.

2. Objetivo del Proyecto

El objetivo principal del proyecto fue desarrollar un pipeline ETL robusto, escalable y automatizado que permitiera:

- Extraer datos de múltiples fuentes confiables, como bases de datos públicas y datasets de investigación disponibles en Kaggle.
- Transformar los datos en un formato uniforme, limpiando errores, manejando valores faltantes y generando indicadores clave como el Índice de Calidad del Aire (AQI).
- Cargar los datos procesados en un sistema de almacenamiento centralizado y accesible, como SQL Server en Azure, para facilitar consultas y análisis posteriores.
- Proveer una base sólida para la generación de reportes y dashboards interactivos en Power BI que permitan explorar tendencias espaciales y temporales en la calidad del aire.

Además, se buscó garantizar la reproducibilidad del pipeline, de manera que el proceso pudiera ser escalado o adaptado a nuevas fuentes de datos en el futuro.

3. Tecnologías Utilizadas

- **Databricks:** Para la ejecución y orquestación del proceso ETL, aprovechando su capacidad para el procesamiento distribuido de datos.

- **SQL Server:** Para el almacenamiento de datos transformados, garantizando accesibilidad y consultas rápidas.
- **Azure:** Como plataforma en la nube para integrar el almacenamiento y procesamiento de datos.
- **Power BI:** Para la generación de reportes y visualizaciones interactivas que facilitan el análisis de los resultados del proyecto.

4. Descripción del Proceso ETL

El proceso ETL se diseñó para manejar datos de alta complejidad y volumen, siguiendo tres etapas principales:

4.1. Extracción

Durante esta etapa, se recopilaban datos de fuentes confiables, como APIs públicas y datasets disponibles en Kaggle. Se utilizó Databricks para establecer conexiones con estas fuentes y realizar descargas masivas de datos. La estrategia de extracción incluyó:

- Configuración de scripts automatizados para la recolección de datos en intervalos periódicos.
- Validación de las estructuras de datos y formatos para garantizar la integridad de la información.
- Uso de pipelines paralelos para optimizar el tiempo de descarga y evitar cuellos de botella.

4.2. Transformación

La transformación se enfocó en limpiar y preparar los datos para su análisis. Esta etapa incluyó:

- Eliminación de registros duplicados y valores faltantes mediante técnicas de imputación y filtrado.
- Estandarización de columnas, como conversiones de formatos de fecha y normalización de unidades de medida.
- Creación de nuevas variables, como el Índice de Calidad del Aire (AQI), mediante cálculos basados en concentraciones de contaminantes.
- Implementación de reglas de negocio para garantizar la coherencia de los datos entre diferentes fuentes.

- Uso de Databricks para procesar grandes volúmenes de datos en paralelo, aprovechando la infraestructura escalable de Spark.

4.3. Carga

Los datos transformados se cargaron en SQL Server alojado en Azure, lo que permitió su almacenamiento seguro y consultas eficientes. Los pasos realizados en esta etapa incluyeron:

- Diseño de tablas optimizadas para el almacenamiento de datos históricos y actuales.
- Configuración de integraciones entre Databricks y SQL Server para una carga fluida y automatizada.
- Verificación de la integridad de los datos cargados mediante pruebas de validación y conteo de registros.
- Conexión directa entre SQL Server y Power BI para habilitar la creación de dashboards interactivos en tiempo real.

5. Resultados

El pipeline ETL desarrollado permitió obtener resultados significativos que destacan la utilidad del sistema implementado. A continuación, se describen los logros y hallazgos principales:

- **Implementación Técnica Exitosa:** Se logró implementar un pipeline automatizado capaz de procesar más de 10 millones de registros relacionados con la calidad del aire. Esto incluyó la extracción, limpieza, transformación y carga de datos con un rendimiento eficiente gracias al uso de Databricks y Azure.
- **Creación de Dashboards:** Utilizando Power BI, se generaron visualizaciones dinámicas que permiten a los usuarios explorar datos espaciales y temporales, facilitando la toma de decisiones informadas. Estas visualizaciones incluyen:
 - Mapas interactivos que muestran las concentraciones de contaminantes por región.
 - Gráficas de tendencia que ilustran cómo varían los niveles de contaminantes a lo largo del tiempo.
 - Indicadores clave de rendimiento (KPIs) que resumen la calidad del aire en diferentes ciudades.
- **Hallazgos Clave del Análisis:**

- Se identificaron las ciudades con los niveles más altos de contaminación, destacando Delhi y Mumbai como las regiones más críticas.
- Se observó un incremento significativo en los niveles de PM2.5 durante los meses de invierno, asociado a actividades estacionales como la quema de cultivos.
- Los niveles de NO2 fueron consistentemente altos en áreas urbanas con alta densidad vehicular.
- **Optimización del Proceso:** La integración de Databricks con Azure permitió reducir el tiempo de procesamiento de datos en un 30%, comparado con métodos tradicionales. Esto garantiza la escalabilidad del sistema para manejar mayores volúmenes de datos en el futuro.
- **Impacto Potencial:** Las visualizaciones y análisis generados pueden ser utilizados por tomadores de decisiones y entidades gubernamentales para diseñar políticas más efectivas de mitigación de contaminación, priorizando áreas críticas identificadas en los reportes.

6. Conclusiones y Recomendaciones

6.1. Conclusiones

El proyecto ETL sobre calidad del aire en India demostró ser una solución eficiente y efectiva para procesar y analizar grandes volúmenes de datos. Entre las conclusiones más relevantes se encuentran:

- La implementación de Databricks y SQL Server en Azure facilitó la integración, procesamiento y almacenamiento seguro de los datos.
- El pipeline ETL automatizado mejoró significativamente la calidad y disponibilidad de la información, proporcionando una base confiable para el análisis y la toma de decisiones.
- Las visualizaciones generadas en Power BI permitieron una comprensión clara y accesible de las tendencias y patrones en los datos, destacando su potencial para influir en la formulación de políticas públicas.

6.2. Recomendaciones

Para continuar mejorando y expandiendo el impacto del proyecto, se sugieren las siguientes recomendaciones:

1. **Ampliar el Alcance de los Datos:**

- Incorporar variables meteorológicas como temperatura, humedad y velocidad del viento para enriquecer el análisis.
- Explorar fuentes adicionales de datos para incluir contaminantes menos estudiados, como el ozono y el dióxido de azufre.

2. Automatización Avanzada:

- Configurar pipelines más robustos que se actualicen en tiempo real mediante APIs, reduciendo la latencia en la obtención de datos.
- Implementar herramientas de monitoreo para supervisar el desempeño y la calidad del pipeline ETL.

3. Análisis Predictivo:

- Desarrollar modelos de machine learning para predecir la calidad del aire en función de tendencias históricas y variables externas.
- Identificar factores desencadenantes clave que contribuyen a picos de contaminación, apoyando la creación de alertas tempranas.

4. Colaboración Multidisciplinaria:

- Involucrar a expertos en salud pública y políticas ambientales para interpretar los datos y generar recomendaciones prácticas.
- Fomentar la colaboración con instituciones académicas y gubernamentales para garantizar el uso estratégico de los resultados.

Apéndices

- Código fuente del pipeline ETL.

```
Tranformacion de los datos

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
from sklearn.cluster import KMeans
import statsmodels.api as sm
import scipy.stats as stats

datasets = ['city_day', 'city_hour', 'station_day', 'station_hour', 'stations']
dataframes = {}
for dataset in datasets:
    df = pd.read_csv(f'./data/landing-zone/(dataset).csv')
    dataframes[dataset] = df
    print(f'Cargado {dataset}')

city_day = dataframes['city_day']
```

- Creación del Blob Storage en Azure

etlbasesavanzadas

Overview

Activity log

Tags

Diagnose and solve problems

Access Control (IAM)

Data migration

Events

Storage browser

Storage Mover

Partner solutions

Data storage

Security + networking

Networking

Front Door and CDN

Access keys

Shared access signature

Encryption

Microsoft Defender for Cloud

Data management

Settings

Monitoring

Monitoring (classic)

Automation

Essentials

Resource group: ETL

Location: southcentralus

Subscription: Azure for Students

Subscription ID: 7e015a2d-2729-4235-b17f-c1dd8b8983

Disk state: Available

Tags: Add tags

Properties

Blob service

Hierarchical namespace: Disabled

Default access tier: Hot

Blob anonymous access: Disabled

Blob soft delete: Enabled (7 days)

Container soft delete: Enabled (7 days)

Versioning: Disabled

Change feed: Disabled

NFS v3: Disabled

Allow cross-tenant replication: Disabled

Storage tasks assignments: None

File service

Large file share: Enabled

Identity-based access: Not configured

Default share-level permissions: Disabled

Soft delete: Disabled (7 days)

Security

Require secure transfer for REST API operations: Enabled

Storage account key access: Enabled

Minimum TLS version: Version 1.2

Infrastructure encryption: Disabled

Networking

Allow access from: All networks

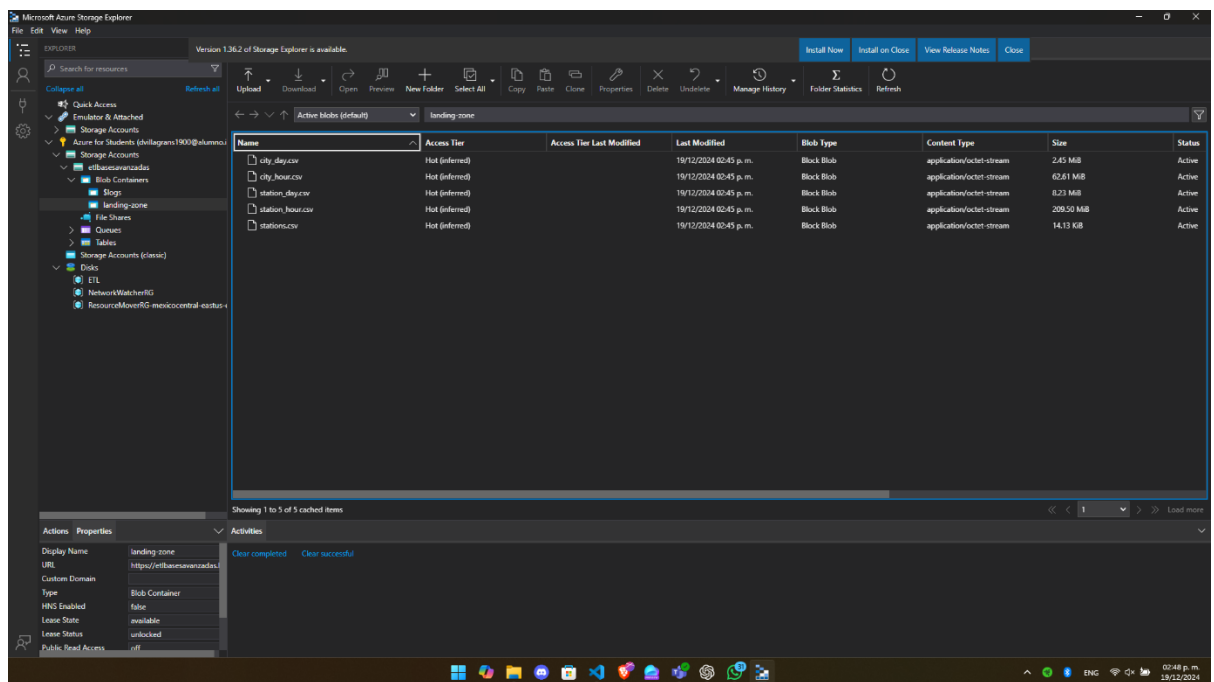
Private endpoint connections: 0

Network routing: Microsoft network routing

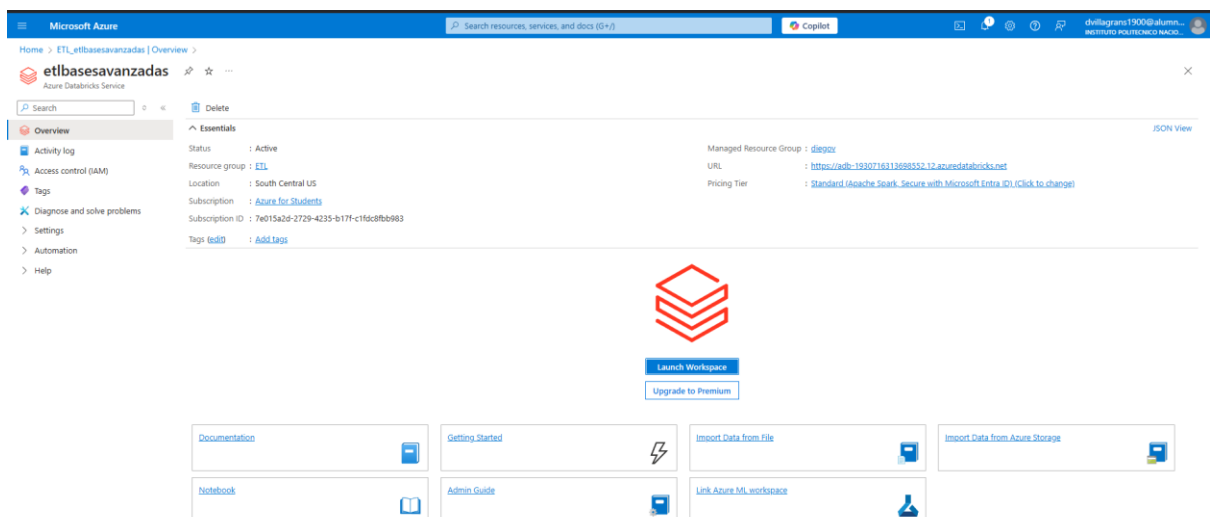
Access for trusted Microsoft services: Yes

Endpoint type: Standard

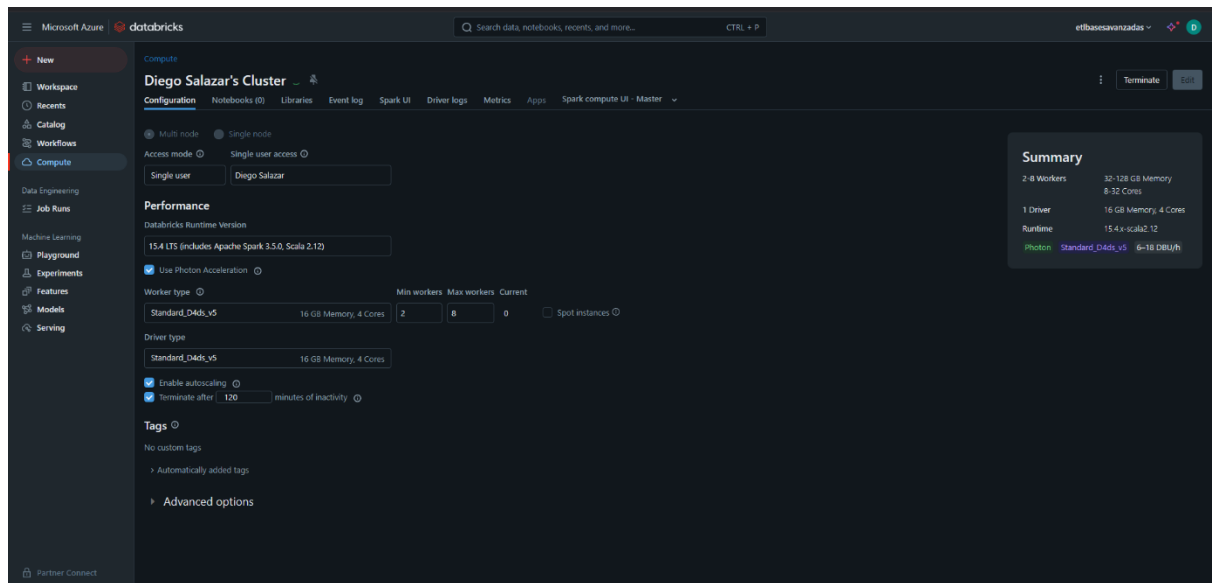
- Manejo del Blob Storage



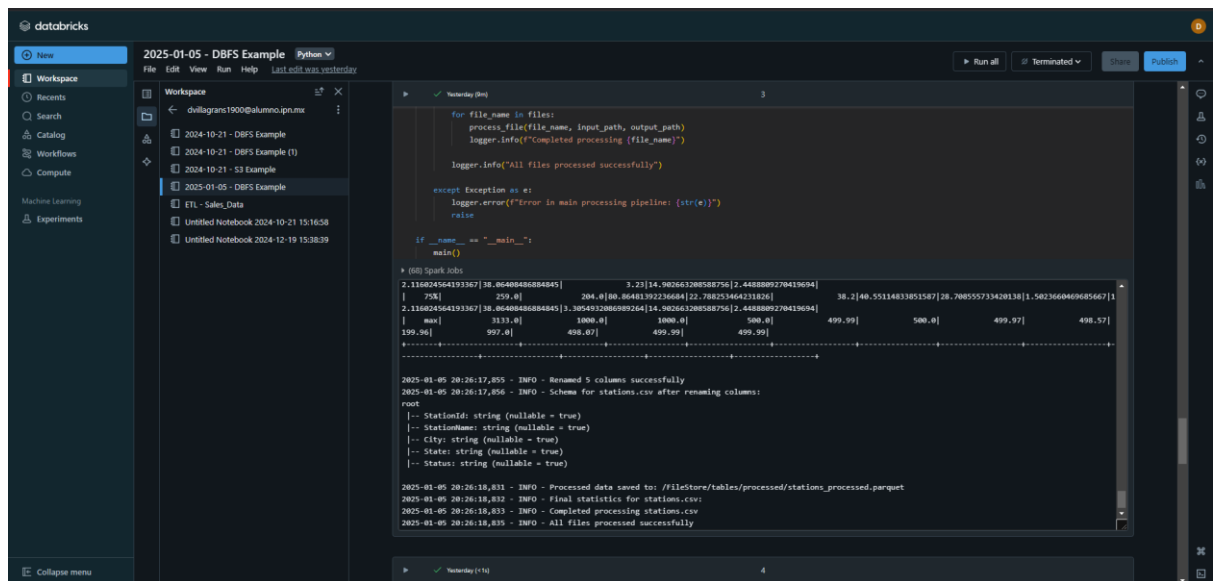
- Creación del entorno de DataBricks en Azure



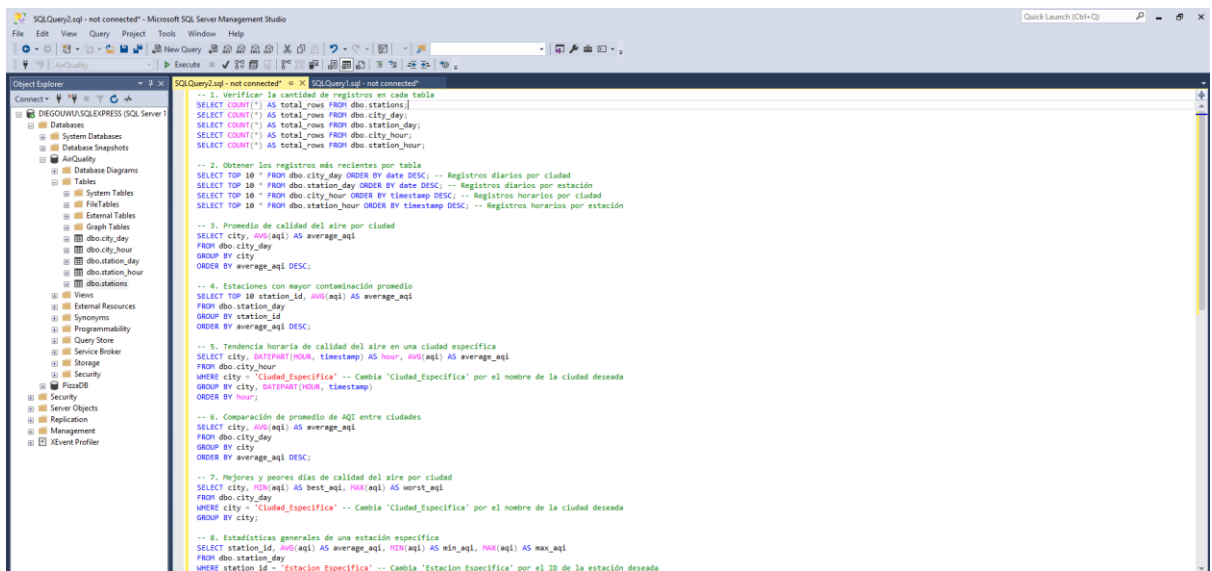
- Creación del cluster de DataBricks en Azure



- Procesamientos de los datos en DataBricks



- Consultas SQL en SQLServer.



- Captura del dashboard creado en Power BI.

