AZURE DATA FACTORY / AZURE DATABRICKS







INTEGRANTES



Carlos Eduardo Denett



Cecilia Marcela Espada



Federico Pfund



Juan Martín Elena



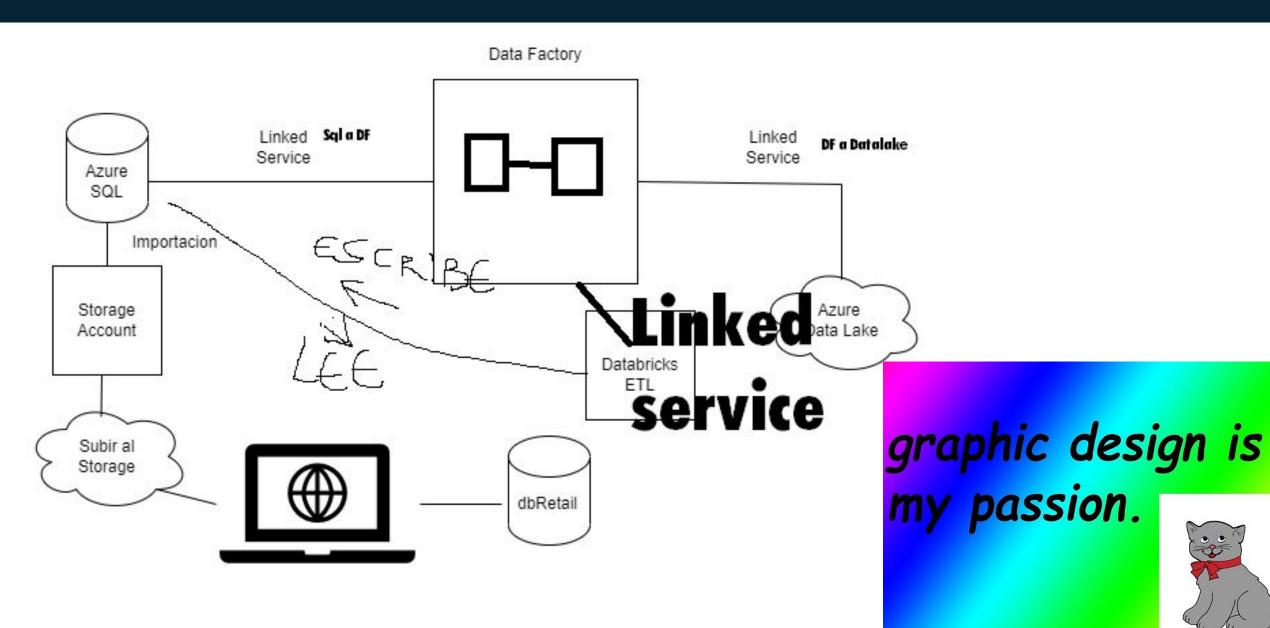
Agustín Fernández



Patricio Perrone

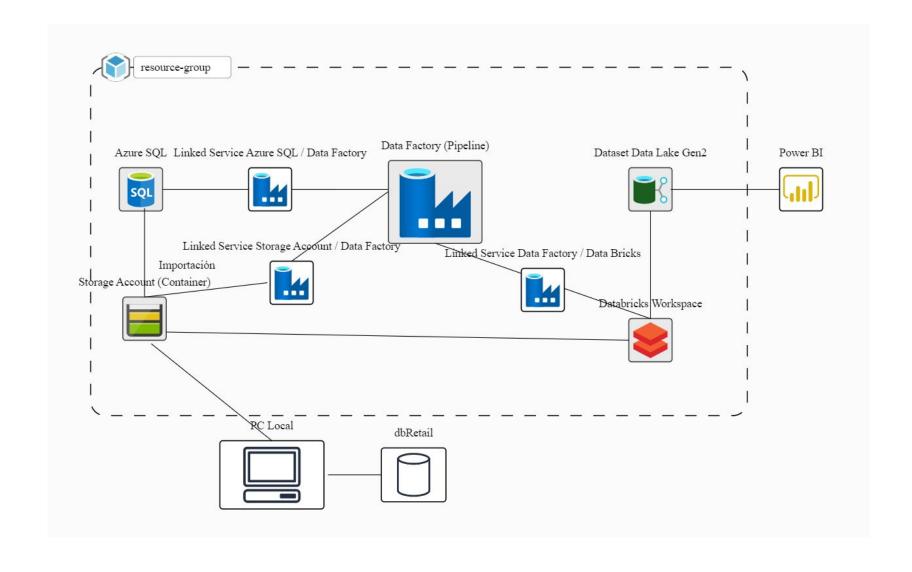
Review Sprint 2 - Grupo 1Proyecto Integrador Final

Cambios de Arquitectura - Antigua



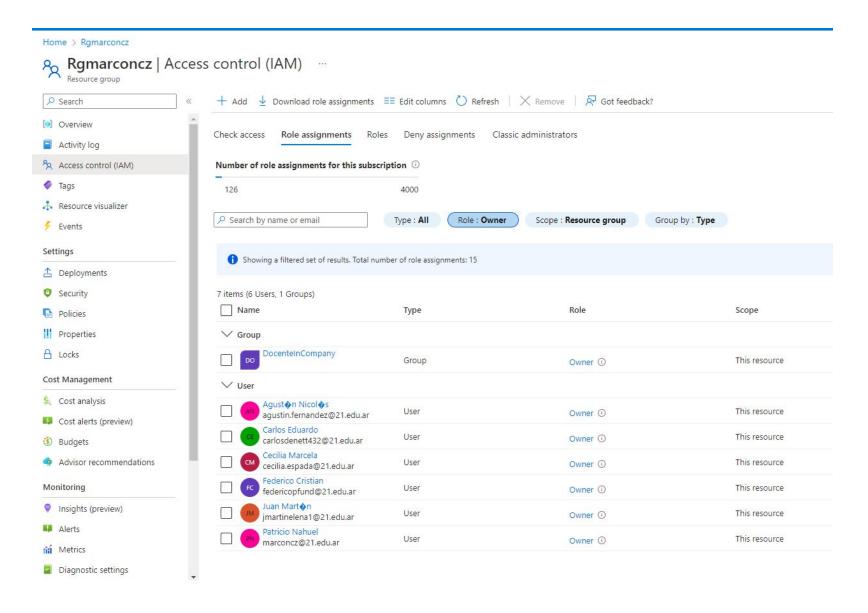
Cambios de Arquitectura

En este segundo sprint, realizamos cambios en la arquitectura para trabajar directamente con archivos csv y evitar el gasto de consumo continuo del recurso Azure SQL. También logramos mejorar los tiempos de ejecución.



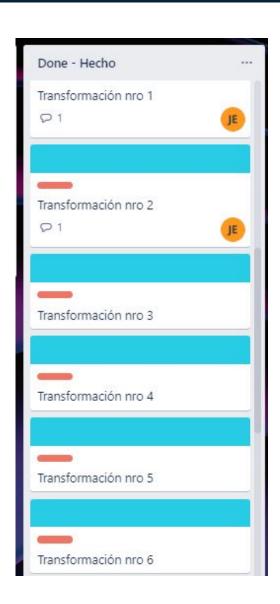
Grupo de Recursos Compartidos

 Decidimos utilizar un sólo grupo de recursos y compartir los accesos con el resto del equipo para ahorrar presupuesto y focalizarnos en una sola solución en común.



Tareas del Sprint 2

- En este sprint nos concentramos en el cambio de la arquitectura, el pipeline con las transformaciones y la presentación de los datos de manera visual
- Trello oficial con todos los grupos.
- Trello del grupo

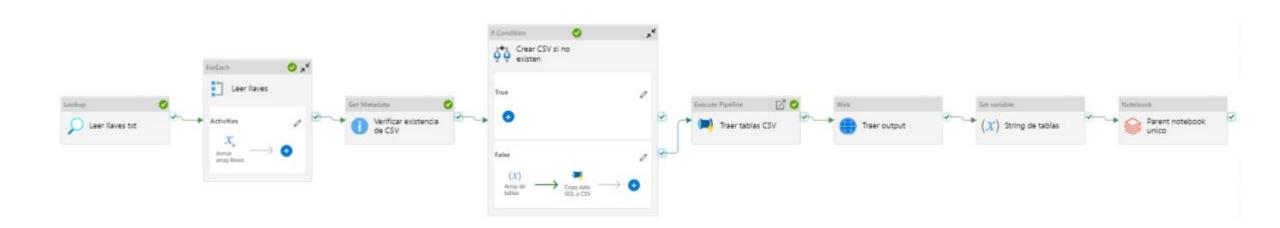




> Pipeline ETL en Data Factory

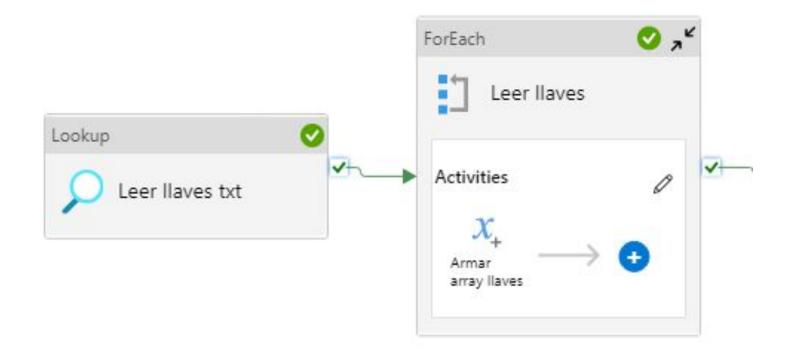
Creación de Pipeline en Data Factory

• Desarrollamos un Pipeline que crea las tablas en CSV si no existen y se las transfiere a Databricks.



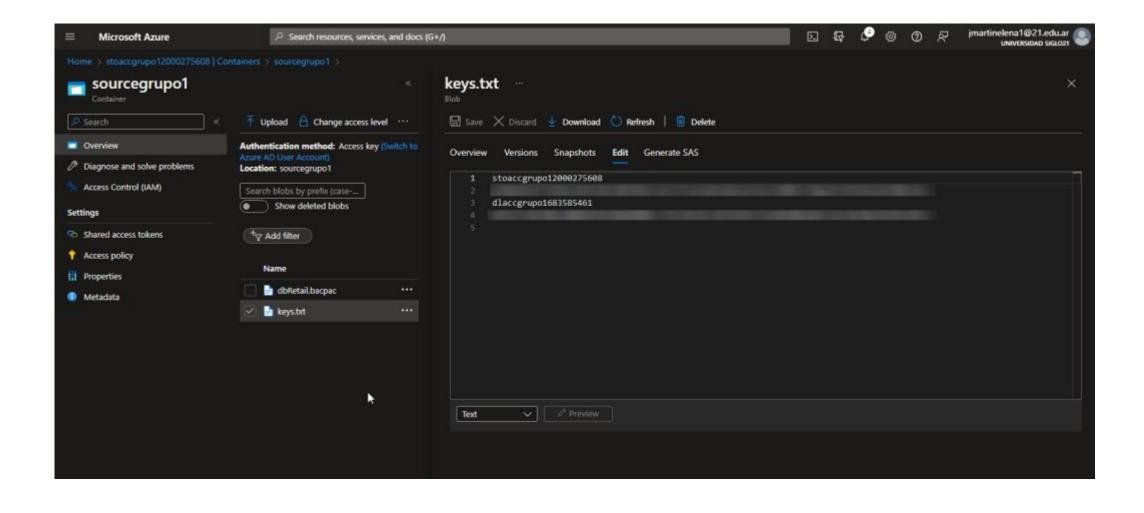
Lectura de datos de acceso

• En este primer paso, se accede a un archivo txt donde se almacenan los datos de acceso que emula un Key Vault, como lo son nombres del storage account y data lake, y sus access keys.



Lectura de datos de acceso

Archivo txt dentro del Blob Storage:

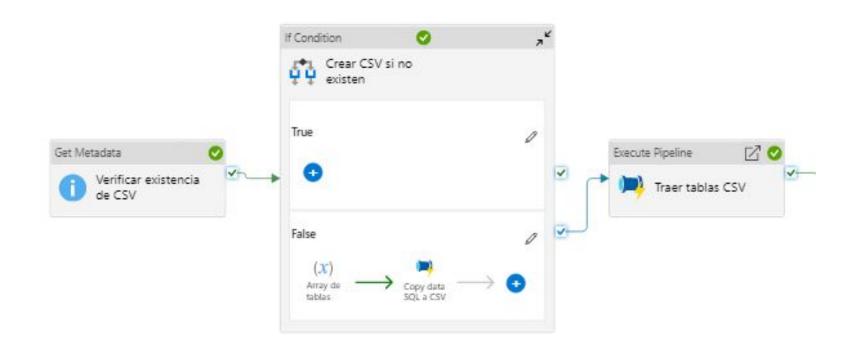


Lectura de datos de acceso

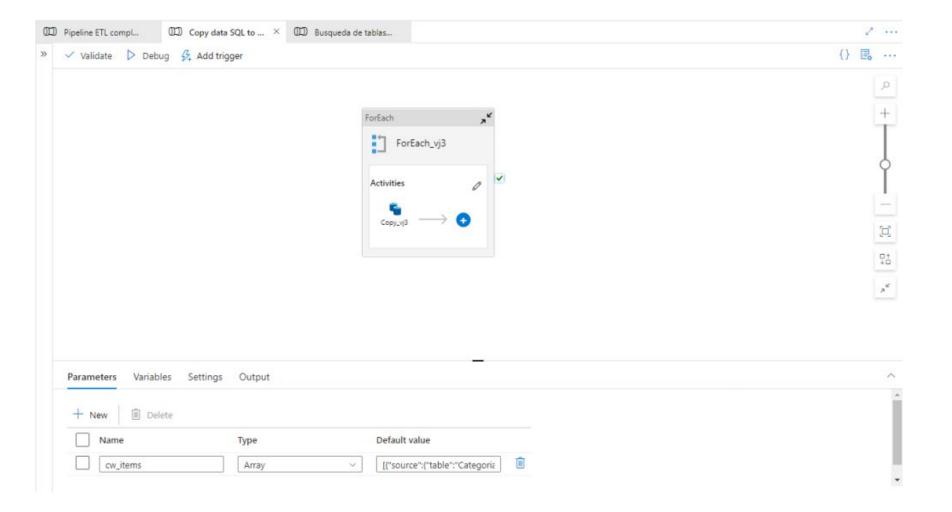
Vista previa del archivo del archivo de texto tras el Lookup activity:



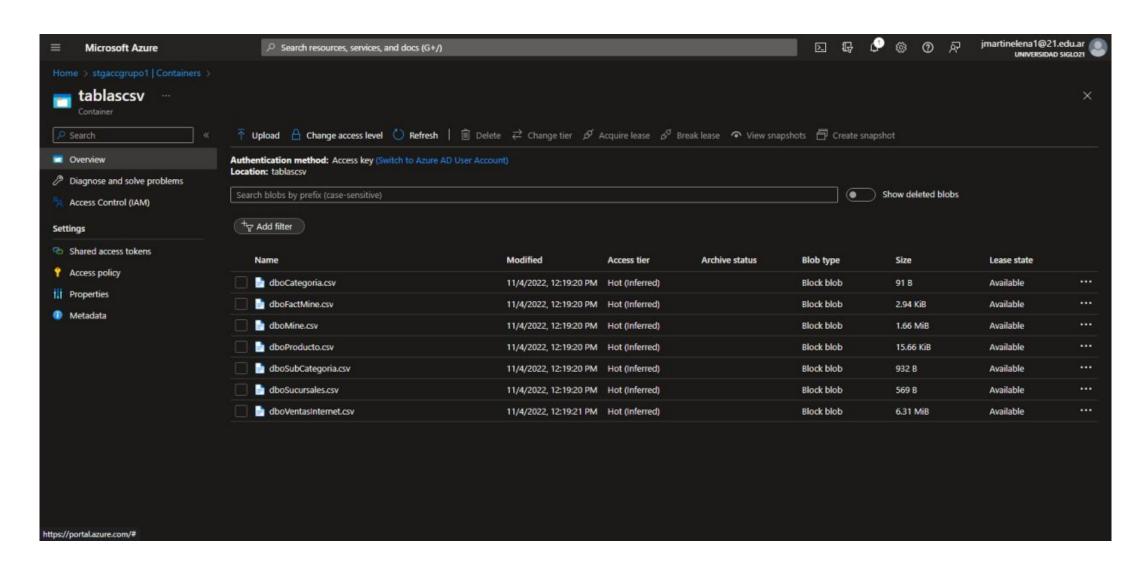
• El segundo paso verifica si las tablas ya existen en el Blob Storage para no ejecutar el Copy Data de no ser necesario.



• Luego, convierte las tablas de la base de datos a tablas en archivos csv (en caso de que no existan). Este ForEach se integra al Pipeline ETL Completo a través de un Pipeline Activity.

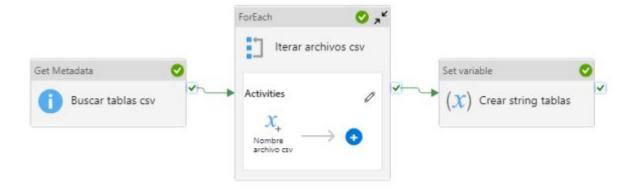


Los archivos csv se depositan en un container en el storage account.



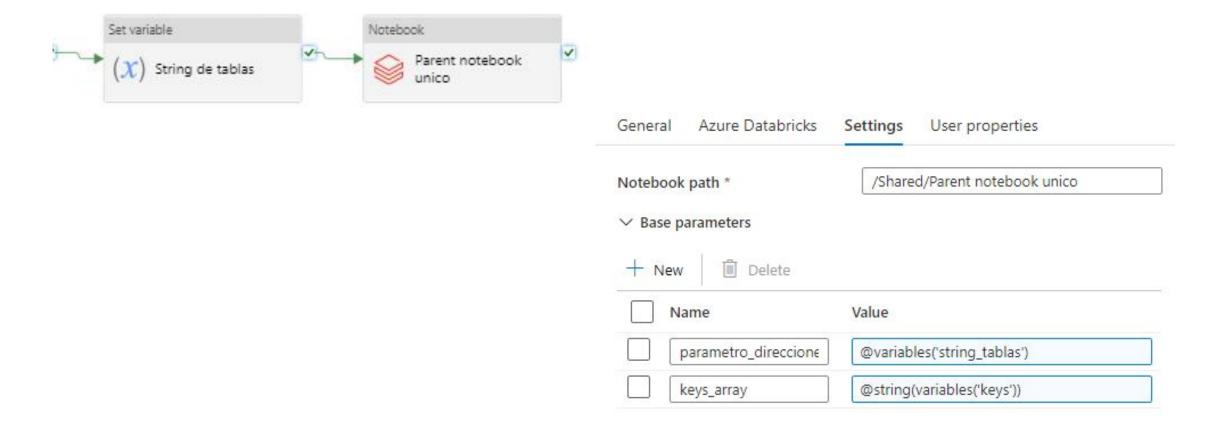
Después, se preparan las tablas para enviar a Databricks como parámetro.





Envío de tablas csv a Databricks

• En este último paso, se envían las tablas y parámetros a una notebook de Databricks.



> Transformaciones en Databricks

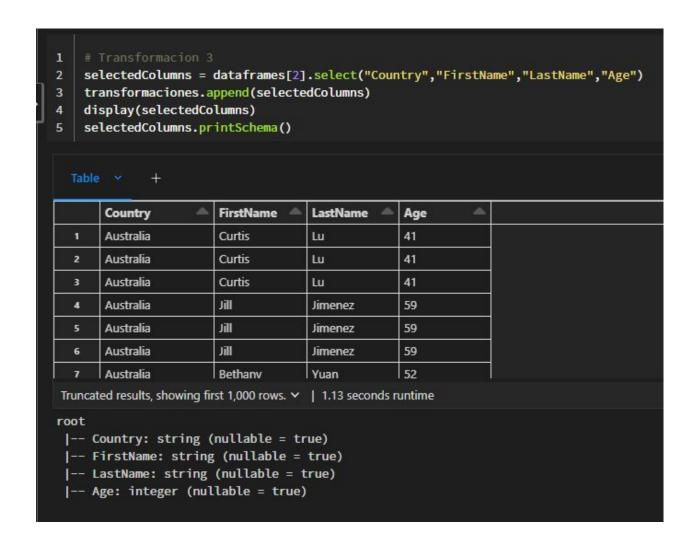
Transformación de tablas en Databricks

 Tomamos los csv y realizamos las transformaciones. Optamos por usar un solo notebook, en lugar de un notebook padre que llama a notebooks hijos, para mejorar los tiempos de ejecución.

```
Microsoft Azure @ databricks
                                                                                                                                                                              ② jmartinelena1@21.edu.ar ~
      Parent notebook unico Python Y
                                                                                                                                                                  ● Connect ∨
                                                                                                                                                                                               Share
       File Edit View Run Help Last edit was 2 hours ago Give feedback
                port pyspark.sql.functions as F
             from pyspark.sql.types import IntegerType
              from pyspark.sql.types import DecimalType
             import pandas as pd
0
       Cmd 3
a
         2 storage_account_name = dbutils.widgets.get('storage_account_name')
a
            storage_account_access_key = dbutils.widgets.get('storage_account_access_key')
            spark.conf.set('fs.azure.account.key.' + storage_account_name + '.blob.core.windows.net', storage_account_access_key)
A
         Command took 0.65 seconds -- by a user at 1/11/2022, 17:07:09 on unknown cluster
             blob_container = 'tablascsv'
             parametro_direcciones = dbutils.widgets.get('parametro_direcciones')
             direcciones = parametro_direcciones.split(',')
            dataframes = []
             for i in range(0, len(direcciones)):
23
                 if (directiones[i] == 'dbolanding_tables.csv' or directiones[i] == 'dboSucursales.csv'):
        12
1/4
        13
        14
                     filePath = "wasbs://" + blob_container + "@" + storage_account_name + f".blob.core.windows.net/{direcciones[i]}"
                     dataframes.append(spark.read.format("csv").load(filePath, inferSchema = True, header = True))
```

Ejemplo de Transformación 3

• A modo de ejemplo, se muestra la Transformación 3 en Databricks:



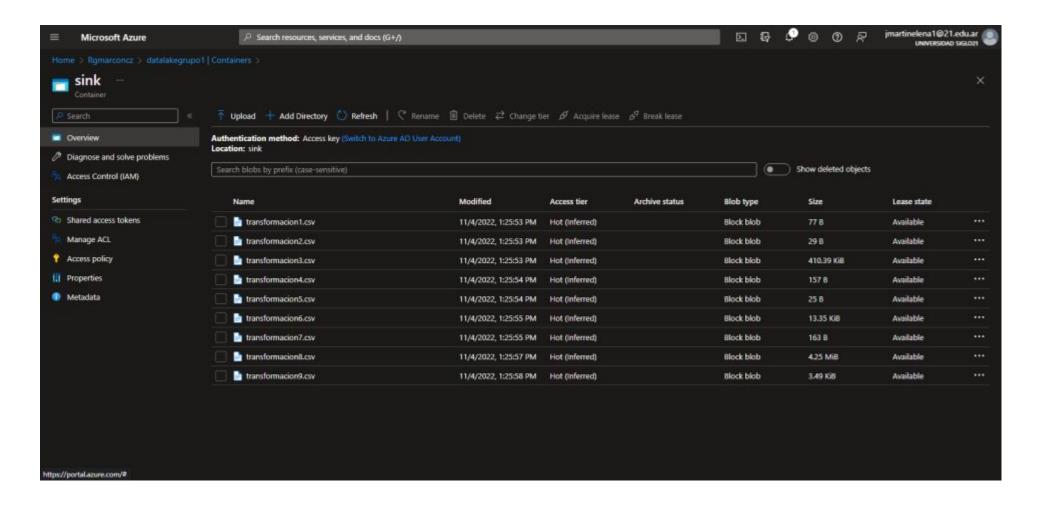
Salida al Data Lake Storage

• Configuramos la salida desde Databricks al Data Lake Storage, y escribimos las tablas en formato CSV.

```
| # Cargar las tablas en formato CSV dentro del data lake
| datalake_container = 'sink' |
| datalake_account_name = dbutils.widgets.get('datalake_account_name') |
| datalake_access_key = dbutils.widgets.get('datalake_access_key') |
| for i in range(len(transformaciones)): |
| transformaciones[i].toPandas().to_csv(f'abfs://{datalake_container}@{datalake_account_name}.dfs.core.windows.net/transformacion{i+1}.csv',storage_options = {'account_key': datalake_access_key} ,index=False) |
```

Salida de Data Lake

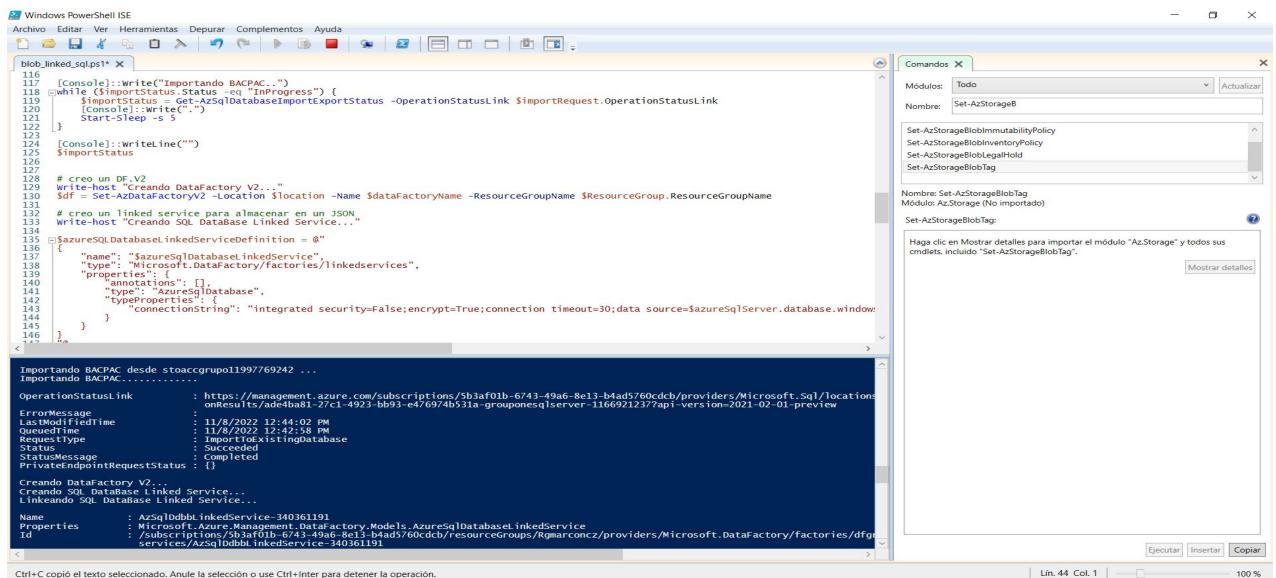
Tablas transformadas en formato CSV en el Data Lake Storage:



Adicionales

Script para automatizar creación de recursos

https://youtu.be/b8gvkS6xk9Y

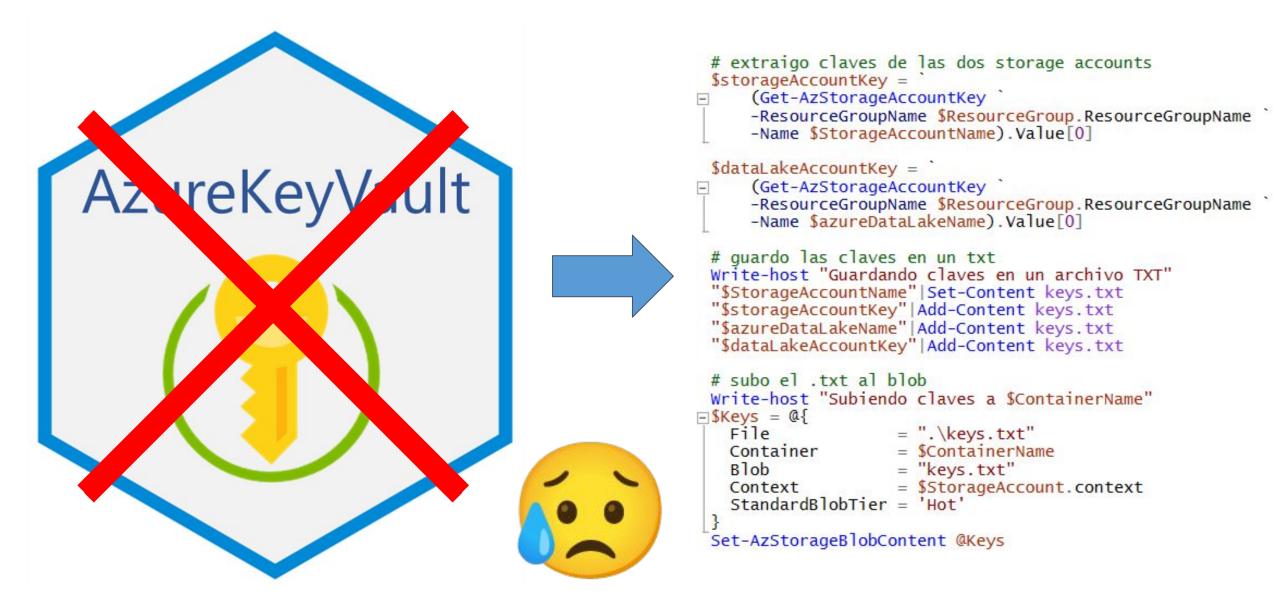


Script para automatizar creación de recursos

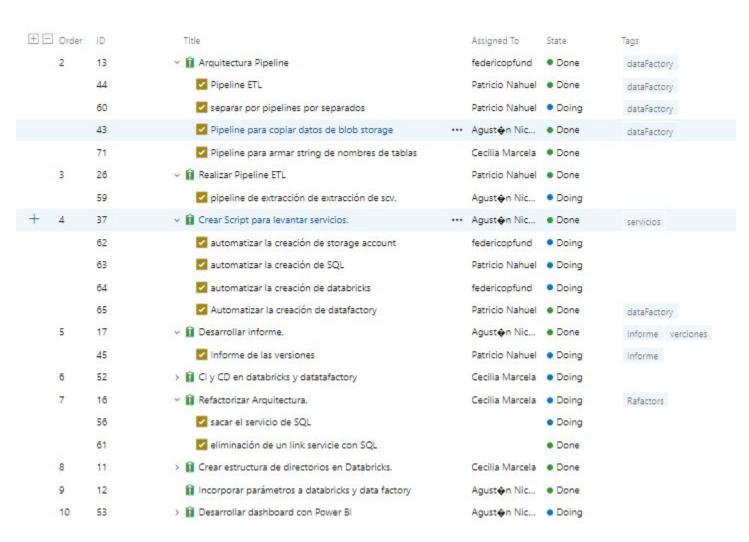
https://youtu.be/b8gvkS6xk9Y

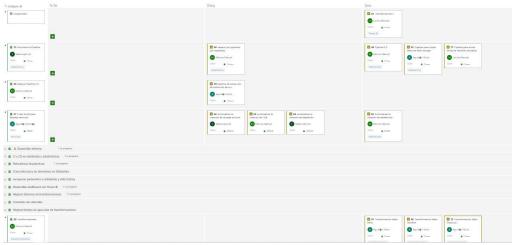
Name ↑↓	Type ↑↓	Location ↑↓
dbricksgrupo11270365620	Azure Databricks Service	West US 3
dfgrupo1siglo21367831347	Data factory (V2)	West US 3
dlaccgrupo11160844581	Storage account	West US 3
grouponesqlserver-331627171	SQL server	West US 3
groupOneSqlServerDatabase (grouponesqlserver-331627171/groupOneSqlServerDatabase)	SQL database	West US 3
stoaccgrupo11803122672	Storage account	West US 3

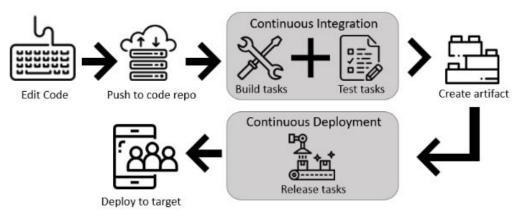
¿Y las claves del almacenamiento? Simulando el Keyvault



Integración con Azure DevOps







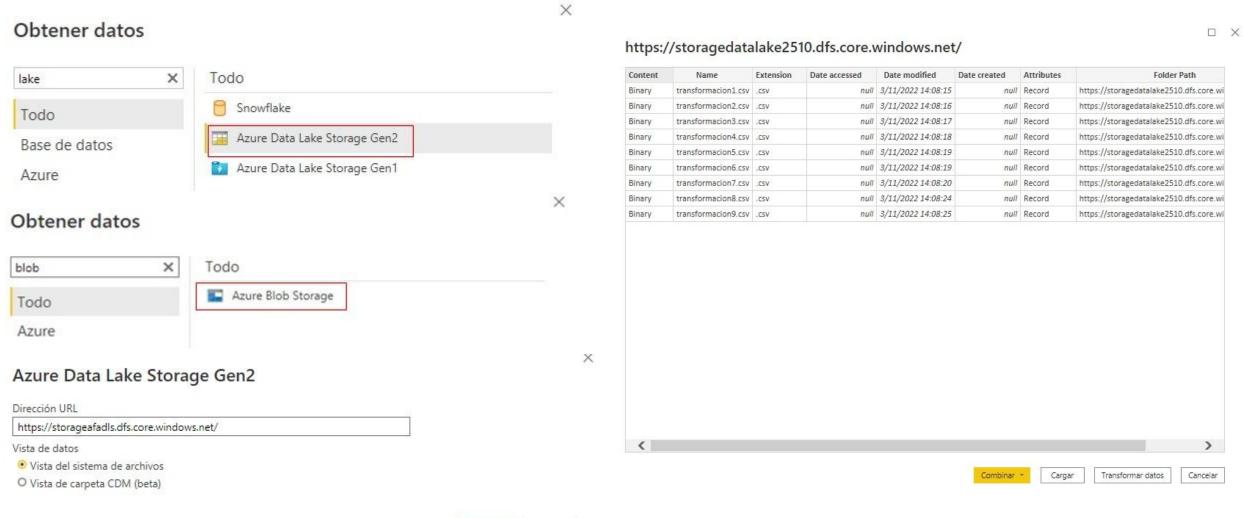
Visualización de Datos con Power BI

Visualización de datos en PowerBi

Conexión de blob storage y Data Lake Storage con PowerBi Desktop

Aceptar

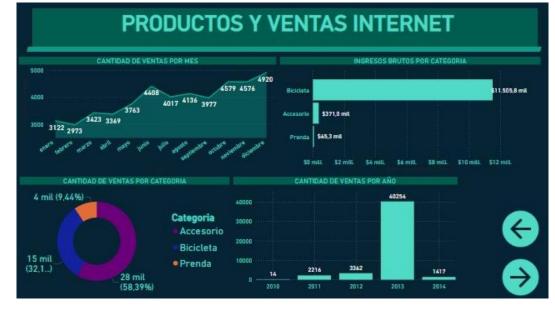
Cancelar



Dashboard Interactivo









MUCHAS GRACIAS M