

# **Assignment #1**

## **Data Warehouse Implementation Project**

**ITMD-526**



# **ILLINOIS TECH**



**ILLINOIS INSTITUTE OF TECHNOLOGY**

**Ignacio Hidalgo Power**

## Summary of the project

The objective of the project was to implement a data warehouse that could connect to an external, azure based, storage unit that held the data sources.

The first step was gathering the data. In order to have reliable, up to date, real world data we used the Bureau of Transport Statistics from the US. Once we had the data and we stored it into a container on the Azure platform we implemented through Databricks an ETL that connected to the Container through a SAS.

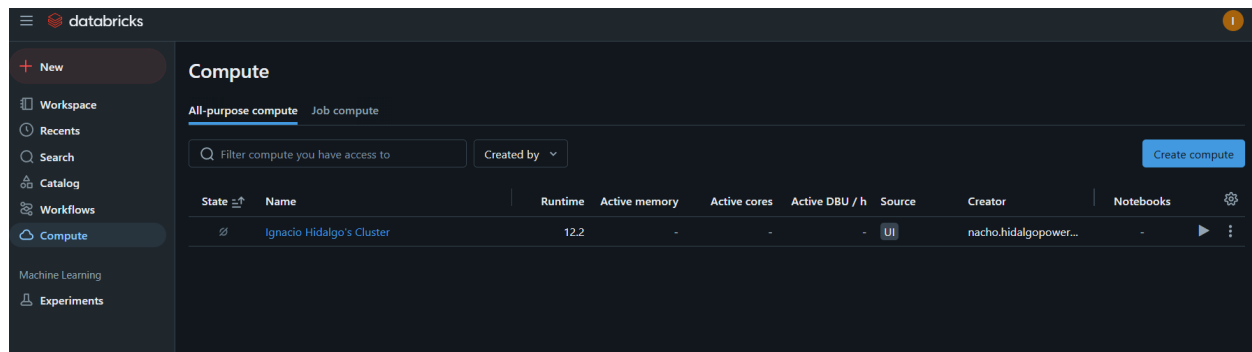
The screenshot shows the 'Containers' page for a storage account named 'datawarehouse10'. The left sidebar contains a navigation menu with options like Overview, Activity log, Tags, Diagnose and solve problems, Access Control (IAM), Data migration, Events, Storage browser, Storage Mover, Partner solutions, Data storage, Containers (selected), File shares, Queues, Tables, Security + networking, Networking, and Front Door and CDN. The main area has a search bar and a table of containers. The table has columns for Name, Last modified, Anonymous access le..., and Lease state. Two containers are listed: '\$logs' and 'demo-container', both created on 3/7/2025.

Name	Last modified	Anonymous access le...	Lease state
\$logs	3/7/2025, 6:57:25 PM	Private	Available
demo-container	3/7/2025, 6:58:34 PM	Container	Available

The screenshot shows the 'demo-container' page. The left sidebar has options like Overview (selected), Diagnose and solve problems, Access Control (IAM), Settings, Shared access tokens, Access policy, Properties, and Metadata. The main area shows the container's details, including the authentication method (Access key) and location (demo-container). Below this is a search bar for blobs and a table of blobs. The table has columns for Name, Modified, Access tier, Archive status, Blob type, Size, and Lease state. Four blobs are listed: 'Border\_Crossing.csv', 'Freight\_Shipments.csv', 'helloworld.txt', and 'Port\_Data.csv', all created on 3/12/2025.

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
Border_Crossing.csv	3/12/2025, 8:21:50 PM	Hot (Inferred)		Block blob	44.05 MiB	Available
Freight_Shipments.csv	3/12/2025, 8:21:47 PM	Hot (Inferred)		Block blob	318.83 KiB	Available
helloworld.txt	3/11/2025, 2:04:01 PM	Hot (Inferred)		Block blob	12 B	Available
Port_Data.csv	3/12/2025, 8:21:47 PM	Hot (Inferred)		Block blob	607.82 KiB	Available

Once the connection was through, we stored the CSV files onto a Spark Data Frame and preprocessed the data on the Spark Cluster to make use of the parallelization that is supported by Spark Data Frames, instead of the pandas one.



When the preprocessing was done, we normalized the data. Our approach, to saving resources and make use of the optimization that databricks allows, consisted on doing the preprocessing in the data frame and then convert it into a SQL Table, stored in our Logistics\_DB. This made use of the functionality that Spark offers of transforming the data frame into a table and not requiring a DDL script.

```
from pyspark.sql.functions import col, regexp_extract

df_shipment_filtered = df_shipment.select(
    col("Mode"),
    col("Statistic"),
    col("Year"),
    col("Value"),
    col("Units"),
    col("Statistic Short Name").alias("StatisticShortName")).filter(
    (col("Mode")=="Water") |
    (col("Mode")=="Rail") |
    (col("Mode")=="Air and truck-air") |
    (col("Mode")=="Air") |
    (col("Mode")=="Pipeline") |
    (col("Mode")=="Highway"))
```

Once we had the tables stored we made the Dimensions and Facts tables. One dimension table for each of the different sources, Ports, Shipments and Borders, and fact tables for the Ports that are at each border and another for trades.

### **Challenges:**

There were many challenges throughout the project. The first was the Azure Blob Storage that we wanted to implement. We first had in mind to use Azure Databricks, thanks to the ease of use, rich resources online, and ease of integration with the storage. However, we were unable to use it due to the nature of our Azure subscription and the scope that was limited to us.

When the Azure Databricks were unable for us, we moved onto Databricks. To improve performance and security we tried to launch Databricks locally to increase the number of resources we could use and to make use of the key vaults that can be used and give extra security to our cluster, but because we were using the community edition, we were also unable to do this.

Thankfully Databricks community edition supports a VM of 15GB of memory and 2 cores for free, so we could use that for our cluster. However, one thing that this VM had was that the way to connect to the Azure storage was through the web using the SAS key anonymously. But because of our lack of experience with Azure we had set the security requirements to enter the container too strict and for some days we were stuck being unable to connect. To fix this we had to create another storage unit with another container which had a less strict security requirement.

Our next challenge came with preprocessing. The data we gathered was quite dirty and dispersed, so we were forced to filter with the data frames functionality all of data, putting limits and taking away rows that shouldn't be there.

Our last challenge was that when deploying the normalized data we had to delete and create a new compute cluster due to the Databricks Community policy. When this happened all the memory databases and tables were removed, but the cache stored the names of the tables. To be able to restart and re run the DDL scripts that generated our tables we made use of the SQL syntaxis 'CREATE OR REPLACE TABLE', that way when the memory saw that the name was already in use it would just replace it.

**Analysis:**

After seeing the implementation of the dimensional model, there are possible improvements to do. For future development we should be able to create more facts tables and further increase the number of dimensions tables.

A goal is to generate tables that specialize in certain queries that can be used for analysis. Along with using those new tables we would also like to improve the optimization of the table through the use of SQL techniques, like indexes and temp tables.

**Future Improvements:**

Our main goals for future improvements are:

- Improve security connection: We want to improve the security of our cloud storage, so that if it were to store sensitive data it would comply with strict security metrics.
- DB Schema: We feel that the schema design can be improved as if to allow for better analysis, faster queries and more robust data quality.
- Report Generation: We want to implement a report generation that automatically gathers the data from the source if there has been an update, preprocess it, analyses the data, generates the report and stores it in a Azure Container in a different storage account.