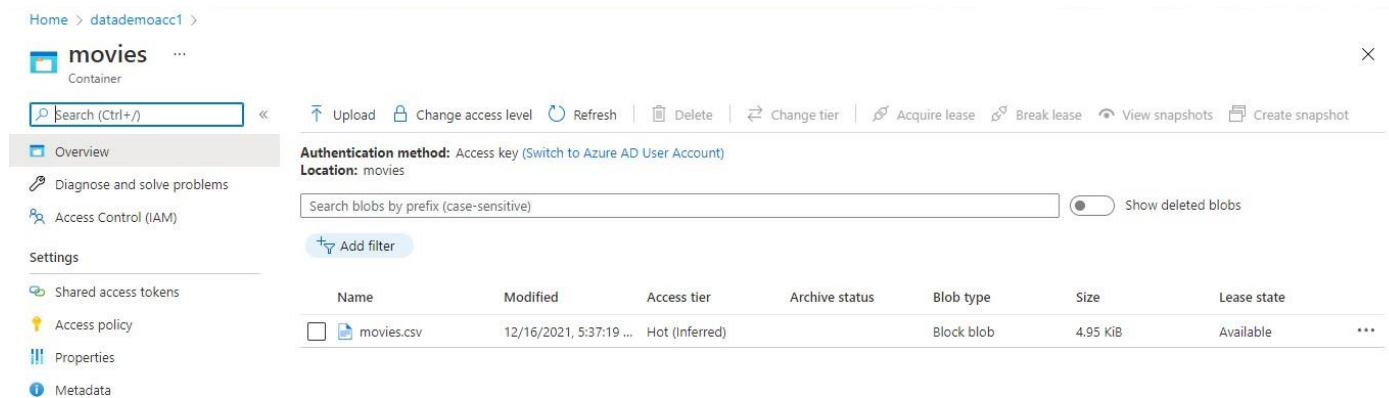# Building an ETL Pipeline using Azure Data Services

## Prerequisites:

**Step 1:** Create a blob storage, then create a container named "movies" and put the dataset (movies.csv) into the container as below:



**Step 2:** Create the SQL database "appdb" and create a table named "movies" and "agg_movies" :

```sql
CREATE TABLE movies
(
    film  varchar(200) NULL,
    genre varchar(200) NULL,
    lead_studio varchar(200) NULL,
    audience_score int NULL,
    profitability real NULL,
    rotten_tomatoes int NULL,
    worldwide_gross varchar(20) NULL,
    year varchar(4) NULL
);

CREATE TABLE agg_movies
(
    film  varchar(200) NULL,
    genre varchar(200) NULL,
    lead_studio varchar(200) NULL,
    audience_score int NULL,
    profitability real NULL,
    rotten_tomatoes int NULL,
    worldwide_gross real NULL,
    year varchar(4) NULL,
    film_count int NULL
);
```

**Step 3:** Create a Data Lake Gen 2 Storage, then create a container named "movies" to store the processed data after cleaning and aggregating as shown below:

**movies** ...
Container

✕

🔍 Search (Ctrl+/)   ≪

⬆ Upload   ➕ Add Directory   🔄 Refresh   |   ⟲ Rename   🗑 Delete   |   ⇄ Change tier   |   🔓 Acquire lease   🔓 Break lease

▢ Overview

**Authentication method:** Access key (Switch to Azure AD User Account)
**Location:** movies

🔗 Diagnose and solve problems

🔐 Access Control (IAM)

Search blobs by prefix (case-sensitive)                                    ⬤ Show deleted objects

**Settings**

| Name | Modified | Access tier | Archive status | Blob type | Size | Lease state |
|------|----------|-------------|----------------|-----------|------|-------------|

☁ Shared access tokens

🔐 Manage ACL

No blobs found.

🔑 Access policy

III Properties

ℹ Metadata

**Step 4:** Create a Synapse Analytics, the create a SQL pool named "appexample1234" and create a table named "movies" using the below query:

SQLQuery2.sql - ap...qladminuser (2165))*   📌 ✕   SQLQuery1.sql - ap...ppdb (sqluser (73))*

```sql
CREATE TABLE movies
(
    film  varchar(200) NULL,
    genre varchar(200) NULL,
    lead_studio varchar(200) NULL,
    audience_score int NULL,
    profitability real NULL,
    rotten_tomatoes int NULL,
    worldwide_gross real NULL,
    year varchar(4) NULL,
    film_count int NULL,
    rank int NULL
);
```

**Step 5:** Create a key vault, then navigate to secrets and create one secret to use in Azure Databricks

Home > demovault1234554

**demovault1234554 | Secrets** ...
Key vault

| Name | Type | Status | Expiration date |
|------|------|--------|-----------------|
| demovault | | ✓ Enabled | |
| demovault1 | | ✓ Enabled | |
| demovault2 | | ✓ Enabled | |

Sidebar:
- Overview
- Activity log
- Access control (IAM)
- Tags
- Diagnose and solve problems
- Events

Settings
- Keys
- Secrets
- Certificates
- Access policies

Toolbar: + Generate/Import  ↻ Refresh  ⤒ Restore Backup  🔗 Manage deleted secrets

Home > demovault1234554 >

# Create a secret ...

| | |
|---|---|
| Upload options | Manual |
| Name * ⓘ | |
| Value * ⓘ | Enter the secret. |
| Content type (optional) | |
| Set activation date ⓘ | ☐ |
| Set expiration date ⓘ | ☐ |
| Enabled | Yes   No |
| Tags | 0 tags |

Create

- Once created, Now launch your Azure databricks and append this to the URI "/#secrets/createScope ", It will open up a window like this



- DNS Name and Resource ID can be copied from key vault properties "Vault URI" and "Resource ID"



- Make sure you remember the scope name created in databricks since it will be used to connect with ADLS in the later part.

## Proceedings:

**Step 1:** Importing dataset from blob to DB



**Step 2:** Create a dataflow to clean and join the table as required

- Get the Input Data from SQL

- Remove '$' from worldwide_gross and make it float using query in source options



**Step 3:** Clean the Genre column which has typo mistakes as required

- Cleaning using expression builder

**Column name** *

genre

**Expression**                                                                          💾 Save

    initCap(case(lower(genre)=="romence","romance",case(lower(genre)=="comdy", "comedy", lower(genre))))

- You can see that genre column is updated



**Step 4:** Find the count of films based on genre of the film using aggregate functions

**Step 5:** Join the film count(step 4) with existing clean data(step 3) using genre column



**Step 6:** Now that we have joined two tables, we'll have two genre columns which is there in both table, we can remove this using select function

**Step 7:** Store the resultant data in "agg_movies" table which we already created



**Step 8:** Grab the dataflow to the existing pipeline, then add a copy data to transfer the table data to ADLS in PARQUET FORMAT

**Step 9:** Create a notebook to rank the data existing in ADLS and store it back in ADLS in the same parquet format

- To do this, you'll need to create a key vault and connect it to the notebook to make the ADLS accessible
- A new folder "finalOutput" will be created inside the movies container in ADLS.
- Scope and key mentioned here are created as part of prerequisites



MoviesAssignment (Python)

Free trial ends in 4 days. Upgrade to Premium in Azure Portal

Schedule ∨

newCluster

Cmd 1

```python
spark.conf.set(
    "fs.azure.account.key.dataforbigdata1.dfs.core.windows.net",
    dbutils.secrets.get(scope="newScope",key="demovault1"))

df = spark.read.format("parquet").option("header","true").load("abfss://movies@dataforbigdata1.dfs.core.windows.net/")

df.registerTempTable("movies")

resultset = spark.sql("SELECT distinct *, DENSE_RANK() OVER(PARTITION BY genre ORDER BY profitability DESC) AS rank from movies order by genre desc,profitability desc,rank;")

resultset.write.format("parquet").option("header","true").save("abfss://movies@dataforbigdata1.dfs.core.windows.net/finalOutput")
```

Shift+Enter to run

**Step 10:** Now grab the notebook created and connect it to the existing pipeline such that the notebook execution happens only in the sequential order of the pipeline (applies to all the steps)



**Step 11:** Grab a copy data which performs copying the dataset in ADLS to SQL POOL in Synapse Analytics

- Data will be inserted to the table "movies" that is already created in SQL POOL

# Result Screenshots:

- **Pipeline**



- **Data Flow**

- **Pipeline Successful**



- **Result data**



```
select * from movies;
```

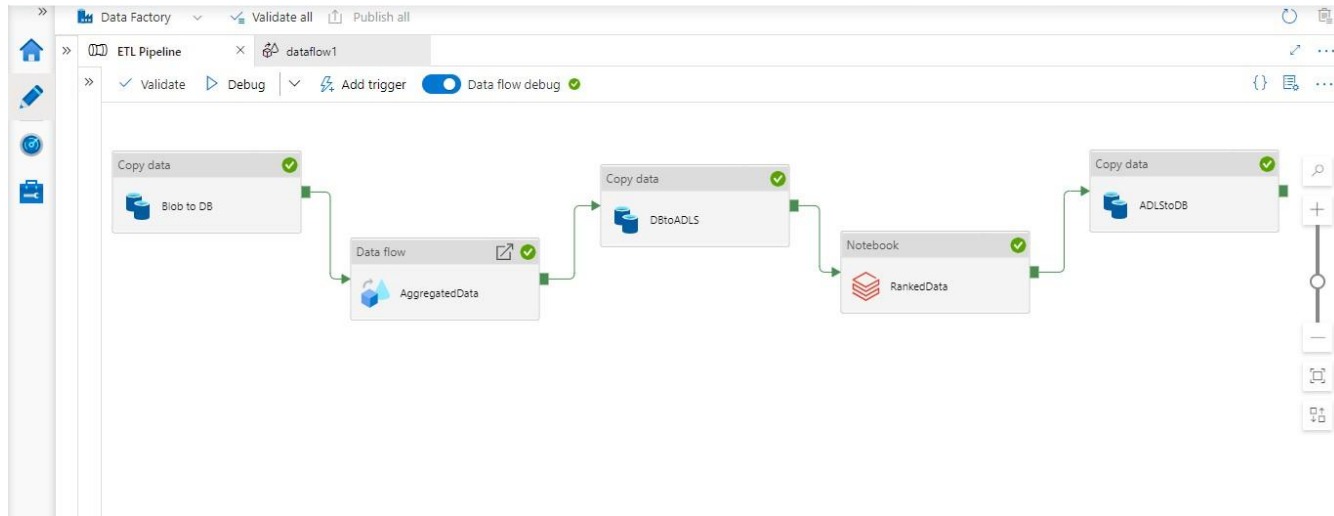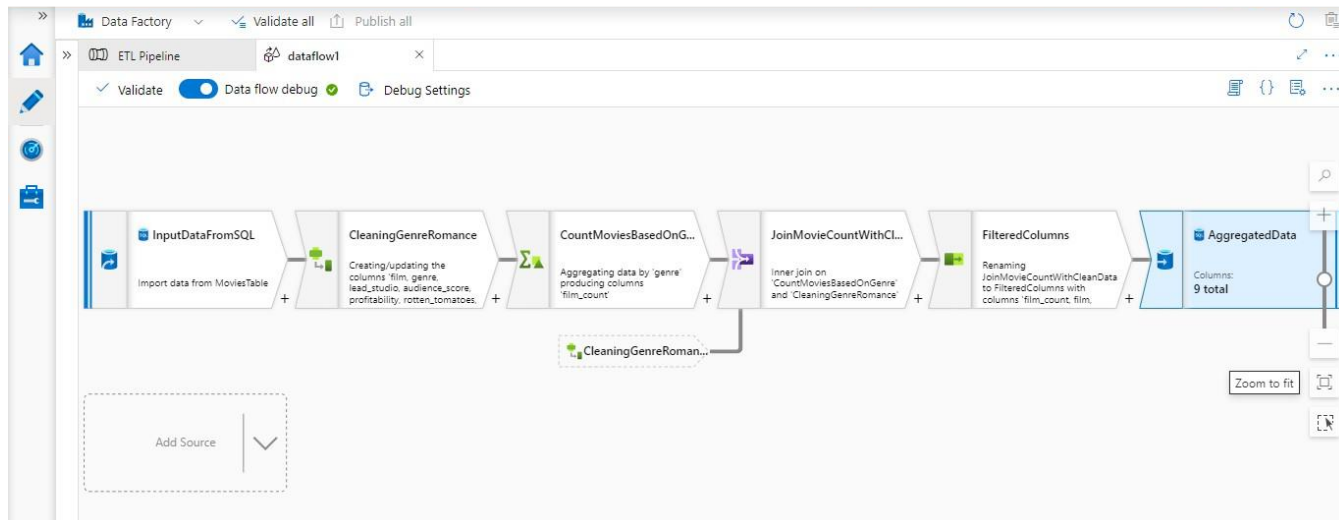| | film | genre | lead_studio | audience_score | profitability | rotten_tomatoes | worldwide_gross | year | film_count | rank |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Waitress | Romance | Independent | 67 | 11.08974 | 89 | 22.18 | 2007 | 15 | 1 |
| 2 | Twilight | Romance | Summit | 82 | 10.18003 | 49 | 376.66 | 2008 | 15 | 2 |
| 3 | Midnight in Paris | Romance | Sony | 84 | 8.744706 | 93 | 148.66 | 2011 | 15 | 3 |
| 4 | Twilight: Breaking Dawn | Romance | Independent | 68 | 6.383364 | 26 | 702.17 | 2011 | 15 | 4 |
| 5 | P.S. I Love You | Romance | Independent | 82 | 5.103117 | 21 | 153.09 | 2007 | 15 | 5 |
| 6 | Tyler Perry's Why Did I get Married | Romance | Independent | 47 | 3.724192 | 46 | 55.86 | 2007 | 15 | 6 |
| 7 | One Day | Romance | Independent | 54 | 3.682733 | 37 | 55.24 | 2011 | 15 | 7 |
| 8 | Music and Lyrics | Romance | Warner Bros. | 70 | 3.647411 | 63 | 145.9 | 2007 | 15 | 8 |
| 9 | New Year's Eve | Romance | Warner Bros. | 48 | 2.536429 | 8 | 142.04 | 2011 | 15 | 9 |
| 10 | Monte Carlo | Romance | 20th Century Fox | 50 | 1.9832 | 38 | 39.66 | 2011 | 15 | 10 |
| 11 | Zack and Miri Make a Porno | Romance | The Weinstein Company | 70 | 1.747542 | 64 | 41.94 | 2008 | 15 | 11 |
| 12 | Something Borrowed | Romance | Independent | 48 | 1.719514 | 15 | 60.18 | 2011 | 15 | 12 |
| 13 | Across the Universe | Romance | Independent | 84 | 0.65260... | 54 | 29.37 | 2007 | 15 | 13 |
| 14 | Waiting For Forever | Romance | Independent | 53 | 0.005 | 6 | 0.03 | 2011 | 15 | 14 |
| 15 | Jane Eyre | Romance | Universal | 77 | 0 | 85 | 30.15 | 2011 | 15 | 15 |
| 16 | The Curious Case of Benjamin B... | Fantasy | Warner Bros. | 81 | 1.783944 | 73 | 285.43 | 2008 | 1 | 1 |
| 17 | Fireproof | Drama | Independent | 51 | 66.934 | 40 | 33.47 | 2008 | 13 | 1 |
| 18 | The Twilight Saga: New Moon | Drama | Summit | 78 | 14.1964 | 27 | 709.82 | 2009 | 13 | 2 |
| 19 | Dear John | Drama | Sony | 66 | 4.5988 | 29 | 114.97 | 2010 | 13 | 3 |
| 20 | A Serious Man | Drama | Universal | 64 | 4.382857 | 89 | 30.68 | 2009 | 13 | 4 |
| 21 | Remember Me | Drama | Summit | 70 | 3.49125 | 28 | 55.86 | 2010 | 13 | 5 |
| 22 | The Duchess | Drama | Paramount | 68 | 3.20785 | 60 | 43.31 | 2008 | 13 | 6 |
| 23 | Water For Elephants | Drama | 20th Century Fox | 72 | 3.081421 | 60 | 117.09 | 2011 | 13 | 7 |
| 24 | The Time Traveler's Wife | Drama | Paramount | 65 | 2.598205 | 38 | 101.33 | 2009 | 13 | 8 |

Query executed successfully.      appexample1234.sql.azuresyn...  sqladminuser (144)  appexampleDedicat