# Pós-graduação em Data Science & Business Analytics Formato Blended 2ª Edição - 2022

ISEG Executive Education

Data Warehousing

Artur Vieira

# Índice

# Azure Data Factory Concepts

## Data Factory Copy Activity

In Azure Data Factory, you can use the Copy activity to copy data among data stores located on-premises and in the cloud. After you copy the data, you can use other activities to further transform and analyze it. You can also use the Copy activity to publish transformation and analysis results for business intelligence (BI) and application consumption.



**The Copy activity is executed on an integration runtime**. You can use different types of integration runtimes for different data copy scenarios:

- When you're copying data between two data stores that are publicly accessible through the internet from any IP, you can use the **Azure integration runtime** for the copy activity. This integration runtime is secure, reliable, scalable, and globally available.

- When you're copying data to and from data stores that are located on-premises or in a network with access control (for example, an Azure virtual network), you need to set up a **self-hosted integration runtime**.
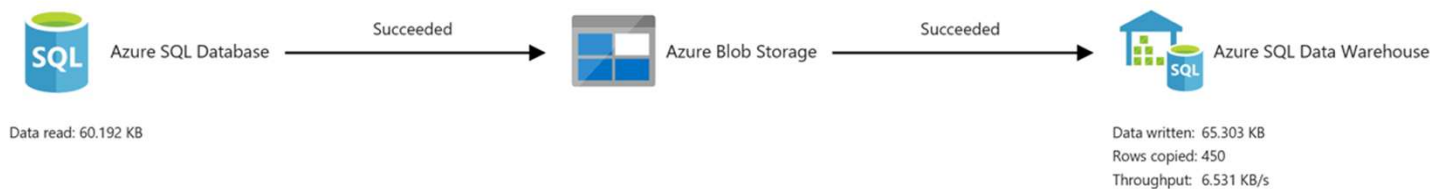
# Azure Data Factory Concepts

## Data Factory Copy Activity

To copy data from a source to a sink, the service that runs the Copy activity performs these steps:

- **Reads data from a source data store.**

- **Performs serialization/deserialization, compression/decompression, column mapping, and so on**. It performs these operations based on the configuration of the input dataset, output dataset, and Copy activity.

- **Writes data to the sink/destination data store.**

# Azure Data Factory Concepts

## Copy Activity - Supported file formats

You can use the Copy activity to copy files as is between two file-based data stores. In this case, the data is copied efficiently without any serialization or deserialization.

Azure Data Factory support the following file formats. Refer to each article on format-based settings.

- Avro format
- Binary format
- Delimited text format
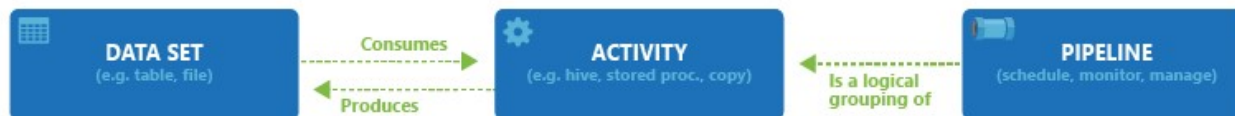- JSON format
- ORC format
- Parquet format

### Supported data stores and formats

| Category | Data store | Supported as a source | Supported as a sink | Supported by Azure IR | Supported by self-hosted IR |
|---|---|---|---|---|---|
| Azure | Azure Blob storage | ✓ | ✓ | ✓ | ✓ |
| | Azure Cosmos DB (SQL API) | ✓ | ✓ | ✓ | ✓ |
| | Azure Cosmos DB's API for MongoDB | ✓ | ✓ | ✓ | ✓ |
| | Azure Data Explorer | ✓ | ✓ | ✓ | ✓ |
| | Azure Data Lake Storage Gen1 | ✓ | ✓ | ✓ | ✓ |
| | Azure Data Lake Storage Gen2 | ✓ | ✓ | ✓ | ✓ |
| | Azure Database for MariaDB | ✓ | | ✓ | ✓ |
| | Azure Database for MySQL | ✓ | ✓ | ✓ | ✓ |
| | Azure Database for PostgreSQL | ✓ | ✓ | ✓ | ✓ |
| | Azure File Storage | ✓ | ✓ | ✓ | ✓ |
| | Azure SQL Database | ✓ | ✓ | ✓ | ✓ |
| | Azure SQL Database | ✓ | ✓ | ✓ | ✓ |

# Azure Data Factory Concepts

## Data Factory Pipeline, Activities Datasets and Linked Services

- A data factory can have one or more pipelines.

- **A pipeline is a logical grouping of activities that together perform a task**. For example, a pipeline could contain a set of activities that ingest and clean log data, and then kick off a mapping data flow to analyze the log data. The pipeline allows you to manage the activities as a set instead of each one individually. You deploy and schedule the pipeline instead of the activities independently.

- **The activities in a pipeline define actions to perform on your data**. For example, you may use a copy activity to copy data from an on-premises SQL Server to an Azure Blob Storage. Then, use a data flow activity or a Databricks Notebook activity to process and transform data from the blob storage to an Azure Synapse Analytics pool on top of which business intelligence reporting solutions are built.

# Azure Data Factory Concepts

**Data Factory Pipeline, Activities Datasets and Linked Services**

- **Data Factory has three groupings of activities: data movement activities, data transformation activities, and control activities**. An activity can take zero or more input datasets and produce one or more output datasets.

    - **Data movement activities:** Copy Activity in Data Factory copies data from a source data store to a sink data store. Data Factory supports the data stores listed in https://docs.microsoft.com/en-us/azure/data-factory/concepts-pipelines-activities. Data from any source can be written to any sink.

    - **Data transformation activities:** Azure Data Factory supports a great set of transformation activities that can be added to pipelines either individually or chained with another activity. (ex. Azure Function, Spark , Databricks Notebook, etc)

    - **Control flow activities:** Azure Data Factory supports a great set of control flow activities that can be added to pipelines (ex. For Each, Get Metadata, etc)

- An **input dataset** represents the input for an activity in the pipeline and an **output dataset** represents the output for the activity. **Datasets identify data within different data stores, such as tables, files, folders, and documents.** After you create a dataset, you can use it with activities in a pipeline.

# Azure Data Factory Concepts

**Data Factory Self-Hosted Integration Runtime**

- The Integration Runtime is a **customer managed data integration infrastructure used by Azure Data Factory to provide data integration capabilities across different network environments**. It was formerly called as Data Management Gateway.

- The integration runtime is **capable of moving data in and out of data stores within private network**, as well as dispatching activities against compute service within private network. **You can install a self-hosted integration runtime on an on-premises machine or a virtual machine inside a private network**. This was formerly called the Data Management Gateway (DMG) and is fully backward compatible. Note: An Integration Runtime instance can be registered with only one of the versions of Azure Data Factory (version 1 -GA or version 2 -GA).

# Azure Portal Resources – DW Resource Group



- Azure DataLake Store gen 2
- Azure Data Factory
- Self Hosted Virtual Machine
  - Disk
  - Network
  - Storage
- SQL Azure Operational Db
- SQL Azure Data Warehouse DB

# Navigate to the Azure Data Factory

**In this task, you will Navigate to the Azure**

**Data Factory.**

1. Go to https://adf.azure.com and select your new

Azure Data Factory:

2. In Azure Active Directory choose **IDEFE SA**

3. In Subscriptions Choose **--subscription--**

4. In Data Factory Name choose **--adf--**

5. Press the **Continue Button**

6. You will see the let's get started page

**Azure Data Factory UI features currently only support the Microsoft Edge e Google Chrome. Make sure you are using one of this browsers**

## Select Data Factory

Microsoft Azure Data Factory is a cloud-based data integration service that automates data movement and transformation. Learn more

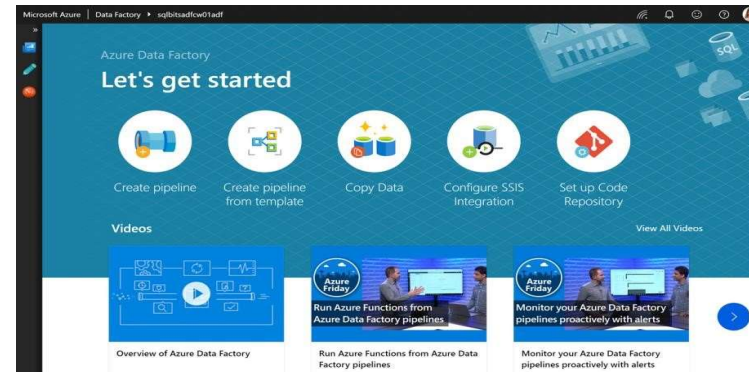**Azure Active Directory**
ISEG (9c184083-de20-4859-968c-96f26db0bfe1)

**Subscription**
ISEG Education Cloud (000a9103-fec4-4ce8-aba5-d7b0908f5aca)

**Data Factory name ***
dwintazdtf02

[Continue]

# Navigate to the Azure Data Factory

**In this task, you will validate that the Self-hosted integration runtime is set and ready to use**

1. Click on the ==Manage== button in the left panel

2. In the Connections tab, click on ==**Integration Runtime**==,

3. Validate that an Azure IR is set and that the actual status is running

if a Self-Hosted IR is also running:
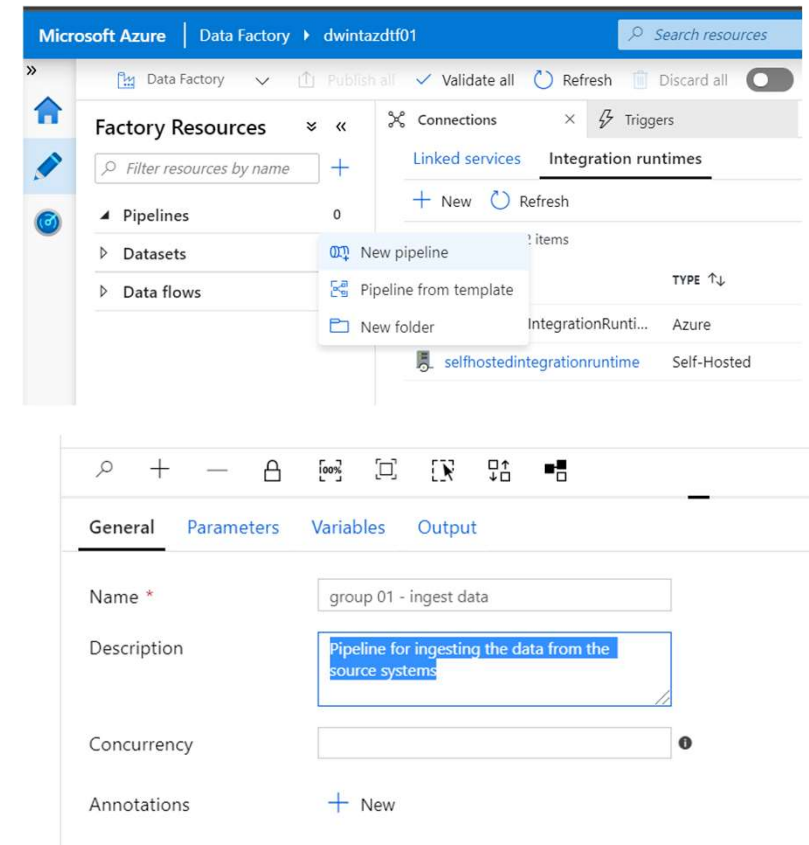
Linked services    Integration runtimes

+ New    ⟳ Refresh

Showing 1 - 2 of 2 items

| NAME ↑↓ | TYPE ↑↓ | SUB-TYPE ↑↓ | STATUS ↑↓ |
|---|---|---|---|
| AutoResolveIntegrationRunti... | Azure | Public | ✅ Running |
| selfhostedintegrationruntime | Self-Hosted | --- | ✅ Running |

# ADF Copy Activity – Configure Pipeline

**In this task, you will create your first Pipeline**

1. On the left pane click on the three dots **...** Right to the pipeline section

and select ==**New Pipeline  - if necessary create your personal working**==

==  **folder to work with /ixxxxxx**==

2. On the General TAB set the pipeline name property

to ==**iXXXXXX - ingest data**==

==**(names on images may be different)**==

3. In the Description field insert the value ==**Pipeline for ingesting the data**==
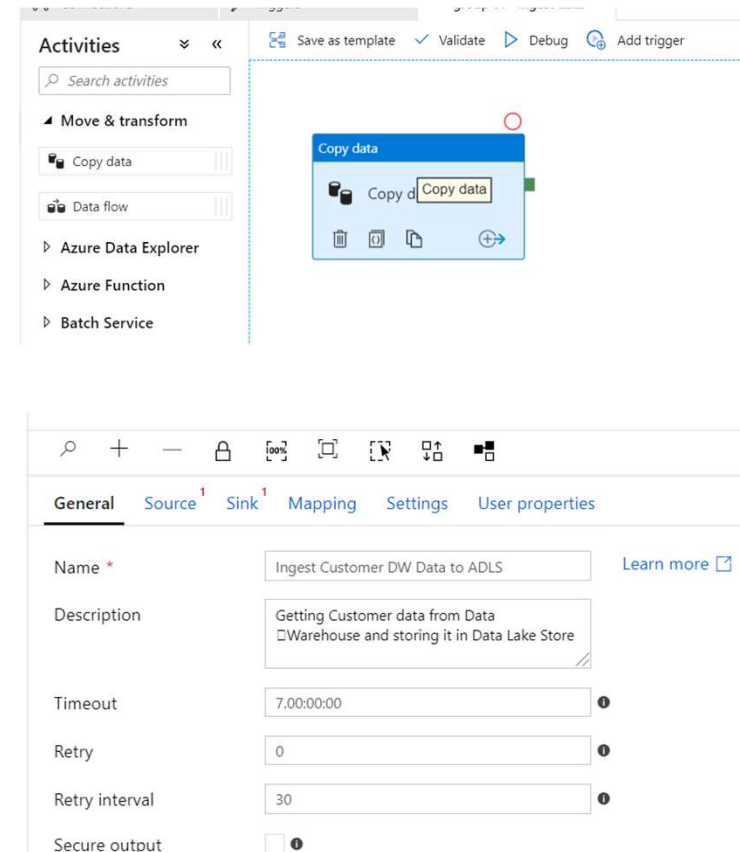
==**from the source systems**==

# ADF Copy Activity – Configure Pipeline

**In this task, you will add activities to the Pipeline created before**

1. In the activities list open the <mark>Move & Transform Section</mark>

2. Drag the <mark>Copy Data Activity</mark> into the <mark>Main Pane on the Right</mark>

3. <mark>Select the Copy data Activity</mark> on the main pane

4. Make sure the **General Tab Window** is selected

5. On the <mark>Name Property</mark> set the value <mark>Ingest Customer DW Data to ADLS</mark>

6. On the <mark>Description Property</mark> set the value <mark>Getting Customer data from Data Warehouse and storing it in Data Lake Store</mark>
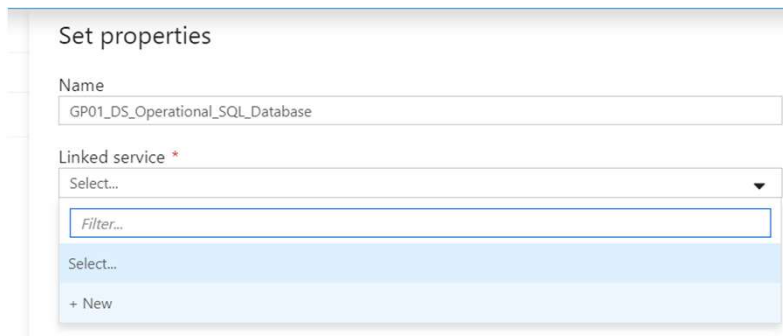
**The Copy activities allow**

# ADF Copy Activity – Configure Source

**In this task, you will set the Source settings for Copy Activity**

1. Choose the **Source TAB Window**

2. Press **New** button next to the Source Dataset Drop Down List

3. In the New Dataset Window select the **Database Tab**

4. Select the **SQL Server** Dataset and Select **Continue**

5. In the Set Properties Window insert **ixxxx_DS_Operational_SQLDB_Customers**

6. In the Linked Service Drop Down List select **New** option

# ADF Copy Activity – Configure Source

**In this task, you will set the Source settings for Copy Activity**

1. In the New linked Service (SQL Server Database) set Name value to

   **ixxxxxx_LS_Operacional_Database**

2. Set Description value to **Linked Service to fetch data from SQL Operational Database**

3. In the Connect Via Integration Runtime select **AutoResolveIntegrationRuntimee**

4. In servername field type **--sqlserver operacional--**

5. In Database field name type **--db operacional--**

6. Make sure the authentication is set to **SQL Authentication**

7. In the field username type **--user db operacional--**

8. Set **--psw db operacional--** As the Password field value

9. Click on the **Create** button

New linked service (Azure SQL Database)

Name *

bb

Description

Connect via integration runtime * ⓘ

AutoResolveIntegrationRuntime

Connection string    Azure Key Vault

Account selection method ⓘ

From Azure subscription    ● Enter manually

Fully qualified domain name *

daesqlsrv01.database.windows.net

Database name *

daesqldtb01

Authentication type *

SQL authentication

User name *

daereader

Password    Azure Key Vault

Password *

••••••••••

Add dynamic content [Alt+P]

Additional connection properties

+ New

**In this task, you will set the Source settings for Copy Activity**

1. In the Table name select dbo.DimCustomer

2. Make sure the option From Connection/store is selected in Import Schema field

3. Press the OK button

4. In the Source Tab press the Preview Data button and ensure data is displayed

5. Close the preview data window

# ADF Copy Activity – Configure Sink

**In this task, you will set the Sink settings for Copy Activity**

1. Choose the Sink TAB Window

2. Press New button next to the Sink Dataset Drop Down List

3. In the New Dataset Window select the Azure Tab

4. Select the Azure Data Lake Storage Gen 2 and Select Continue

5. In the select format Window select delimited

6. Press the Continue button

Tópico

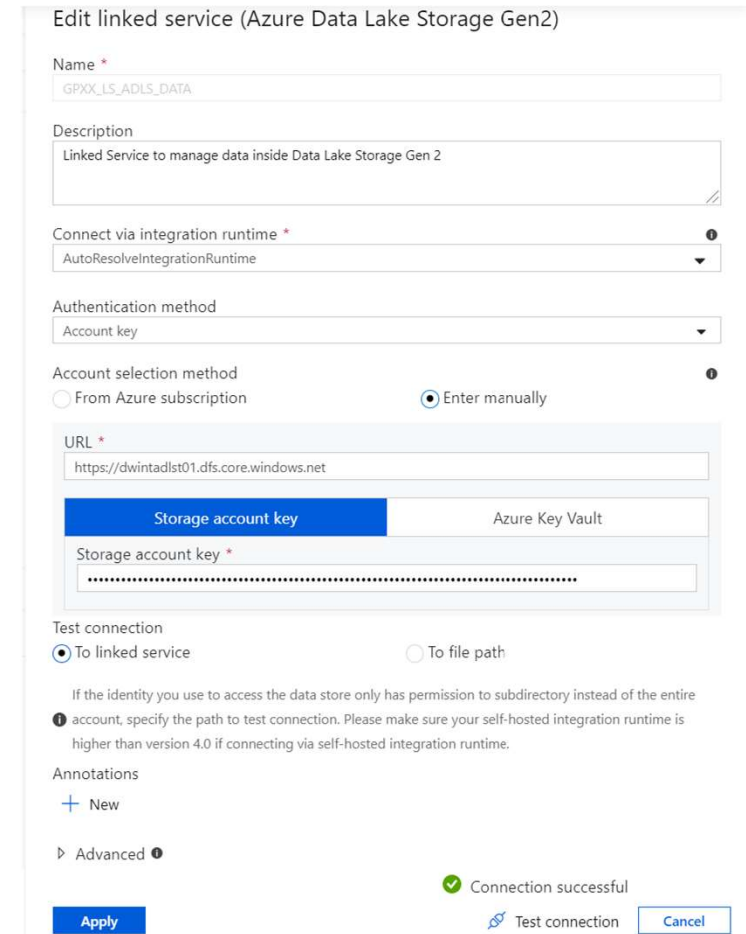**In this task, you will set the Sinc settings for Copy Activity**

1. In the set properties window set the name to **ixxxxxx_DS_RAW_ADLS_Customers**

2. In the Linked Service Drop Down List select **New** option

3. In the New linked Service (Azure Data Lake Storage Gen2) set Name value to **ixxxxxx_LS_ADLS_DATA**

4. Set Description value to **Linked Service to manage data inside Data Lake Storage Gen 2**

5. In the Connect Via Integration Runtime select **AutoResolveIntegrationRuntime**

6. In Authentication Method select **Account Key**

7. In Account Selection Method select **Enter manually**

8. In the URL set the value to **--adls--**

9. In the Storage Account Key insert **--adls key--**

10. Ensure Test Connection is set to **To Linked Service**

11. Press **Test Connection**

12. Press **Create** button

Edit linked service (Azure Data Lake Storage Gen2)

Name *
GPXX_LS_ADLS_DATA

Description
Linked Service to manage data inside Data Lake Storage Gen 2

Connect via integration runtime *
AutoResolveIntegrationRuntime

Authentication method
Account key

Account selection method
○ From Azure subscription    ● Enter manually

URL *
https://dwintadlst01.dfs.core.windows.net

Storage account key | Azure Key Vault

Storage account key *
••••••••••••••••••••••••••••••••••••••••••••••••

Test connection
● To linked service    ○ To file path

ⓘ If the identity you use to access the data store only has permission to subdirectory instead of the entire account, specify the path to test connection. Please make sure your self-hosted integration runtime is higher than version 4.0 if connecting via self-hosted integration runtime.
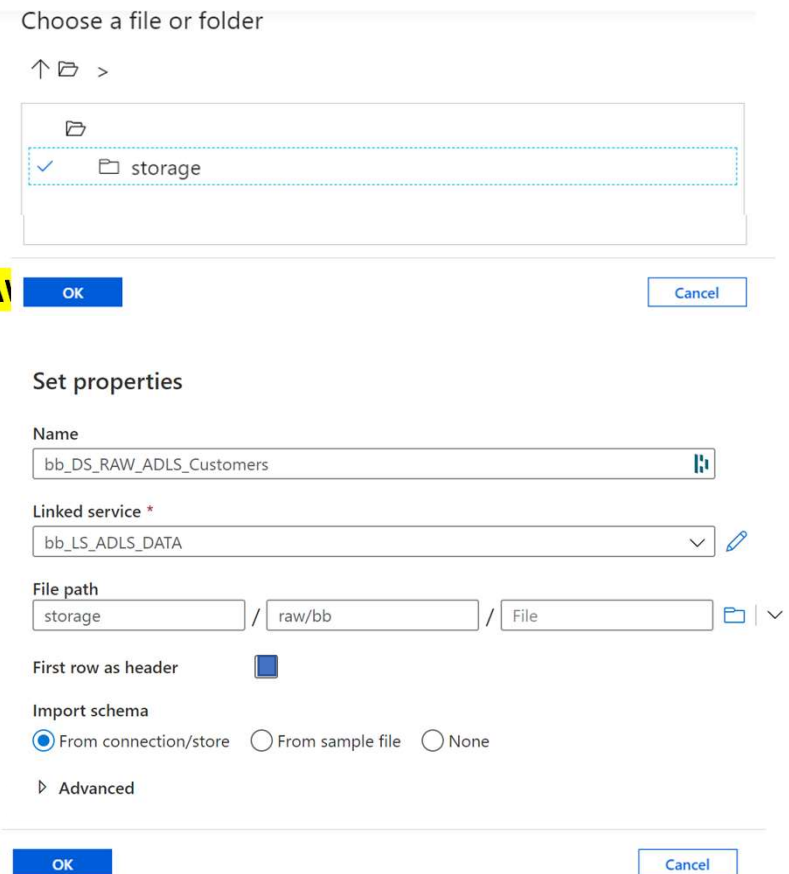
Annotations
+ New

▷ Advanced ⓘ

✓ Connection successful

Apply                    Test connection    Cancel

**In this task, you will set the Sink settings for Copy Activity**

1. In the set properties press the **Browse** button

2. Select the **shared-storage-dw** Path  - in the image (raw/bb) –

     path must be pre-created

3. Press the **OK** button

4. In the Directory field make sure it has the **shared-storage-dw/iXXXXXX/RAW**

5. Let the **File field** empty

6. Select the **First row as header** option

7. In the Import Schema select **From Connection/store**

8. Press **OK**

9. Inside the Sink TAB choose **Open**

10. In the Connection TAB set the File Field to **dimCustomer.csv**
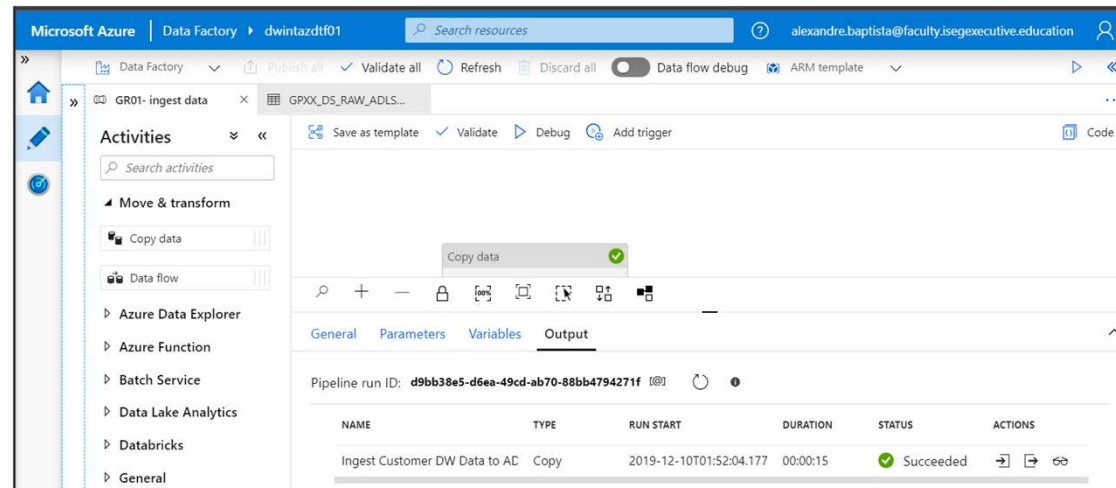
11. Press **Publish All** button

# Execute the Azure Data Factory pipeline

**In this task, you will run the pipeline and**

**Validate the outcome**

1. Open the Pipeline from the left panel

2. Press the <mark>Debug</mark> button and run the package



3. Make sure the <mark>pipeline runs successfully</mark>

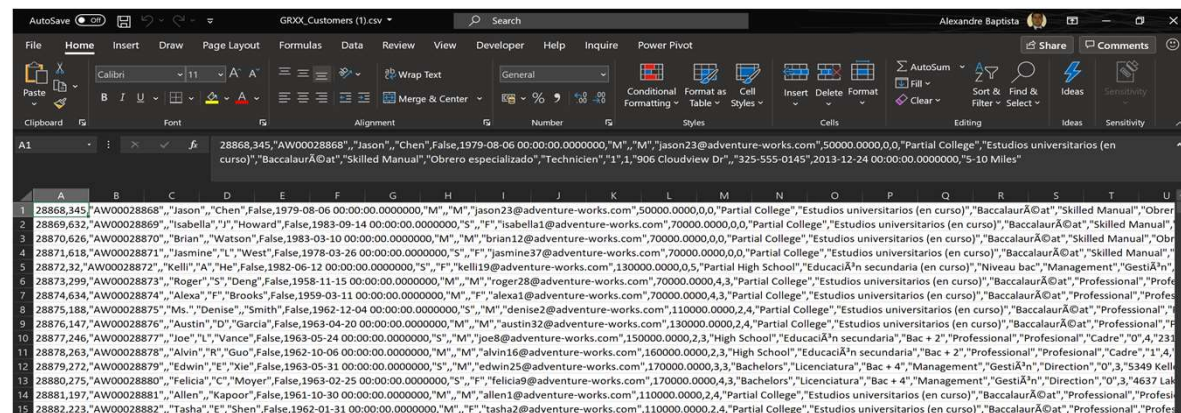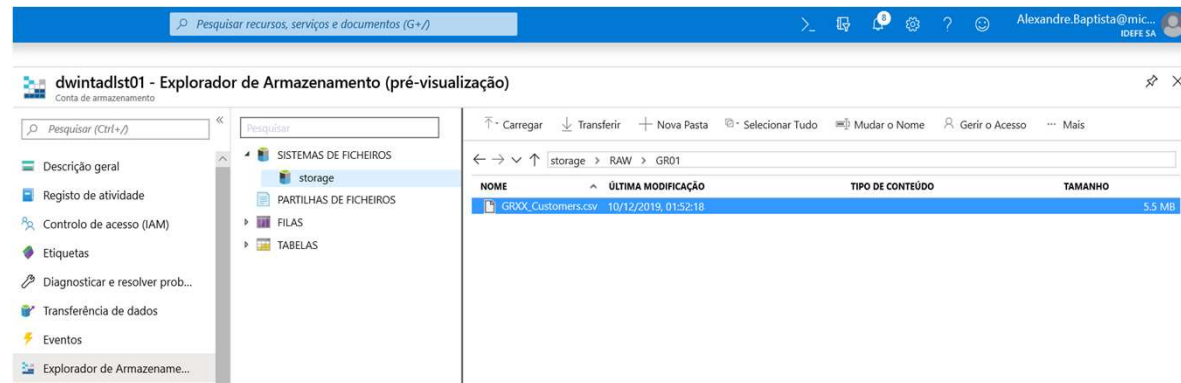4. Open the Actions buttons and
   explore its content



> If the pipeline execution fails, use the error action button to understand and correct the problem that occurred. If no error occurred use ghe pipeline details button mentioned above to understand pipeline's statistics (how long did it took, number of lines that were processed, etc.

# Validate pipeline execution results

**In this task, you will access the datalake storage to ensure data was saved**

1. In Azure Portal open the DW Resource Group

2. Choose the Data Lake Store dsbafb2sharedstor

3. On the left Panel select Storage Explorer (preview)

4. Open FILE SYSTEM

5. Open Storage

6. Open directory

7. Open the dimCustomers.csv file, download it and and check its content using Excel

# Validate pipeline execution results

**In this task, you will repeat the bellow steps for all the other Tables**

1. Perform the same steps for the following tables:

   - **dbo.DimCurrency**

   - **dbo.DimGeography**

   - **dbo.DimDate**

   - **dbo.DimProduct**

   - **dbo,DimProductCategory**

   - **dbo,DimProductSubCategory**

   - **dbo.FactInternetSales**

2. Validate the sucess of the ingesting these tables