# Data Engineering Batch 1


# PROJECT 2


# Data Processing Pipeline

**Akilesh K**

k.akilesh123@gmail.com

## Project statement

Implement a serverless data processing pipeline where Azure Data Factory orchestrates data workflows, and Azure Databricks is used as a serverless processing engine for on-demand analytics and transformations.

## Project Overview

This project establishes a serverless code analysis pipeline in Azure. By leveraging Azure Data Factory (ADF), code seamlessly moves from a GitHub repository to secure Azure Blob Storage. Azure Databricks, empowered by PySpark, then analyses the data, offering valuable insights into functionality, quality, and potential issues. This automated process streamlines data analysis, reduces manual intervention, and fosters data-driven decision-making for developers.

## About the Project

### Dataset:

This dataset, named "organisations", contains information about 10,000 organisations. It was downloaded from the website Datablist: https://www.datablist.com/ and is formatted as a comma-separated values (CSV) file with the following schema:

### Attributes:

Index: An integer representing the unique identifier for each organisation within the dataset.

Organization Id: A unique identifier for each organisation, potentially in a non-human-readable format.
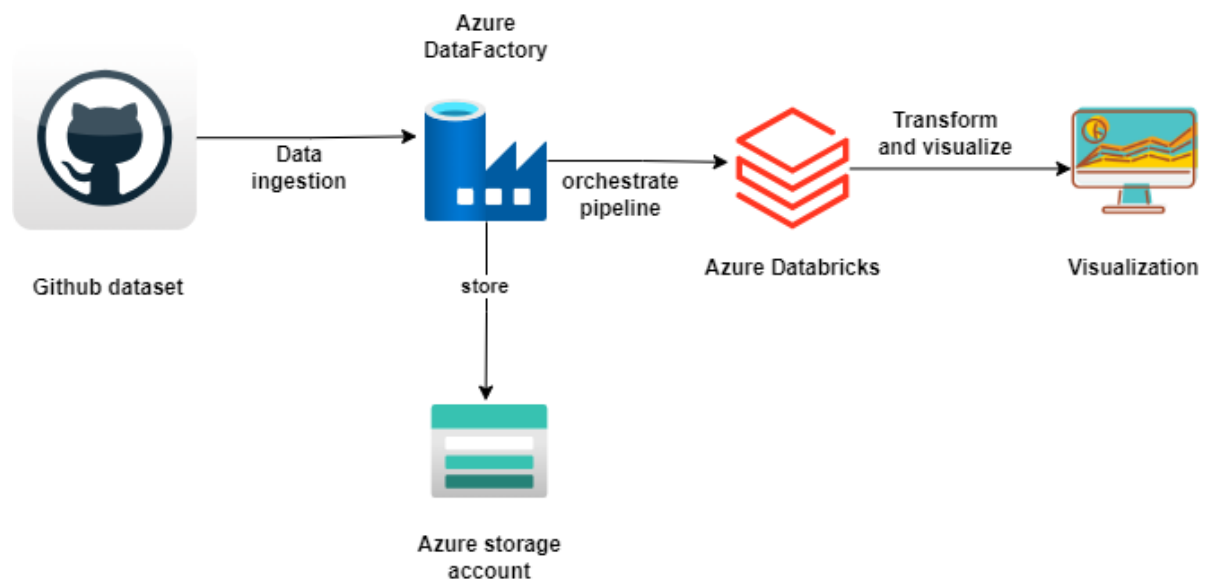
Name: The official name of the organisation.

Website: The organisation's website URL (if available).

Country: The country in which the organisation is headquartered.

Description: A brief description of the organisation's activities or mission.

Founded: The year in which the organisation was founded.

Industry: The primary industry sector in which the organisation operates.

Number of employees: The approximate number of employees working for the organisation (may not be entirely accurate)

| Index | Organization id | Name | Website | Country | Description | Founded | Industry | Number of employees |
|---|---|---|---|---|---|---|---|---|
| 1 | 522816eF8f | Mckinney P | http://soto | Sri Lanka | Synergized | 1988 | Dairy | 3930 |
| 2 | 70C7FBD7e | Cunninghar | http://hard | Namibia | Team-orier | 2018 | Library | 7871 |
| 3 | 428B397eA | Ruiz-Walls | http://www | Iran | Re-context | 2003 | Hospital / H | 3095 |
| 4 | 9D234Ae8C | Parrish, Osl | http://sala | British Indi | Fully-config | 1989 | Supermarke | 5422 |
| 5 | 6CDCcdE3D | Diaz, Roble | https://ww | Botswana | Inverse inta | 2013 | Nanotechn | 3135 |
| 6 | cdAD9BBF2 | Keith PLC | http://www | Ecuador | Cross-grouy | 1978 | Online Pub | 7233 |
| 7 | 0fe6F8Dd1( | Humphrey- | https://fau | Sierra Leon | Adaptive di | 2005 | Publishing | 6022 |
| 8 | ECC0FBd0d | Castaneda- | http://www | Zimbabwe | Front-line v | 2015 | Defense / S | 4580 |
| 9 | e0E6cfAE68 | Santos-Bow | https://ww | Ecuador | Multi-layer | 1979 | Computer I | 3245 |
| 10 | A7DdBb239 | Valdez-Este | http://mel: | Timor-Leste | Diverse ney | 1985 | Plastics | 1785 |
| 11 | 6a74D1bF2I | Young-Zava | http://www | Ukraine | Sharable m | 1972 | Music | 2985 |
| 12 | Cc5293Dbc: | Shaffer Inc | https://ww | United King | Sharable cc | 1997 | Automotive | 839 |
| 13 | 849CEAb2f( | Gaines-Van | http://www | Thailand | Compatible | 2021 | Online Pub | 3135 |
| 14 | DcFCcD6B1 | Larsen-Garr | https://wal | Turkey | Polarized o | 2014 | Building Ma | 7261 |
| 15 | a8dc16ba1t | Taylor LLC | https://ww | Kiribati | Re-context | 1971 | Information | 2427 |
| 16 | 2B6c7b5bEI | Lindsey Ltd | https://log: | Canada | Profit-focus | 1981 | Chemicals | 6477 |
| 17 | fD4BD3662I | Rich-Kelly | https://ww | Colombia | Triple-buffe | 1994 | Investment | 4263 |
| 18 | 3F2A9d9c44 | Solis PLC | https://gon | Suriname | Stand-alone | 2004 | E - Learning | 2546 |
| 19 | 4eBE2067f1 | David-Sum | http://levir | Burundi | Monitored | 2019 | Machinery | 219 |
| 20 | e8F94361bt | Orr-Stanley | http://www | Kazakhstan | Triple-buffe | 1979 | Photograph | 581 |
| 21 | aeCF5B129( | Watkins-Gi | https://kell | Malta | Pre-emptiv | 2004 | Online Pub | 6796 |
| 22 | 20c0DFDE6I | Yoder LLC | http://www | Andorra | Cross-group | 2015 | Mining / M | 7859 |
| 23 | EDb2DAef8 | Ramsey Ltd | https://ww | Yemen | Universal w | 1998 | Investment | 3670 |
| 24 | 220c1D09F5 | Hampton Ir | https://ww | Central Afri | Upgradable | 2019 | Civil Engine | 236 |
| 25 | d4B566cF72 | Schaefer, B | http://mar: | Swaziland | Optimized | 1987 | Design | 8044 |
| 26 | 1Cec2bE2f8 | Sosa-Lynn | https://wal | Brazil | Monitored | 1994 | Architectur | 4905 |
| 27 | 92cAd11Ec4 | Nichols, Ro | https://ols: | Bangladesh | Ameliorate | 2014 | Commercia | 1996 |
| 28 | 1ADB796E3 | Landry, Per | http://matl | Ghana | Optional ne | 2014 | Textiles | 4214 |
| 29 | aaFDd69B4( | George Ltd | https://pitt | Isle of Man | Realigned a | 1986 | Commercia | 9958 |
| 30 | 65bB9Baa2: | Greer-Watt | https://ww | Timor-Leste | Reverse-en | 2017 | Wholesale | 1884 |
| 31 | da750fccbe | Archer, Gec | http://www | Russian Fec | Innovative | 1994 | Environme | 6377 |
| 32 | 907DE64cA | Zhang-Esco | https://hoc | Iraq | Innovative | 1979 | Internet | 358 |
| 33 | BB4fecbcEL | Kelley-Luce | https://ww | Kenya | Versatile fu | 2001 | Mental Hea | 5154 |
| 34 | F8cAa7DA8 | Gallegos, P | http://www | Belarus | Virtual clie | 1970 | Photograph | 6153 |
| 35 | 144FBe6c3E | Benjamin, ( | http://www | Moldova | Adaptive in | 1976 | Venture Ca | 3528 |
| 36 | AdC821b25 | Hatfield PL | http://www | Malta | Distributed | 1985 | Translation | 554 |
| 37 | eabCabc94I | Richards PL | http://todc | Albania | Persevering | 1988 | Public Safe | 7554 |
| 38 | E5E75Eb6cF | Brandt, Vac | http://riggs | Andorra | Proactive n | 2019 | Cosmetics | 7372 |

organizations-10000

## Architecture



The architectural diagram depicts a systematic approach to data processing and analysis within the project. It begins with the acquisition of code from a GitHub repository, facilitated by Azure Data Factory (ADF). ADF orchestrates the transfer of code to Azure Blob Storage, where it is securely stored for further processing. Azure Databricks workspace is then provisioned to handle data processing tasks efficiently. Within the workspace, a Notebook activity configured within ADF executes PySpark code responsible for data analysis. This PySpark code loads the code data from Blob Storage into Spark DataFrame, conducts necessary transformations, and performs analysis using PySpark libraries. The results of this analysis are visualized using tools available within Azure Databricks, with options for generating visualizations such as pie charts to represent the code analysis findings. Additionally, if required, the processed data or analysis results can be stored back to Blob Storage for future reference or sharing.

## Azure Resources Used for This Project

- **Azure Data Factory (ADF):**

An ADF instance will be created to orchestrate and automate the data transfer process between GitHub and Azure Blob Storage.

Linked Services: Configured to connect to GitHub using an appropriate authentication method (e.g., personal access token, OAuth) and Azure Blob Storage.

Datasets: Defined to represent the schema of the code files (e.g., text format) in the GitHub repository.

Pipelines: Created to define the sequence of activities for copying code from the GitHub repository to a specific blob storage container.

- **Azure Blob Storage:**

A blob storage container will be used as the destination for storing code transferred from the GitHub repository.

This stored code will be accessed by Azure Databricks for data processing tasks.

- **Azure Databricks:**

A Databricks workspace will be provisioned to perform data processing tasks using PySpark.

Databricks Notebooks: Utilized to write and execute PySpark code for code analysis, including tasks like:

Extracting data from the code files (e.g., lines of code, function names, comments).

Performing data transformations and cleaning (e.g., removing irrelevant content, formatting code).

## Project Requirements

- **Data Source:**

Type: Identify the specific source of your code, such as a GitHub repository, Azure DevOps repository, or a local file system.

Connectivity: Ensure you have the necessary credentials and access methods to securely connect to this source from the Azure environment.

- **Azure Data Factory (ADF) Configuration:**

Instance: Create an Azure Data Factory (ADF) instance within your Azure subscription.

Linked Services:

Configure a linked service to connect to the GitHub repository using an appropriate authentication method (e.g., personal access token, OAuth).

Configure a linked service to connect to Azure Blob Storage for storing the transferred code.

Datasets: Define a dataset in ADF to represent the schema of the code files (e.g., text format) from the source location.

Pipelines: Create a pipeline in ADF to define the data transfer process:

Use a Copy Activity to copy code from the source location (e.g., GitHub repository) to a specific blob storage container.

Configure the copy activity with source and destination details, scheduling options (e.g., daily), and error handling mechanisms.

- **Data Storage:**

Azure Blob Storage: Utilize a blob storage container as the destination for storing the transferred code retrieved through ADF. This stored code will be accessed by Azure Databricks for further processing.

- **Data Processing with Azure Databricks:**

Provision an Azure Databricks workspace in your Azure subscription.

Define and implement data processing tasks using PySpark within the Databricks environment. PySpark provides a powerful framework for distributed data processing, enabling tasks such as data cleaning, transformation, aggregation, and analysis.

Use Databricks notebooks to write and execute PySpark code interactively, leveraging the scalability and performance of the Databricks runtime.

## Tasks performed:

- Acquire data: Retrieve code from a GitHub repository.
- Transfer data: Utilize Azure Data Factory (ADF) to orchestrate the transfer of code from GitHub to Azure Blob Storage.
- Store data: Utilize Azure Blob Storage to securely store the transferred data.
- Process data:
- Provision an Azure Databricks workspace.
- Within ADF, add a Notebook activity: Configure the activity to use your Databricks workspace and specify the PySpark notebook responsible for code analysis.
- Within the PySpark notebook:
- Load data from Blob Storage into a Spark DataFrame.
- Implement data transformations and analysis using PySpark libraries.
- Visualize data: Generate visualizations (e.g., pie charts) depicting the results of code analysis using libraries within Databricks.

# Implementation

- ## CSV data on github



- ## In azure ,create a storage account mentioning the location

● **create a new blob container**



● **create a azure data factory account**

- **select ingest data to create a new data pipeline**



- **For Source click on new connection and type http**

- **Copy and paste the raw url of the csv file**



- **Create a connection for destination blob storage**

- **Select the container**



- **summary of the deployment of from SQL server to azure blob storage**

- **The csv file is stored in the container**



- **Create a Azure databricks and launch it**



**Mount blob storage to azure databricks**

- **Add source line of the blob storage**
- **enter the mount points to store in databricks**
- **Configure extra_configs by mentioning the access key and pasting the key of the blob storage.**
- **list of the files mounted on Databricks**
- **display the data frame**

- **Performing data preprocessing activity such trimming blank spaces from columns**

● **Drop rows with Null values**



```
df.na.drop().show()
```
▸ (2) Spark Jobs
▸ ▤ df: pyspark.sql.dataframe.DataFrame = [Index: integer, Organization Id: string ... 7 more fields]

```
|  314|D06f4A7F5Cd6ecB|Luna, Mcmillan an...|https://www.dixon...|          SOMALIA|Intuitive asynchr...| 1987|   Online Publishing|       7649|
|  434|8aEa9A19068f77E|          Mccall Ltd|http://www.glover...|            TONGA|Implemented heuri...| 2008|      Other Industry|       3498|
|  578|EbC8B7e60C503cd|           Smith Ltd|http://lam-robins...|RUSSIAN FEDERATION|Triple-buffered e...| 2009|   Apparel / Fashion|       3398|
|  643|02d6AecA05DDaaC|       Gutierrez Ltd|https://www.mathi...|   CAYMAN ISLANDS|Multi-lateral dis...| 1999|     Computer Games|       2746|
|  646|Bd8057f80A250aB|          Weber Ltd|http://cisneros-e...|         BOTSWANA|Focused asynchron...| 2000|Oil / Energy / So...|       5440|
| 1174|5cCe4dBAceffc40|Cabrera, Yoder an...|https://stanley.com/|DOMINICAN REPUBLIC|Right-sized dedic...| 2021|Capital Markets /...|        943|
| 1189|Eea2a2212D9BA1e|        Gonzales PLC|http://www.nunez-...|       BANGLADESH|Robust maximized ...| 2002|     Transportation|       4883|
| 1363|57af1f2D8F62B2c|           Heath Inc|http://blankenshi...|          COMOROS|Self-enabling exp...| 2021|           Plastics|       7024|
| 1370|aDe2ab1dD8aa71D|        Cain-Hensley|   https://west.org/|       MONTSERRAT|Reduced leadinged...| 1998|       Shipbuilding|       3664|
| 1946|Da69C89370A0bC3|        Levy and Sons|https://fuentes-v...|FALKLAND ISLANDS ...|Cross-platform ex...| 1972|     Wine / Spirits|       9104|
| 2091|F3F7C0ffd09c065|    Reynolds-Goodman|https://www.chapm...|             PERU|Self-enabling mul...| 1973|      Other Industry|        944|
| 2214|842f1DEcdA80dD5|           Mills Inc|https://www.lower...|           MEXICO|Progressive didac...| 1974|    Leisure / Travel|       8409|
| 2244|86ba8881d21ed8E|        Potter-Hines|      http://ball.com/|         CAMEROON|Open-architected ...| 1977|Logistics / Procu...|       1990|
| 2310|83eb230Ddf4A0fb|    Livingston-Oliver|   http://walls.info/|   CAYMAN ISLANDS|Optimized directi...| 1989|    Law Enforcement|        361|
| 2378|3FdFE1982be12FE|Phelps, Lutz and ...|https://www.hendr...|WALLIS AND FUTUNA|Integrated compos...| 2003|        Accounting|       2331|
| 2490|9Ad1149F4569Cc8|Faulkner, Nash an...|http://www.greer....|        MAURITIUS|Universal system-...| 1997|   Research Industry|       8540|
| 2652|e710d22B6Fae7b2|Hutchinson, Walsh...|https://www.juare...|    NEW CALEDONIA|Digitized mission...| 1998|Staffing / Recrui...|       3364|
| 3207|37139c7ad4B48EC|Woods, Wheeler an...|   http://decker.com/|  FRENCH POLYNESIA|Cloned responsive...| 2020|   Research Industry|       5027|
+-----+---------------+--------------------+--------------------+-----------------+--------------------+-----+--------------------+-----------+
only showing top 20 rows
```

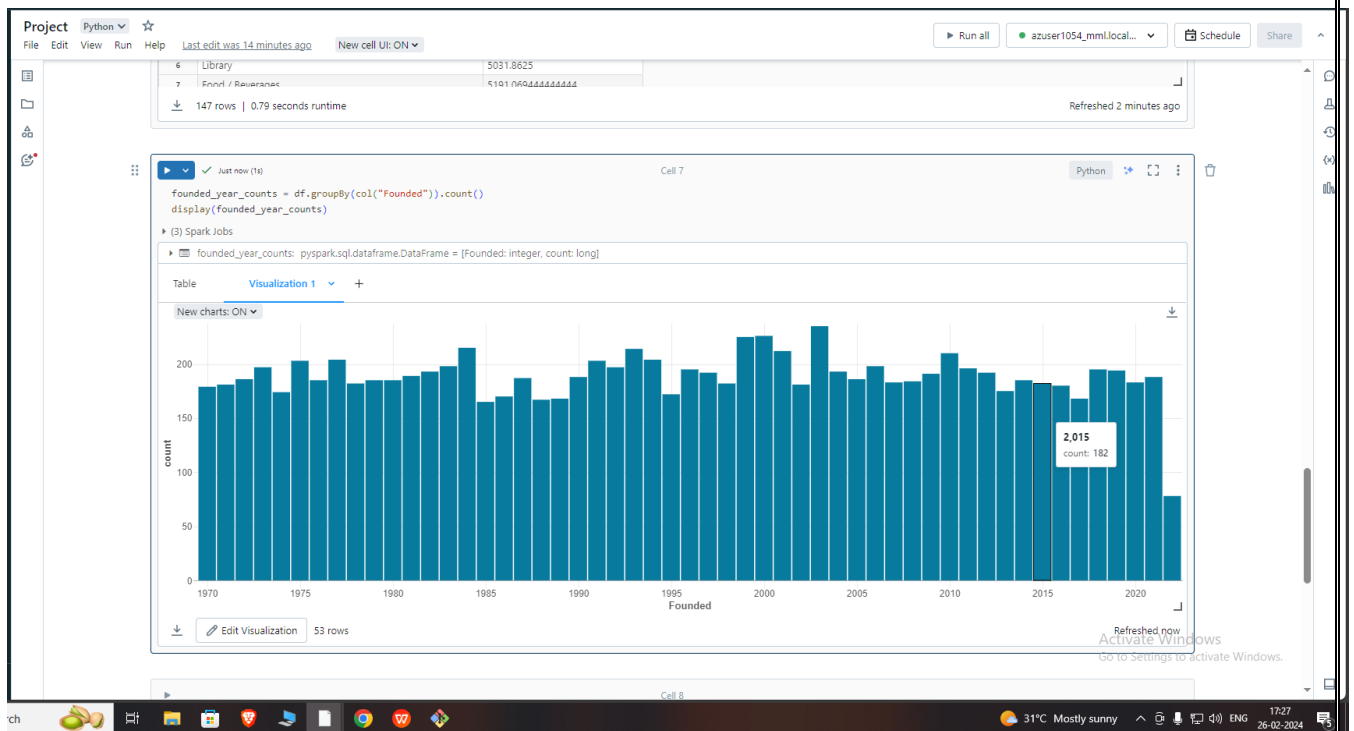● **Count of Organisations by Country**

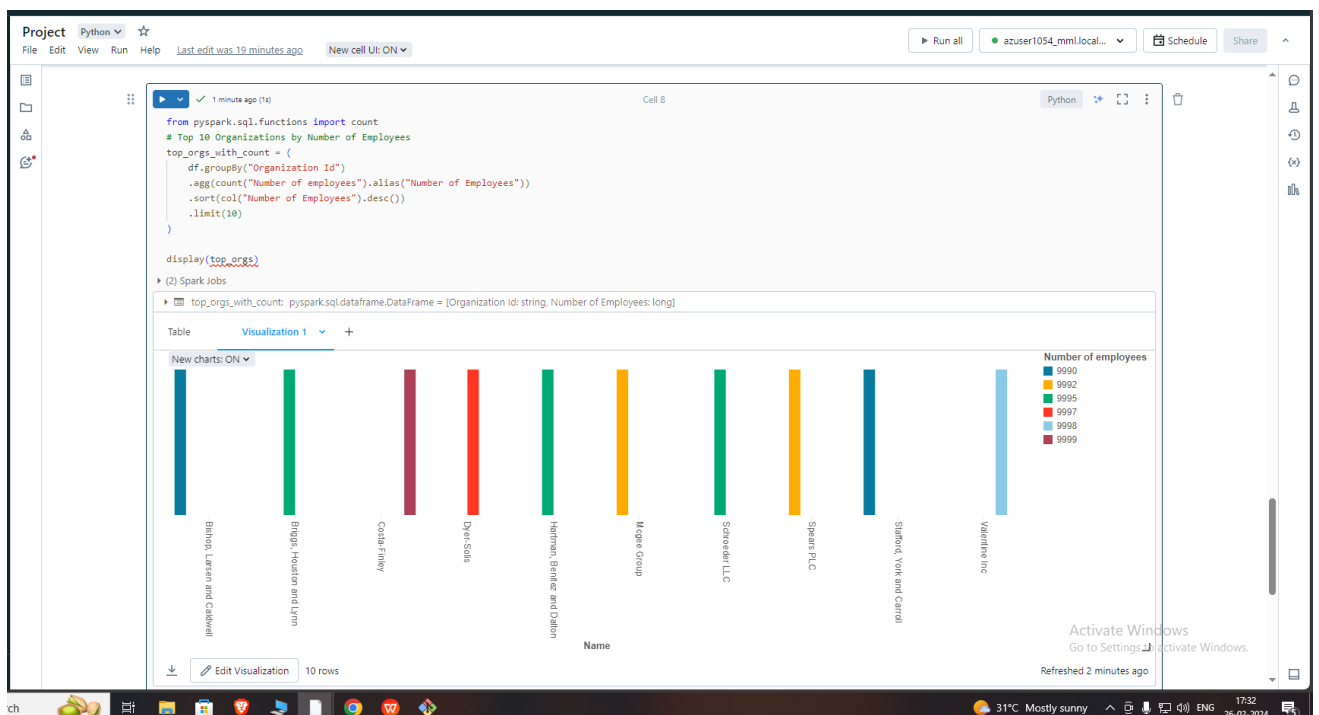- **Visualising count of Organisation by Country using a pie chart**



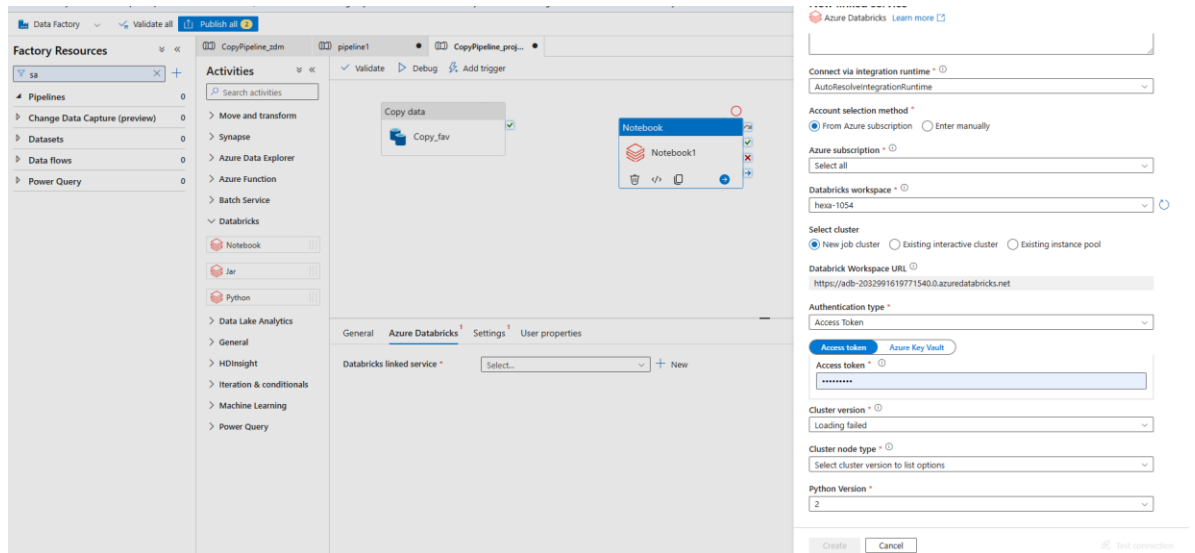- **Average number of Employees per industry**

- **Visualising count of country grouped by the year in which they were founded using a bar graph**
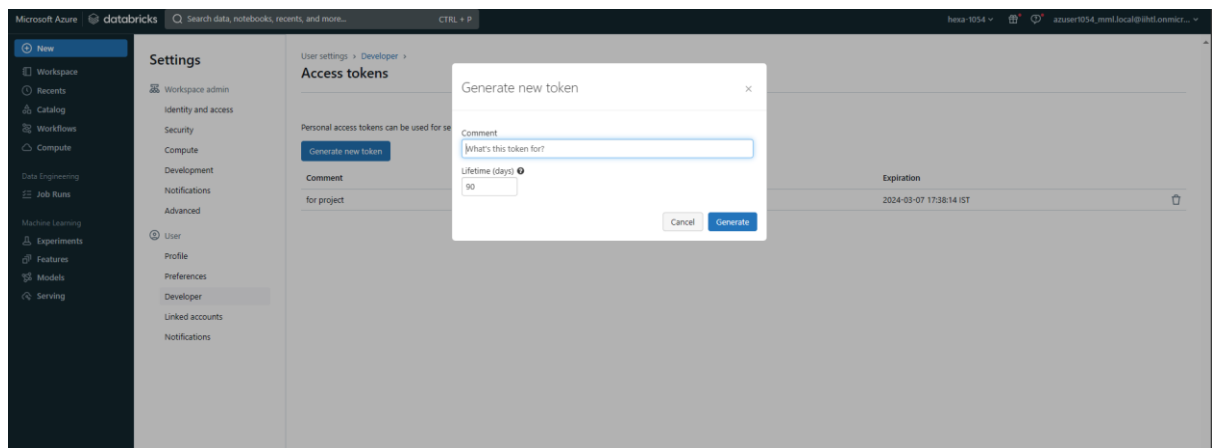


- **Top 10 Organisations sorted by the number of employees**

- **In Azure Datafactory drag and drop the Notebook activity and link it to the databricks Notebook**



- **Copy the access token from Databricks Notebook by clicking on User settings -> Developer -> Generate new Token**

- **In Datafactory Validate the pipeline and then click on debug**



- **The pipeline is running successfully**

## Conclusion

In conclusion, this project successfully demonstrated the integration of various Azure services to create an end-to-end data processing pipeline. Beginning with data acquisition from GitHub, we utilized Azure Data Factory to orchestrate data movement to Azure Blob Storage. Subsequently, Azure Databricks was employed for data processing and analysis, including preprocessing tasks and generating insightful visualizations. The seamless integration of these Azure services facilitated efficient data handling and analysis, showcasing the power of cloud-based data solutions in modern data engineering workflows. This project highlights the effectiveness of Azure's ecosystem in enabling scalable and efficient data processing pipelines for diverse use cases.