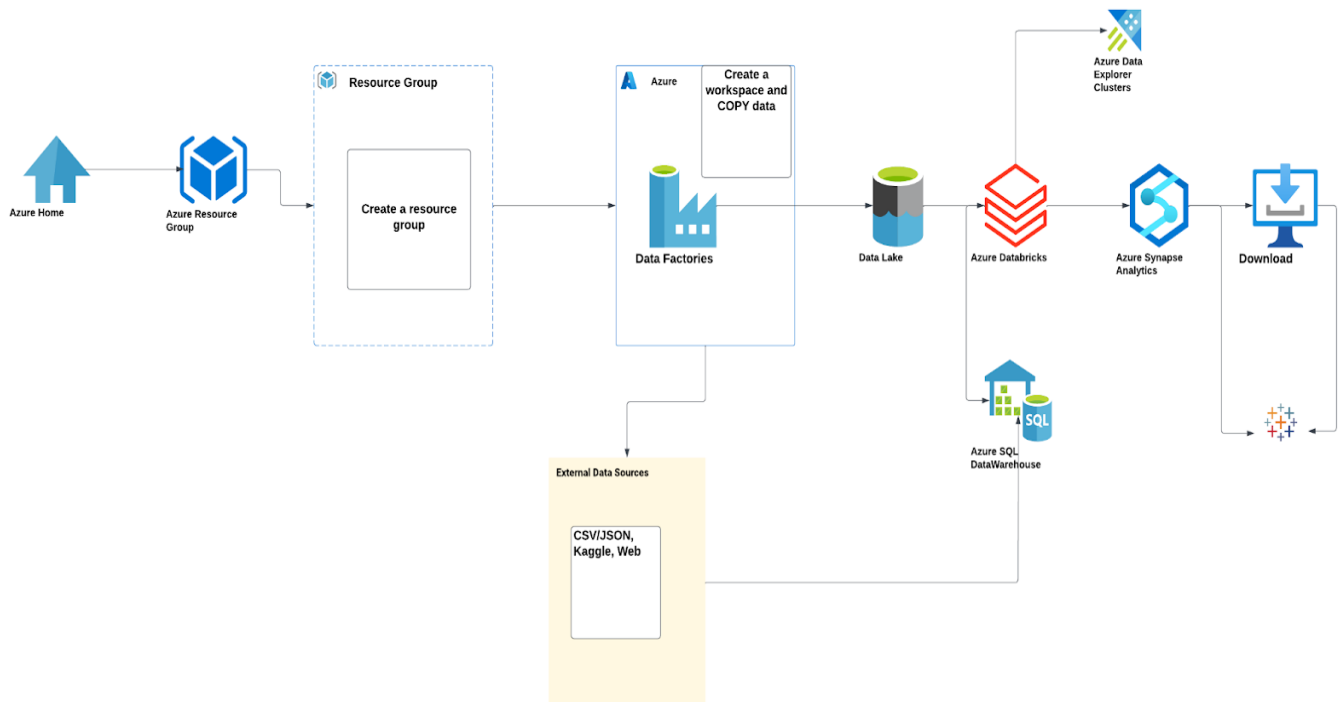# Analyzing Salary data

**Data set: [Click to checkout Data](#)**

**About the Data:** This dataset offers an extensive compilation of salary data spanning diverse industries and global regions. Gathered from reputable employment platforms and surveys, it encompasses information on job titles, compensation, industry sectors, geographical locations, and more. We can utilize this dataset for a thorough examination of trends in the job market, a comparative analysis of salaries across various professions, and to make well-informed decisions regarding career paths or hiring strategies. The dataset has been meticulously cleaned and preprocessed to facilitate analysis, and it is accessible under an open license for research and data analysis endeavors.

**Tools used for Project:**
- Azure Services - Data Factory, Databricks, Synapse Analytics
- Tableau
- Lucid Chart

**Process Flow:**

# Analyzing Salary data

Implementing Azure Analytics with Tableau involves a series of well-defined steps for seamless integration and effective data analysis. Here is an elaboration of the key steps:

**Azure Resource Group Creation:**
Start by establishing an Azure Resource Group. This serves as the foundational container for organizing and managing resources related to your analytics project on the Azure platform.

**Azure Data Factory Workspace Setup:**
Create an Azure Data Factory workspace, which acts as the orchestration hub for managing and orchestrating data workflows.
Within this workspace, set up a container with two directories, "raw data" and "transformed data," providing structured storage for data at different stages of processing.

**Data Ingestion with Azure Data Factory:**
Utilize the "COPY data" feature in Azure Data Factory to connect to the raw data source, typically hosted on a Git repository.
Direct the data to the desired destination, often referred to as the "sink," and store it in a CSV format after undergoing necessary transformations.

**Azure Databricks Workspace Creation:**
Establish an Azure Databricks workspace, a collaborative environment for big data analytics, and create a new App registration within it.
Extract essential credentials such as Client ID, Tenant IP, and secret key during the App registration process. These credentials are vital for seamless data access between Azure Data Factory and Databricks.

**Direct Data Ingestion into Azure Databricks:**
As an alternative to using Azure Data Factory, the data file can be ingested directly into Azure Databricks, eliminating the need for an intermediate step. This involves reading the data into a tabular format suitable for analysis.

# Analyzing Salary data

**Data Exploration and Analysis in Azure Databricks:**
Leverage SQL and/or Spark within the Azure Databricks workspace to access and manipulate the tabular data.
Create a new notebook in the workspace to load necessary Python libraries, inspect the data for inconsistencies, and perform statistical analysis.

**Data Visualization and Analysis with Tableau and Azure Synapse Analytics:**
Utilize Azure Synapse Analytics for efficient data warehousing and analysis.
Develop visualizations and dashboards using Tableau, connecting to the processed data in Azure Synapse Analytics to gain meaningful insights.
Leverage the combined capabilities of Tableau and Azure Synapse Analytics for advanced data visualization, analysis, and reporting.

## Tableau Integration:

Download Tableau Desktop and the required drivers to establish connectivity with Azure Databricks.
Alternatively, download the data from Azure Databricks and use it directly integration to Tableau for creating visualizations and dashboards



**Tableau visualizations**

## Screenshots Below show the Process followed after creating DataFactory Workspace

# Analyzing Salary data

- Create a workspace in Data Factory Studio(Note: Resource group should be created before this)



- **We can see the details here data factory is attached to resource group and at left we can see containers -select to create**



- **Create two directories in containers**

# Analyzing Salary data



- **Create a pipeline and utilize COPY Data to load data then validate, Debug and published**



- **Add access to Azure Blob container**

# Analyzing Salary data



- **Create DataBricks workspace**

# Analyzing Salary data





- **Create a new App registration to mount data factory in databricks**

# Analyzing Salary data

# Analyzing Salary data

- **Else add data through files/ various integrations**



- **Create a new notebook >Python and ensure your cluster is selected. Then load data used PySpark**

- **Using Spark Dataframe we performed Analysis**



**Tableau:**

The Tableau process undertaken involved the utilization of a dataset in CSV format. To initiate the analysis, the dataset was uploaded onto Tableau using the Spatial file upload option. This method allows for the seamless integration of spatial data, providing a rich foundation for exploration and visualization.

Once the dataset was successfully uploaded, the next step involved extracting the data into the Tableau workspace. Extraction allows for faster data processing within Tableau, enhancing the overall efficiency of the analysis.

Subsequently, the analytical exploration of the dataset took place. This was accomplished by creating separate sheets within Tableau, each focusing on specific features or aspects of the data. Breaking down the analysis into individual sheets enables a more granular examination of each variable or category, facilitating a deeper understanding of the dataset.

Concurrently, a key aspect of the process involved the creation of a dashboard. The dashboard serves as a consolidated and visually intuitive representation of the

analyzed data. It acts as a user-friendly interface that brings together insights from various sheets, making it easy to communicate and comprehend the findings.

The emphasis on creating an easily explainable dashboard suggests a focus on clarity and accessibility. By designing the dashboard with a user-friendly approach, the intent is to make the insights readily understandable to a diverse audience, ensuring that the analytical results are effectively communicated.

**Tableau Viz:**



## Let's break down each visualization

**Bar Graph - Salary Range by Job Title:**
- The bar graph illustrates the salary ranges associated with different job titles.
- Notably, the 'CEO' position stands out with the highest remuneration, followed by positions such as CTO, CDO, Director, and VP.
- A trend emerges where executive-level titles consistently command higher salaries on average.

# Analyzing Salary data

**Map - Average Salary of Selected Job Titles Worldwide:**
- The map provides a global perspective, displaying the average salary for selected job titles across different countries, with USD as the base currency.
- The color-coded map facilitates easy differentiation between countries.
- The observation indicates that the USA generally offers more competitive salaries compared to other countries.

**Bar Graph - Salary vs. Experience:**
- The bar graph explores the relationship between work experience and salary.
- Contrary to expectations, it reveals that salaries vary irrespective of experience levels. Some individuals with less experience earn higher salaries, while those with more experience may earn lower salaries.
- The range of 10-15 years of experience emerges as a period with more competitive salaries.

**Pie Charts - Gender Disparity in Salary:**
- The pie charts aim to understand salary disparities between male and female employees.
- The observation indicates that, on average, females earn a comparatively lower salary than males.
- However, the difference, while present, is not overwhelmingly significant, suggesting a potential for equalization in the future.

**Bar Chart - Salary vs. Race:**
- The bar chart examines salary discrepancies based on race.
- Clear trends emerge, indicating that individuals belonging to the white and Asian races generally earn higher salaries compared to other racial groups.

In conclusion, the data visualizations present a comprehensive snapshot of salary dynamics across job titles, global locations, experience levels, gender, and race. Executives consistently earn higher salaries, the USA stands out for competitive wages, experience alone does not linearly determine earnings, gender pay gaps exist but are not overwhelming, and disparities based on race are evident, with white and Asian individuals earning higher salaries on average. These findings provide valuable insights for understanding and addressing various aspects of salary inequality.