# CA675 CLOUD TECHNOLOGIES- Assignment 2

| | |
|---|---|
| Student's Name(s): | Alen George Shibu, Sarath Mukundan Adavakkat, Suhara Fathima Shahul Hameed |
| Student Number(s): | 23271013, 23267908, 23266257 |
| Student Mail ID(s): | alen.georgeshibu7@mail.dcu.ie, sarath.mukundanadavakkat2@mail.dcu.ie , suhara.shahulhameed2@mail.dcu.ie |
| Group Number: | 34 |
| Program: | Master of Science in Computing (Data Analytics) |
| Project Title: | Loan Defaulters' Demographic Visualizer |
| Module Code: | CA675 |
| Lecturer: | Dr. Michael Scriney |
| Assignment Submission Date: | 19th December 2023 |
| GitLab Repository: | https://gitlab.computing.dcu.ie/georgea7/loan-defaulters-demographic-visualiser |
| Video Link: | https://drive.google.com/file/d/1TxjfonAQfwKhxW9sZhYix8H_NiM3RFzr/view?usp=drive_link |

# Table of Contents

# List Of Figures

# INTRODUCTION

In the constantly evolving financial services industry, understanding the intricacies of loan defaulters is critical for risk assessment and proactive decision-making. It provides an extensive collection of information, uncovering significant patterns and weaknesses among those who default on loans.

Visualizing age, income, occupation type and gender data specifically for loan defaulters is an important analytical tool in identifying the risk factors linked with loan repayment. By focusing on this specific subset of persons who have defaulted on loans, we can acquire useful insights into the demographic and socioeconomic aspects that may contribute to defaulting behaviour. Understanding the demographic profile of defaulters by age, income, and gender enables financial institutions to estimate risk more properly. It aids in the refinement of lending strategies and risk assessment models by identifying trends or groups that may be more prone to loan defaults. We hope that by visualizing this, we can improve knowledge and guide solutions for mitigating default risk while increasing financial inclusion and stability.

## Choice of Technology

The technology used for visualizing demographic data of loan defaulters - Amazon Web Services (AWS), PySpark, and Tableau—provides a strong and complete approach.

Amazon EMR (Elastic MapReduce) is a big data cloud platform provided by Amazon Web Services (AWS). It is a fantastic solution for efficiently processing large-scale datasets. It makes use of Apache Spark, among other frameworks, to enable distributed data processing. PySpark, a Python API for Spark, makes it easier to manipulate and analyze data. Using Amazon EMR with PySpark allows for the handling of huge volumes of demographic data, computations, and the extraction of valuable insights.

PySpark is a strong data processing tool, particularly when dealing with enormous datasets. Its interface with Amazon EMR enables distributed computing, making data translation, cleaning, and analysis more efficient.

Tableau, on the other hand, specializes in data visualization and interactive dashboard creation. It can connect to data sources such as Amazon EMR and display processed and analyzed demographic insights. Tableau has an easy-to-use interface for making visually appealing charts, graphs, and dashboards. It enables stakeholders to study and comprehend complicated demographic patterns among loan defaulters interactively.

# RELATED WORK

1.Demographic Analysis 2022 [link]

The Central Bank of Ireland's "Demographics of the Financial Sector" report is likely to provide a detailed analysis of the composition and diversity of the financial industry workforce. This study sheds light on the demographic makeup, including age, gender, ethnicity, and maybe other diversity-related criteria. It examines how different demographic groups are represented at various career levels, noting potential trends or changes over time. The report concentrates on gender diversity, stressing the presence of women in positions of leadership, and it may also discuss ethnic or racial diversity within the sector. It could also highlight problems in reaching diversity targets, actions to encourage inclusion and the positive effects of a diverse workforce on corporate performance.

2.German Loan Default Prediction Project[link]

The focus of a German loan default prediction project is anticipated to be on using machine learning techniques to forecast the risk of loan defaults. Analysing a dataset containing borrower information, loan attributes, and previous default records is required. To prepare the dataset for modelling, the project may include data pre-treatment activities such as cleaning and feature engineering. To develop prediction models, many machine learning techniques, such as logistic regression or decision trees, may be used. The accuracy and efficacy of these algorithms in forecasting loan defaults are assessed during their evaluation. The project's findings frequently provide insights into major elements impacting default probabilities and aid in the improvement of credit risk assessment methods.
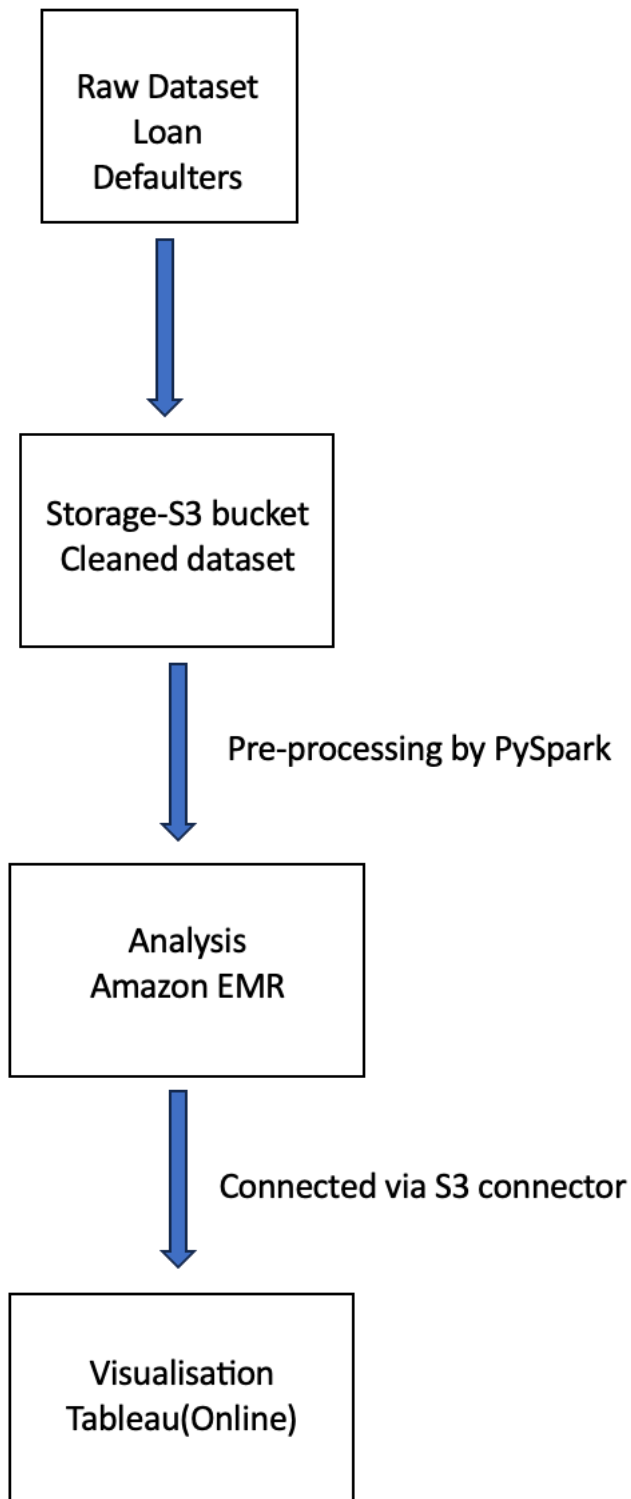
3.The Effect of Demographic Characteristics on Loan Performance of Commercial Banks In Kenya[link]

This study investigates how demographic characteristics influence commercial bank loan performance in Kenya. The study explores the impact of demographic variables such as age, income, education, and work status on borrowers' repayment behavior using data from various demographic variables. Statistical analysis results show substantial relationships between specific demographic variables and loan performance, providing insights that could guide more focused lending strategies or risk assessment processes in Kenya's commercial banking sector.

Our work focuses on demographic factors such as age, gender, income, and occupation of loan defaulters. We visualize these factors in a dashboard to understand these trends based on the gender of loan defaulters.

# SYSTEM ARCHITECTURE

The system architecture of our visualization is given below:

```
          ┌─────────────────┐
          │   Raw Dataset    │
          │      Loan        │
          │   Defaulters     │
          └─────────────────┘
                   │
                   ▼
          ┌─────────────────┐
          │ Storage-S3 bucket│
          │  Cleaned dataset │
          └─────────────────┘
                   │
                   │   Pre-processing by PySpark
                   ▼
          ┌─────────────────┐
          │     Analysis     │
          │   Amazon EMR     │
          └─────────────────┘
                   │
                   │   Connected via S3 connector
                   ▼
          ┌─────────────────┐
          │  Visualisation   │
          │  Tableau(Online) │
          └─────────────────┘
```

# DATASET

## Source of the data

The data was collected from the Kaggle website.
Link:https://www.kaggle.com/datasets/gauravduttakiit/loandefaulter/data?select=application_data.csv
The original dataset consists of 122 columns and 307512 rows which included loan default status as 0 and 1. Target variable (1 - client with payment difficulties: he/she had a late payment, 0 – no payment difficulties). Out of 122 columns we required only 11 columns for our aim.

## Data pre-processing

Link to source code repository: https://gitlab.computing.dcu.ie/georgea7/loan-defaulters-demographic-visualiser

The dataset preparation and cleaning are done in PySpark in the following steps :



Figure 1:EMR-cloud2

Initially, we selected rows with target = 1 i.e., loan defaulters as we are focussing on these data closely. Selection of required columns is performed which are SK_ID_CURR", "TARGET", "NAME_CONTRACT_TYPE","CODE_GENDER","AMT_INCOME_TOTAL","AMT_CREDIT","AMT_ANNUITY","NAME_INCOME_TYPE","NAME_EDUCATION_TYPE","DAYS_BIRTH", "OCCUPATION_TYPE".
Next, the DAYS_BIRTH column consisted of negative values which were converted to age by dividing the negative value by 365 and taking the absolute floor value as the age. Row with null values were eliminated. Checked for outliers in the tableau by box plotting and removed irrelevant data outliers.
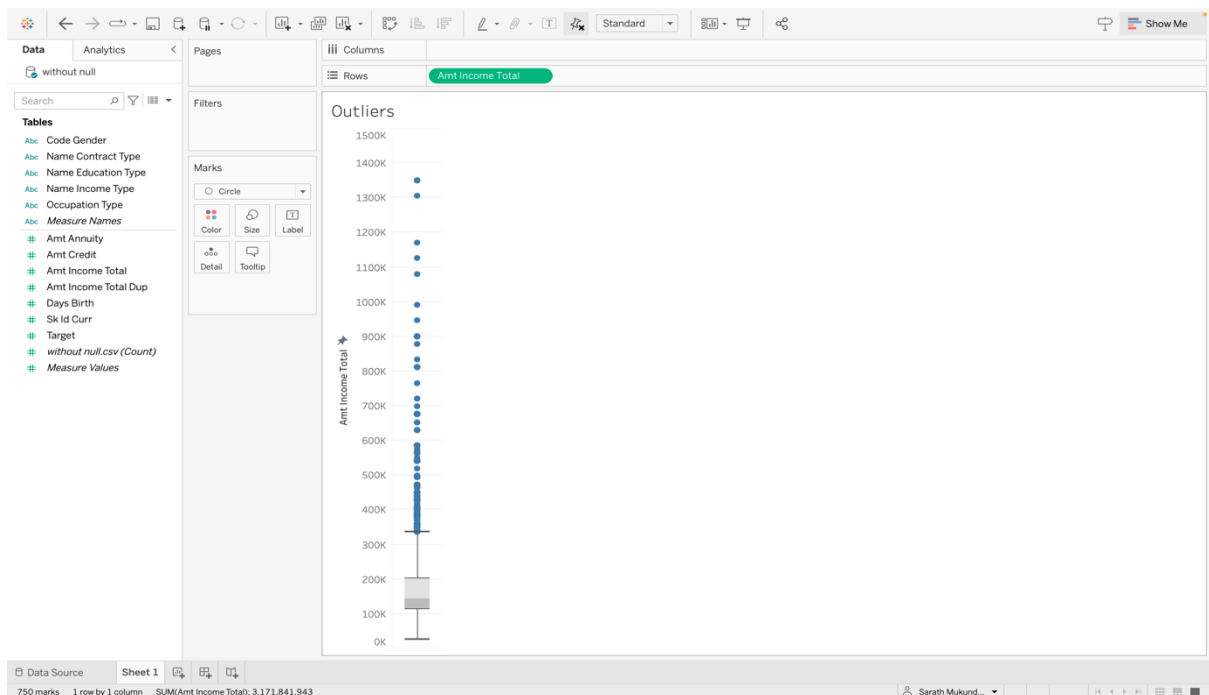
Figure 2: Outliers in data

# DATA PROCESSING

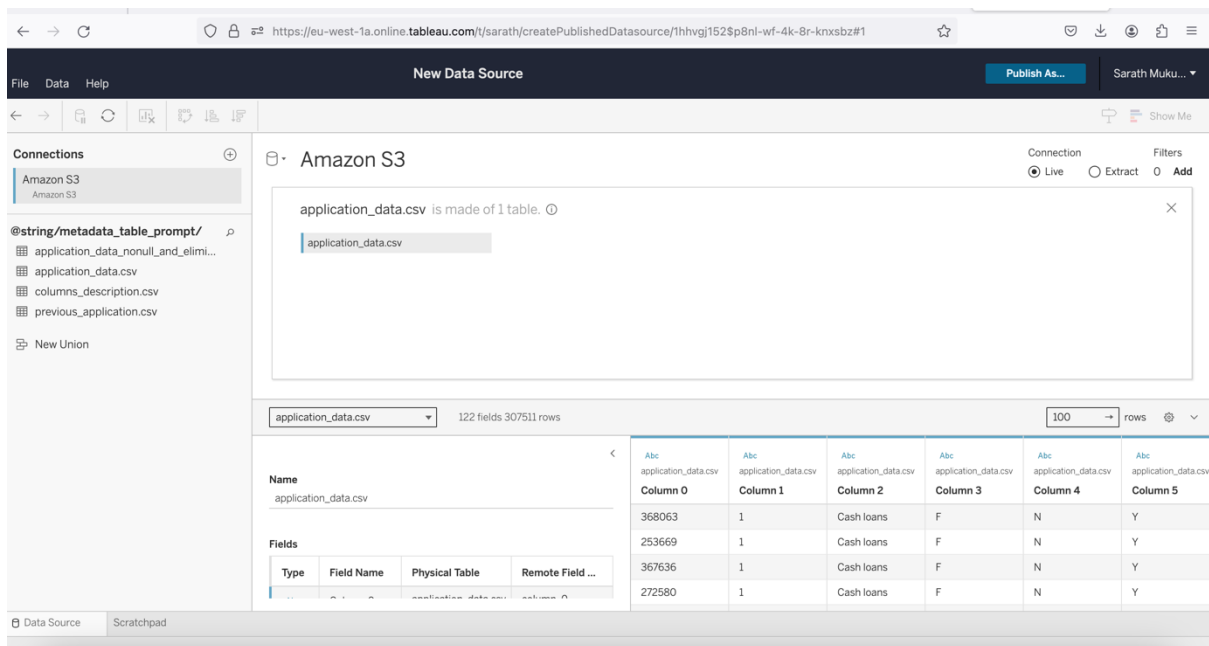The cleaned dataset from S3bucket is imported to Tableau online by using S3 connector.



Figure 3: S3 bucket

The distribution of age, income and occupation of loan defaulters based on gender is visualised.

## A.Distribution of Age

We have visualised the distribution for the age of loan defaulters based on gender which can be seen below.
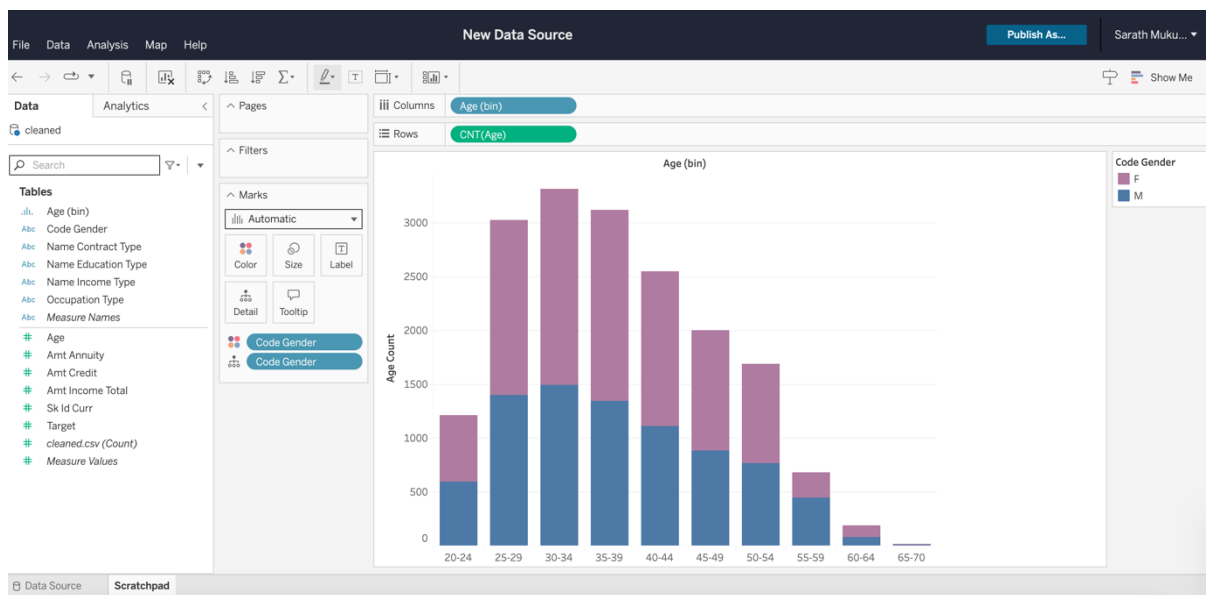


Figure 4:Age distribution

Insights from graph:

- The graph shows the distribution of age for males and females, with males in blue and females in purple colour. The x-axis shows the age interval, and the y-axis shows the number of people in each interval.
- The graph depicts the count of people in different age categories, from 20 to 70.
- Higher number of women in each age category compared to men.
- A slight peak in the 30-34 age category can be seen which says more people fall under this age category for loan defaulters.
- There is a gradual decline after the age category 30-34 i.e., fewer loan defaulters after this age group.

## B.Distribution of Income

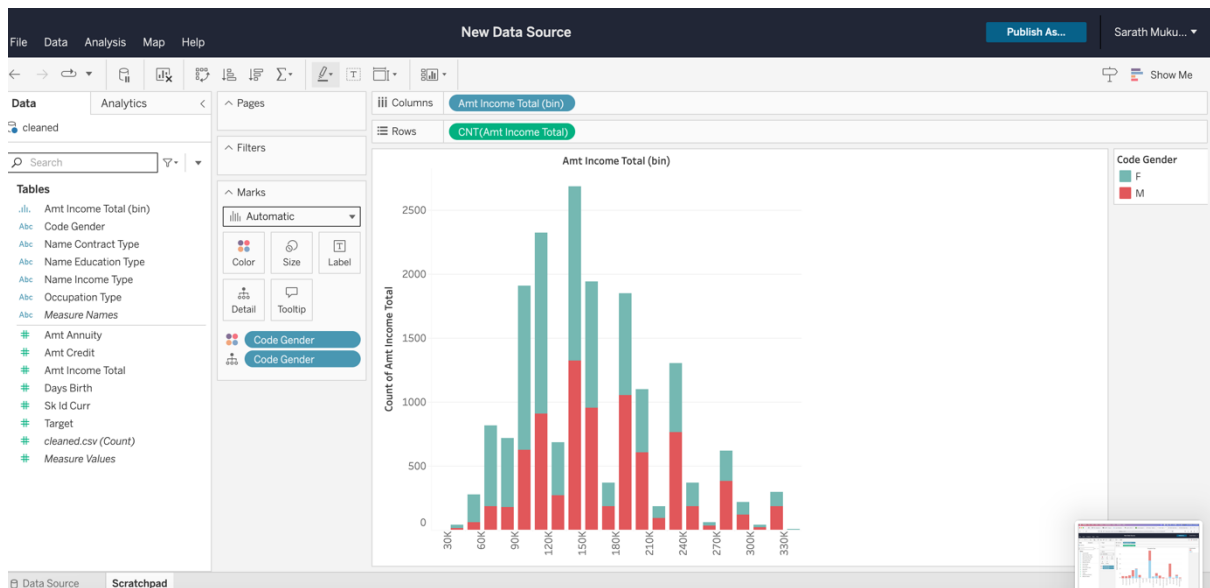The below graph shows the distribution of loan defaulters income based on gender:

Figure 5:Income Distribution

Insights from the graph:

- The graph shows the distribution of income for males and females, with males in red and females in green. The x-axis shows the income section, and the y-axis shows the number of people in each section.
- More number of people reside in the range of €150k with almost the same for males and females.
- In the high-income range, we can see comparatively fewer females than male loan defaulters.
- There can be seen high no. of defaulters in the mid-income range.

## C.Distribution of Occupation Type

The below graph represents the distribution of occupation type based on the gender of loan defaulters:
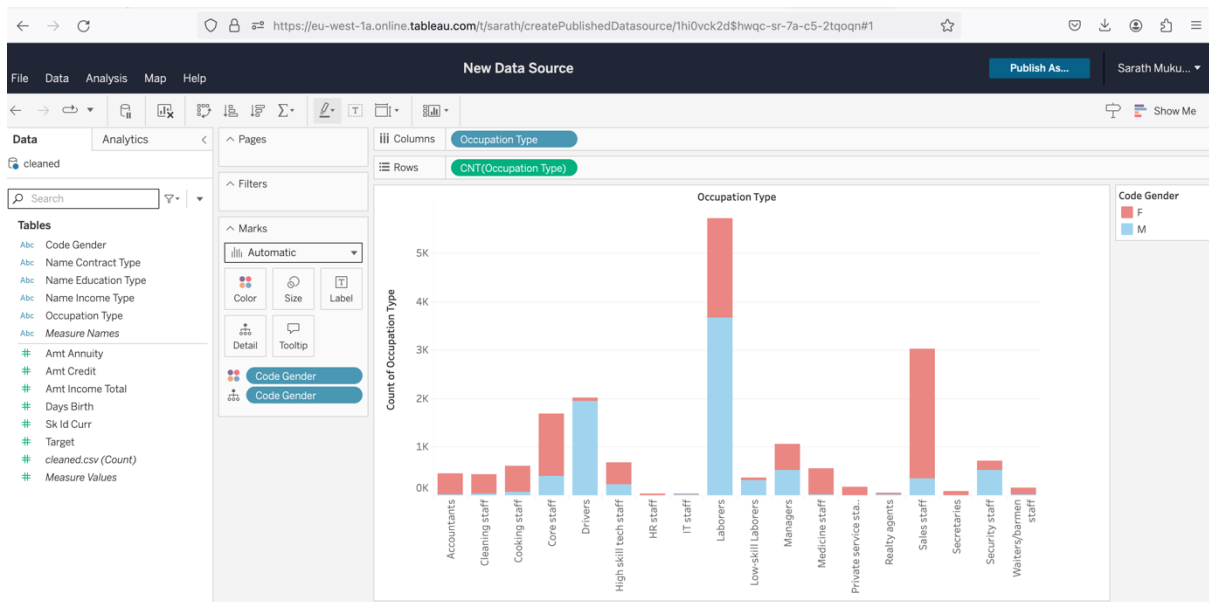
Figure 6: Occupation type distribution

Insights from the graph:

- The graph shows the distribution of occupation type for males and females, with males in blue and females in peach colour. The x-axis shows the occupation type, and the y-axis shows the number of people in each occupation.
- It can be seen more loan defaulters fall under the labourers category and the least in HR staff, realty agents and IT staff.
- More male loan defaulters in drivers and labourers occupations. In the case of female loan defaulters, more no. in labourers and sales staff.

# VISUALISATION

Tableau (Online) was employed to visualise a dashboard with the distribution of age, income and occupation type based on gender for loan defaulters which can seen below:
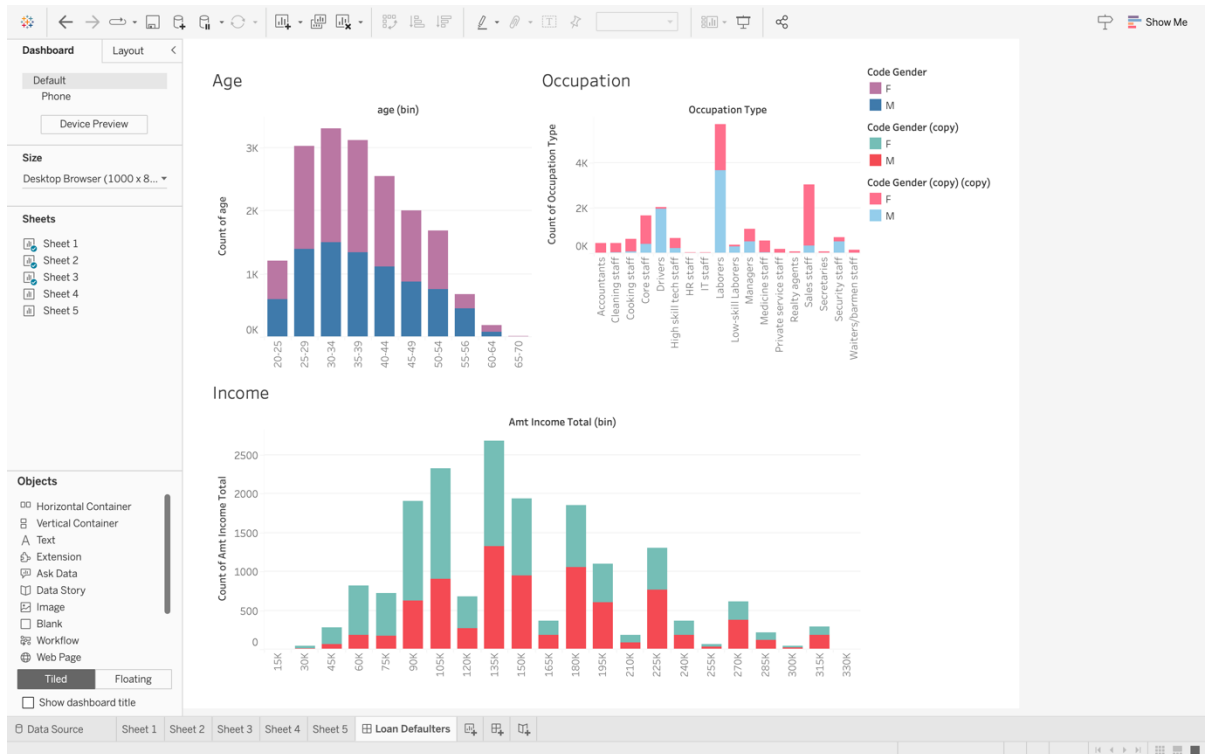


Figure 7: Dashboard

- The visualization shows a constant trend of more female loan defaulters throughout all age groups, with a noteworthy peak in the 30-34 age group. This knowledge could be critical for targeted financial initiatives or risk assessment aimed at reducing loan defaults, particularly among younger groups.

- A remarkable income distribution among male and female loan defaulters can be seen, emphasizing similarities in the middle and discrepancies in the upper-income groups. This knowledge can be useful in establishing targeted financial strategies or risk assessment models, particularly for high-income male defaulters and middle-income people of both genders.

- Loan defaulters are prevalent in the "Labourers" category, showing that this occupation type has the highest number of defaulters, regardless of gender.

- The graph, on the other hand, shows significantly lower counts of defaulters in the occupation categories of HR workers, real estate agents, and IT staff, implying that these occupations had fewer cases of loan defaults.

# CHALLENGES

- There were no challenges or difficulties in data processing and cleaning in PySpark.

- The configurations of different AWS were quite difficult to grasp and keep track of, it was a time-consuming process. Later, after understanding the concept we found it to be really simple.

- There was some issue while connecting the AWS S3 bucket to Tableau online but later after some research, it was resolved.

# RESPONSIBILITY SECTION

Regular communication and collaboration among all team members were done. To track work, share code snippets, and record any issues discovered, we used shared documentation like Google Docs. Search for the dataset was performed by everyone.

Alen:

- In charge of data input into AWS via PySpark.
- Created Python scripts to clean and pre-process raw data.
- Handled difficulties with missing values, outliers, and data formatting.
- Worked with Sarath to ensure data quality and consistency.

Sarath:

- Collaborated with Alen on data quality checks and pre-processing.
- Used Tableau to perform exploratory data analysis (EDA) to gain insights about age, income, employment, and loan defaulters based on gender.
- Performs feature engineering to prepare data for modelling and visualization.
- Documented the analysis procedure and results.

Suhara:

- Based on the preprocessed data, used Tableau to build insightful visualizations.
- Created dashboards and visual representations of age, income, occupation, and loan defaulters based on gender.
- Collaborated seriously with Alen and Sarath to fully comprehend the data and insights obtained.
- Wrote a detailed report outlining the important facts and insights.

# RELEVANT SCREENSHOTS

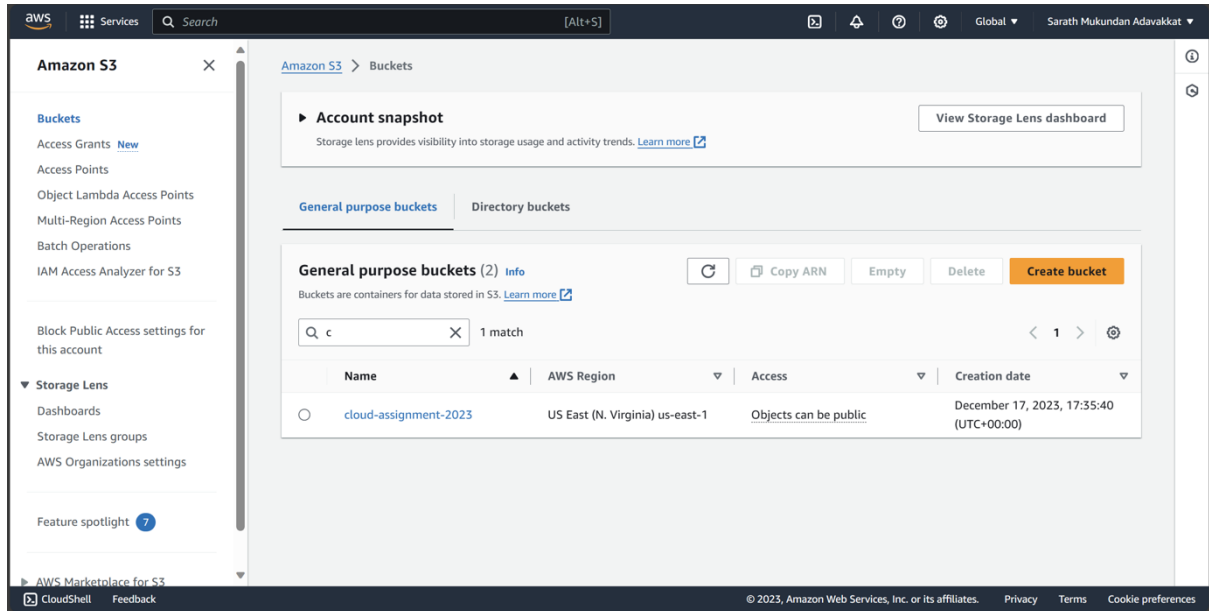1.Created S3 bucket named cloud-assignment-2023



Figure 8: S3 bucket cloud-assignment-S3
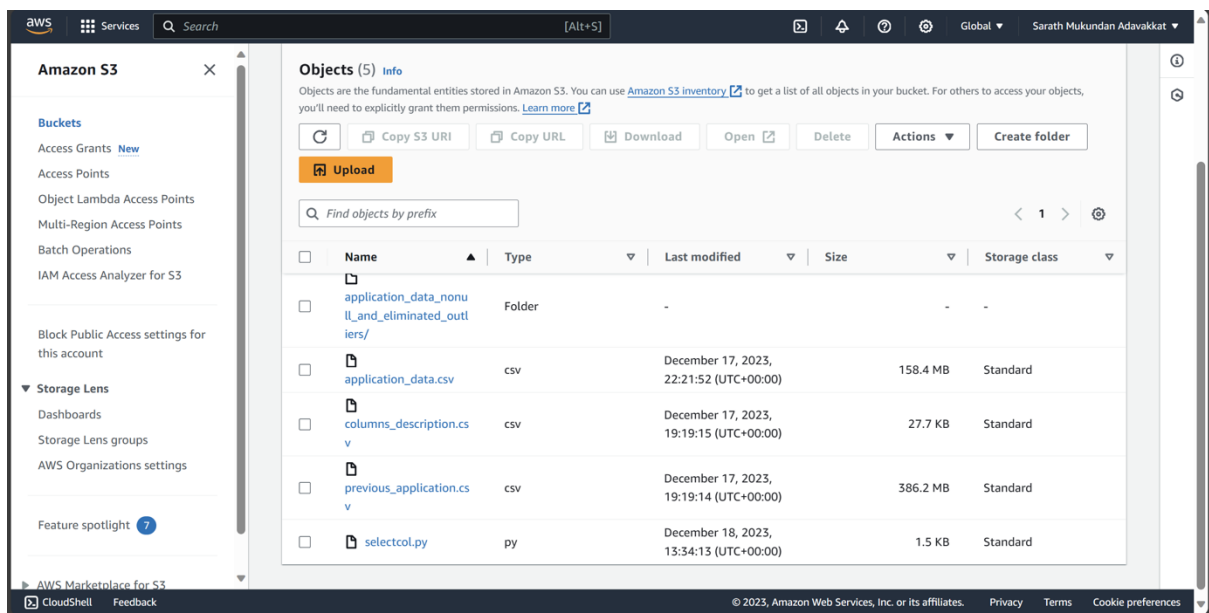
2. Stored our raw dataset and PySpark code.



Figure 9: Code and dataset stored

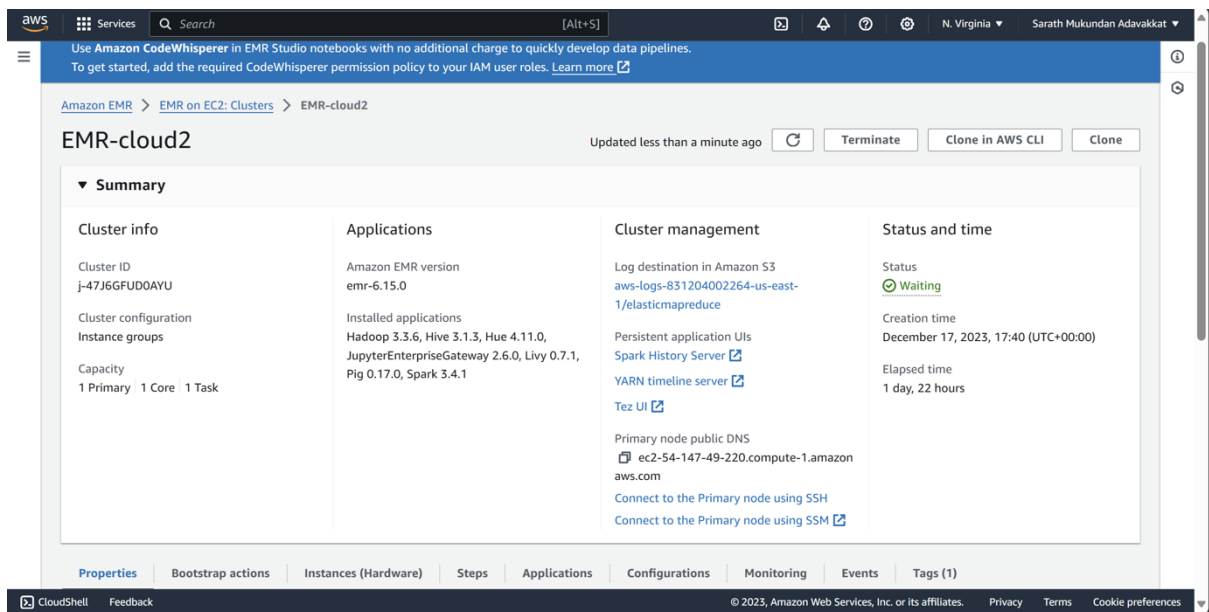## 3.Configuration of cluster created EMR-cloud2 using AWS EMR



Figure 10: Cluster EMR cloud-2
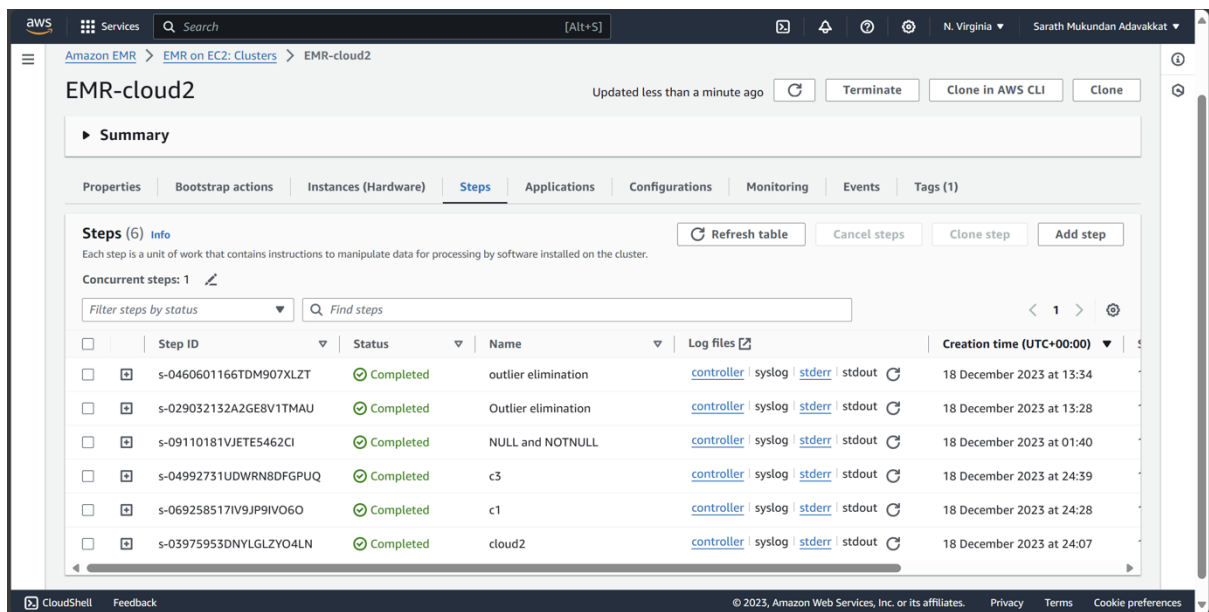
## 4.Executed PySpark codes using steps in EMR



Figure 11: Executed PySpark code

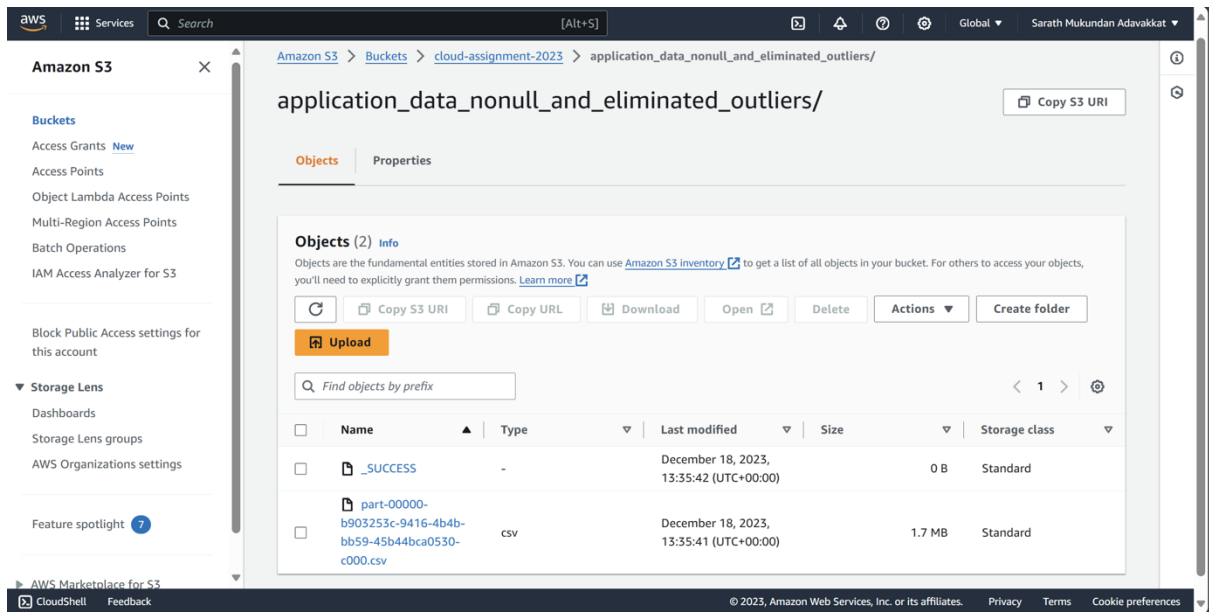5. The output of PySpark code executed is stored in the same S3 bucket under the folder application_data_nonull_and_eliminated_outliers



Figure 12: Output stored in S3 bucket