# Futurense

## Democratizing Tech Talent to deliver impact at scale

The
Godfather
of Talent | Futurense

# Credit Card Analysis

The
ZOOKEEPERS

GUIDE : Venkat
Sir
DE Expert
(Sensei)

Akhil
Shubham
Ashpak
Karthik
Kanika
Gaurav

# Business Problem/Overview

In this project, we have 2 problems – Transactions problem & Defaulters problem.

## Customer Spending Behavior Analysis

The primary objective of this project is to leverage big data technologies to perform an in-depth analysis of credit card transactions and credit card defaulters datasets.

By applying advanced data processing techniques in cloud, this project aims to uncover valuable insights and patterns that can assist in making informed decisions to mitigate credit card default risks and improve overall financial strategies.

- Gain insights into customer spending patterns and behaviors.
- Identify trends and patterns in transaction data.
- Optimize marketing strategies and tailor promotions based on transaction history.
- Enhance customer experiences by understanding their preferences and behaviors.

## Credit Card Defaulter Risk Analysis

The outcomes of credit card defaulter analysis empower credit issuers to make informed decisions, manage risks effectively, enhance customer relationships, and optimize their overall business strategies in the dynamic landscape of credit lending. The major goal is to reduce the risk of credit card defaults.

# Customer Spending Behavior Analysis

The primary objective of this project is to leverage big data technologies to perform an in-depth analysis of credit card transactions and credit card defaulters datasets.

By applying advanced data processing techniques in cloud, this project aims to uncover valuable insights and patterns that can assist in making informed decisions to mitigate credit card default risks and improve overall financial strategies.

- Gain insights into customer spending patterns and behaviors.
- Identify trends and patterns in transaction data.
- Optimize marketing strategies and tailor promotions based on transaction history.
- Enhance customer experiences by understanding their preferences and behaviors.

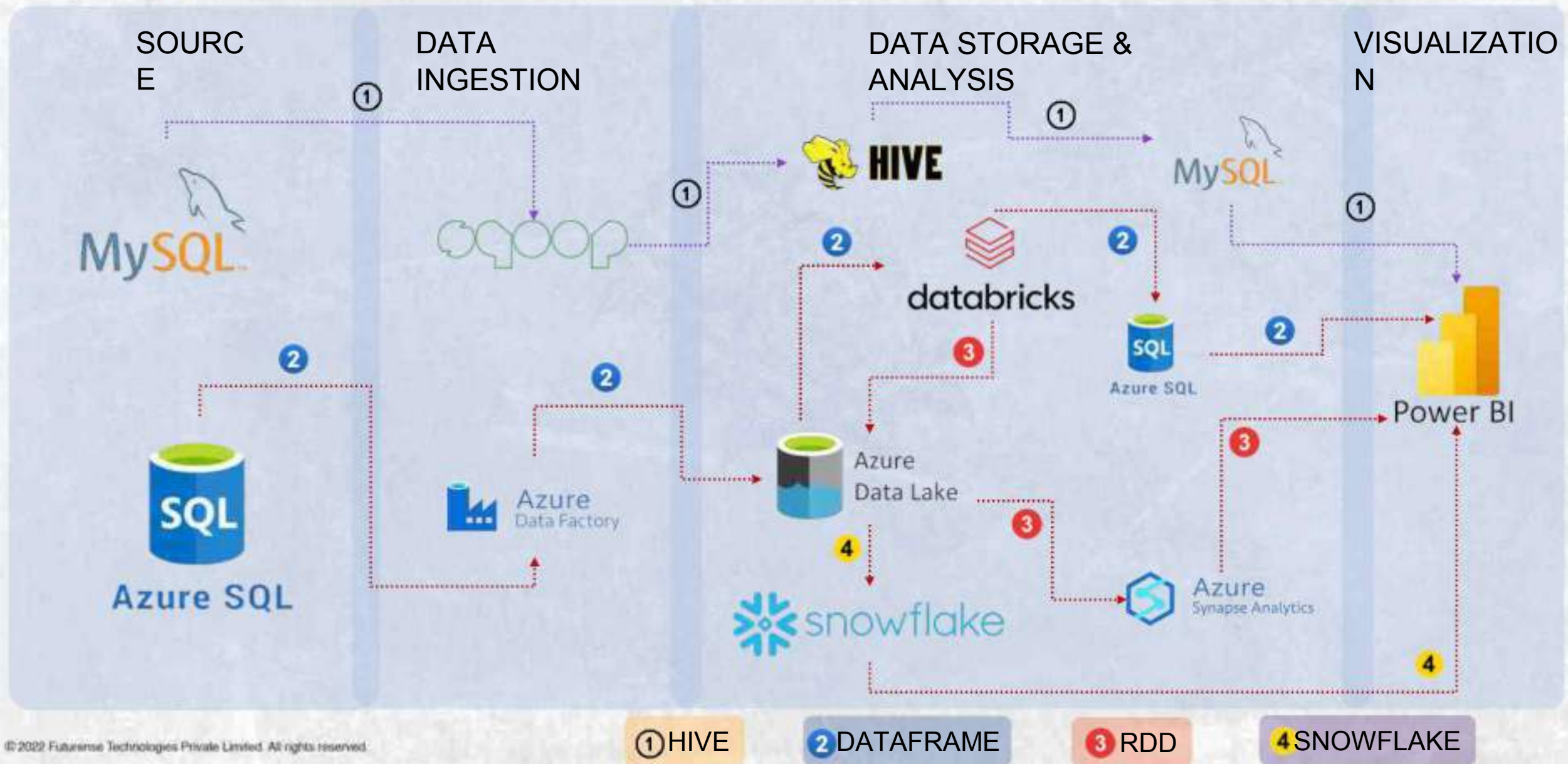# Credit Card Defaulter Risk Analysis

The outcomes of credit card defaulter analysis empower credit issuers to make informed decisions, manage risks effectively, enhance customer relationships, and optimize their overall business strategies in the dynamic landscape of credit lending. The major goal is to reduce the risk of credit card defaults.

# Architecture of the solution

# Data Representation

For the given problems, we have 2 different datasets – one for the transactions problem and another for Defaulters problem.

## Transactions Data

The transaction problem has the dataset as following features:

**Index:** A unique identifier for each record in the dataset.
**City:** The city where transaction was done.
**Date:** The date on which transaction occurred.
**Card Type:** Indicates the card type used for transactions – Silver, Gold, Platinum, Signature.
**Exp Type:** Indicates the expense type for which the card was used – Bills, Entertainment, Food, Fuel, Grocery Travel.
**Gender:** Denotes the gender of the cardholder – Male, Female.
**Amount:** The amount of transaction done by the customer.

Out of the given data, Index can be categorized as **Numerical identifier.**
City, Card Type, Exp Type and Gender as **Categorical variables**
Date as **Date variable** & Amount as **Numerical variable.**

Total Records: 26052

## Defaulters Data

The fraud detection problem has the dataset as following attributes:

**CustID:** A unique identifier for each customer.
**Limit_Bal:** Maximum spending limit assigned to the customer.
**Sex:** Gender of the customer – 1 (Male) or 2 (Female).
**Education:** Education level of the customer – 1 (Graduate), 2 (University), 3 (High school), 4 (Others).
**Marriage:** Marital status of the customer - 1 (Single), 2 (Married), 3 (Others).
**Age:** Age of the customer.
**PAY_1 to PAY_6:** Repayment status of the customer for the last 6 months. The values indicate the number of months of delayed for payment.
**BILL_AMT1 to BILL_AMT6:** Bill amount for each of the last six months.
**PAY_AMT1 to PAY_AMT6:** Actual amount customer paid for each of the last six months.
**DEFAULTED:** Whether the customer is defaulted or not on their credit card payment – 0 (not defaulted), 1 (defaulted).

Total Records: 1002

# Transactions Data

The transaction problem has the dataset as following features:

**Index**: A unique identifier for each record in the dataset.
**City**: The city where transaction was done.
**Date**: The date on which transaction occurred.
**Card Type**: Indicates the card type used for transactions – Silver, Gold, Platinum, Signature.
**Exp Type**: Indicates the expense type for which the card was used – **Bills, Entertainment, Food, Fuel, Grocery Travel.**
**Gender**: Denotes the gender of the cardholder – **Male, Female.**
**Amount**: The amount of transaction done by the customer.

Out of the given data, Index can be categorized as **Numerical identifier**,
City, Card Type, Exp Type and Gender as **Categorical variables**
Date as **Date variable** & Amount as **Numerical variable**.

Total Records

26052

# Defaulters Data

The fraud detection problem has the dataset as following attributes:

**CustID:** A unique identifier for each customer.

**Limit_Bal:** Maximum spending limit assigned to the customer.

**Sex:** Gender of the customer – **1** (Male) or **2** (Female).

**Education:** Education level of the customer – **1** (Graduate), **2** (University), **3** (High school), **4** (Others).

**Marriage:** Marital status of the customer - **1** (Single), **2** (Married), **3** (Others).

**Age:** Age of the customer.

**PAY_1 to PAY_6:** Repayment status of the customer for the last 6 months. The values indicate the number of months of delayed for payment.

**BILL_AMT1 to BILL_AMT6:** Bill amount for each of the last six months.

**PAY_AMT1 to PAY_AMT6:** Actual amount customer paid for each of the last six months.

**DEFAULTED:** Whether the customer is defaulted or not on their credit card payment – **0** (not defaulted), **1** (defaulted).

Total Records

1002

# Batch Processing using 🐝 HIVE

Creating the table

```
mysql> CREATE TABLE CCD (CUSTID INT,LIMIT_BAL DECIMAL(10, 2),SEX int,EDUCATION int,MARRIAGE int,AGE INT,PAY_1 INT,PA
Y_2 INT,PAY_3 INT,PAY_4 INT,PAY_5 INT,PAY_6 INT,BILL_AMT1 DECIMAL(10, 2),BILL_AMT2 DECIMAL(10, 2),BILL_AMT3 DECIMAL(
10, 2),BILL_AMT4 DECIMAL(10, 2),BILL_AMT5 DECIMAL(10, 2),BILL_AMT6 DECIMAL(10, 2),PAY_AMT1 DECIMAL(10, 2),PAY_AMT2 D
ECIMAL(10, 2),PAY_AMT3 DECIMAL(10, 2),PAY_AMT4 DECIMAL(10, 2),PAY_AMT5 DECIMAL(10, 2),PAY_AMT6 DECIMAL(10, 2),DEFAUL
TED INT);
```

Activate Windows

Using the following command we will load the data in the table created

```
LOAD DATA INFILE '/home/cloudera/CCD_ProcessedData.csv' INTO TABLE CreditCardData FIELDS
TERMINATED BY ',' (@var1,@var2,@var3,@var4,@var5,@var6,@var7) SET Date=STR_TO_DATE(@var1,'%d
-%m-
%Y'),Low=@var2,Open=@var3,Volume=@var4,high=@var5,Close=@var6,Adjusted _close=@var7;
```

# Batch Processing using HIVE

Transfer data using Sqoop

```
[cloudera@quickstart ~]$ sqoop import --connect jdbc:mysql://localhost:3306/project --username root --password clou
dera --table CCD --target-dir /user/cloudera/Shubz/CCD.txt -m 1
```
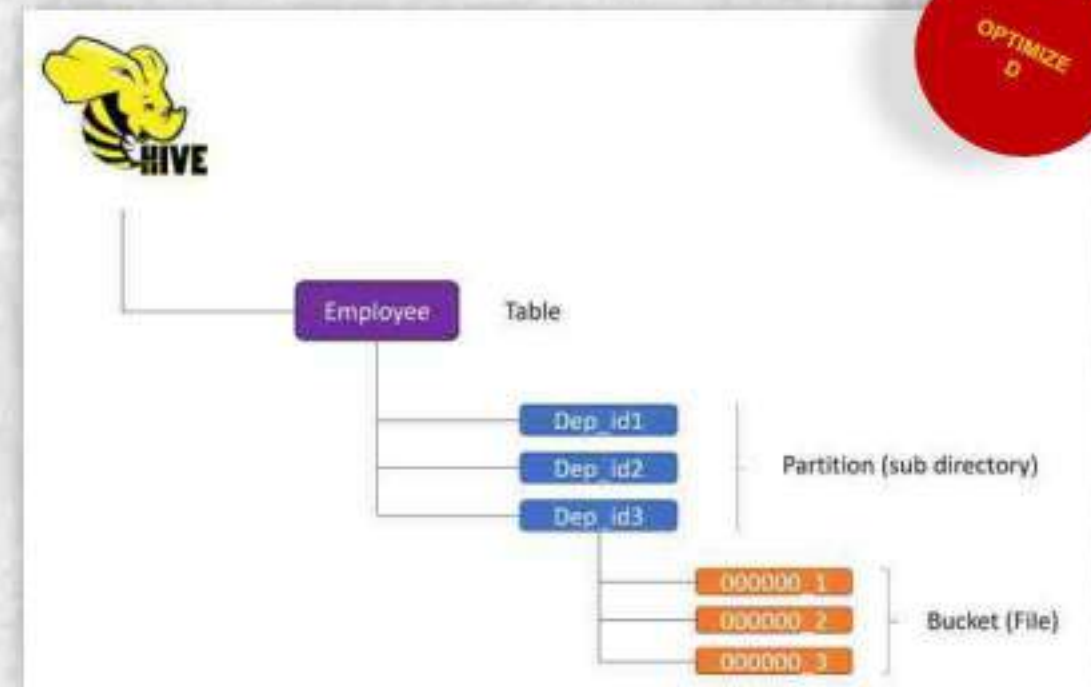
Create the schema in Hive

```
hive> create table CCD(CUSTID INT,LIMIT_BAL DECIMAL(10, 2),SEX int,EDUCATION int,MARRIAGE int,AGE INT,PAY_1 INT,PAY_
2 INT,PAY_3 INT,PAY_4 INT,PAY_5 INT,PAY_6 INT,BILL_AMT1 DECIMAL(10, 2),BILL_AMT2 DECIMAL(10, 2),BILL_AMT3 DECIMAL(10
, 2),BILL_AMT4 DECIMAL(10, 2),BILL_AMT5 DECIMAL(10, 2),BILL_AMT6 DECIMAL(10, 2),PAY_AMT1 DECIMAL(10, 2),PAY_AMT2 DEC
IMAL(10, 2),PAY_AMT3 DECIMAL(10, 2),PAY_AMT4 DECIMAL(10, 2),PAY_AMT5 DECIMAL(10, 2),PAY_AMT6 DECIMAL(10, 2),DEFAULTE
D INT) row format delimited fields terminated by ',';
```

# Batch Processing using 🐝 HIVE

Create partitioning on Hive table

```
hive> create table CCDP(CUSTID INT,LIMIT_BAL DECIMAL(10, 2),SEX int,EDUCATION int,MARRIAGE int,AGE INT,PAY_1 INT,PAY
_2 INT,PAY_3 INT,PAY_4 INT,PAY_5 INT,PAY_6 INT,BILL_AMT1 DECIMAL(10, 2),BILL_AMT2 DECIMAL(10, 2),BILL_AMT3 DECIMAL(1
0, 2),BILL_AMT4 DECIMAL(10, 2),BILL_AMT5 DECIMAL(10, 2),BILL_AMT6 DECIMAL(10, 2),PAY_AMT1 DECIMAL(10, 2),PAY_AMT2 DE
CIMAL(10, 2),PAY_AMT3 DECIMAL(10, 2),PAY_AMT4 DECIMAL(10, 2),PAY_AMT5 DECIMAL(10, 2),PAY_AMT6 DECIMAL(10, 2)) partit
ioned by (DEFAULTED INT) row format delimited fields terminated by
```

Partitioning is one of the optimization techniques used in Hive to improve the performance of the query

OPTIMIZED



| | |
|---|---|
| Employee | Table |
| Dep_id1 | |
| Dep_id2 | Partition (sub directory) |
| Dep_id3 | |
| 000000_1 | |
| 000000_2 | Bucket (File) |
| 000000_3 | |

# Batch Processing using HIVE

**1. Write a SQL query to determine the count of defaulted (1) and non-defaulted (0) records for both males and females in the Credit Card Data table.**

```
SELECT SEX,DEFAULTED,COUNT(*) AS COUNT FROM CCDP
GROUP BY SEX, DEFAULTED;
```

```
mysql> select * from statement6;

+------+-----------+-------------+
| Sex  | DEFAULTED | Total_count |
+------+-----------+-------------+
|    2 |         0 |         373 |
|    2 |         1 |         218 |
|    0 |         0 |           1 |
|    1 |         1 |         185 |
|    1 |         0 |         224 |
+------+-----------+-------------+
5 rows in set (0.01 sec)
```

**2. average billed Amount & average pay amount**

```
SELECT CUSTID, SEX, EDUCATION, MARRIAGE, AGE,
AVG((BILL_AMT1 + BILL_AMT2 + BILL_AMT3 + BILL_AMT4 +
BILL_AMT5 + BILL_AMT6) / 6) AS AVERAGE_BILLED_AMOUNT,
AVG((PAY_AMT1 + PAY_AMT2 + PAY_AMT3 + PAY_AMT4 +
PAY_AMT5 + PAY_AMT6) / 6) AS AVERAGE_PAY_AMOUNT FROM
CCDP GROUP BY CUSTID,SEX,EDUCATION,MARRIAGE,AGE;
```

```
mysql> select * from statement78 limit 10;

+--------+-----+-----------+----------+-----+-----------------------+--------------------+
| CUSTID | SEX | EDUCATION | MARRIAGE | AGE | AVERAGE BILLED AMOUNT | AVERAGE PAY AMOUNT |
+--------+-----+-----------+----------+-----+-----------------------+--------------------+
|    582 |   2 |         2 |        1 |  40 |              54306.50 |            3233.33 |
|    255 |   2 |         2 |        2 |  30 |              10500.00 |            4166.67 |
|    750 |   2 |         2 |        2 |  30 |               1699.67 |            1636.33 |
|      0 |   0 |         0 |        0 |   0 |                  0.00 |               0.00 |
|    256 |   1 |         2 |        1 |  30 |               4673.17 |            4654.50 |
|    257 |   2 |         2 |        1 |  50 |              67273.00 |            2201.67 |
|    258 |   2 |         2 |        1 |  30 |              32185.83 |            2862.83 |
|    259 |   2 |         3 |        1 |  40 |              60358.67 |            2300.00 |
|    260 |   1 |         1 |        1 |  50 |              39685.67 |            1295.33 |
|    261 |   2 |         1 |        2 |  30 |              56231.83 |           22051.83 |
+--------+-----+-----------+----------+-----+-----------------------+--------------------+
10 rows in set (0.01 sec)
```

# Risk Analysis using snowflake



**Data Pre-Processing**

**Data Ingestion**

Mounting

**Data Visualization**

**Data Storage**

**Data Analysis**

# Risk Analysis using snowflake

## ADLS MOUNTING TO DATABRICKS :

Connecting ADLS with Databricks

After completing the thorough data cleaning and preprocessing procedures, the processed and refined data is transferred to the Azure Data Lake Factory

Data Pre-processing

# Risk Analysis using snowflake

**Data Pre-Processing :** In the initial five scenarios, the data cleaning and preprocessing tasks are executed within the Databricks environment

Load the file into RDD

```
data = sc.textFile('/mnt/mounted_SA3/credit-card-default-1000.csv')
```

Seperating Header

```
headers= data.first()
```

```
headers1 = [name for name in headers.split(',')]
```

Convert to DataFrame

```
ccard_df = cleaned_rdd.toDF(headers1)
```

Removing header and unwanted inwerted comma , splitting columns

```
#removing header and unwanted inverted comma , splitting columns
rdd_split = data.filter(lambda x : x!=headers).map(lambda x : x.replace('"','')).map(lambda x : x.split(','))
```

Removing Lines that are not CSV

```
rdd_partial = rdd_split.filter(lambda x : x[0].isdigit())
```

Normalize sex to only 1 and 2

```
cleaned_rdd = rdd_partial.map(lambda x : [int(i.replace('M','1').replace('F','2')) for i in x])
```

# Risk Analysis using snowflake

**Data Ingestion :** Data Ingestion Pipeline from Data Lake Storage to Snowflake via External Stage and COPY INTO. From ADLS, the data is then seamlessly moved to Snowflake, a data warehouse platform, in order to facilitate subsequent in-depth analysis and exploration

**Connecting with SnowSQL**

**Create File Format**

# Risk Analysis using snowflake

**Create External Stage**

```
KANIKAAGG#COMPUTE_WH...DATABASE.ZOOSCHEMA>CREATE OR REPLACE STAGE zoo_azure_stage
                                          URL = 'azure://zooadlsaccount_blob_core.windows.net/cont1/CCdefaulter_csv/part-00000-tid-739541815169323954-5a2e6ce4-b86e
                                          -4cc6-082337fa0c96-11-1-c0000.csv'
                                          CREDENTIALS = (AZURE_SAS_TOKEN = '?sv=2022-11-02&ss=bfqt&srt=sco&sp=rwdlacupx&se=2023-08-15T22:43:44Z&st=2023-08-15T14
                                          :pr=https&sig=PIXE32PGyRq2aP6iNyslZCwnQ=ndb=%U48Ph77E1Q1AVgXXD' )
                                          FILE_FORMAT = zoo_csv_format;
+--------------------------------------------+
| status                                     |
+--------------------------------------------+
| Stage area ZOO_AZURE_STAGE successfully created. |
+--------------------------------------------+
1 Row(s) produced. Time Elapsed: 0.247s
```

**Copy Into Snowflake Internal Table**

```
KANIKAAGG#COMPUTE_WH@ZOODATABASE.ZOOSCHEMA>COPY INTO cc_default
                                           FROM @zoo_azure_stage
                                           FILE_FORMAT = zoo_csv_format;
+------+
| file
rsed | rows_loaded | error_limit | errors_seen | first_error | first_error_line | first_error_character | first_error_column_name |
+------+
| azure://zooadlsaccount.blob.core.windows.net/cont1/CCdefaulter_csv/part-00000-tid-739541815169323954-5a2e6ce4-b86e-4cc6-ae2d-082337fa0c
1000 |        1000 |           1 |           0 | NULL        |             NULL |                  NULL | NULL                    |
+------+
1 Row(s) produced. Time Elapsed: 2.623s
```

# Risk Analysis using snowflake

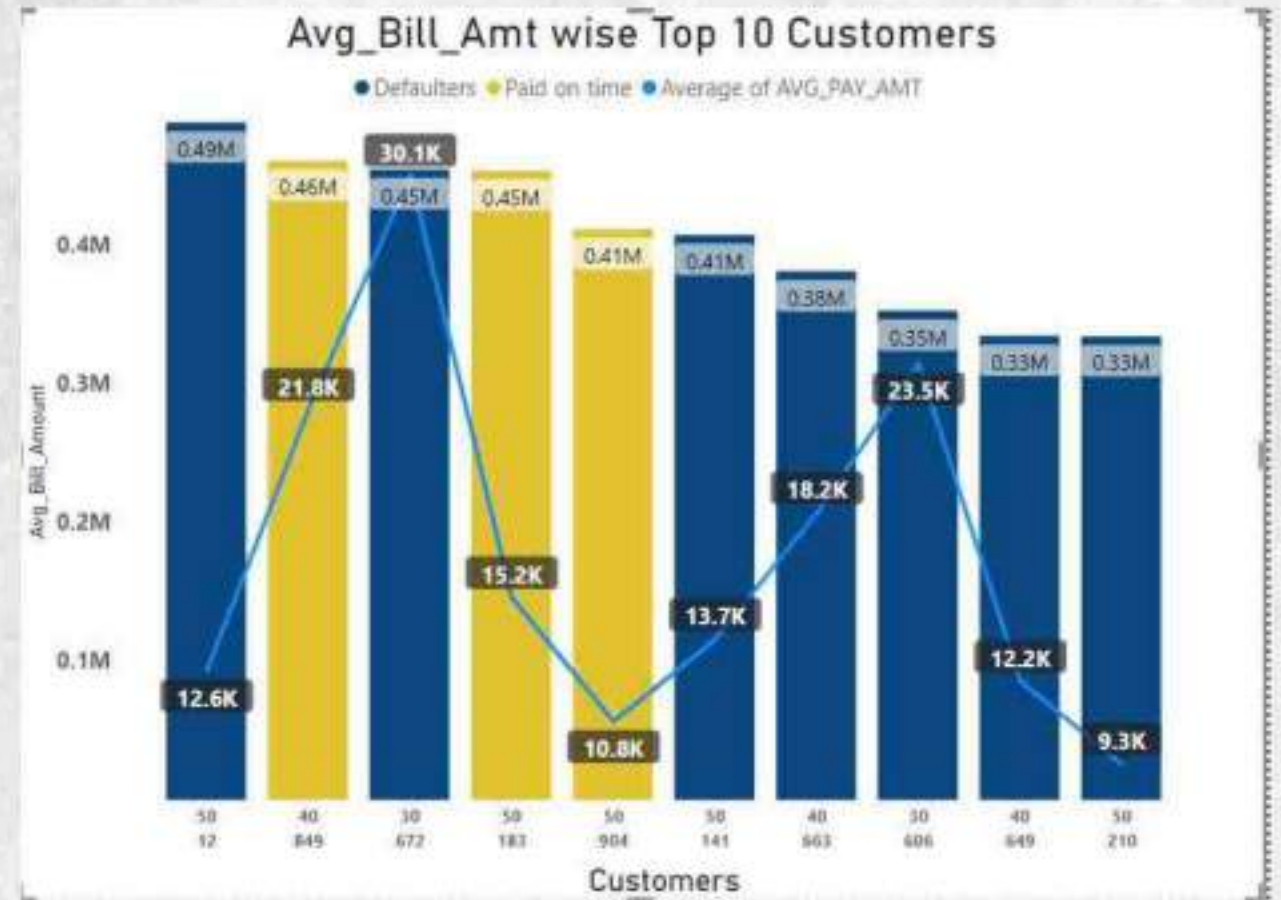**Data Processing/ Analysis :** Rounding of age to range of 10s and add SEXNAME to the data using SQL Joins.

**Create Gender table and insert values**
Create table gender (

sex_code int,

sex_name varchar(50) );


select *, FLOOR((Age + 5) / 10) * 10 AS RoundedAge

from cc_default c

inner join gender g on c.sex=g.sex_code;



Gender
● Female ● Male
409 (40.9%)
591 (59.1%)



CREDIT CARD DEFAULTERS
● Paid on time ● Defaulters
403 (40.3%)
597 (59.7%)



AGE RANGE WISE DEFAULTERS
● Defaulters ● Paid on time
Total Customers
Rounded_age

# Risk Analysis using snowflake

**Data Processing/ Analysis :** Rounding of age to range of 10s and add SEXNAME to the data using SQL Joins.

**Data Processing/ Analysis :** Find Average billed Amount and average pay amount for each customer.

```
select *, ROUND(
(CASE WHEN BILL_AMT1 > 0 THEN BILL_AMT1 ELSE 0 END
+

CASE WHEN BILL_AMT2 > 0 THEN BILL_AMT2 ELSE 0 END
+

CASE WHEN BILL_AMT3 > 0 THEN BILL_AMT3 ELSE 0  END
+

 CASE WHEN BILL_AMT4 > 0 THEN BILL_AMT4 ELSE 0 END
+

 CASE WHEN BILL_AMT5 > 0 THEN BILL_AMT5 ELSE 0 END
+

  CASE WHEN BILL_AMT6 > 0 THEN BILL_AMT6 ELSE 0
END )/6,2)

AS AVG_BILL_AMT,

round((PAY_AMT1+PAY_AMT2+PAY_AMT3+PAY_AMT4+PAY_AMT5+PAY_AMT6)/6,2) AS AVG_PAY_AMT

FROM CC_DEFAULT;
```



Avg_Bill_Amt wise Top 10 Customers

# Risk Analysis using snowflake

**Data Processing/ Analysis :** Find average pay duration. Make sure numbers are rounded and negative values are eliminated.

```
SELECT *, ROUND ((

CASE WHEN PAY_1 > 0 THEN PAY_1ELSE 0 END +

 CASE WHEN PAY_2 > 0 THEN PAY_2ELSE 0 END +

 CASE WHEN PAY_3 > 0 THEN PAY_3ELSE 0 END +

 CASE WHEN PAY_4 > 0 THEN PAY_4ELSE 0 END +

 CASE WHEN PAY_5 > 0 THEN PAY_5ELSE 0 END +

 CASE WHEN PAY_6 > 0 THEN PAY_6ELSE 0 END

   )/6, 0) AS AVG_PAY_DURATION

FROM

   CC_DEFAULT;
```

| 0.00 | 739 |
|------|-----|
| Min of AVGPAYDURATION | Count of CUSTID |
| **3.00** | **6** |
| Max of AVGPAYDURATION | Count of CUSTID |

# Batch Processing using RDD

# Cluster Configuration

# Importing Data



The source for us is MS SQL Server. We are importing the data from the SQL Server into the Azure cloud, using the Azure Data Factory pipeline.

# Batch Processing using RDD



## Mounting ADLS

Cmd 3

```python
SAS_Token = "?sv=2022-11-02&ss=bfqt&srt=co&sp=rwdlacupyx&se=2023-08-20T12:48:32Z&st=2023-08-09T04:48:32Z&spr=https&sig=pZN3UpPsvxFn8Fz276yGvAdSu2x8GOl%2B7KxAjl7Ifjc%3D"
```

Command took 0.06 seconds -- by theashpak_5@live.com at 9/8/2023, 10:20:59 am on Ashpak Sheikh's Cluster

Cmd 4

```python
1
2    dbutils.fs.mount(
3        source = 'wasbs://project@meraaccountdeletekarega.blob.core.windows.net',
4        mount_point = '/mnt/mounted_SAS',
5        extra_configs={
6            "fs.azure.sas.project.meraaccountdeletekarega.blob.core.windows.net":SAS_Token
7        }
8    )
```

Out[9]: True

Command took 10.60 seconds -- by theashpak_5@live.com at 9/8/2023, 10:21:07 am on Ashpak Sheikh's Cluster

# Batch Processing using RDD

Write a query to print top 5 cities with highest spends and their percentage contribution of total credit card spends

Cmd 11

```python
1   # Mapping the data to (City, Amount) pairs
2   city_amount_rdd = rdd_split.map(lambda row: (row[1], float(row[7])))
3   # Total spends per city
4   city_total_spends_rdd = city_amount_rdd.reduceByKey(lambda a, b: a + b)
5
6   # Total spends across all cities
7   total_spends = city_total_spends_rdd.values().sum()
8
9   # Total spends in descending order
10  sorted_cities = city_total_spends_rdd.sortBy(lambda x: x[1], ascending=False)
11
12  # Top 5 cities
13  top_cities = sorted_cities.take(5)
14
15  # Percentage contribution of each city's spends
16  city_spends_percentage = [(city, (amount / total_spends) * 100, amount) for city, amount in top_cities]
17  # Printing the results
18  for city,percentage,amount in city_spends_percentage:
19      print("City: {}, Amount: {}, Spends: {:.2f}%".format(city, amount,percentage))
20
```

▶ (4) Spark Jobs

```
City: Greater Mumbai, Amount: 576751476.0, Spends: 14.15%
City: Bengaluru, Amount: 572326739.0, Spends: 14.05%
City: Ahmedabad, Amount: 567794310.0, Spends: 13.93%
City: Delhi, Amount: 556929212.0, Spends: 13.67%
City: Kolkata, Amount: 115466943.0, Spends: 2.83%
```

Command took 1.13 seconds -- by theashpak.5@live.com at 17/9/2023, 11:40:22 am on Ashpak Sheikh's Cluster

# Batch Processing using RDD



Write a query to print highest spend month and amount spent in that month for each card type

Cmd 13

```python
1    #'0', 'Delhi', ' India', '29-Oct-14', 'Gold', 'Bills', 'F', '82475'
2    # Splitting data monthwise
3    card_month_amount_rdd = rdd_split.map(lambda x: (x[3].split('-')[1], float(x[7])))
4
5    # Calculating Monthly Spent
6    monthly_sum = card_month_amount_rdd.reduceByKey(lambda x,y : x+y).max(lambda x : x[1])
7
8    #Getting all the data from month having maximum spent
9    card_month_amount_rdd = rdd_split.filter(lambda x : x[3].split('-')[1]==monthly_sum[0]).map(lambda x: ((x[4], x[3].split('-')[1]),
     float(x[7])))
10
11   # calculating spent by each card type
12   monthly_grouped = card_month_amount_rdd.reduceByKey(lambda x,y:x+y)
13
14   for card_type, max_month in monthly_grouped.collect():
15       print(f"Card Type : {card_type}, Spent : {max_month}")
```

▶ (2) Spark Jobs

Card Type : ('Platinum', 'Jan'), Spent : 112784373.0
Card Type : ('Signature', 'Jan'), Spent : 98919381.0
Card Type : ('Silver', 'Jan'), Spent : 109359598.0
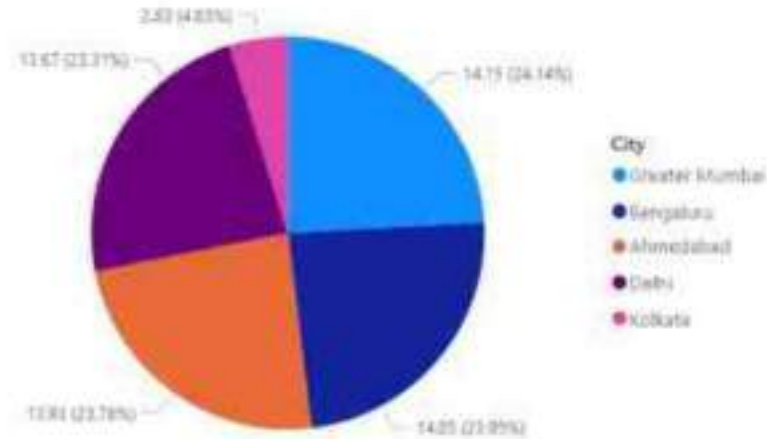Card Type : ('Gold', 'Jan'), Spent : 110146204.0

Command took 0.38 seconds -- by theashpak_5@live.com at 9/8/2023, 5:36:34 pm on Ashpak Sheikh's Cluster

# Optimization using RDD

**VISUALIZATION:**

**Scenario 6 - Write a query to find percentage contribution of spends by females for each expense type**

```
Cmd 7

1   female_rdd = rdd2.filter(lambda x: x[6] == "F")
2   # Map the RDD to (expense_type, amount)
3   expense_type_amount_rdd = female_rdd.map(lambda x: (x[5], float(x[7])))
4   # Reduce by key to calculate total spend by females for each expense type
5   total_spend_by_expense_type = expense_type_amount_rdd.reduceByKey(lambda a, b: a + b)
6   # Collect the total spend by expense type as a dictionary for easy lookup
7   total_spend_dict = dict(total_spend_by_expense_type.collect())
8   # Calculate the total spend by all genders for each expense type
9   total_spend_all_rdd = rdd2.map(lambda x: (x[5], float(x[7]))).reduceByKey(lambda a, b: a + b)
10  # Calculate the percentage contribution of spends by females for each expense type
11  percentage_contribution_rdd = total_spend_by_expense_type.join(total_spend_all_rdd).mapValues(lambda x: (x[0] / x[1]) * 100)
12  # Collect and print the result
13  result = percentage_contribution_rdd.collect()
14  for expense_type, percentage in result:
15      print(f"Expense Type: {expense_type}, Percentage Contribution: {percentage:.2f}%")
```

▶ (2) Spark Jobs

```
Expense Type: Bills, Percentage Contribution: 63.95%
Expense Type: Entertainment, Percentage Contribution: 49.37%
Expense Type: Grocery, Percentage Contribution: 50.91%
Expense Type: Fuel, Percentage Contribution: 49.71%
Expense Type: Food, Percentage Contribution: 54.91%
Expense Type: Travel, Percentage Contribution: 51.13%
```

# Batch Processing using RDD

**VISUALIZATION:**



Sum of Female_Spend_Percentage by Exp_Type

# Batch Processing using RDD

**Scenario 7 - Which card and expense type combination saw highest month over month growth in Jan -2014**

Cmd 12

```python
1   # Filter data for January 2014
2   jan_2014_data_rdd = rddy.filter(lambda x: x[3].split('-',1)[1] ==('Jan-14'))
3   Dec_2013_data_rdd = rddy.filter(lambda x: x[3].split('-',1)[1] ==('Dec-13'))
4   # Map the RDD to ((card_type, expense_type), amount)
5   card_expense_amount_rdd_jan = jan_2014_data_rdd.map(lambda x: ((x[4], x[5]), float(x[7])))
6   card_expense_amount_rdd_Dec = Dec_2013_data_rdd.map(lambda x: ((x[4], x[5]), float(x[7])))
7   # Reduce by key to calculate total amount spent for each card and expense type combination
8   total_amount_by_card_expense_jan = card_expense_amount_rdd_jan.reduceByKey(lambda a, b: a + b)
9   total_amount_by_card_expense_Dec = card_expense_amount_rdd_Dec.reduceByKey(lambda a, b: a + b)
10  final_total_amountbycardexpense=total_amount_by_card_expense_jan.join(total_amount_by_card_expense_Dec)
11  # Calculate the growth from the previous month (December 2013) for each card and expense type combination
12  growth_rdd = final_total_amountbycardexpense.map(lambda x : (x[0],100*(x[1][0]-x[1][1])/x[1][1]))
13  # Find the combination with the highest month-over-month growth
14  max_growth_combination = growth_rdd.max(lambda x: x[1])
15  for i in max_growth_combination:
16      print(i)
```

▶ (1) Spark Jobs

('Gold', 'Travel')
87.92008147034576

Command took 0.91 seconds — by niharikajs03321@gmail.com at 8/15/2023, 11:54:49 PM on Niharika J 5's Cluster

# Interactive Processing using DataFrame

# Connecting to Azure SQL Database & Reading the data

**Azure SQL Database Configuration**

```
Cmd 3

1   jdbcHostname = "ccaserver.database.windows.net"
2   jdbcPort = "1433"
3   jdbcDatabase = "credit_analysis"
4
5   url = f"jdbc:sqlserver://{jdbcHostname}:{jdbcPort};database={jdbcDatabase}"
```

▶ ▦ data: pyspark.sql.dataframe.DataFrame = [index: short, City: string ...5 more fields]

Command took 0.16 seconds -- by niharika[ja0321]@gmail.com at 6/15/2023, 3:41:22 PM on Niharika J S's Cluster

**Reading Table from Azure SQL Database as DataFrame using Spark Read API**

```
Cmd 5

1   credit_card = spark.read.format("jdbc") \
2                        .option("driver", "com.microsoft.sqlserver.jdbc.SQLServerDriver") \
3                        .option("url", url) \
4                        .option("dbtable", "credit_card_transactions") \
5                        .option("user", "sqladmin") \
6                        .option("password", "Password@123") \
7                        .load()
```

▶ ▦ credit_card: pyspark.sql.dataframe.DataFrame = [index: short, City: string ...5 more fields]

Command took 0.18 seconds -- by niharika[ja0321]@gmail.com at 6/15/2023, 3:51:55 PM on Niharika J S's Cluster

# Cleaning Data and Creating Temporary Table

## Removing country name from City column

Cmd 11

```
1    from pyspark.sql.functions import *
2    # Creating an UDF for removing the country name
3    func = udf(lambda x: str(x)[:-7])
4
5    # Applying UDF function on the column
6    credit_card = credit_card.withColumn('City', func(col('City')))
```

▸ ▦ credit_card: pyspark.sql.dataframe.DataFrame = [index: integer, City: string ... 5 more fields]

Command took 0.21 seconds -- by niharikaju8321@gmail.com at 3/15/2023, 7:05:05 PM on Niharika J 5's Cluster

## Creating temporary table for data frame

Cmd 11

```
1    credit_card.createOrReplaceTempView('credit_tbl')
```

Command took 0.18 seconds -- by niharikaju8321@gmail.com at 3/15/2023, 1:11:03 PM on Niharika J 5's Cluster

# Problem Statement

**Write a query to print 3 columns: city, highest_expense_type, lowest_expense_type**

Scenario 5 - Write a query to print 3 columns: city, highest_expense_type, lowest_expense_type

Cmd 38

```
1   sc5_out = spark.sql(""" WITH cte AS
2                         (SELECT City, Exp_Type, SUM(Amount)
3                          FROM credit_tbl
4                          GROUP BY City, EXP_TYPE
5                          ORDER BY 1, 3 DESC
6                         )
7                         SELECT DISTINCT City,
8                              FIRST_VALUE(Exp_Type) OVER(PARTITION BY City) AS Highest_Expense_Type,
9                              LAST_VALUE(Exp_Type) OVER(PARTITION BY City) AS Lowest_Expense_Type
10                        FROM cte
11                        """)
```

▶ ▥ sc5_out: pyspark.sql.dataframe.DataFrame = [City: string, Highest_Expense_Type: string ... 1 more field]

Command took 1.17 seconds — by niharikajs8321@gmail.com at 8/15/2023, 7:06:06 PM on Niharika J S's Cluster

# Problem Statement

**Output:**



```
1    sc5_out.show()
```

▸ (5) Spark Jobs

```
+----------+-------------------+------------------+
|      City|Highest_Expense_Type|Lowest_Expense_Type|
+----------+-------------------+------------------+
|  Achalpur|            Grocery|     Entertainment|
|  Adilabad|              Bills|              Food|
| Adityapur|               Food|           Grocery|
|     Adoni|              Bills|     Entertainment|
|     Adoor|               Fuel|             Bills|
|  Afzalpur|               Fuel|              Food|
|  Agartala|            Grocery|              Food|
|      Agra|              Bills|           Grocery|
| Ahmedabad|              Bills|           Grocery|
|Ahmednagar|               Fuel|           Grocery|
|    Aizawl|               Food|           Grocery|
|     Ajmer|      Entertainment|              Fuel|
|     Akola|              Bills|              Fuel|
|      Akot|               Fuel|     Entertainment|
| Alappuzha|               Food|     Entertainment|
|   Aligarh|              Bills|     Entertainment|
|Alipurduar|               Food|     Entertainment|
| Alirajpur|      Entertainment|     Entertainment|
```

Command took 4.18 seconds — by niharikaju0321@gmail.com at 6/15/2023, 7:06:00 PM on Niharika 2.5's Cluster

# Catalyst Optimizer

# Exporting result to Azure SQL Database

**Exporting result DataFrame to Azure SQL Server Table**

Cmd 37

```
1   sc5_out.write.format("jdbc") \
2           .option("driver", "com.microsoft.sqlserver.jdbc.SQLServerDriver") \
3           .option("url", url) \
4           .option("dbtable", "Scenario_5") \
5           .option("user", "sqladmin") \
6           .option("password", "Password@123") \
7           .save()
```
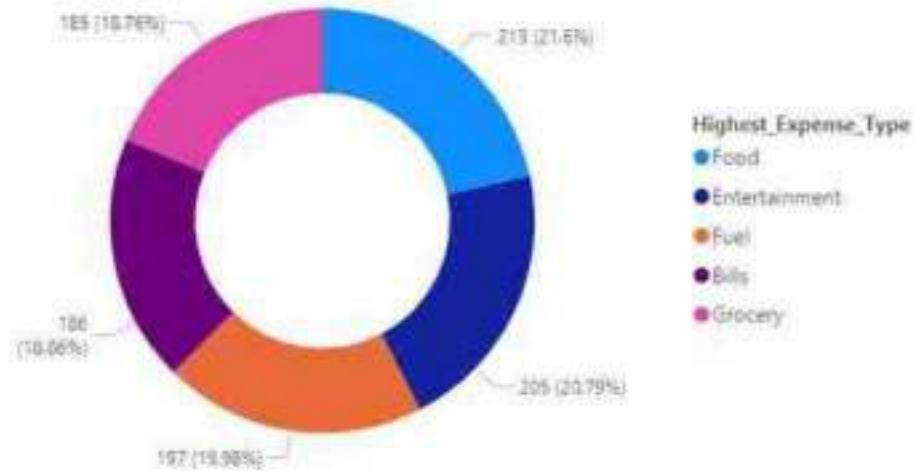
▶ (5) Spark Jobs

Command took 4.59 seconds -- by niharikajs03210gmail.com at 8/10/2023, 11:39:28 AM on Niharika J S's Cluster
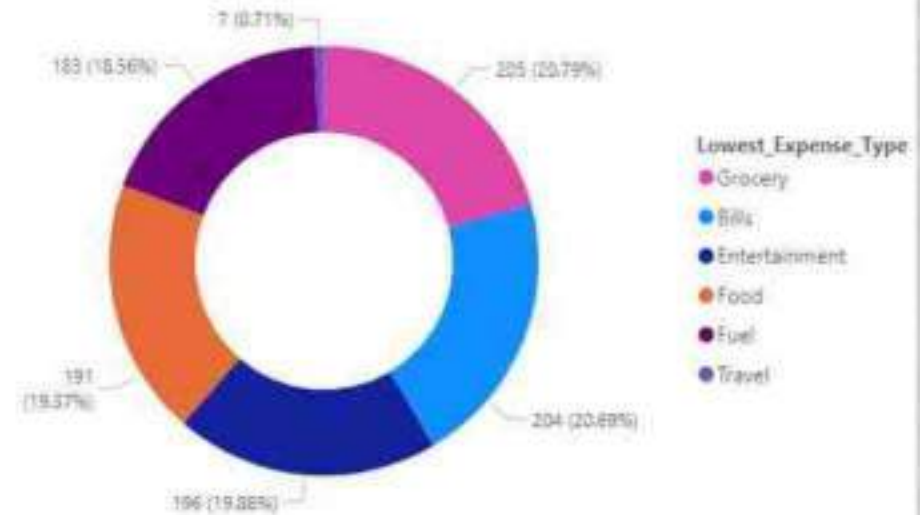
# Problem Statement

**Which city took least number of days to reach its 500th transaction**



Scenario 9 - Which city took least number of days to reach its 500th transaction after the first transaction in that city

```
sc9_out = spark.sql(""" WITH cte1 AS
                (SELECT City, Date, ROW_NUMBER() OVER(PARTITION BY City ORDER BY Date) AS RN
                 FROM credit_tbl
                ), cte2 AS
                (SELECT *, LAST_VALUE(RN) OVER(PARTITION BY City) AS Low
                 FROM cte1
                 WHERE RN<=500
                ), cte3 AS
                (SELECT DISTINCT City, DATEDIFF(MAX(Date) OVER(PARTITION BY City), Date) AS Difference
                 FROM cte2 WHERE Low=500
                )
                SELECT City, MAX(Difference) Days_Took
                FROM cte3
                GROUP BY City
                ORDER BY 2 LIMIT 1
                """)
```

sc9_out: pyspark.sql.dataframe.DataFrame = [City: string, Days_Took: integer]

```
sc9_out.show()
```

```
+---------+---------+
|     City|Days_Took|
+---------+---------+
|Bengaluru|       81|
+---------+---------+
```

# Scheduling Databricks Notebook in ADF for Periodical Analysis

# Roadblocks

- Faced issues with data understanding and handling irrelevant data.

**Resolution**: Consulting with domain experts who can help you determine which features are important for credit card analysis.

- The automated data ingestion Snowflake pipeline from ADLS to Snowflake encountered challenges due to a region mismatch between the two platforms and unavailability of an enterprise version of Snowflake.

**Resolution:** Manually loaded the data using COPY INTO and External stage for data ingestion.

# Conclusion

- **Customer Insights:** Through in-depth analysis, we gained valuable insights into how customers use their credit cards, revealing spending patterns and preferences.

- **Risk Identification:** Our analysis helped identify potential risks, enabling us to proactively manage and minimize the chances of defaults.

- **Smart Decisions:** Armed with data-driven insights, we can make informed decisions that benefit both customers and the business.

- **Financial Well-being:** By analyzing behavior, we contribute to customers' financial well-being by offering suitable credit limits and advice.

Thank You