



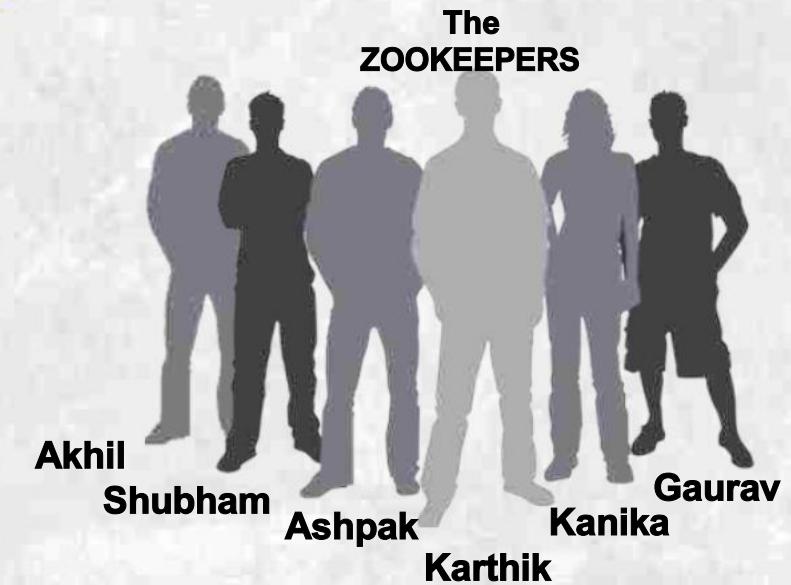
Futurenese

Democratizing Tech Talent
to deliver impact at scale

Credit Card Analysis



GUIDE : Venkat
Sir
DE Expert
(Sensei)



Business Problem/Overview



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

Customer Spending Behavior Analysis

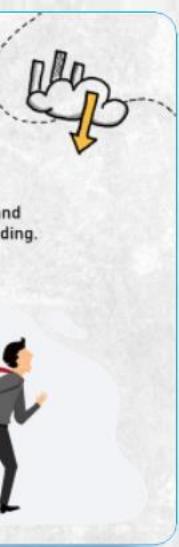


The primary objective of this project is to leverage big data technologies to perform an in-depth analysis of credit card transactions and credit card defaulters datasets. By applying advanced data processing techniques in cloud, this project aims to uncover valuable insights and patterns that can assist 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.

© 2022 Futurense Technologies Private Limited. All rights reserved.

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.

© 2022 Futurense Technologies Private Limited. All rights reserved.

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.



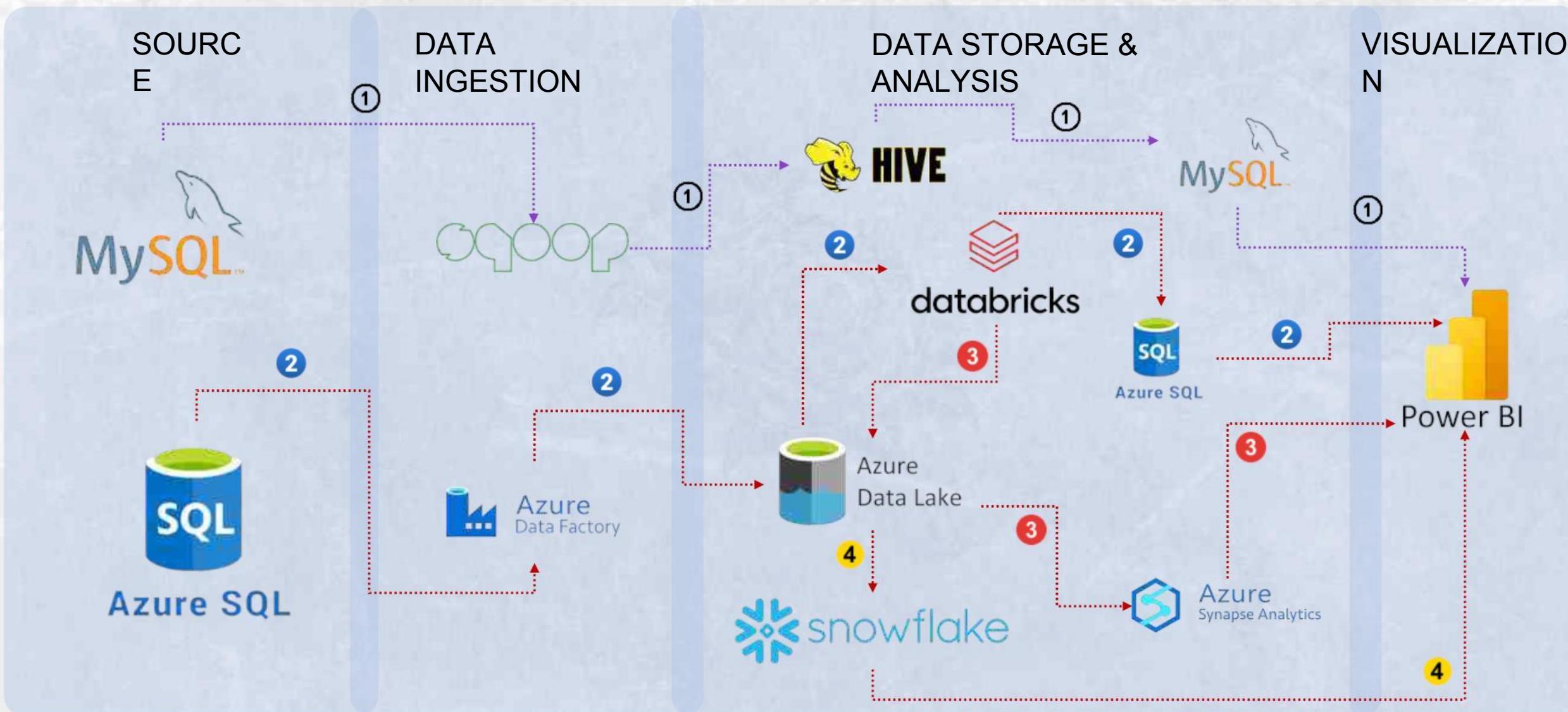
Credit Card Defaulter Risk Analysis



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 DECIMAL(10, 2),PAY_AMT3 DECIMAL(10, 2),PAY_AMT4 DECIMAL(10, 2),PAY_AMT5 DECIMAL(10, 2),PAY_AMT6 DECIMAL(10, 2),DEFALTED INT);
```

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 cloudera --table CCD --target-dir /user/cloudera/Shubz/CCD.txt -m 1
```

Activate Windows

Go to Settings to activate Windows.

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 DECIMAL(10, 2),PAY_AMT3 DECIMAL(10, 2),PAY_AMT4 DECIMAL(10, 2),PAY_AMT5 DECIMAL(10, 2),PAY_AMT6 DECIMAL(10, 2),DEFAULTED INT) row format delimited fields terminated by ',';
```

Activate Windows

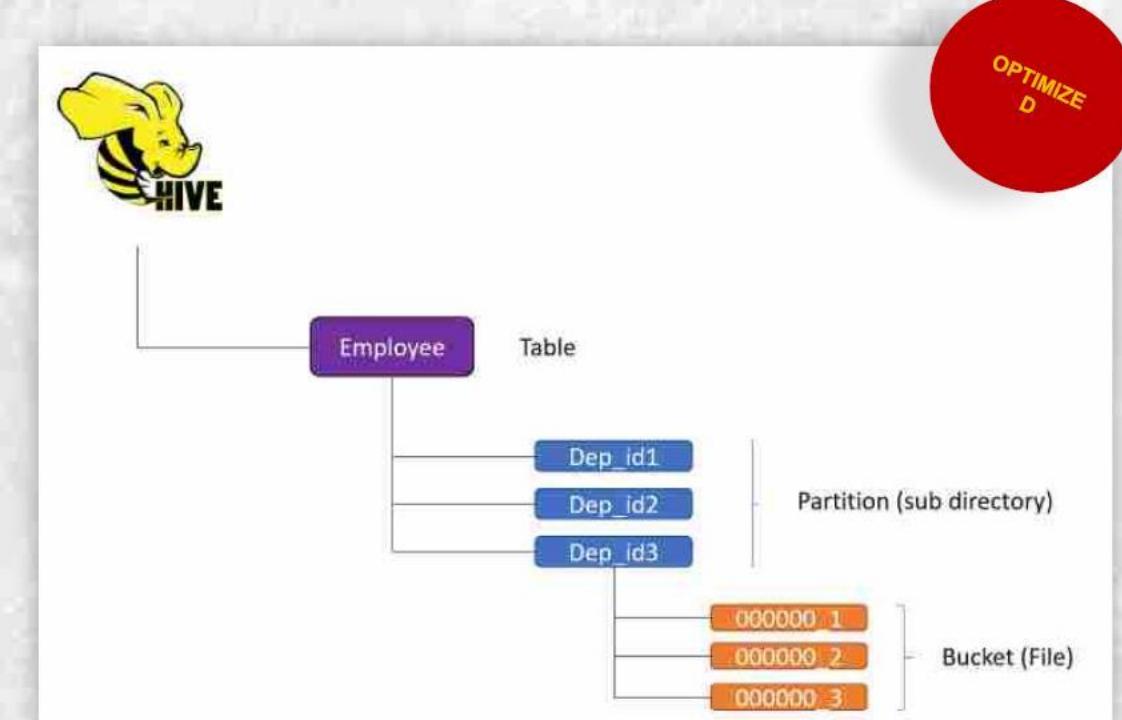


Batch Processing using HIVE

Create partitioning on Hive table

```
hive> create table CCPD(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 DECIMAL(10, 2),PAY_AMT3 DECIMAL(10, 2),PAY_AMT4 DECIMAL(10, 2),PAY_AMT5 DECIMAL(10, 2),PAY_AMT6 DECIMAL(10, 2)) partitioned by (DEFAULTED INT) row format delimited fields terminated by '\t' stored as textfile
```

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



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)
```

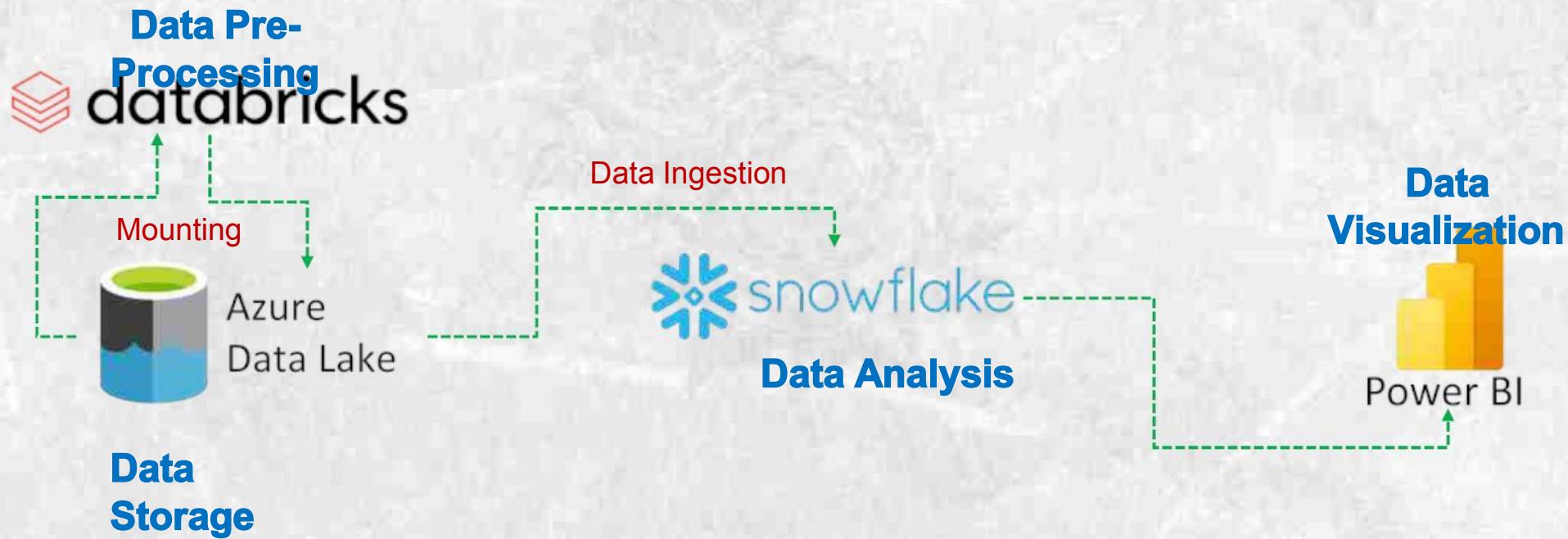
2. average billed Amount & average pay amount

```
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 |
+---+---+---+---+---+---+---+
| 502 | 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



Risk Analysis using

ADLS MOUNTING TO DATABRICKS :

Connecting ADLS with Databricks

```
Generate SAS Token from ADLS ACCOUNT and Paste it here

SAS_Token = "sv=2022-11-02&ss=b&tkt=sk0spwvdlacupy&se=2023-08-09T16:57:17Z&t=

Connect with ADLS account via mounting to load the data

dbutils.fs.mount(
  source = "wasbs://cont1.zoadlstorage.blob.core.windows.net",
  mount_point = "/mnt/mounted_SAS",
  extra_configs = {
    "fs.azure.sas.cont1.zoadlstorage.blob.core.windows.net":SAS_Token
  }
)
```

Data Pre-processing



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

```
Merge two part files in 1 using repartition

df3=ccard_df.repartition(1)

Save the processed data in ADLS

df3.write.format("csv").mode("overwrite").save("/mnt/mounted_SAS/CCdefaulter_csv",header=True)
```

```
dbutils.fs.unmount('/mnt/mounted_SAS')

/mnt/mounted_SAS has been unmounted.
Out[14]: True
```



Risk Analysis using

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_SAS/credit-card-default-1000.csv')
```

Separating Header

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

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

Convert to DataFrame

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

Removing header and unwanted inverted 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

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

```
C:\Users\ Miles>snowsql -a gfsysnv-ep39977 -u KANIKAAGG
Password:
* SnowSQL * v1.2.28
Type SQL statements or !help
KANIKAAGG#COMPUTE_WH@(no database).(no schema)>use ZOODATABASE;
+-----+
| status
+-----+
| Statement executed successfully.
+-----+
1 Row(s) produced. Time Elapsed: 18.938s
KANIKAAGG#COMPUTE_WH@ZOODATABASE.PUBLIC>use ZOODATABASE.ZOOSCHEMA;
+-----+
| status
+-----+
| Statement executed successfully.
+-----+
1 Row(s) produced. Time Elapsed: 0.624s
```

Create File Format

```
KANIKAAGG#COMPUTE_WH@ZOODATABASE.ZOOSCHEMA>CREATE OR REPLACE File Format zoo_csv_format
  type=csv
  field_delimiter=','
  skip_header=1
  null_if=('NULL','null')
  empty_field_as_null=true;
+-----+
| status
+-----+
| File format ZOO_CSV_FORMAT successfully created.
+-----+
1 Row(s) produced. Time Elapsed: 0.270s
```

Risk Analysis using snowflake

Create External Stage

```
KANIKAAGG#COMPUTE_WH@ZOOADATA#>CREATE OR REPLACE STAGE zoo_azure_stage
    URL = 'azure://zooadlsaccount.blob.core.windows.net/cont1/CCdefaulter_csv/part-00000-tid-739541815169323954-5a2e6ce4-b86e-4cc6-ae2d-082337fa0c86-11-1-c098.csv'
    CREDENTIALS = (AZURE_SAS_TOKEN = '?sv=2022-11-02&ss=t&st=2023-08-15T22:43:44Z&se=2023-08-15T14:48:59Z&sr=https&sig=FIXEX2F6ylg2s96zNys2CwvQHm8dvrM49Eh77E1Q1AVg%20')
    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

```
KANIKAAGG#COMPUTE_WH@ZOOADATA#>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-082337fa0c86-11-1-c098.csv | 1000 | 1000 | 1 | 0 | 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.

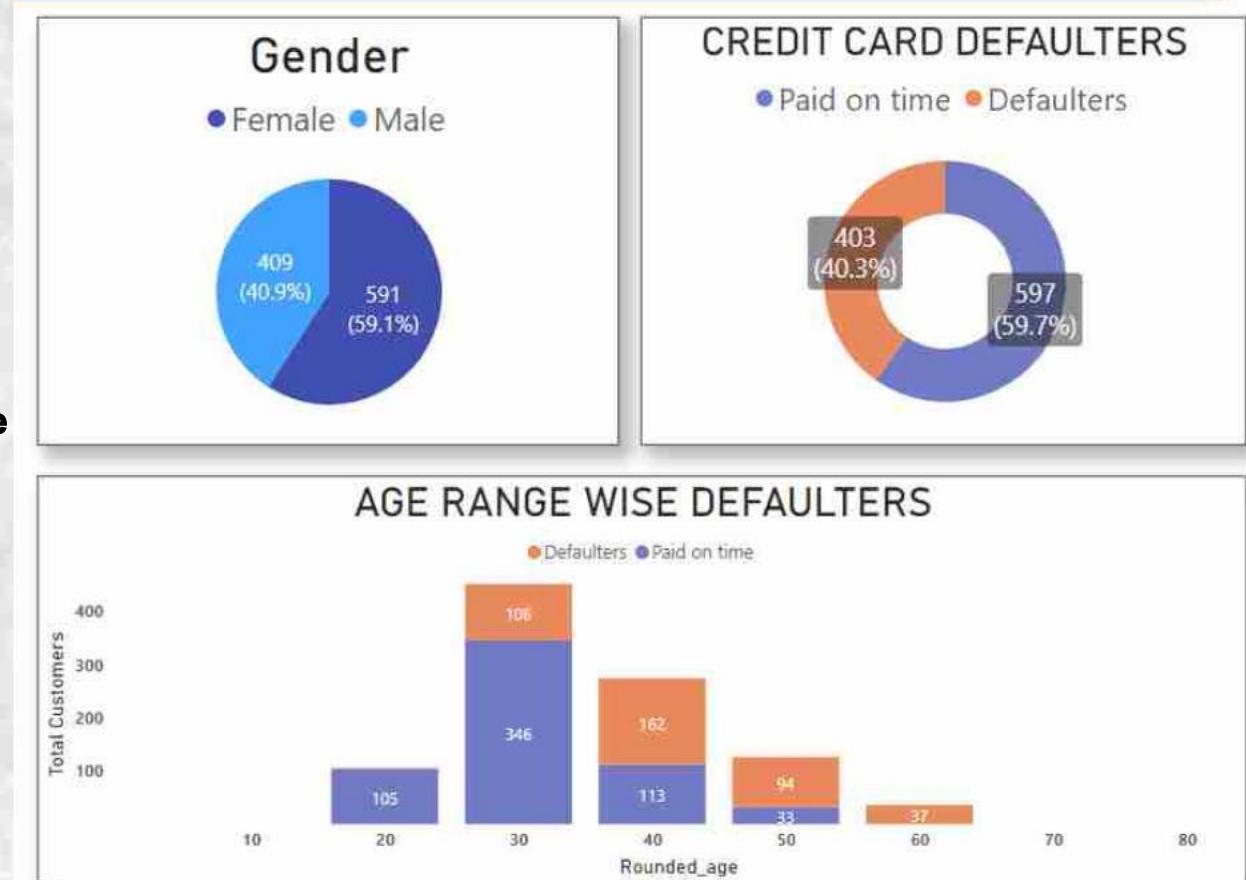
Create Gender table and insert values

```
Create table gender (
```

```
sex_code int,
```

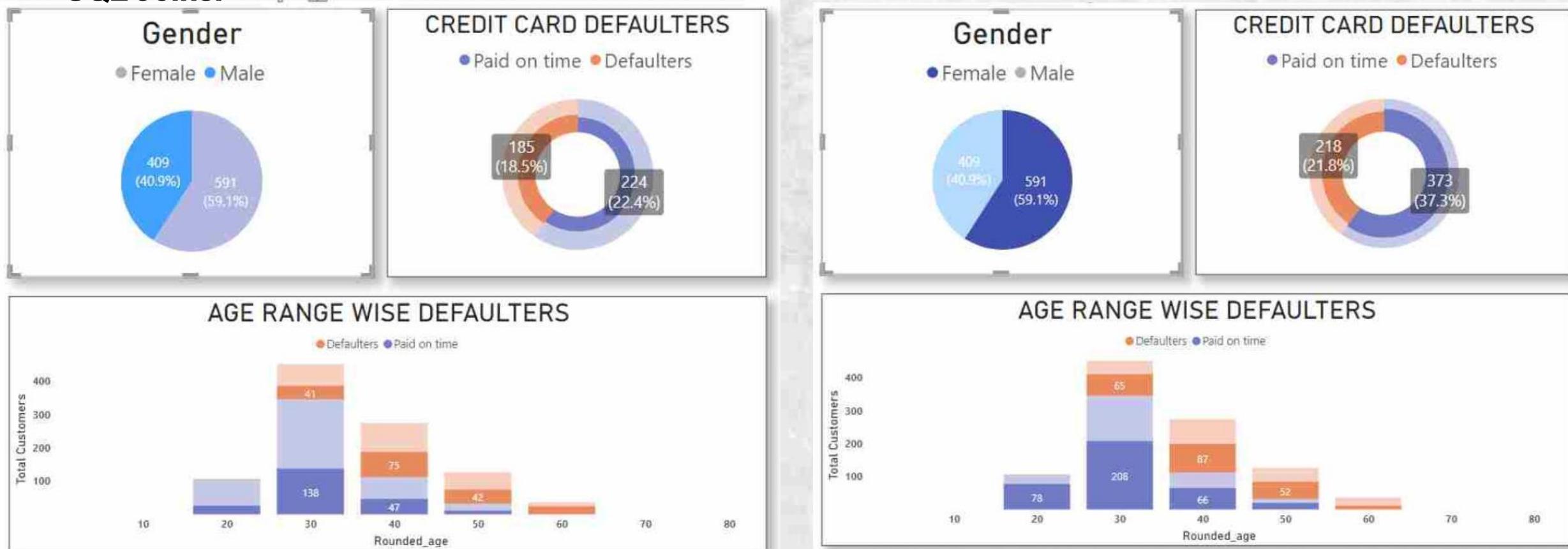
```
sex_name varchar(50) );
```

```
select *, FLOOR((Age + 5) / 10) * 10 AS RoundedAge
from cc_default c
inner join gender g on c.sex=g.sex_code;
```



Risk Analysis using

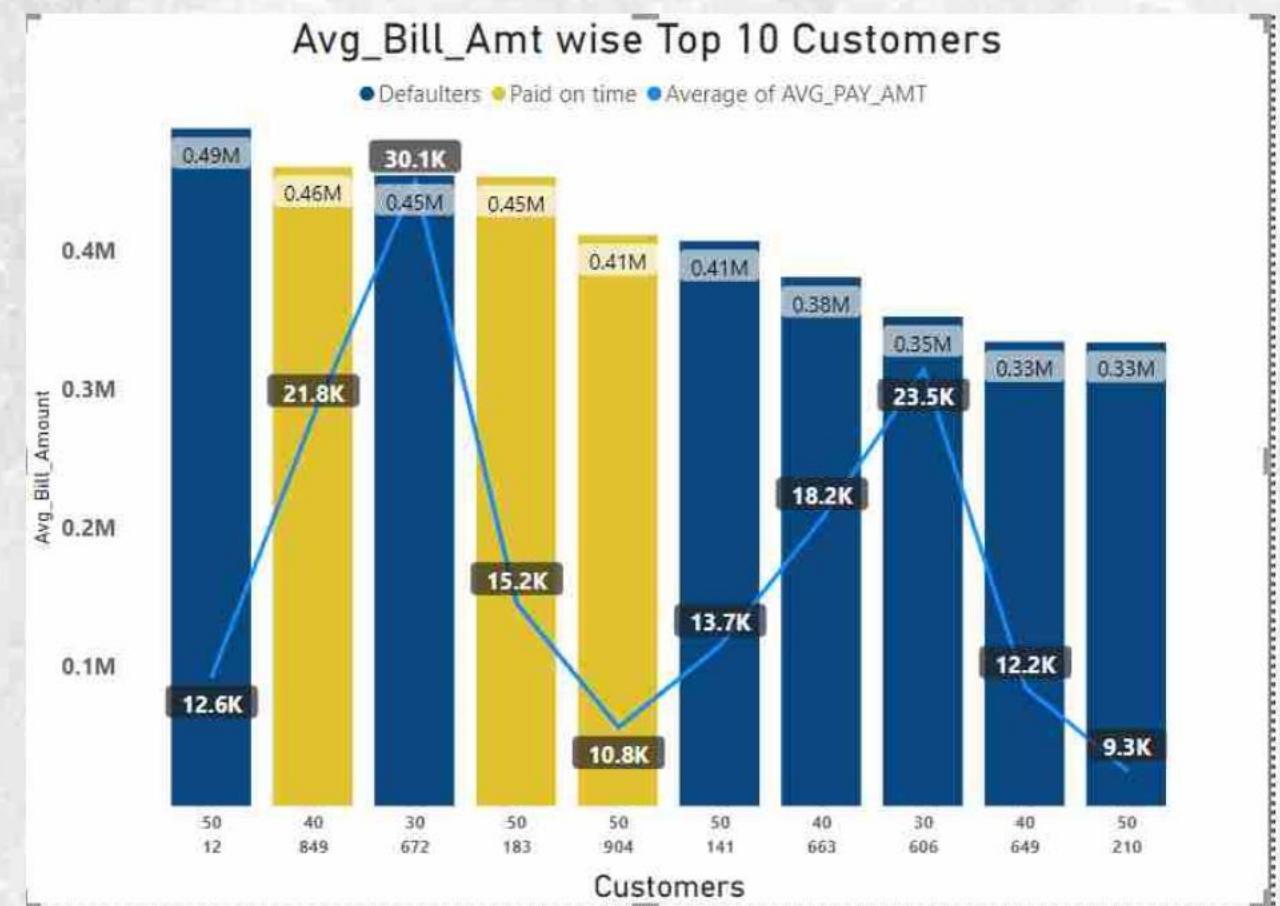
Data Processing/ Analysis : Rounding of age to range of 10s and add SEXNAME to the data using SQL Joins.



Risk Analysis using snowflake

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+PA
Y_AMT5+PAY_AMT6)/6,2) AS AVG_PAY_AMT
FROM CC_DEFAULT;
```

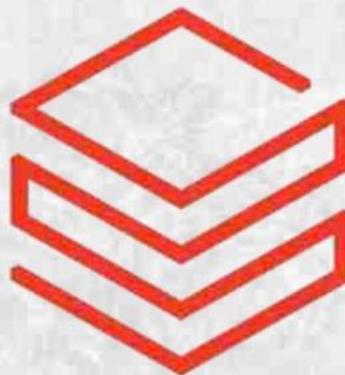


Risk Analysis using snowflake

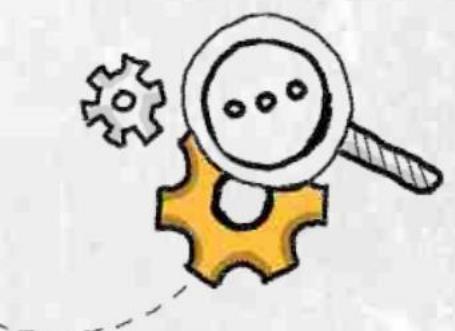
Data Processing/ Analysis : Find average pay duration. Make sure numbers are rounded and negative values are eliminated.

```
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;
```

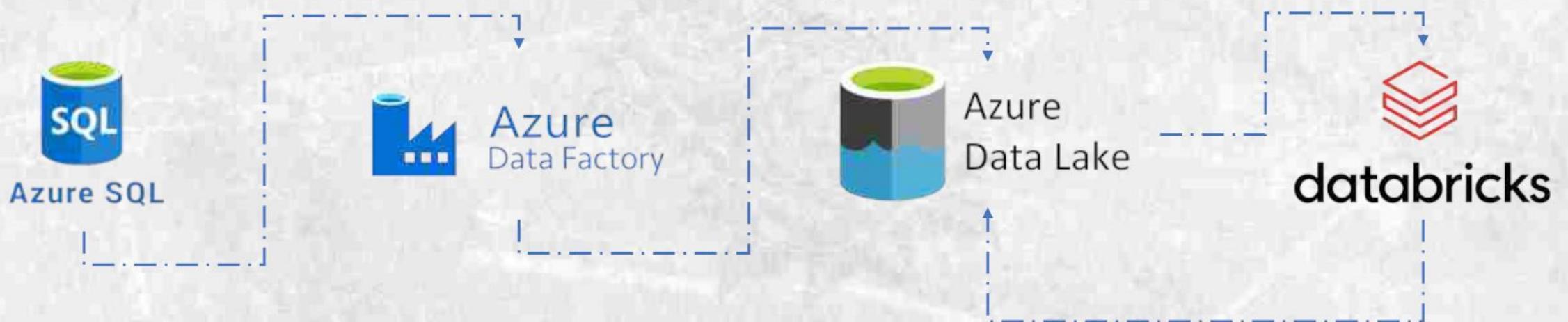




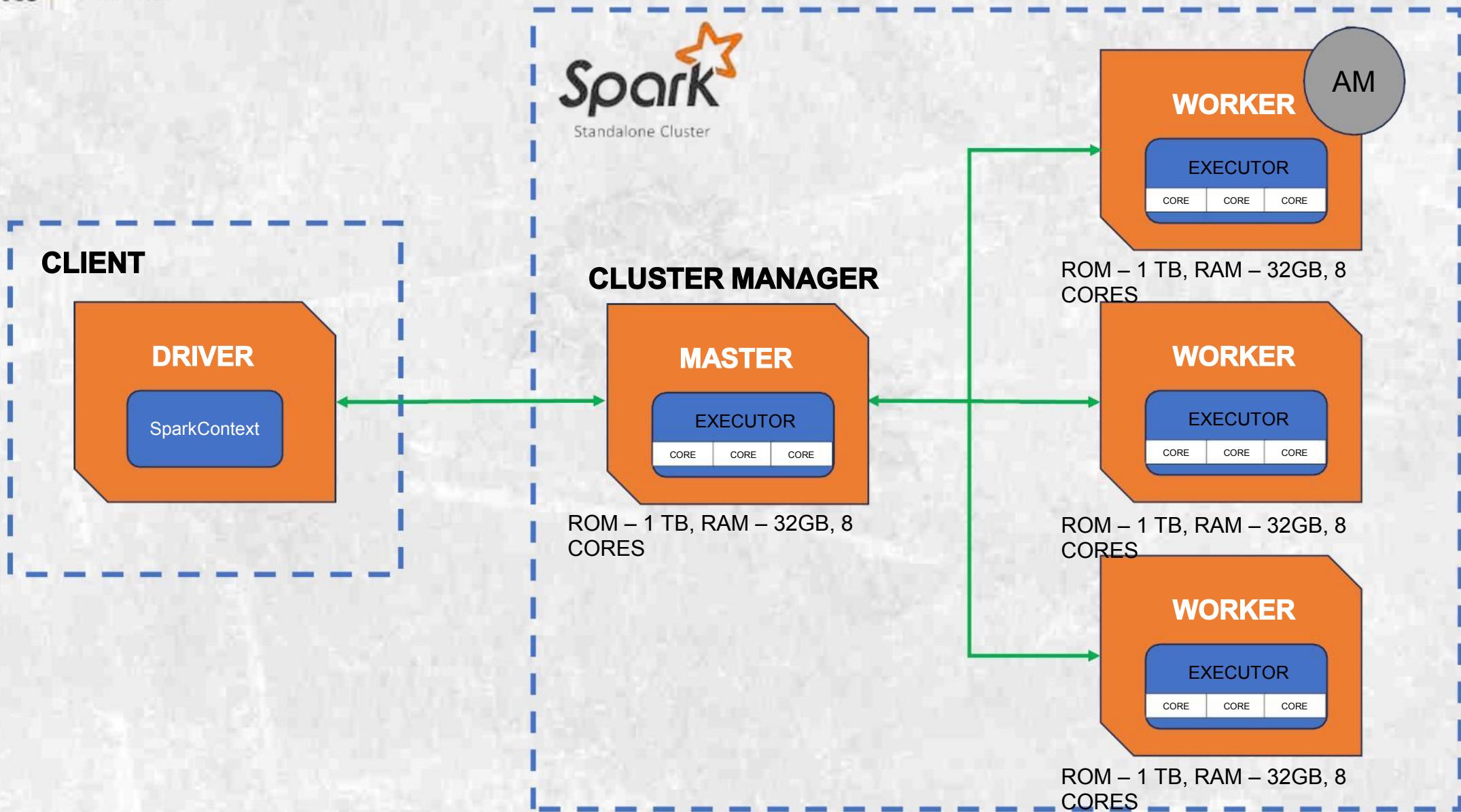
databricks



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.

The screenshot shows the Microsoft Azure Data Factory interface. The left sidebar lists 'Factory Resources' including 'Pipelines', 'Datasets', and 'Data flows'. The 'Pipelines' section shows three pipelines: 'SQL_To_ADLs' (selected), 'Change Data Capture (preview)', and 'ADLS_Output_Dataset'. The main workspace displays the 'Copy data' activity for the 'SQL_To_ADLs' pipeline. The 'Source' tab is selected, showing the 'Source dataset' dropdown set to 'SQL_Input_Dataset'. Other settings visible include 'Query timeout (minutes)' at 120 and 'Isolation level' at 'None'. The top navigation bar includes tabs for 'Data Factory', 'Validate all', and 'Published'.

Batch Processing using RDD

Mounting ADLS

Cmd 3

Python ► ▶ ▾ ✎

```
1 SAS_Token = "?sv=2022-11-02&ss=bfqt&srt=co&sp=rwdlacupyx&se=2023-08-28T12:48:32Z&st=2023-08-09T04:48:32Z&spr=https&sig=pZN3UpPsvxFn8Fz276y6vAdSu2x8G01t2B7KxAj17Ifjc%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

```
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

```
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_58@live.com at 17/11/2023, 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 ► ▾ ▷ - X

```
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]),
10   Float(x[7])))
11
12  # calculating spent by each card type
13  monthly_grouped = card_month_amount_rdd.reduceByKey(lambda x,y:x+y)
14
15  for card_type, max_month in monthly_grouped.collect():
16      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 : 110146294.0
```

```
Command took 0.89 seconds -- by theashpak_5@live.com at 9/8/2023, 5:36:34 pm on Ashpak Sheikh's Cluster
```

Optimization using RDD



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

```
Code: 11
1 # Mapping the data to (city, amount) pairs.
2 rdd_split_cached = rdd_split.cache()
3 city_amount_rdd = rdd_split_cached.map(lambda row: (row[1], float(row[0])))
4 # Total spends per city
5 city_total_spends_rdd = city_amount_rdd.reduceByKey(lambda a, b: a + b)
6
7 # Total spends across all cities
8 total_spends = city_total_spends_rdd.values().sum()
9
10 # Total spends in descending order
11 sorted_cities = city_total_spends_rdd.sortBy(lambda x: x[1], ascending=False)
12
13 # Top 5 cities
14 top_cities = sorted_cities.take(5)
15
16 # Percentage contribution of each city's spends
17 city_spends_percentage = [(city, (amount / total_spends) * 100, amount) for city, amount in top_cities]
18 # Printing the results
19 for city,percentage,amount in city_spends_percentage:
20     print("City: {}, Amount: {}, Spends: {:.2f}%,".format(city, amount,percentage))
21
```

```
* [0] Spark Job
City: Greater Mumbai, Amount: 576751476.0, Spends: 14.19%
City: Bengaluru, Amount: 572326730.0, Spends: 14.03%
City: Ahmedabad, Amount: 567794310.0, Spends: 13.93%
City: Delhi, Amount: 556929212.0, Spends: 13.47%
City: Kolkata, Amount: 115466943.0, Spends: 2.82%
```

Command took 0.44 seconds -- by theashish_80@live.com at 17/12/2013, 11:43:06 am on Adonis: Shalin's Cluster

► (4) Spark Jobs

```
City: Greater Mumbai, Amount: 576751476.0
City: Bengaluru, Amount: 572326730.0
City: Ahmedabad, Amount: 567794310.0
City: Delhi, Amount: 556929212.0
City: Kolkata, Amount: 115466943.0
```

Command took 0.66 seconds -- by theashish_80@live.com at 17/12/2013, 11:43:16 am on Adonis: Shalin's Cluster

► (2) Spark Jobs

```
Card Type : ('Platinum', 'Jan'), Spent : 112784373.0
Card Type : ('Signature', 'Jan'), Spent : 98819381.0
Card Type : ('Silver', 'Jan'), Spent : 109359590.0
Card Type : ('Gold', 'Jan'), Spent : 118146284.0
```

Command took 0.83 seconds -- by theashish_80@live.com at 17/12/2013, 11:49:11 am on Adonis: Shalin's Cluster

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

```
Code: 12
1 #@('Palhi', 'India', '29-Dec-14', 'Gold', 'Hills', 'F', '82425'
2 # splitting data monthwise
3 rdd_split_cached = rdd_split.cache()
4 card_month_amount_rdd = rdd_split_cached.map(lambda x: [x[3].split('-')[1], float(x[7])])
5
6 # calculating Monthly Spend
7 monthly_sum = card_month_amount_rdd.reduceByKey(lambda x,y: x+y).map(lambda x: x[1])
8
9 #Getting all the data from month having maximum spent
10 card_month_amount_rdd = rdd_split_cached.filter(lambda x: x[3].split('-')[1]==monthly_sum[0]).map(lambda x: ((x[4], x[3].split('-')[1]), float(x[7])))
11
12 # calculating spent by each card type
13 monthly_grouped = card_month_amount_rdd.reduceByKey(lambda x,y:x+y)
14
15 for card_type, max_month in monthly_grouped.collect():
16     print("Card Type : ({}, Jan), Spent : {}".format(card_type, max_month))
```

* [1] Spark Job

```
Card Type : ('Platinum', 'Jan'), Spent : 112784373.0
Card Type : ('Signature', 'Jan'), Spent : 98819381.0
Card Type : ('Silver', 'Jan'), Spent : 109359590.0
Card Type : ('Gold', 'Jan'), Spent : 118146284.0
```

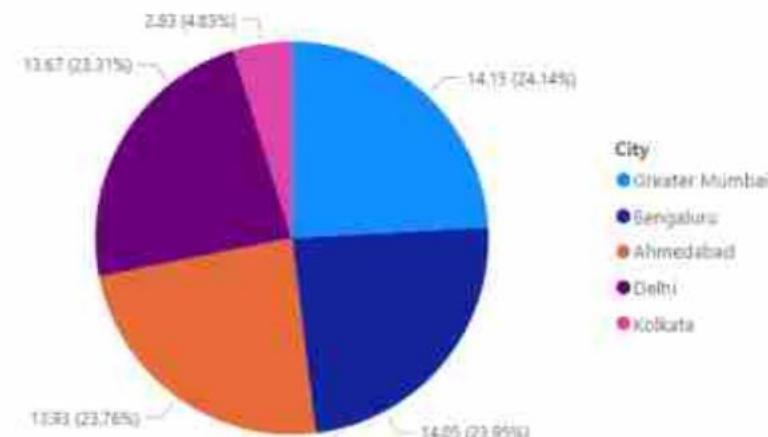
Command took 0.83 seconds -- by theashish_80@live.com at 17/12/2013, 11:49:11 am on Adonis: Shalin's Cluster

Batch Processing using RDD

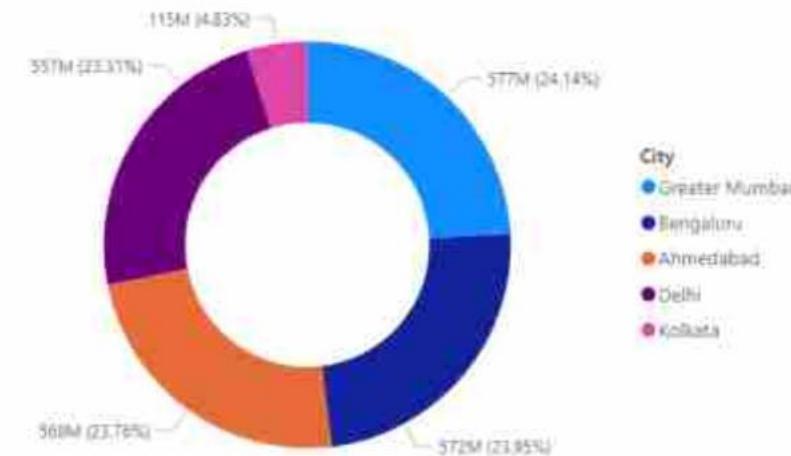
VISUALIZATION

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

Percentage Contribution by City

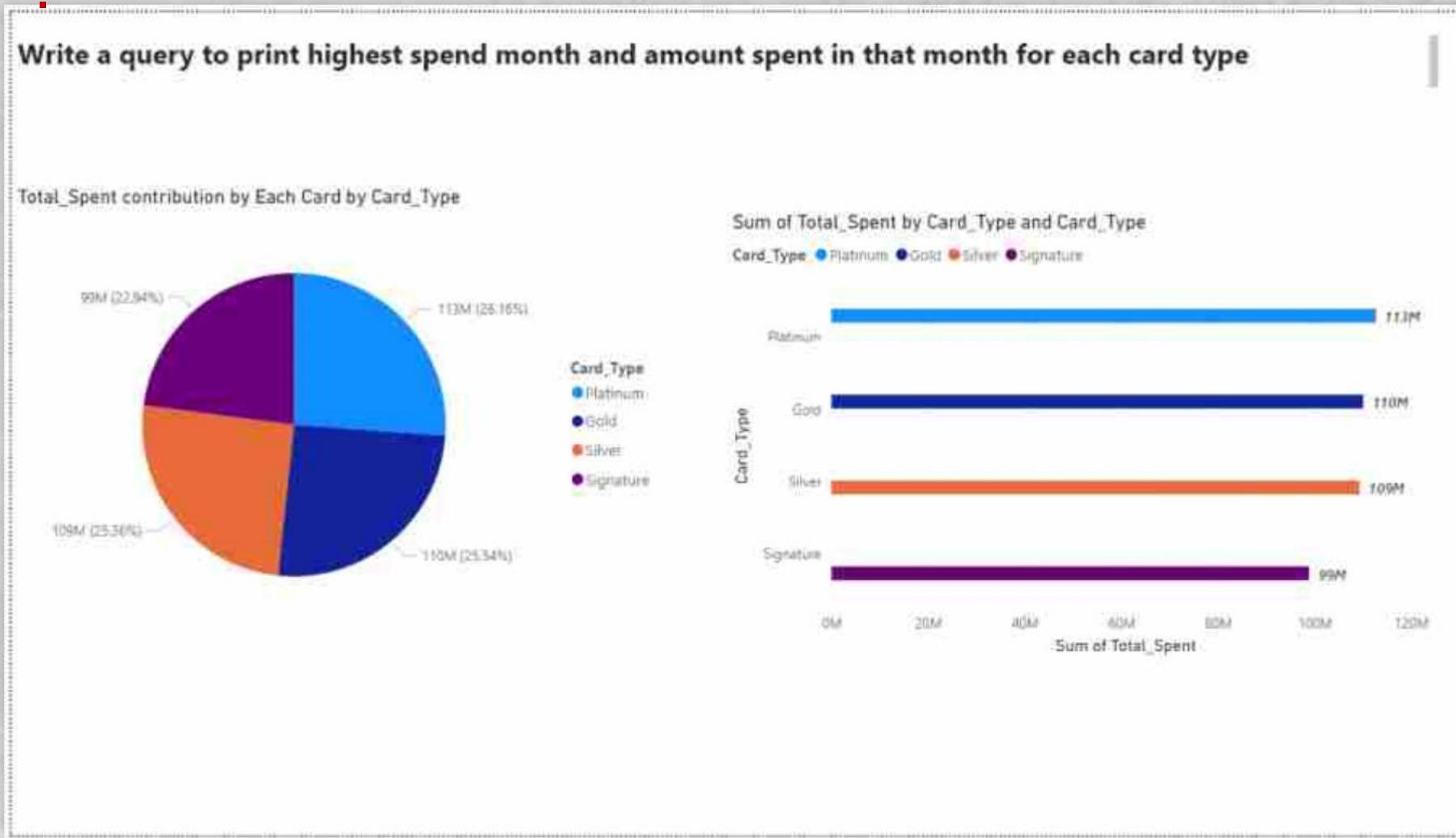


Total Spent By City by City



Batch Processing using RDD

VISUALIZATION



Batch Processing using RDD

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}%")

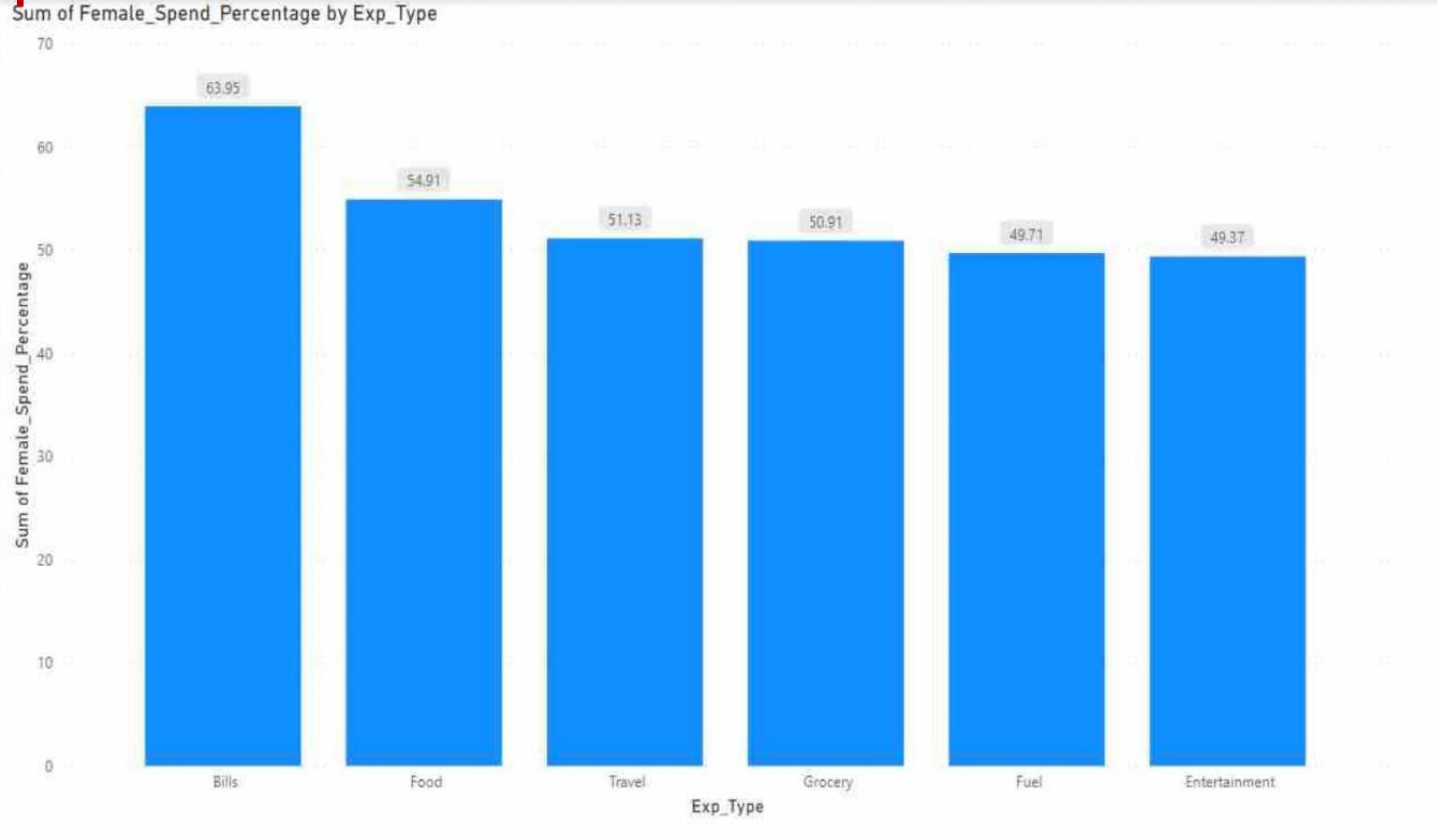
▶ (2) Spark Jobs

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



Batch Processing using RDD

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

```
Cmd 12

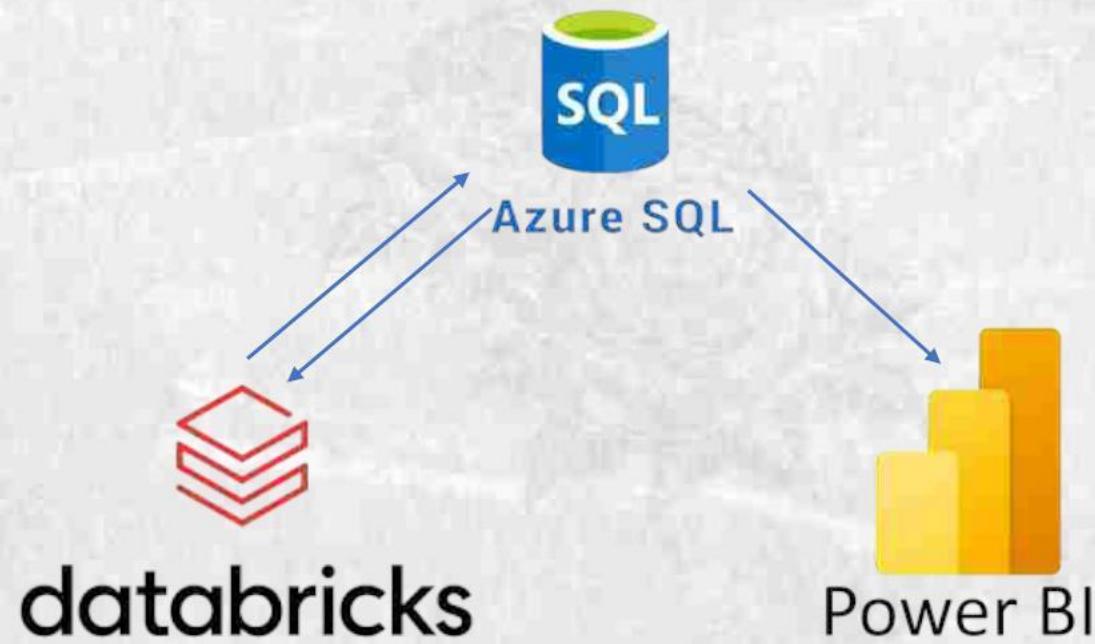
1 # Filter data for January 2014
2 jan_2014_data_rdd = rddy.filter(lambda x: x[3].split('-',1)[1] ==('Jan-14'))
3 Dec_2013_data_rdd = rddy.filter(lambda x: x[3].split('-',1)[1] ==('Dec-13'))
4 # Map the RDD to ((card_type, expense_type), amount)
5 card_expense_amount_rdd_jan = jan_2014_data_rdd.map(lambda 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 niharikajs0321@gmail.com at 8/15/2023, 11:54:49 PM on Niharika J S'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}}"

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

Command took 0.10 seconds -- by niharika.joshi@gmail.com at 8/15/2023, 9:41:22 PM on Niharika J S's Cluster
```

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

```
Cmd 3

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()

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

Command took 0.28 seconds -- by niharika.joshi@gmail.com at 8/15/2023, 9:51:55 PM on Niharika J S's Cluster
```

Cleaning Data and Creating Temporary Table

Removing country name from City column

```
Cell 13

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 mihirika.jaiswal@gmail.com at 8/15/2023, 7:05:55 PM on Mihirika's Cluster
```

Creating temporary table for data frame

```
Cell 14

1 credit_card.createOrReplaceTempView('credit_tbl')

Command took 0.16 seconds -- by mihirika.jaiswal@gmail.com at 8/15/2023, 7:06:03 PM on Mihirika'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 39

```
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 |      )
8 |      SELECT DISTINCT City,
9 |              FIRST_VALUE(Exp_Type) OVER(PARTITION BY City) AS Highest_Expense_Type,
10 |              LAST_VALUE(Exp_Type) OVER(PARTITION BY City) AS Lowest_Expense_Type
11 |              FROM cte
12 |          """)
```

▶ sc5_out: pyspark.sql.dataframe.DataFrame = [City: string, Highest Expense Type: string ... 1 more field]

Command took 1.17 seconds -- by niharikais0321@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|
| Adon1|        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|
| Alot|        Fuel| Entertainment|
| Alappuzha|        Food| Entertainment|
| Aligarh|        Bills| Entertainment|
| Alipurduar|        Food| Entertainment|
| Alirajpur| Entertainment| Entertainment|
```

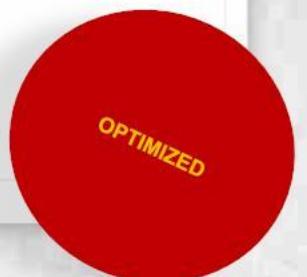
Command took 4.18 seconds -- by niharika.jaiswal@gmail.com at 8/15/2023, 7:06:08 PM on Hadoop 3.5's Cluster

Catalyst Optimizer

```
1 sc5_out.explain(mode="extended")

== Parsed Logical Plan ==
CTE [cte]
: +- 'SubqueryAlias cte
:   +- 'Sort [1 ASC NULLS FIRST, 3 DESC NULLS LAST], true
:     +- 'Aggregate ['City, 'EXP_TYPE], ['City, 'Exp_Type, unresolvedalias('SUM('Amount), None)]
:       +- 'UnresolvedRelation [credit_tbl], [], false
+- 'Distinct
  +- 'Project ['City, 'FIRST_VALUE('Exp_Type) windowspecdefinition('City, unspecifiedframe$()) AS Highest_Expense_Type#119, 'LAST_VALUE('Exp_Type) windowspecdefinition('City, unspecifiedframe$()) AS Lowest_Expense_Type#120]
    +- 'UnresolvedRelation [cte], [], false

== Analyzed Logical Plan ==
City: string, Highest_Expense_Type: string, Lowest_Expense_Type: string
WithCTE
:- CTERelationDef 0, false
: +- SubqueryAlias cte
:   +- Sort [City#111 ASC NULLS FIRST, sum(Amount)#122L DESC NULLS LAST], true
:     +- Aggregate [City#111, EXP_TYPE#63], [City#111, Exp_Type#63, sum(Amount#65) AS sum(Amount)#122L]
:       +- SubqueryAlias credit_tbl
:         +- View ('credit_tbl', [index#59,City#111,Date#61,Card_Type#62,Exp_Type#63,Gender#64,Amount#65])
:           +- Project [index#59, <lambda>(City#60)#110 AS City#111, Date#61, Card_Type#62, Exp_Type#63, Gender#64, Amount#65]
Command took 0.10 seconds -- by niharikajs0321@gmail.com at 8/16/2023, 11:36:09 PM on Niharika J S's Cluster
```



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

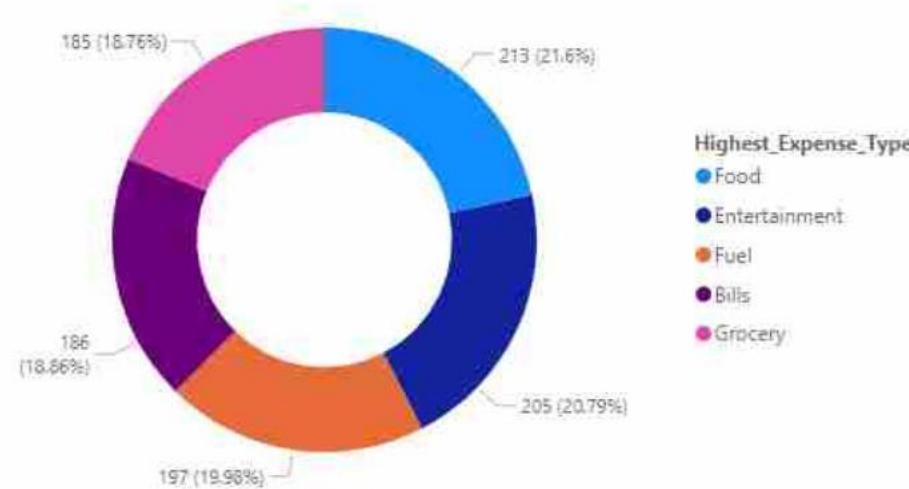
Command took 4.39 seconds -- by niharikajs6321@gmail.com at 8/10/2023, 11:39:28 AM on Niharika J S's Cluster



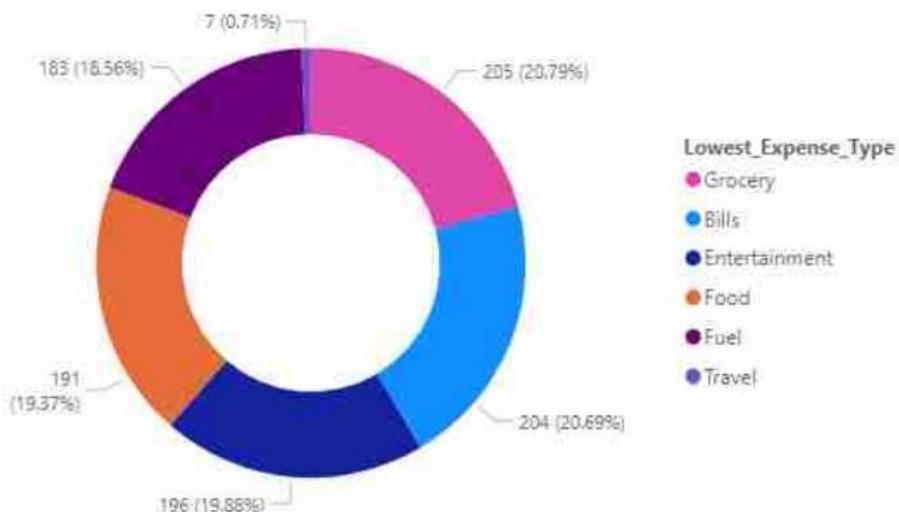
Visualization

Write a query to print 3 columns: city, highest expense type, lowest expense type

Count of City by Highest_Expense_Type



Count of City by Lowest_Expense_Type



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

Code 47

```
1 sc9_out = spark.sql(""" WITH cte1 AS
2     (SELECT City, Date, ROW_NUMBER() OVER(PARTITION BY City ORDER BY Date) AS RN
3      FROM credit_tbl
4    ), cte2 AS
5     (SELECT *, LAST_VALUE(RN) OVER(PARTITION BY City) AS Low
6      FROM cte1
7     WHERE RN<=500
8    ), cte3 AS
9     (SELECT DISTINCT City, DATEDIFF(MAX(Date) OVER(PARTITION BY City), Date) AS Difference
10       FROM cte2 WHERE Low=500
11    )
12     SELECT City, MAX(Difference) Days_Took
13     FROM cte3
14     GROUP BY City
15     ORDER BY 2 LIMIT 1
16   """)

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

Command took 8.11 seconds -- by rithikajoshi32@gmail.com at 8/18/2023, 11:48:35 AM on Rithika-J's Cluster

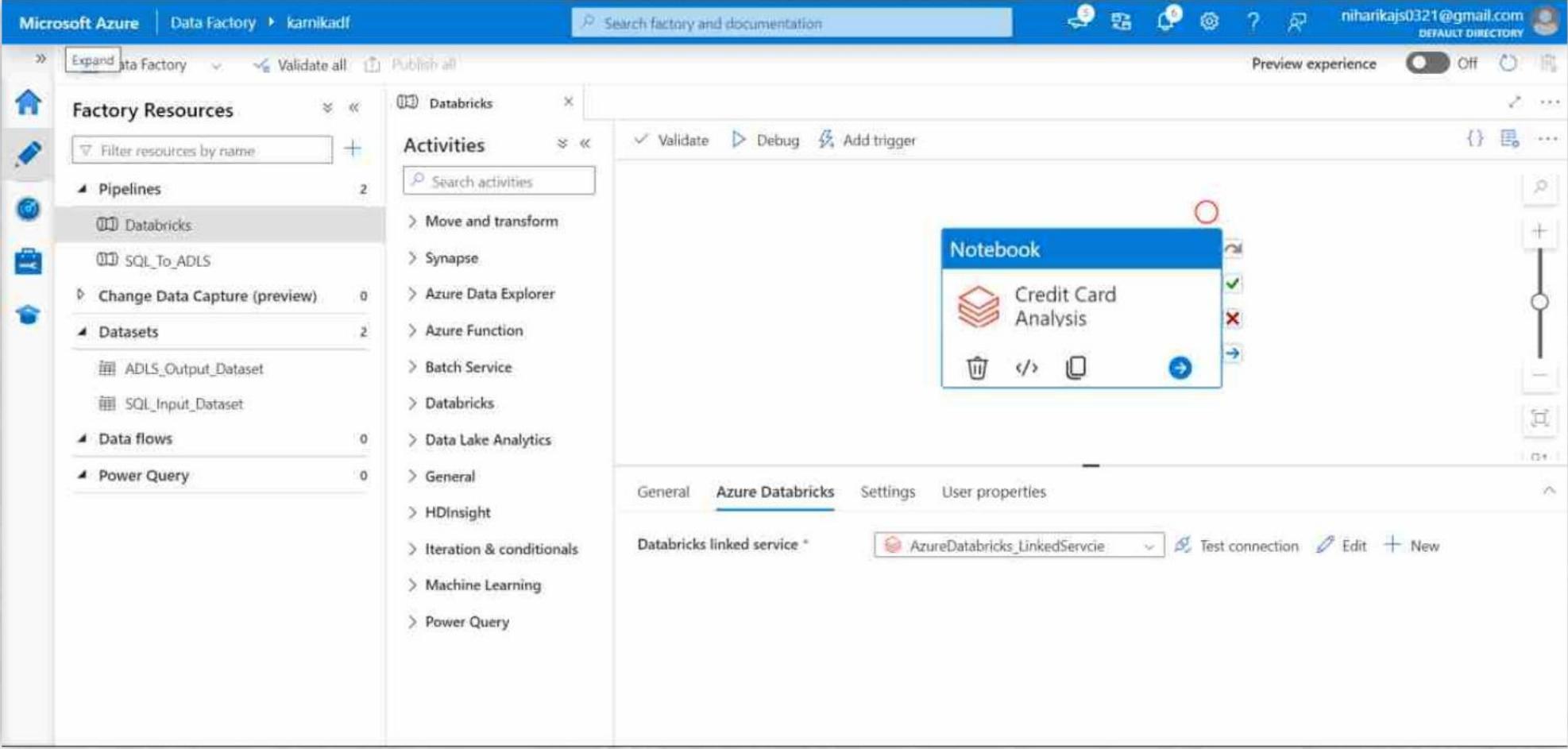
```
1 sc9_out.show()
```

▶ (2) Spark Jobs

City(Days_Took)
[Bengaluru] 81

Command took 1.17 seconds -- by rithikajoshi32@gmail.com at 8/18/2023, 11:48:36 AM on Rithika-J's Cluster

Scheduling Databricks Notebook in ADF for Periodical Analysis



The screenshot shows the Microsoft Azure Data Factory interface. On the left, the 'Factory Resources' sidebar lists various components: Pipelines, Datasets, Data flows, and Power Query. Under Pipelines, 'Databricks' is selected, and under it, 'SQL_To_ADLS' is listed. The main workspace displays a 'Databricks' activity configuration. The 'Activities' pane on the right lists options like Move and transform, Synapse, Azure Data Explorer, Azure Function, Batch Service, Databricks, Data Lake Analytics, General, HDInsight, Iteration & conditionals, Machine Learning, and Power Query. The 'Notebook' section shows a configuration for a 'Credit Card Analysis' notebook, linked to an 'AzureDatabricks_LinkedService'. The 'Azure Databricks' tab is selected. At the bottom, there are tabs for General, Azure Databricks, Settings, and User properties, along with fields for Databricks linked service (set to 'AzureDatabricks_LinkedService') and a 'Test connection' button.

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