

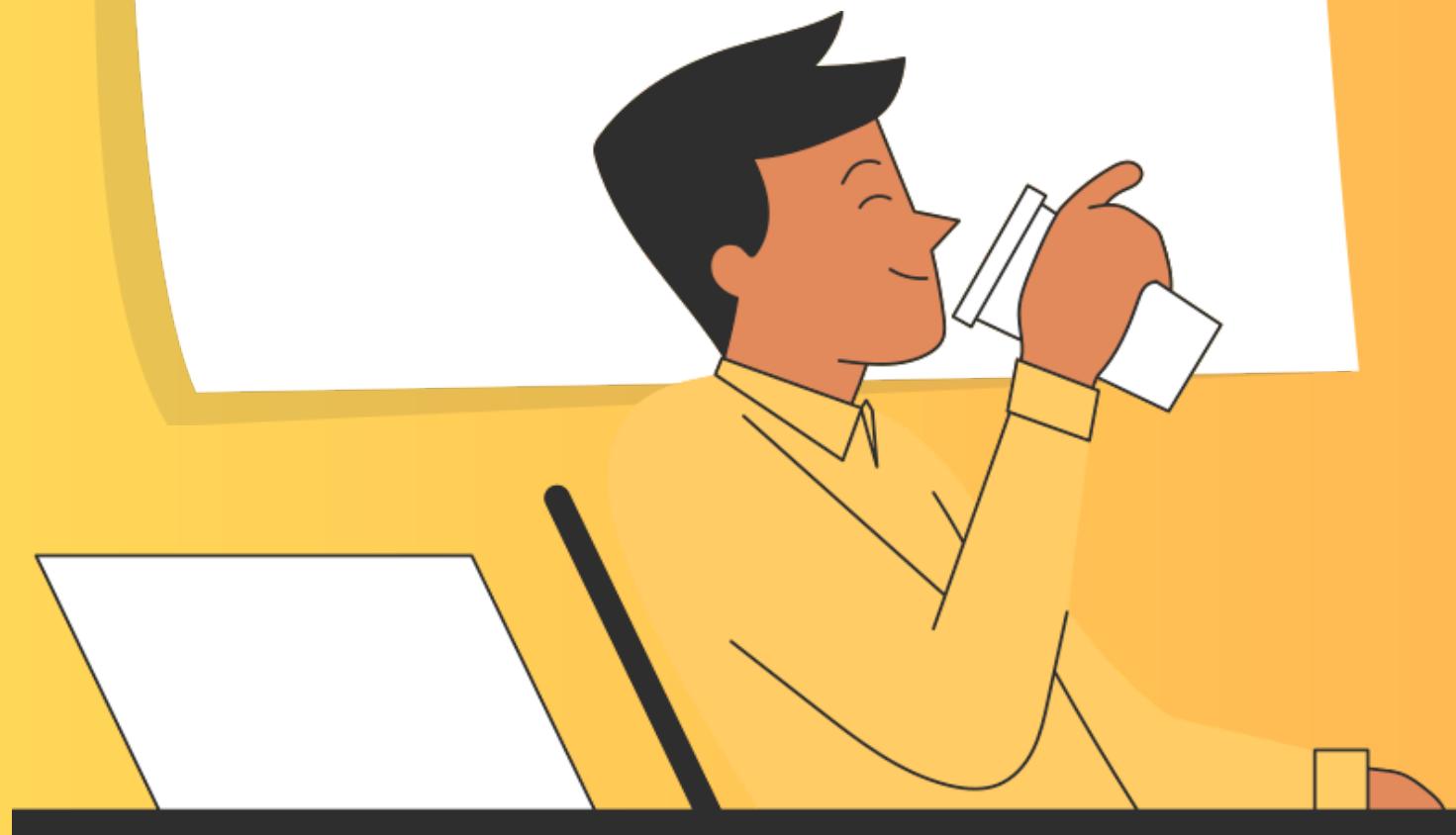


# CREDIT CARD TRANSACTION ANALYSIS

Team Members

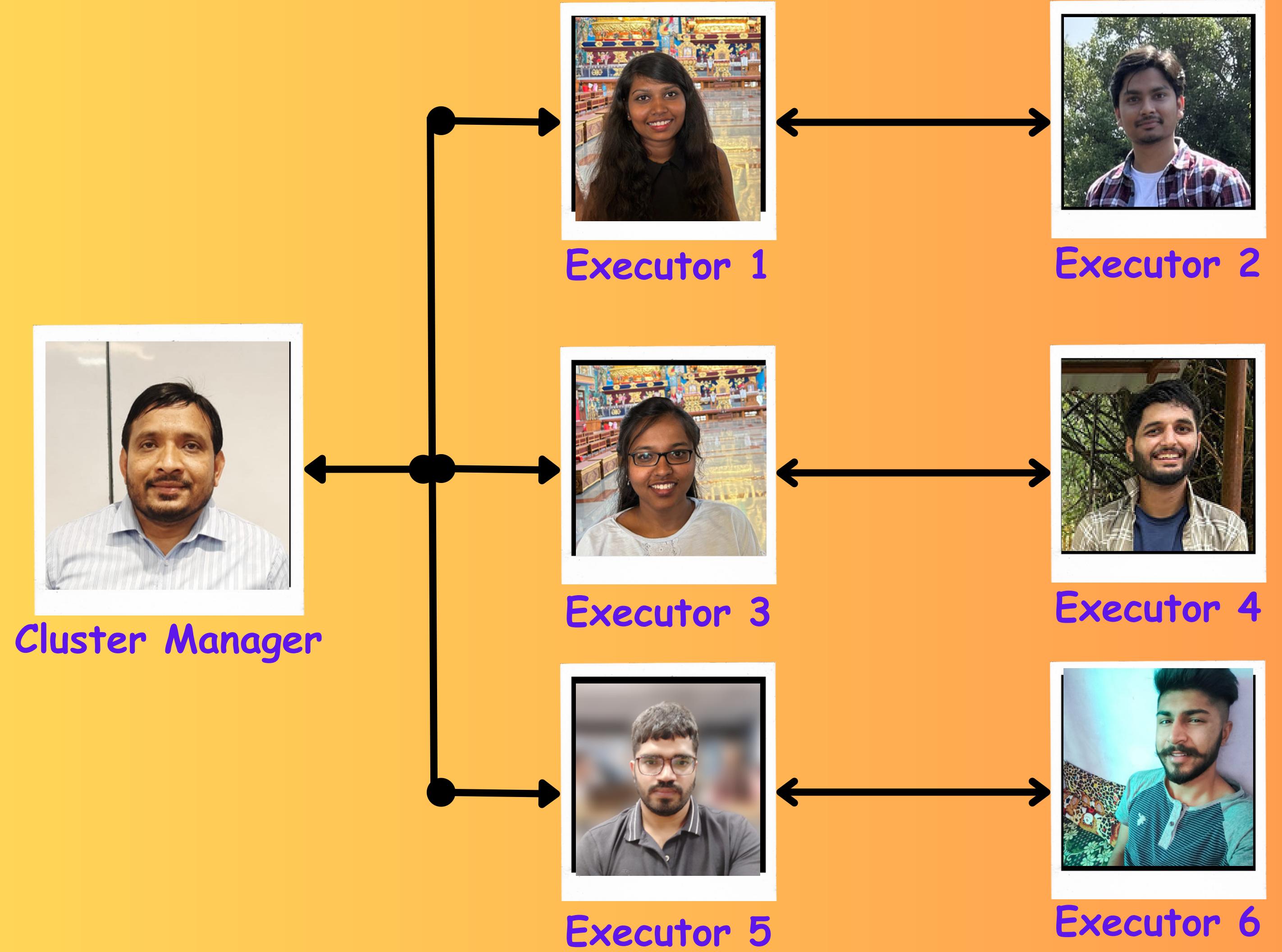
Abhishek Yadav  
Anurag Yadav  
Himanshu Dhiman  
Kaviya R  
Nikita Nannaware  
Shreyash Gupta

# CONTENTS

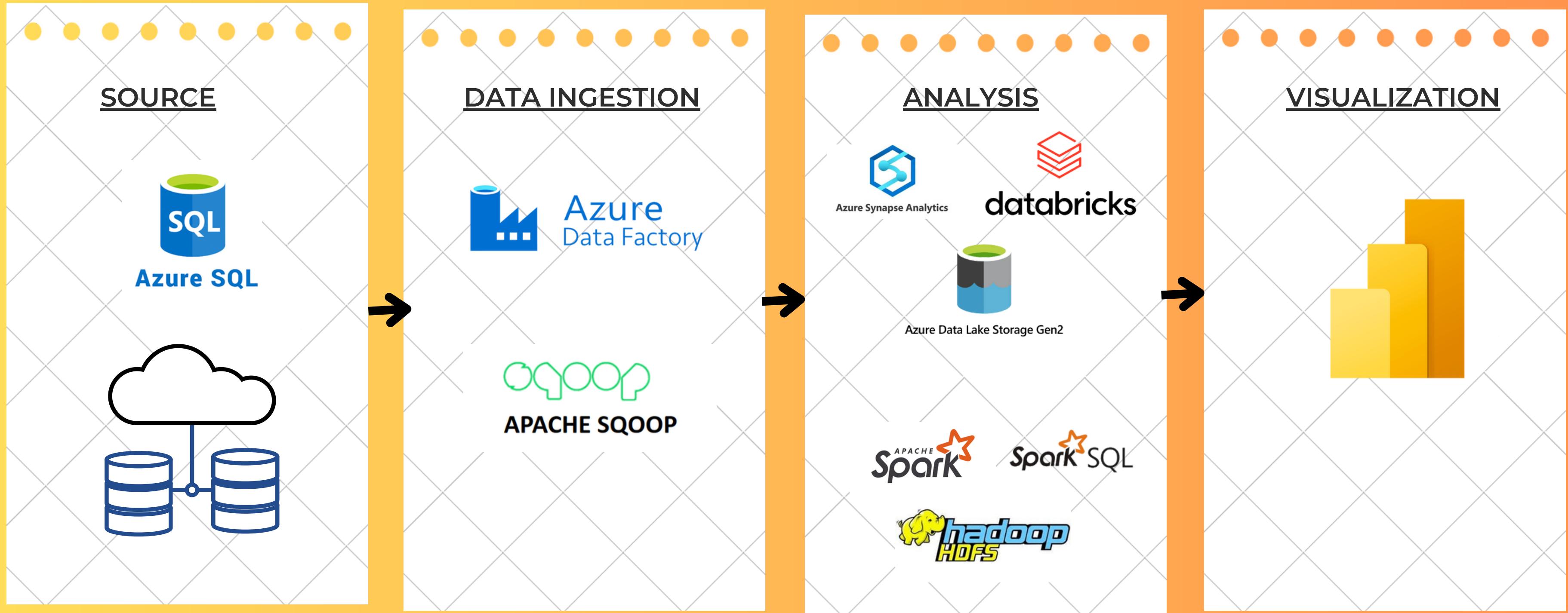


- 1** Architecture Diagram
- 2** Offline Analysis
- 3** Cloud Setup
- 4** Cloud Analysis
- 5** Visualizations

**TEAM  
CLUSTER**



# ARCHITECTURE DIAGRAM



# DATA DESCRIPTION

## CREDIT CARD DEFAULTED. CSV

- NO OF ROWS: 1002
- SIZE: 92.2KB

- CUSTID: UNIQUE CUSTOMER ID
- LIMIT\_BAL: MAXIMUM SPENDING LIMIT FOR THE CUSTOMER
- SEX: SEX OF THE CUSTOMER. 1 ( MALE) AND 2 (FEMALE)
- EDUCATION: EDUCATION LEVEL OF THE CUSTOMER. 1 (GRADUATE), 2 (UNIVERSITY), 3 (HIGH SCHOOL) AND 4 (OTHERS)
- MARRIAGE: MARITAL STATUS OF THE CUSTOMER. 1 (SINGLE), 2 (MARRIED) AND 3 ( OTHERS)
- AGE: AGE OF CUSTOMER
- PAY\_1 TO PAY\_6: PAYMENT STATUS FOR THE LAST 6 MONTHS
- BILL\_AMT1 TO BILL\_AMT6: THE BILLED AMOUNT FOR CREDIT CARD FOR EACH OF THE LAST 6 MONTHS.
- PAY\_AMT1 TO PAY\_AMT6: THE ACTUAL AMOUNT THE CUSTOMER PAID FOR EACH OF THE LAST 6 MONTHS
- DEFAULTED: WHETHER THE CUSTOMER DEFAULTED OR NOT ON THE 7TH MONTH. THE VALUES ARE 0 (DID NOT DEFAULT) AND 1 (DEFAULTED)

## CREDIT CARD. CSV

- NO OF ROWS: 284808
- SIZE: 144MB

- TIME: TIME DIFFERENCE BETWEEN TRANSACTION TIME AND FIRST ENTRY IN DATABASE
- V1 TO V28: RESULT OF PCA TRANSFORMATION OF CUSTOMER INFORMATION
- AMOUNT: AMOUNT PAID BY CUSTOMER
- CLASS: WHETHER CUSTOMER IS DEFaulTER OR NOT

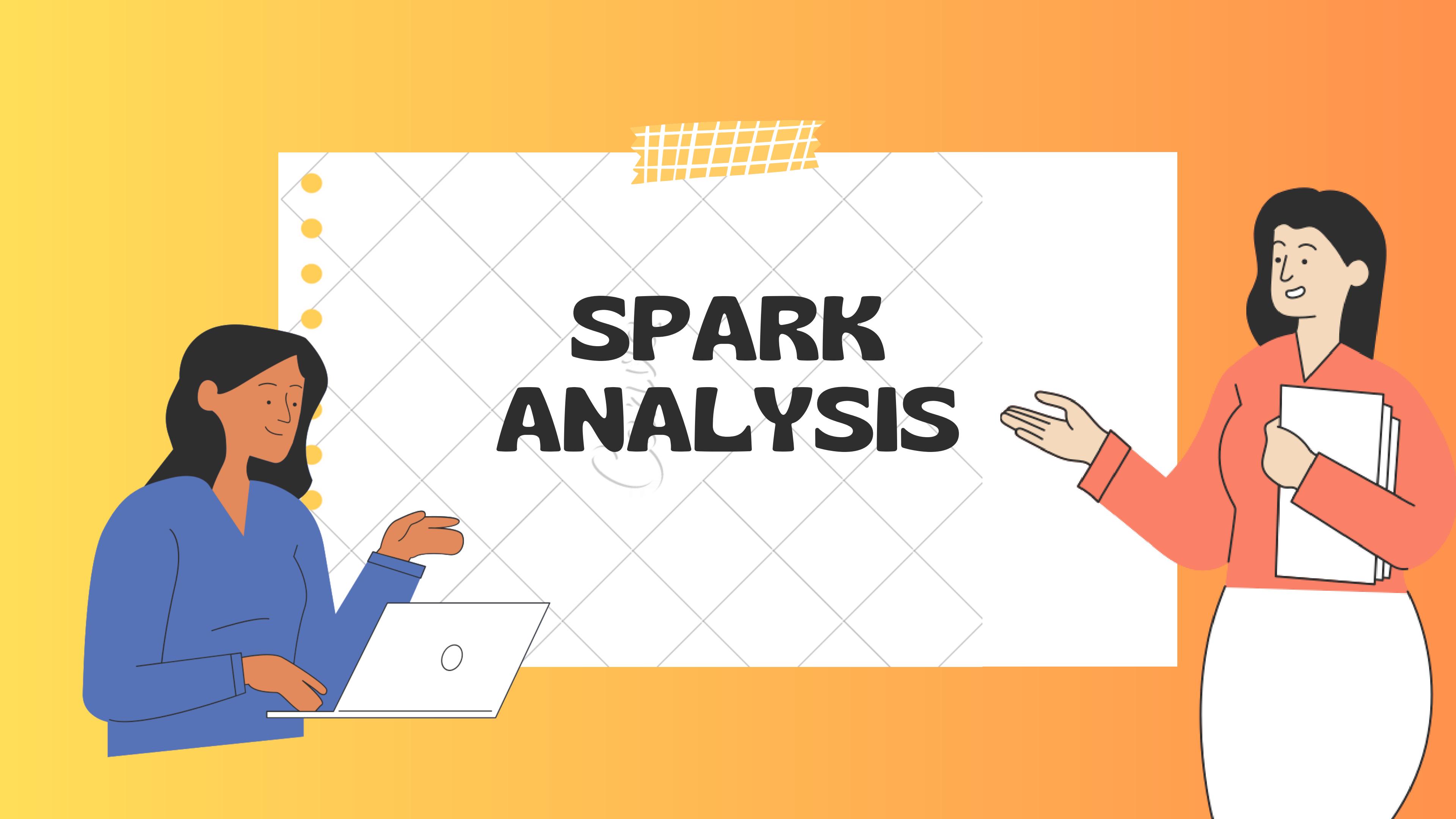


# DATA IMPORT USING SQOOP



```
SQOOP IMPORT
--CONNECT JDBC:SQLSERVER://ABHI23426.DATABASE.WINDOWS.NET:1433/CREDIT
--USERNAME=ABHI23426
--PASSWORD= *WHY_SO_SERIOUS_ABOUT_PASSWORDS*
--TABLE TRANSACTION_DATA
--TARGET-DIR /USER/CLOUDERA/CREDIT_CARD
```





# SPARK ANALYSIS



```
1 import findspark
2 findspark.init()
3 from pyspark import SparkContext
4 from pyspark.sql import SparkSession

1 spark = SparkSession.builder.appName('credit card').getOrCreate()

1 sc = spark.sparkContext

1 rdd = sc.textFile('creditcard_csv.csv')

1 rdd.take(2)

['Time,V1,V2,V3,V4,V5,V6,V7,V8,V9,V10,V11,V12,V13,V14,V15,V16,V17,V18,V19,V20,V21,V22,V23,V24,V25,V26,V27,V28,Amount,Class',
 "0,-1.3598071336738,-0.0727811733098497,2.53634673796914,1.37815522427443,-0.338320769942518,0.46238777762292,0.2395985540612
57,0.0986979012610507,0.363786969611213,0.0907941719789316,-0.55159533260813,-0.617800855762348,-0.991389847235408,-0.31116935
3699879,1.46817697209427,-0.470400525259478,0.207971241929242,0.0257905801985591,0.403992960255733,0.251412098239705,-0.0183067
77944153,0.277837575558899,-0.110473910188767,0.0669280749146731,0.128539358273528,-0.189114843888824,0.133558376740387,-0.0210
530534538215,149.62,'0'"]
```

Splitting data,converting time into proper datetime format and selecting only required columns

Importing required Modules and packages for Analysis.  
Also Reading File as text file and creating rdd out of it

```
1 rdd2 = credit.map(lambda x: x.split(','))

1 import datetime as dt

1 # as we don't have any info about column v1 to v28 we are not selecting them
2
3 ccard = rdd2.map(lambda x: [str(dt.datetime.strptime('01-01-2023 00:00:00', '%d-%m-%Y %H:%M:%S'))+\\
4                             dt.timedelta(seconds = float(x[0])),\
5                             float(x[-2]), int(x[-1][1])])

1 ccard.take(10)

[['2023-01-01 00:00:00', 149.62, 0],
 ['2023-01-01 00:00:00', 2.69, 0],
 ['2023-01-01 00:00:01', 378.66, 0],
 ['2023-01-01 00:00:01', 123.5, 0],
 ['2023-01-01 00:00:02', 69.99, 0],
 ['2023-01-01 00:00:02', 3.67, 0],
 ['2023-01-01 00:00:04', 4.99, 0],
 ['2023-01-01 00:00:07', 40.8, 0],
 ['2023-01-01 00:00:07', 93.2, 0],
 ['2023-01-01 00:00:09', 3.68, 0]]
```

# 1. WHAT IS THE AVERAGE AMOUNT BY CLASS ?



## 2. WHAT IS THE HIGHEST AMOUNT BY CLASS AND TIME ?

```
: 1 # grouping based on class and hour
: 2 classtimegroup = ccard.map(lambda x: [(x[2], x[8][11:13]), x[1]] )

: 1 classtimegroup.keys().distinct().take(6)
: [(1, '00'), (0, '01'), (1, '03'), (1, '04'), (1, '07'), (1, '08')]

: 1 high_amt = classtimegroup.reduceByKey(max).sortBy(lambda x: x[1], ascending = False)

: 1 high_amt.take(5)
: [((0, '22'), 25691.16),
 ((0, '13'), 19656.53),
 ((0, '02'), 18910.8),
 ((0, '11'), 12910.93),
 ((0, '12'), 11898.09)]
```

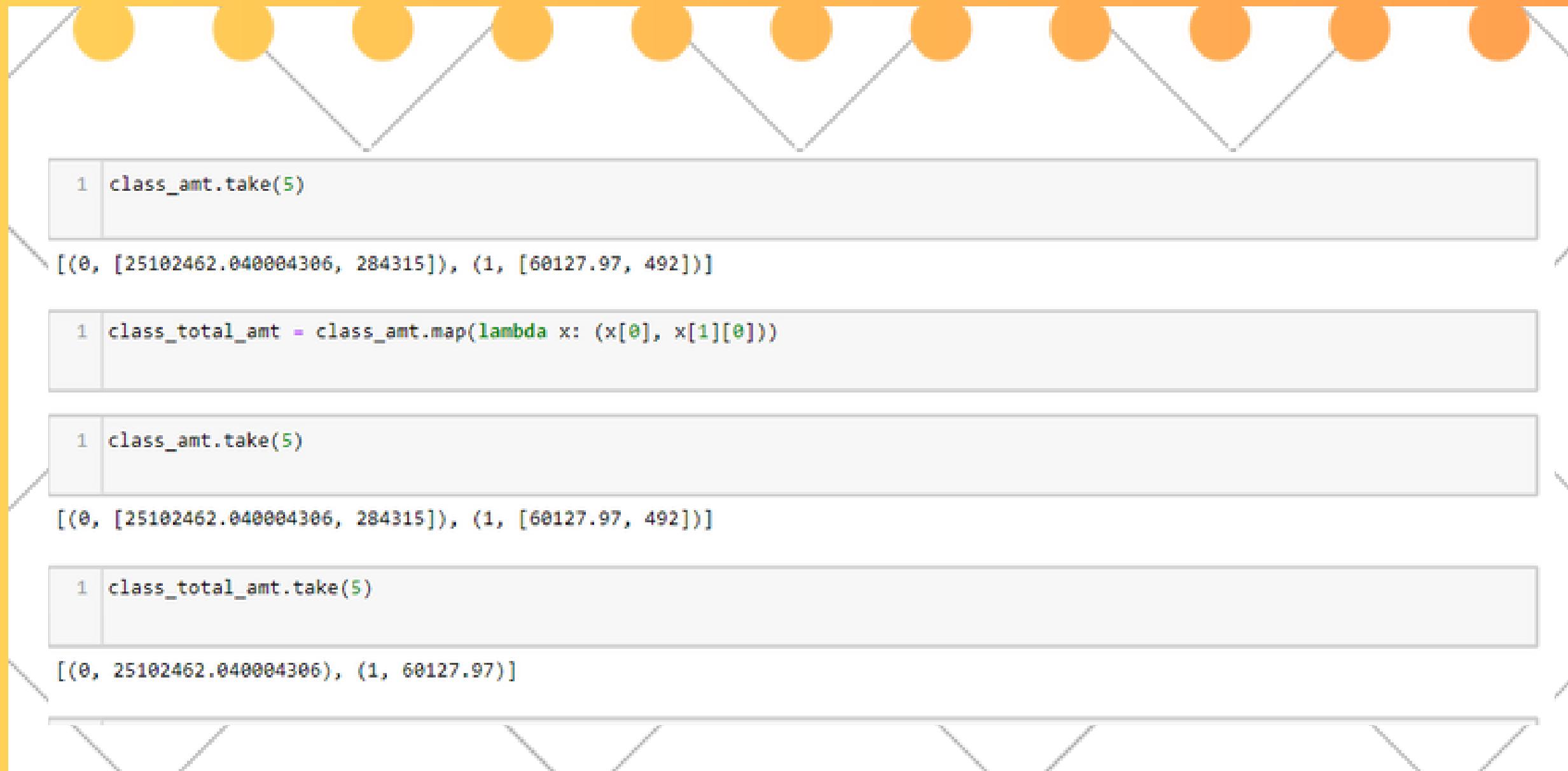
## 3. WHAT IS THE LOWEST AMOUNT BY TIME?

```
: 1 timegroup = ccard.map(lambda x: [ x[8][11:13], x[1]] )

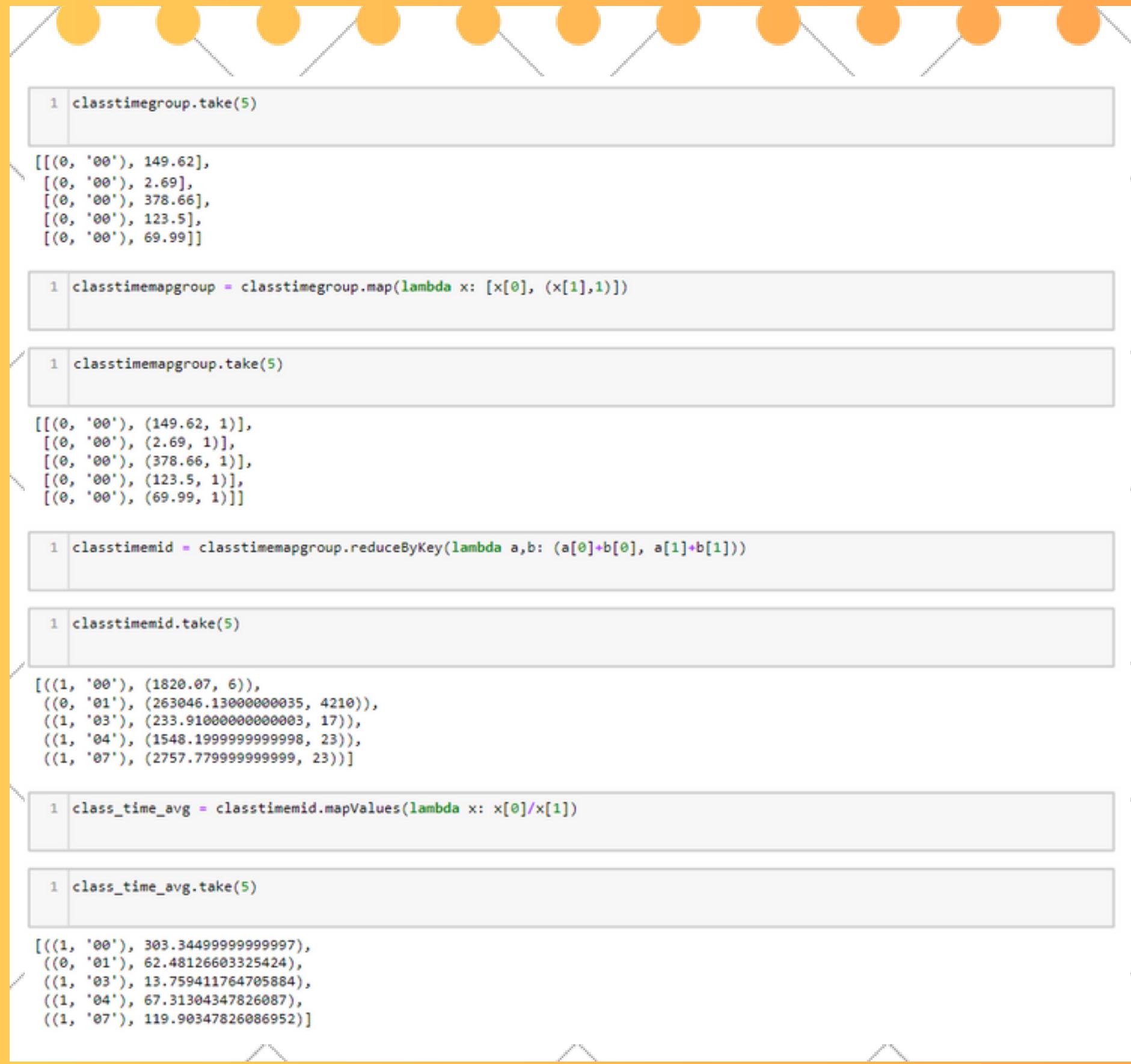
: 1 low_amt = timegroup.reduceByKey(min).sortBy(lambda x: x[1])

: 1 low_amt.take(5)
[('00', 0.0), ('02', 0.0), ('10', 0.0), ('01', 0.0), ('06', 0.0)]
```

## 4. WHAT IS THE TOTAL AMOUNT BY CLASS ?



## 5. WHAT IS THE AVERAGE AMOUNT BY CLASS AND TIME ?



“

**Believe Us. You're  
halfway there. Be with  
Us to the finish line.**



# DATA INGESTION USING AZURE DATA FACTORY

**Linked services**

Linked service defines the connection information to a data store or compute. [Learn more](#)

+ New

Filter by name Annotations : Any

Showing 1 - 2 of 2 items

Name ↑	Type ↑
AzureDataLakeStorage1	Azure Data Lake Storage Gen2
AzureSqlDatabase1	Azure SQL Database

**Factory Resources**

Filter resources by name +

- Pipelines 0
- Change Data Capture (preview) 0
- Datasets 2
  - source
  - target
- Data flows 0
- Power Query 0

target source

Azure SQL Database source

Connection Schema Parameters

Linked service \* AzureSqlDatabase1

Test connection Edit + New Learn more

Table dbo.credit\_card\_default

Move & transform

- Copy data
- Data flow

Synapse

Azure Data Explorer

Azure Function

Batch Service

Databricks

Data Lake Analytics

General

HDInsight

Iteration & conditionals

Machine Learning

Power Query

Copy data

Copy data1

General Source Sink Mapping Settings User properties

Type conversion settings

Import schemas Preview source

New mapping Clear Reset Delete

Source	Type	Destination	Type
CUSTID	12s smallint	Column 1	abc String
LIMIT_BAL	123 int	Column 2	abc String
SEX	abc nchar	Column 3	abc String

Activity runs

Pipeline run ID cf574013-4e4c-4f2e-b218-58f1f8700e9b

All status

Showing 1 - 1 items

Activity name ↑	Status ↑	Activity type ↑	Run start ↑	Duration ↑	Log	Integration runtime ↑	User prop
Copy data1	Succeeded	Copy data	6/5/2023, 10:12:54 PM	00:00:09		AutoResolveIntegrationRu	

# MOUNTING ADLS TO DATABRICKS

```
Code 2
1 dbutils.fs.mount(
2   source="wasbs://capstone/myfirstacc18.blob.core.windows.net",
3   mount_point="/mnt/mount_point",
4   extra_configs={"fs.azure.account.key.myfirstacc18.blob.core.windows.net": "48e+oeEDr80fffff1CHm8H0eGqHWWr/ETEDo/MOBK6IYL7uZGAewftTdrM28abe27rnk5byDME+AStx6t0nw=="})
5

Out[8]: True
Command took 11.34 seconds -- by trinanthudhrisan.dhruv@gmail.com at 6/12/2023, 10:05:46 PM on Rishabh Dhruvan's Cluster

Code 3
1 dbutils.fs.ls("/mnt/mount_point")

Out[10]: [FileInfo(path='dbfs:/mnt/mount_point/credit-card-default-1000.csv', name='credit-card-default-1000.csv', size=94292, modificationTime=1686553962000)]
Command took 0.27 seconds -- by trinanthudhrisan.dhruv@gmail.com at 6/12/2023, 10:05:48 PM on Rishabh Dhruvan's Cluster

Code 4
1 data = sc.textFile("dbfs:/mnt/mount_point/credit-card-default-1000.csv")

Command took 0.11 seconds -- by trinanthudhrisan.dhruv@gmail.com at 6/12/2023, 10:05:50 PM on Rishabh Dhruvan's Cluster

Code 5
1 data.take(5)

▶ (5) SparkJobs
Out[14]: [(1,CUSTID,LIMIT_BAL,SEX,EDUCATION,MARRIAGE,AGE,PAY_1,PAY_2,PAY_3,PAY_4,PAY_5,PAY_6,BILL_AMT1,BILL_AMT2,BILL_AMT3,BILL_AMT4,BILL_AMT5,BILL_AMT6,PAY_AMT1,PAY_AMT2,PAY_AMT3,PAY_AMT4,PAY_AMT5,PAY_AMT6,DEFAUTED),
  '530,20000,2,2,2,21,-1,-1,2,2,-2,-2,0,0,0,0,0,0,0,0,0,162000,0,0',
  '38,60000,2,2,2,22,0,0,0,0,-2,-2,0,0,0,0,0,0,0,0,0,1576,0',
  '40,10000,1,2,2,22,0,0,0,0,-2,-2,0,0,0,0,0,0,0,0,0,1500,0',
  '47,20000,2,1,2,22,0,0,2,-1,0,-1,1131,291,582,291,0,291,582,0,0,130291,651,0']

Command took 1.23 seconds -- by trinanthudhrisan.dhruv@gmail.com at 6/12/2023, 10:06:45 PM on Rishabh Dhruvan's Cluster
```

# DATA CLEANING

```
Cmd 7
1 header = data.first()

+ (1) Spark jobs

Command took 0.29 seconds -- by himanshudhiman.dhiman@gmail.com at 6/12/2023, 1:09:00 PM on Himanshu Dhiman's Cluster

Cmd 8
1 data = data.filter(lambda x:x!=header)

+ (1) Spark jobs

Command took 0.11 seconds -- by himanshudhiman.dhiman@gmail.com at 6/12/2023, 1:09:04 PM on Himanshu Dhiman's Cluster

Cmd 9
1 data.take(5)

+ (1) Spark jobs

Out[17]: ['530,20000,2,2,2,21,-1,-1,2,2,-2,-2,0,0,0,0,0,0,0,0,0,0,0,162000,0,0',
'38,60000,2,2,2,22,0,0,0,-2,-2,0,0,0,0,0,0,0,0,0,0,0,1576,0',
'43,10000,1,2,2,22,0,0,0,-2,-2,0,0,0,0,0,0,0,0,0,0,1500,0',
'47,20000,2,1,2,22,0,0,2,-1,0,-1,1131,291,582,291,0,291,291,582,0,0,136291,651,0',
'70,20000,1,4,2,22,2,0,0,0,-1,-1,1692,13250,433,1831,0,2891,13250,433,1831,0,2891,153504,0']

Command took 0.33 seconds -- by himanshudhiman.dhiman@gmail.com at 6/12/2023, 1:09:46 PM on Himanshu Dhiman's Cluster

Cmd 10
1 data.count()

+ (1) Spark jobs

Out[18]: 1000

Command took 0.12 seconds -- by himanshudhiman.dhiman@gmail.com at 6/12/2023, 1:09:53 PM on Himanshu Dhiman's Cluster
```

```
1   ## 3.Cleanup data. Remove lines that are not "CSV"

Command took 0.11 seconds -- by himanshudhiman.dhiman@gmail.com at 6/12/2023, 1:10:38 PM on Himanshu Dhiman's Cluster

Cmd 12

1   data = data.map(lambda x:x.split(','))

Command took 0.06 seconds -- by himanshudhiman.dhiman@gmail.com at 6/12/2023, 1:09:57 PM on Himanshu Dhiman's Cluster

Cmd 13

1   data = data.filter(lambda x:all(x))

Command took 0.07 seconds -- by himanshudhiman.dhiman@gmail.com at 6/12/2023, 1:10:53 PM on Himanshu Dhiman's Cluster

Cmd 14

1   data.count()

▶ (1) Spark Jobs

Out[23]: 1000

Command took 0.33 seconds -- by himanshudhiman.dhiman@gmail.com at 6/12/2023, 1:10:57 PM on Himanshu Dhiman's Cluster

Cmd 15

1   data.take(1)

▶ (1) Spark Jobs

Out[24]: [['530',
'20000',
'2',
'2',
'2',
'21',
'-1',
'-1',
'2',
'2',
'-2',
'-2',
'0',
'0',
'0',
'0',
'0',
'0',
'0',
'0',
'0',
'0',
'0',
'0']]

Command took 0.30 seconds -- by himanshudhiman.dhiman@gmail.com at 6/12/2023, 1:11:03 PM on Himanshu Dhiman's Cluster
```

# ROUNDING OF AGE TO RANGE OF 10S

```
Cmd 23
1   # or rounding of age to range of 10s.
Command took 0.00 seconds -- by himanshudivman.dhiman@gmail.com at 6/12/2023, 2:17:28 PM on Himanshu Divman's Cluster

Cmd 24
1   from pyspark.sql.functions import *
Command took 0.01 seconds -- by himanshudivman.dhiman@gmail.com at 6/12/2023, 2:17:32 PM on Himanshu Divman's Cluster

Cmd 25
1   def func(x):
2       x = floor(x/10)
3       x = x*10
4       if      return (x/10)*10
5       return x
6
7
8   df = df.withColumn('range_of_10s', func(col('AGE')))

P [df: pyspark.sql.DataFrame = [CUSTID: long, LIMIT_BAL: long ... 24 more fields]
Command took 0.13 seconds -- by himanshudivman.dhiman@gmail.com at 6/12/2023, 2:17:57 PM on Himanshu Divman's Cluster

Cmd 26
1   df.show(5)

P [(1) Spark jobs
+-----+
|CUSTID|LIMIT_BAL|SEX|EDUCATION|MARRIAGE|AGE|PAY_1|PAY_2|PAY_3|PAY_4|PAY_5|PAY_6|BILL_AMT1|BILL_AMT2|BILL_AMT3|BILL_AMT4|BILL_AMT5|BILL_AMT6|PAY_AMT1|PAY_AMT2|PAY_AMT3|PAY_AMT4|PAY_AMT5|PAY_AMT6|DEFAULTED|range_of_10s|
+-----+
| 529| 20000| 2| 2| 2| 21| -1| -1| 2| 2| -2| -2| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 20|
| 38| 60000| 2| 2| 2| 22| 0| 0| 0| 0| -2| -2| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 20|
| 41| 10000| 1| 2| 2| 0| 0| 0| 0| -2| -2| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 0| 20|
| 47| 20000| 2| 1| 2| 22| 0| 0| 2| -1| 0| -1| 1131| 291| 562| 291| 0| 291| 291| 562| 0| 0| 130291| 651| 0| 20|
| 79| 20000| 1| 4| 2| 22| 2| 0| 0| 0| -1| -1| 1692| 13250| 433| 1831| 0| 2891| 13250| 433| 1831| 0| 2891| 153504| 0| 20|
+-----+
only showing top 5 rows

Command took 0.06 seconds -- by himanshudivman.dhiman@gmail.com at 6/12/2023, 2:18:16 PM on Himanshu Divman's Cluster
```

# **NORMALIZE SEX TO ONLY 1 AND 2**

# AVERAGE BILLED AMOUNT

Cmd 31

```
1. %%SQL Average billed amount
```

Command took 0.04 seconds -- by himanshudhiman.dhiman@gmail.com at 4/12/2023, 2:19:41 PM on Himanshu Dhiman's Cluster

Cmd 32

```
1. df.select(avg(col("BILL_AMT1") + col("BILL_AMT2") + col("BILL_AMT3") + col("BILL_AMT4") + col("BILL_AMT5") + col("BILL_AMT6")) .alias("avg_bill_amount")) .show()
```

\* (2) Spark Jobs

```
+-----+  
| avg_bill_amount |  
+-----+  
| 259790.697 |  
+-----+
```

Command took 1.53 seconds -- by himanshudhiman.dhiman@gmail.com at 4/12/2023, 2:19:45 PM on Himanshu Dhiman's Cluster

Cmd 33

# AVERAGE PAY DURATION

Code 34

```
1 df.select(avg(col("PAY_1")+col("PAY_2")+col("PAY_3")+col("PAY_4")+col("PAY_5")+col("PAY_6")).alias("avg_pay_amount")).show()
```

▶ (2) Spark Jobs

```
+-----+  
|avg_pay_amount|  
+-----+  
|      -1.185|  
+-----+
```

Command took 0.58 seconds -- by himanshudhiman.dhiman@gmail.com at 6/12/2023, 2:19:18 PM on Himanshu Dhiman's Cluster

Code 35

# ADD SEXNAME TO THE DATA USING JOINS

```
Cmd 38
1 sex_name = sc.parallelize([('1','Male'), ('2', "Female")])
Command took 0.10 seconds -- by himanshudhiman.dhiman@gmail.com at 6/12/2023, 2:20:16 PM on Himanshu Dhiman's Cluster

Cmd 39
1 sex_name_df = sex_name.toDF(["Id","SEX_NAME"])
Command took 0.07 seconds -- by himanshudhiman.dhiman@gmail.com at 6/12/2023, 2:20:20 PM on Himanshu Dhiman's Cluster

Cmd 40
1 df1 = df.join(sex_name_df,df.SEX == sex_name_df.Id , 'inner').select("*")
Command took 0.09 seconds -- by himanshudhiman.dhiman@gmail.com at 6/12/2023, 2:20:28 PM on Himanshu Dhiman's Cluster

Cmd 41
1 df1 = df.join(sex_name_df,df.SEX == sex_name_df.Id , 'inner').select("*")
Command took 0.16 seconds -- by himanshudhiman.dhiman@gmail.com at 6/12/2023, 2:20:36 PM on Himanshu Dhiman's Cluster
```

# CREATING TEMPORARY VIEW TO DO SQL QUERIES



```
Ques 44
1 app II.

Command took 0.12 seconds. -- by himanshu.bhiman.chirayangmail.com at 6/12/2023, 2:21:30 PM on Himanshu Bhiman's Cluster

Ques 45
1 Creating temporary table for sql query
2 df1.createOrReplaceTempView("CDDATA")

Command took 0.17 seconds. -- by himanshu.bhiman.chirayangmail.com at 6/12/2023, 2:21:40 PM on Himanshu Bhiman's Cluster
```



# PERCENTAGE OF DEFALTERES BY MARRIAGE AND EDUCATION STATUS

```
Cmd 58

spark.sql("SELECT MARRIAGE, EDUCATION, count(*) as Total, SUM(DEFAULTED) as Defaults, \
           ROUND(SUM(DEFAULTED) * 100 / count(*)) as PER_DEFAULT FROM CCDATA GROUP BY MARRIAGE,EDUCATION ORDER BY 1,2").show()

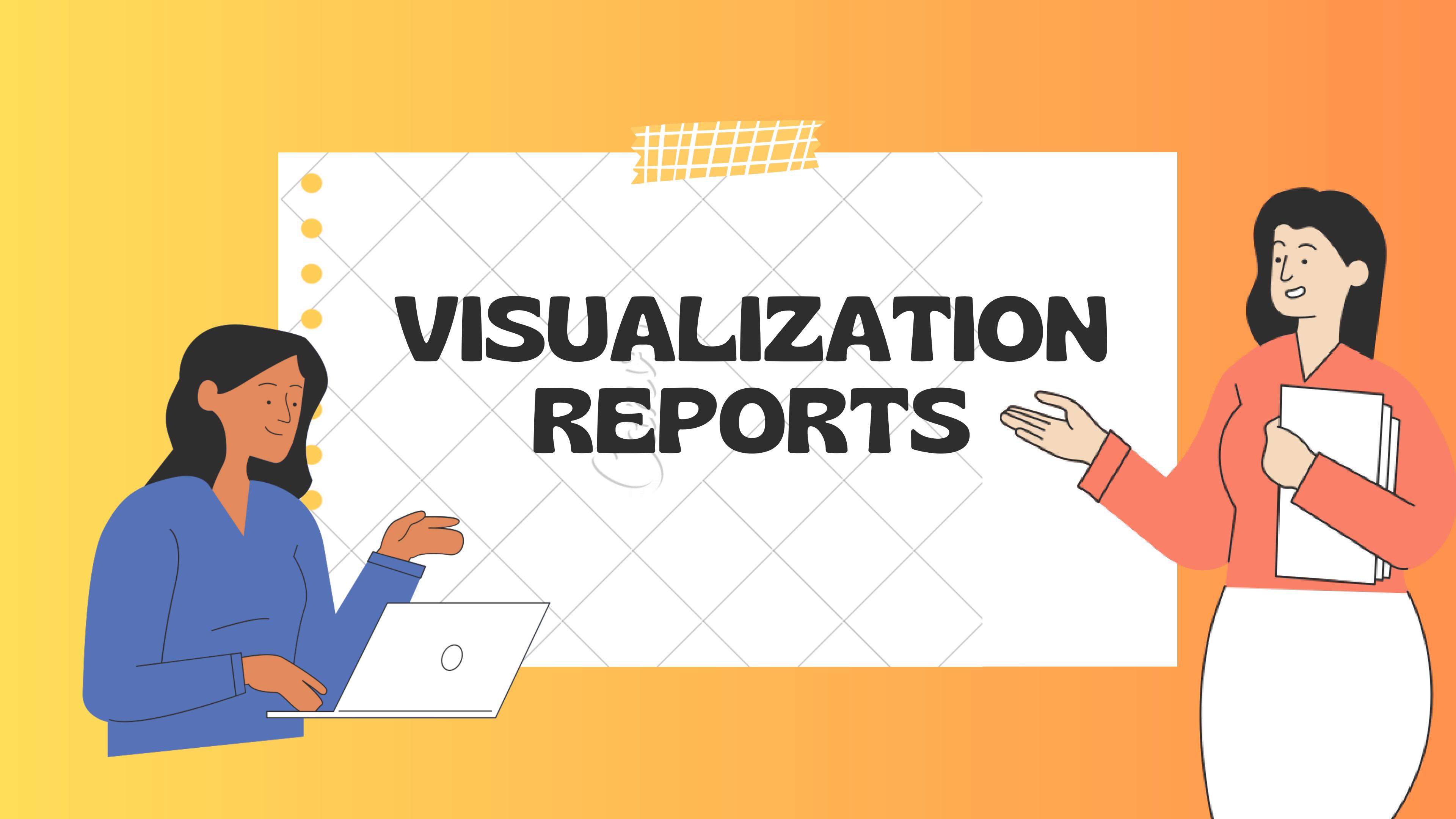
▶ (7) Spark Jobs

+-----+-----+-----+-----+
|MARRIAGE|EDUCATION|Total|Defaults|PER_DEFAULT|
+-----+-----+-----+-----+
|      1|       1|   123|     71|      58.0|
|      1|       2|   198|    105|      53.0|
|      1|       3|    87|     52|      60.0|
|      1|       4|     3|     2|      67.0|
|      2|       1|   268|     69|      26.0|
|      2|       2|   243|     65|      27.0|
|      2|       3|    55|     24|      44.0|
|      2|       4|     4|     2|      50.0|
|      3|       1|     4|     4|     100.0|
|      3|       2|     7|     3|      43.0|
|      3|       3|     8|     6|      75.0|
+-----+-----+-----+-----+
Command took 1.74 seconds -- by himanshudhiman.dhiman@gmail.com at 6/13/2023, 10:17:21 AM on Himanshu Dhiman's Cluster
```

## PERCENTAGE OF DEFALTERES BY GENDER

```
Cmd 47
spark.sql("SELECT SEX_NAME,count(*) as Total,SUM(DEFAULTED) as Defaults, \
           ROUND(SUM(DEFAULTED)*100/count(*)) as PER_DEFAULT FROM CCDATA GROUP BY SEX_NAME").show()

+-----+-----+-----+-----+
|SEX_NAME|Total|Defaults|PER_DEFAULT|
+-----+-----+-----+-----+
| Female|  591|    218|      37.0|
|  Male|  409|    185|      45.0|
+-----+-----+-----+-----+
Command took 2.81 seconds -- by himanshudhiman.dhiman@gmail.com at 6/13/2023, 10:17:19 AM on Himanshu Dhiman's Cluster
Final All
```



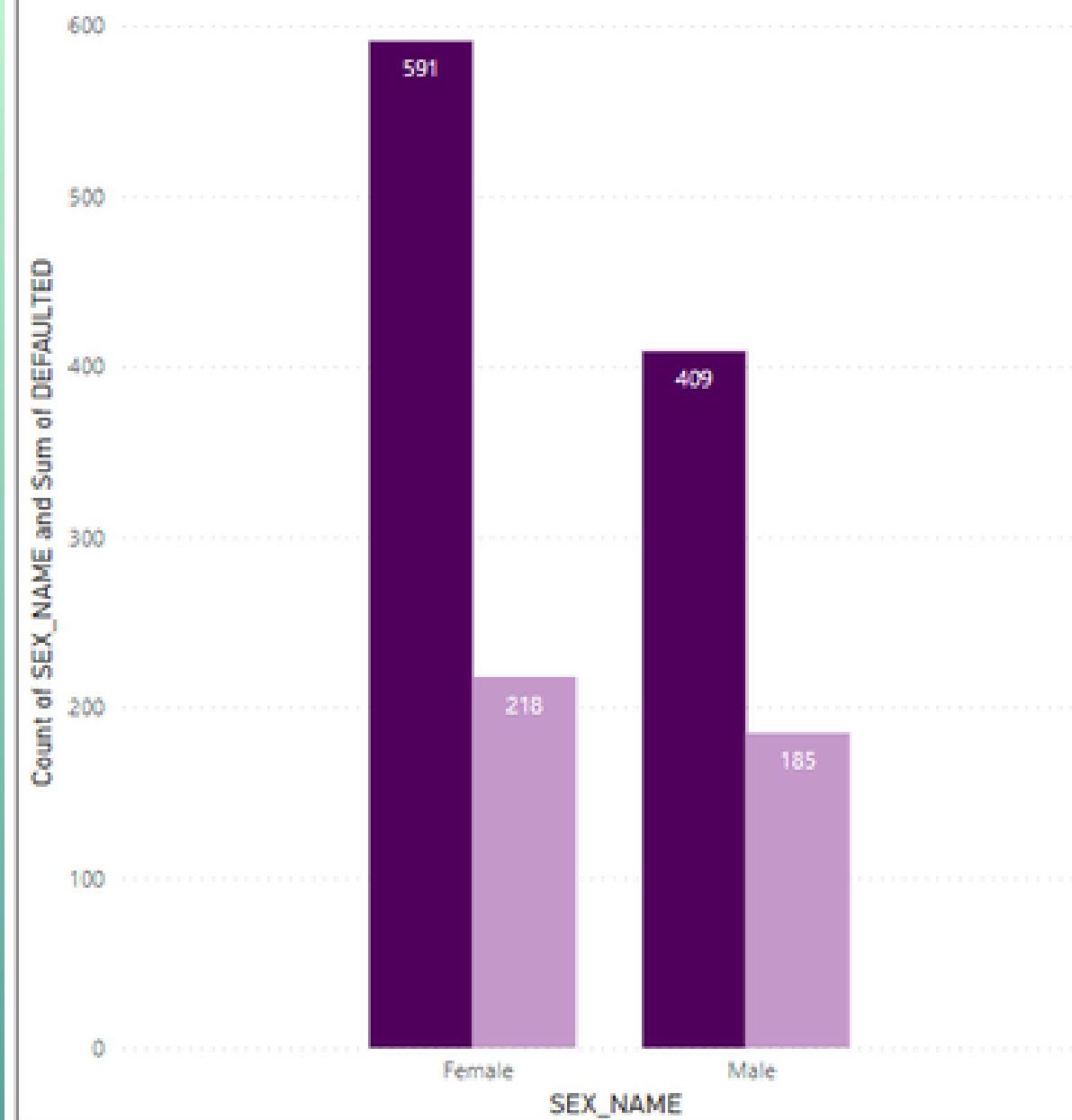
# VISUALIZATION REPORTS

The background features a large white rectangular area with a grid pattern. In the top right corner of this area, there is a yellow decorative element consisting of a grid of squares. To the left of the main title, a woman with dark hair, wearing a blue jacket, is seated at a desk, looking at a laptop screen. To the right, another person in an orange jacket is standing and gesturing with their hands while holding a stack of papers.

# VISUALIZATION

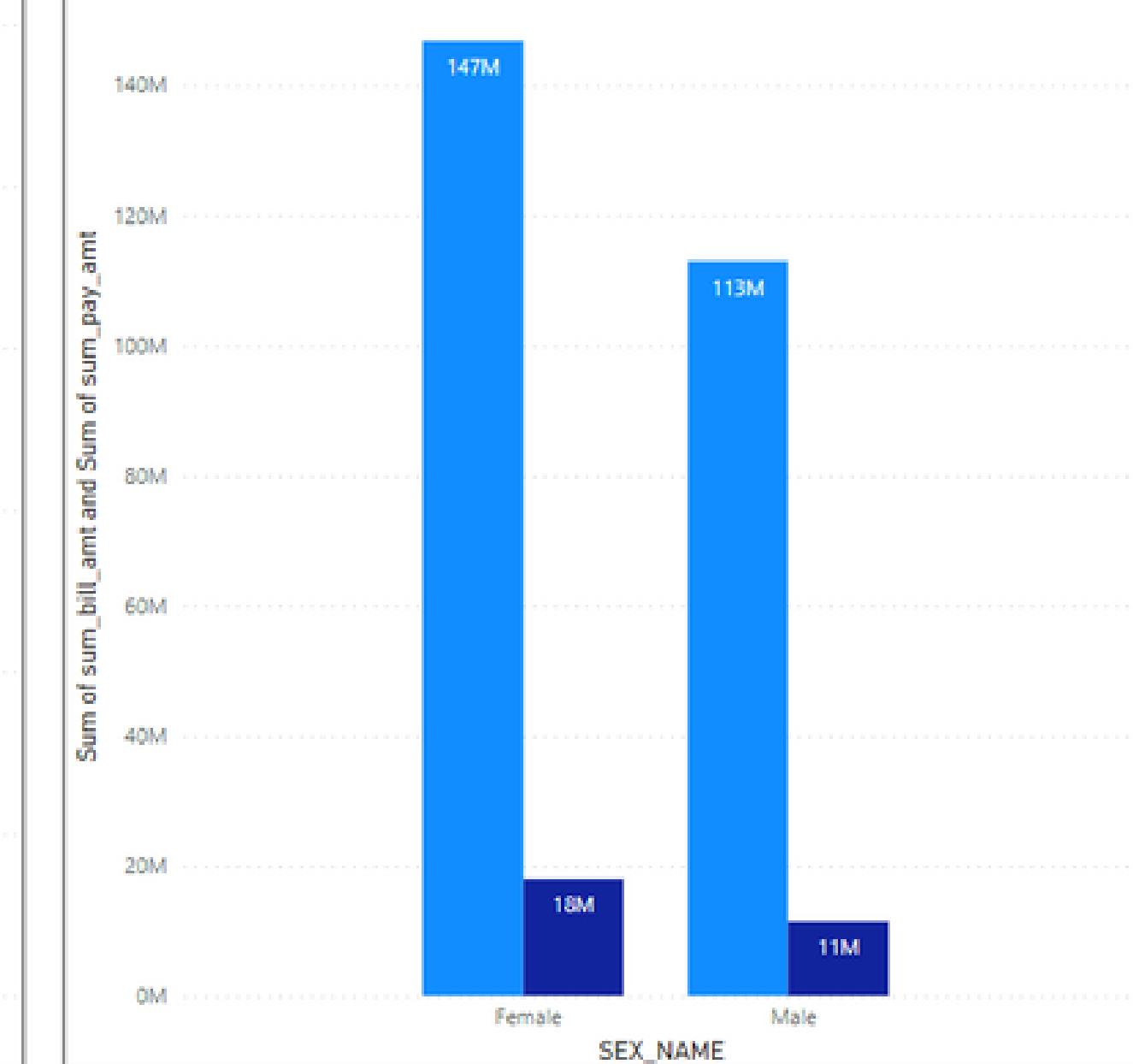
Count of gender Vs Number of defaulter

● Count of SEX\_NAME ● Sum of DEFULTED



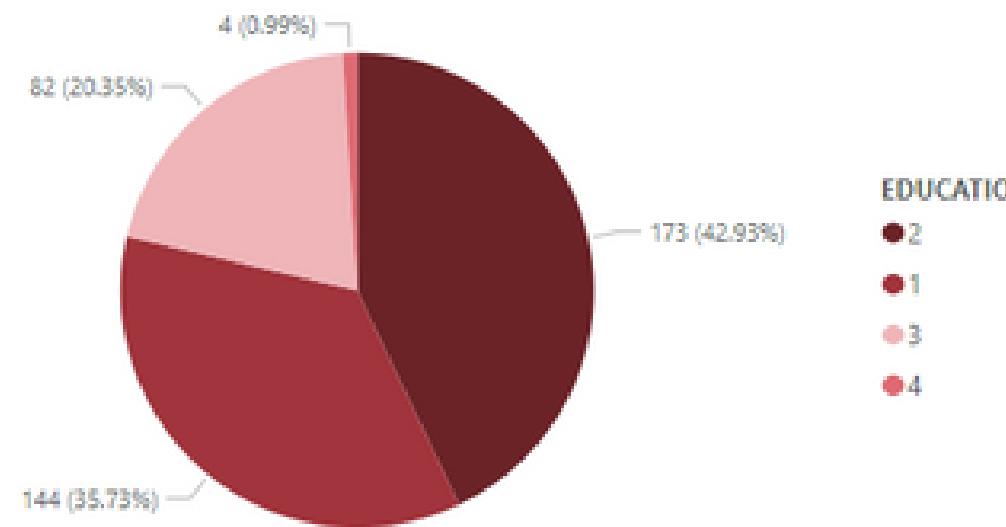
Bill amount VS Amount payed

● Sum of sum\_bill\_amt ● Sum of sum\_pay\_amt

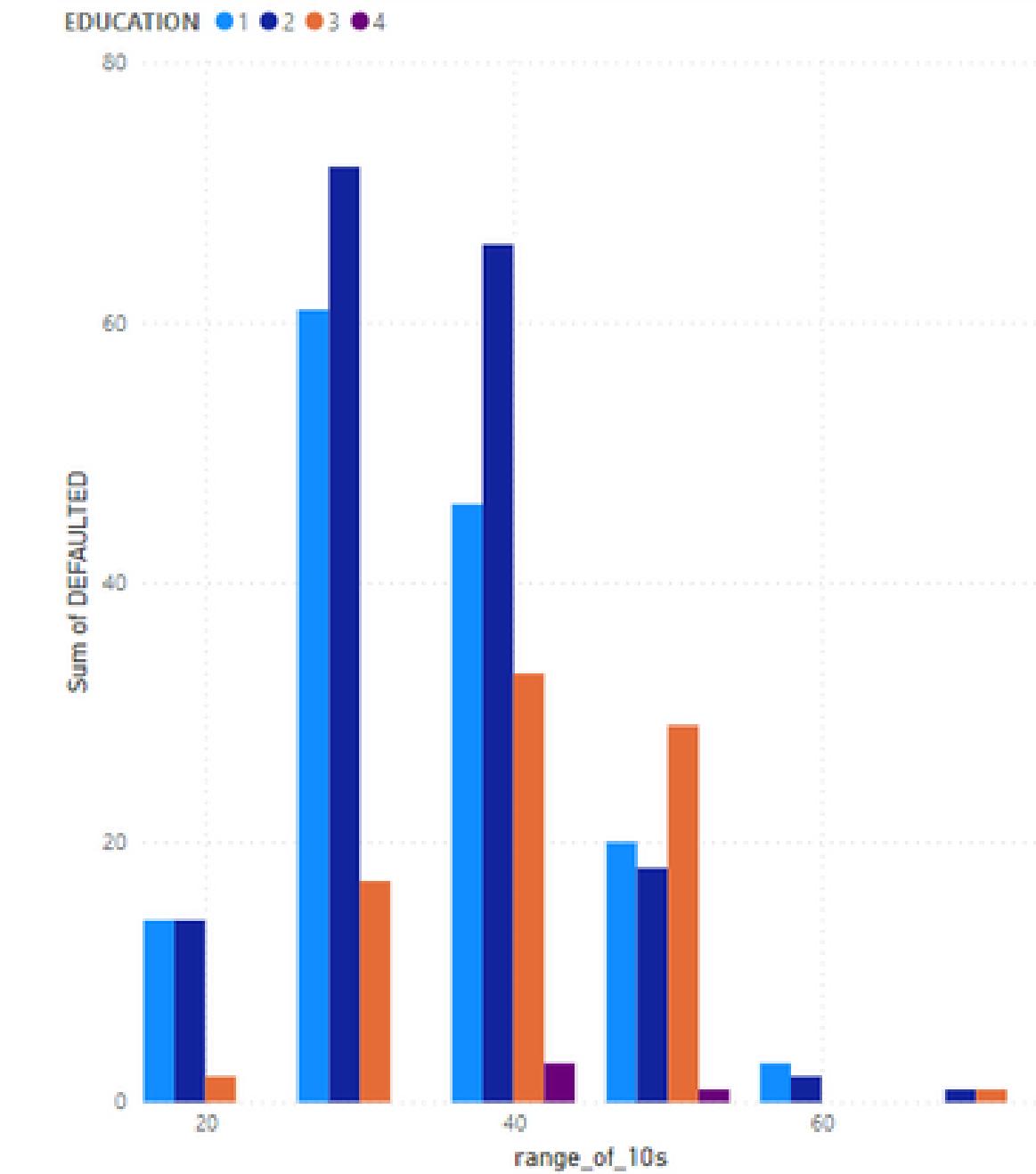


# VISUALIZATION

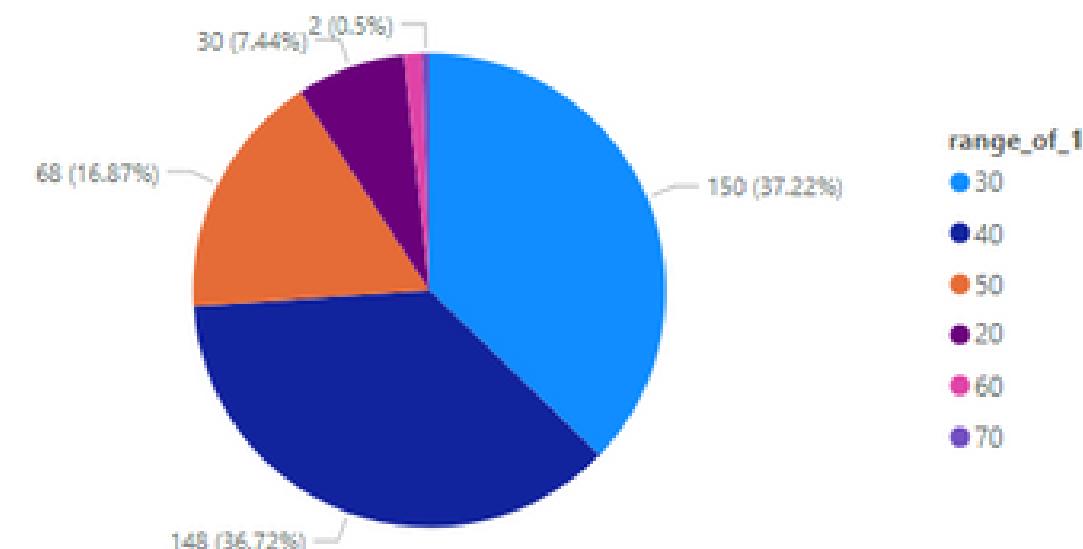
Education wise Defaulter percentage



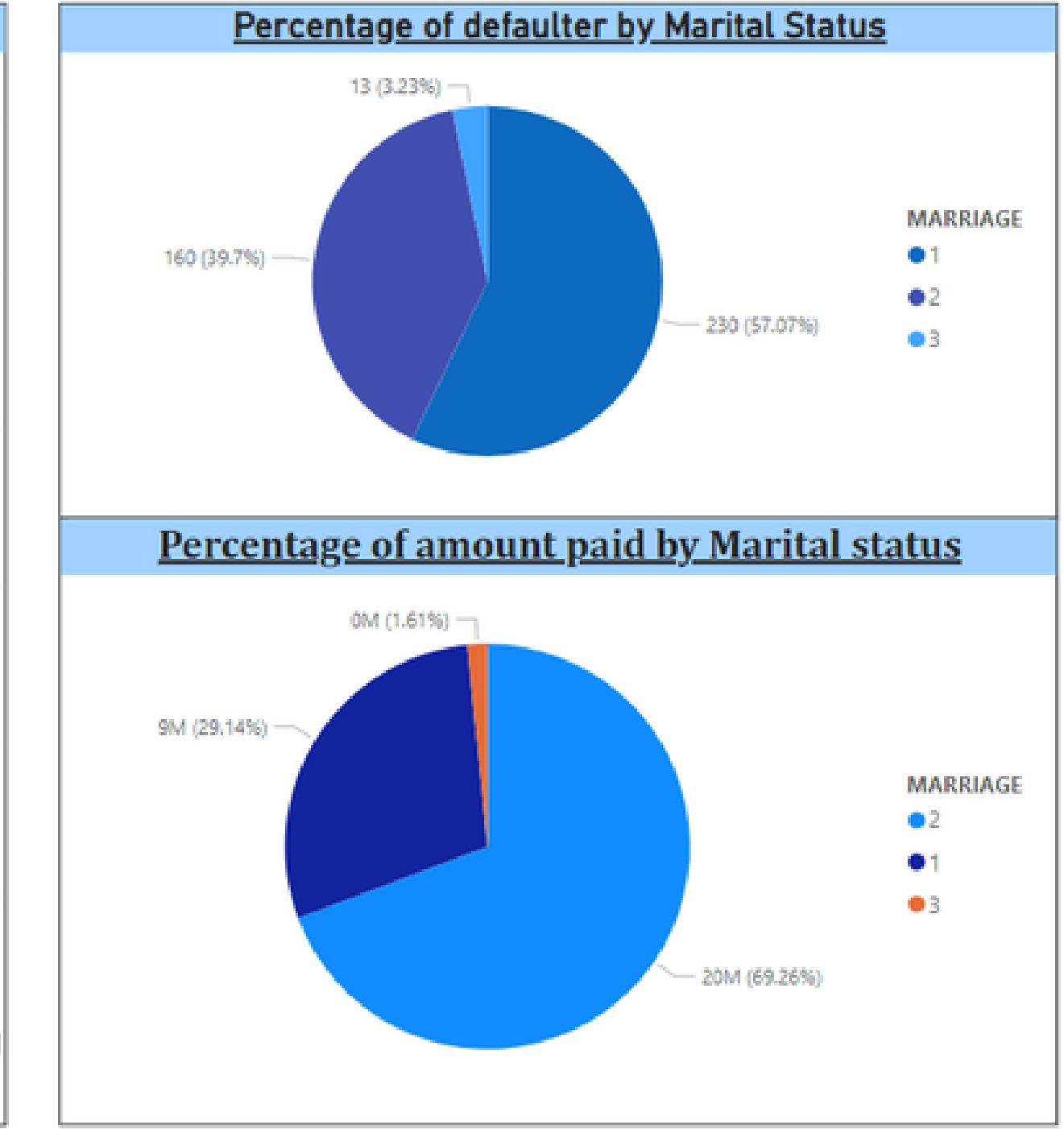
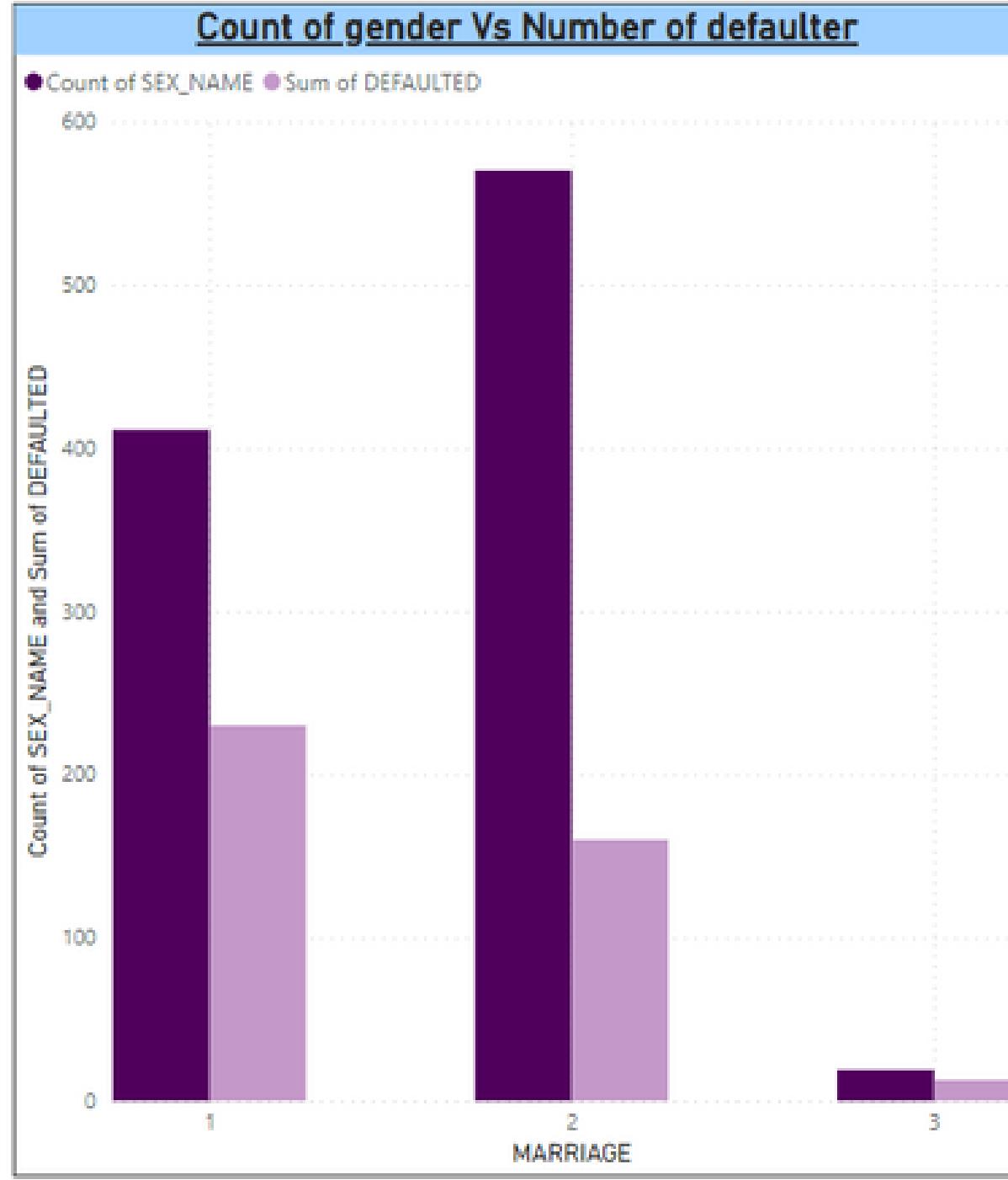
Number of defaulter age range wise



Age range wise Defaulter percentage



# VISUALIZATION





# CONCLUSION

- PAYMENT HISTORY
- CREDIT SCORE
- DEBT-TO-INCOME RATIO
- EMPLOYMENT AND INCOME STABILITY
- PREVIOUS DEFAULT HISTORY

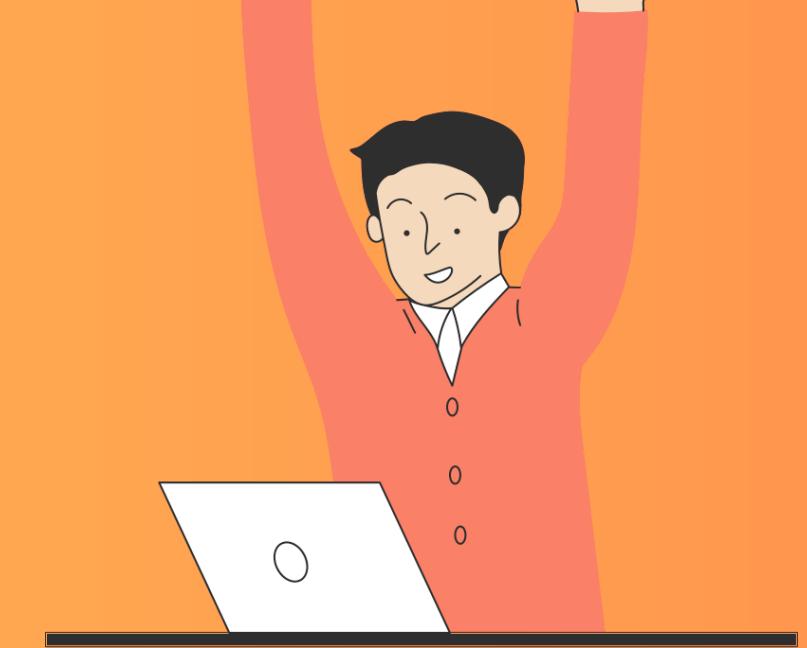
“

## **Let's Take a moment to thank our mentor Venkat sir.**

Thank You so much sir for giving us wonderful training. Being as a good mentor, thank you for giving industrial knowledge as well as sharing your experience, not only technical stuff but also real time team work experience.

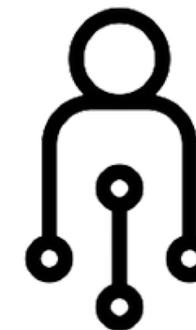
We whole heartedly thank your sir for your irreplaceable and immeasurable work that created an impact on our growth and success

”



“

# Thank You..!



**Futurense**

”





Have a  
great day  
ahead.



6

**ANY**

**QUESTIONS ?**