Module:

**BDM 1024 - Data Technology Solutions 02**

Project Title:

**Fraud Detection in the Banking Sector (Final Submission)**

Intake:

**Spring 2023 | Term One**

Submitted By: **Group A**

Koshish Aryal

Punam Bhattarai

Saaz Neupane

Bipin Pandey

Shweta Laljibhai Thummar

Manan Tushar Kapadia

Robert Thapa

Dipti Baral

# Introduction

This project aims to provide a reliable fraud detection system for the banking industry. Fraudulent transactions have the potential to cause large financial losses and harm a bank's standing. To identify and stop fraudulent activity in real time, an efficient data model in conjunction with data analytics is essential.

# Problem Definition

A GlobalData survey conducted in September 2022 stated that 25% of Canadians had experienced fraud in the past three years. Among them, 35% were concerned about identity theft where as 26% worried fraudsters would steal their card details online. According to the new survey by the TD Bank Group, around 62% of Canadians feel more targeted now than ever by the financial fraud. Along with that according to TD, nearly 72% of Canadians said they were targeted by email/text message fraud, while 66% reported being targeted over the phone. In order to study the claim and further looking into the problem, we decided to looking into the banking data and possible fraud transactions that has been occurring.
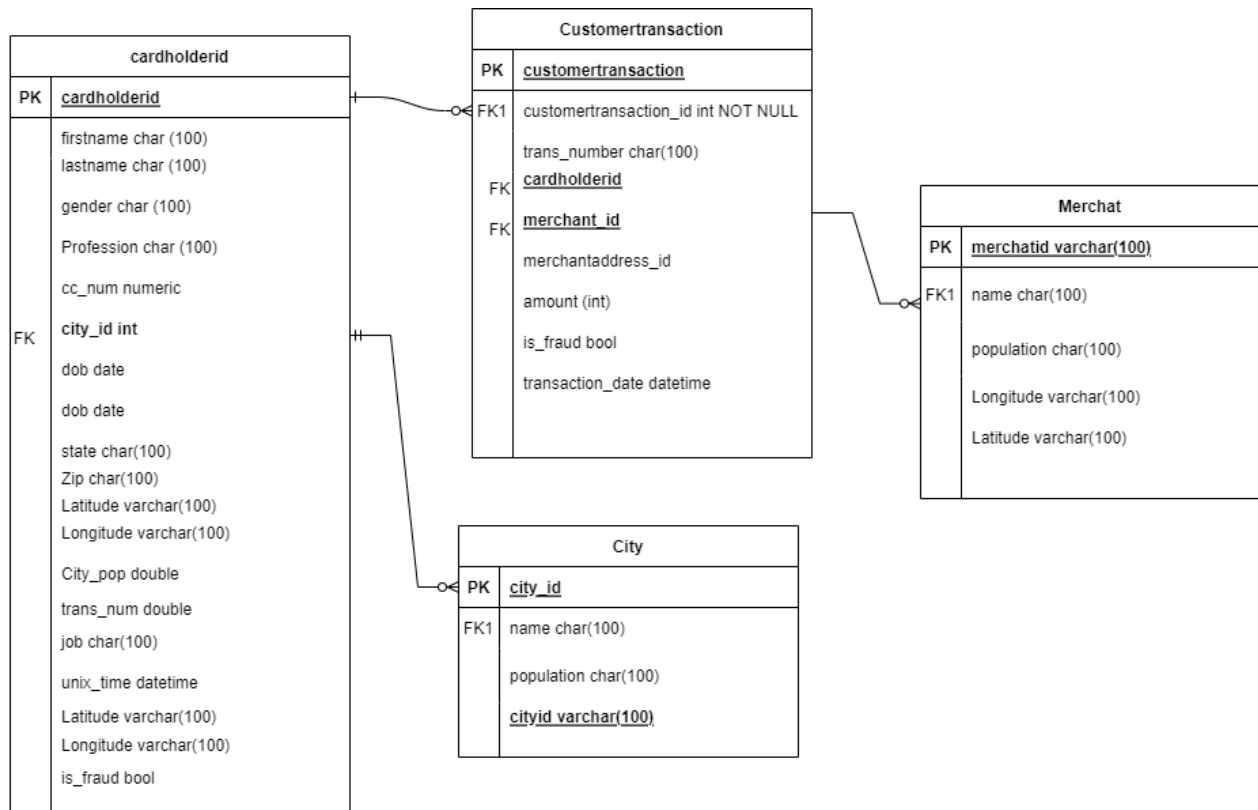
# Data Design Pattern

The system analyzes a fraud dataset using two patterns: Role (Cardholder and Merchant) and Context (transaction details, cities, and merchant addresses). Role Pattern captures specific attributes for each role, while Context Pattern focuses on spatial and temporal aspects of transactions, helping detect anomalies based on locations, dates, and amounts for better fraud detection.

# Data Model Design

The data model is designed to capture essential information about cardholders, merchants, and customer transactions.

**Entity-Relationship Diagram (ERD):**

**cardholderid**

| PK | cardholderid |
|----|----|
| | firstname char (100) |
| | lastname char (100) |
| | gender char (100) |
| | Profession char (100) |
| | cc_num numeric |
| FK | city_id int |
| | dob date |
| | dob date |
| | state char(100) |
| | Zip char(100) |
| | Latitude varchar(100) |
| | Longitude varchar(100) |
| | City_pop double |
| | trans_num double |
| | job char(100) |
| | unix_time datetime |
| | Latitude varchar(100) |
| | Longitude varchar(100) |
| | is_fraud bool |

**Customertransaction**

| PK | customertransaction |
|----|----|
| FK1 | customertransaction_id int NOT NULL |
| | trans_number char(100) |
| FK | cardholderid |
| FK | merchant_id |
| | merchantaddress_id |
| | amount (int) |
| | is_fraud bool |
| | transaction_date datetime |

**Merchat**

| PK | merchatid varchar(100) |
|----|----|
| FK1 | name char(100) |
| | population char(100) |
| | Longitude varchar(100) |
| | Latitude varchar(100) |

**City**

| PK | city_id |
|----|----|
| FK1 | name char(100) |
| | population char(100) |
| | cityid varchar(100) |

**Entity Descriptions:**

- **Cardholder:** Represents banking customers, with attributes like CardholderID, FirstName, LastName, Gender, Profession, Cc_num, Dob, CityID, Street, State, Zip, Lat, and Long.
- **City:** Stores information about cities, including CityID, Name, and Population.
- **Merchant:** Holds data about merchants, with attributes like MerchantID, Name, and Category.
- **Customer Transaction:** Represents individual transactions made by cardholders. It includes CustomerTransactionID, Trans_num, CardholderID, MerchantID, Amount, Is_fraud, and Transaction_date.

# SQL Queries and Insights

SQL (Structured Query Language) is a programming language used for managing and manipulating relational databases. It is widely used in the field of data management and is considered the standard language for interacting with databases. We utilized SQL queries to extract valuable insights from the dataset and analyze transaction patterns.

```
Query    Query History
1  SELECT
2     c.merchantid,
3     COUNT(*) AS transaction_count,
4     SUM(amount) AS total_amount
5  FROM
6     customertransaction c join merchant m on c.merchantid = m.merchantid
7  GROUP BY
8     c.merchantid
9  HAVING
10    SUM(amount) > (SELECT AVG(amount) * 2 FROM customertransaction)
11 ORDER BY
12    total_amount DESC;
13
```

**Insight:** This query calculates the total number of transactions and the total transaction amount for each merchant. It then filters the results to include only those merchants whose total transaction amount is greater than twice the average transaction amount across all transactions. The final result set is sorted in descending order based on the total transaction amount, displaying merchants with the highest total amounts first.

```
Query    Query History
1  SELECT m.Category, AVG(ct.Amount) as AvgTransactionAmount
2  FROM CustomerTransaction ct
3  JOIN Merchant m ON ct.MerchantID = m.MerchantID
4  GROUP BY m.Category
5  ORDER BY AvgTransactionAmount DESC;
6
```

**Insight:** This query calculates the average transaction amount for each merchant category and presents the results in descending order based on the average transaction amount. The result set will show the merchant categories with the highest average transaction amounts first, providing insights into the most profitable or popular categories based on transaction data.

```
Query    Query History
 1   WITH suspicious_transactions AS (
 2     SELECT
 3       custtransactionid, amount, cardholderid, transaction_date,
 4       COUNT(*) OVER (PARTITION BY cardholderid) AS num_transactions,
 5       AVG(amount) OVER (PARTITION BY cardholderid) AS avg_transaction_amount,
 6       STDDEV(amount) OVER (PARTITION BY cardholderid) AS std_dev_transaction_amount,
 7       LAG(amount) OVER (PARTITION BY cardholderid ORDER BY transaction_date) AS prev_transaction_amount
 8     FROM
 9       customertransaction
10   ),
11   potential_fraud AS (
12     SELECT
13       custtransactionid, amount, cardholderid, transaction_date, num_transactions, avg_transaction_amount,
14       std_dev_transaction_amount, prev_transaction_amount,
15       CASE
16         WHEN amount > (avg_transaction_amount + (3 * std_dev_transaction_amount)) THEN 'High Amount'
17         WHEN amount > (1.5 * prev_transaction_amount) THEN 'Significant Increase'
18         WHEN num_transactions > 10 THEN 'High Frequency'
19         ELSE 'None'
20       END AS fraud_type
21     FROM
22       suspicious_transactions
23   )
24   SELECT *
25   FROM potential_fraud
26   WHERE fraud_type <> 'None';
27
```

**Insight:** This query uses CTEs and window functions to analyze transaction data, calculate transaction statistics per cardholder, and identify potential fraudulent transactions based on specific criteria. The result is a list of suspicious transactions with their respective fraud types.

```
Query    Query History
 1   SELECT
 2      custtransactionid,
 3      amount
 4   FROM
 5      customertransaction
 6   WHERE
 7      ABS(
 8          amount - (SELECT AVG(amount) FROM customertransaction))
 9          > 3 * (SELECT STDDEV(amount) FROM customertransaction);
10
```

**Insight:** This query identifies transactions that significantly deviate from the average transaction amount. It does so by comparing the absolute difference between each transaction amount and the average transaction amount to three times the standard deviation. The result set contains the customertransactionid and amount of the transactions that meet the specified condition. These transactions may be considered as outliers or potential anomalies in the dataset.

```
Query  Query History
 1   SELECT
 2     ch.cc_num,concat(ch.firstname,' ',ch.lastname) as name,
 3     COUNT(*) AS transaction_count,
 4     MIN(transaction_date) AS first_transaction_date,
 5     MAX(transaction_date) AS last_transaction_date
 6   FROM
 7     customertransaction ct join cardholder ch on ct.cardholderid = ch.cardholderid
 8   GROUP BY
 9     ch.cc_num,name
10   HAVING
11     COUNT(*) > 100
12     AND (EXTRACT(DAY FROM MAX(transaction_date)) - EXTRACT(DAY FROM MIN(transaction_date))) < 30;
13
```

**Insight:** This query identifies cardholders who have made more than 100 transactions within a 30day period. It retrieves their credit card numbers, full names, total transaction counts, first transaction dates, and last transaction dates. The result set contains the information of cardholders who meet the specified conditions. These cardholders may be considered as active users or potential targets for further analysis based on their transaction activity.

```
Query  Query History
 1   SELECT
 2       ch.Gender,
 3       COUNT(*) AS FraudulentTransactions
 4   FROM
 5       CustomerTransaction ct
 6   JOIN
 7       Cardholder ch ON ct.CardholderID = ch.CardholderID
 8   WHERE
 9       ct.is_fraud = true
10   GROUP BY
11       ch.Gender;
12
```

**Insight:** This query groups fraudulent transactions by gender of the cardholder. It can help identify if there are any gender-specific patterns or vulnerabilities related to fraud, aiding in fraud prevention strategies.

```
Query    Query History
 1   SELECT
 2       m.Category,
 3       AVG(ct.Amount) AS AvgTransactionAmount
 4   FROM
 5       CustomerTransaction ct
 6   JOIN
 7       Merchant m ON ct.MerchantID = m.MerchantID
 8   GROUP BY
 9       m.Category
10   ORDER BY
11       AvgTransactionAmount DESC;
12
```

**Insight:** This query calculates the average transaction amount for each merchant category. It helps identify which categories generate the highest transaction values, guiding marketing strategies and business focus.

```
Query    Query History
 1   SELECT
 2       CASE
 3           WHEN DATE_PART('YEAR', NOW()) - DATE_PART('YEAR', ch.Dob) < 18 THEN 'Under 18'
 4           WHEN DATE_PART('YEAR', NOW()) - DATE_PART('YEAR', ch.Dob) BETWEEN 18 AND 30 THEN '18-30'
 5           WHEN DATE_PART('YEAR', NOW()) - DATE_PART('YEAR', ch.Dob) BETWEEN 31 AND 50 THEN '31-50'
 6           ELSE 'Over 50'
 7       END AS AgeGroup,
 8       COUNT(*) AS TotalTransactions
 9   FROM
10       CustomerTransaction ct
11   JOIN
12       Cardholder ch ON ct.CardholderID = ch.CardholderID
13   GROUP BY
14       AgeGroup
15   ORDER BY
16       AgeGroup;
17
```

**Insight:** This query groups transactions into different age groups based on the cardholder's date of birth. It can reveal the transaction behaviors of different age segments, aiding in targeted marketing and tailored product offerings.

# Data Visualization

Data visualization is all about transforming data and information into visual representations that are easy for humans to understand. It includes creating visual stories using charts, graphs, maps, or infographics by using colors, shapes, and patterns, data visualization simplifies complex data sets, making them more accessible and digestible for people. It helps us to see patterns, trends, and insights that might otherwise be hidden in rows and columns of numbers.



Here, the Sum of the transaction amount that has happen is in the different states is calculate and show in the figure. Here, size of the bubble gets bigger and bigger with the increase in the amount of transaction that has been occurred. By showing this, we have tried to show the regional trend where the most transaction could occur.

Sum of amount by state and category

category ● gas_transport ● grocery_pos ● shopping_pos

Here again we have used the total sum based on top categories where the fraud could take place.

We have taken top three in order to show what would it look like for multiple categories.



Sum of amt by category

In this graph we tend to show the trend of the possible fraud transaction of that could occur in different categories. Here, we could see the in groceries highest chances of occurring fraud transaction.

**Average of num_transactions by fraud_type**



1.72K (32.04%)    1.83K (34.13%)

**fraud_type**
- Significant Increase
- High Frequency
- High Amount

1.81K (33.82%)

It is used to represent data in a visually appealing and concise manner. Here, the donut chart is used to visualize the average number of transactions that has occurred in the specific fraud type. Here, we could average number of transactions is quite similar for the three types of fraud.

**Sum of averagetransactionamount by category**

Bar graph is used to visually display and compare the values of different categories or groups. Here, the graph shows that average amount of transaction that has been occurred in different categories.



Here, fraud transaction based on the population density is represented using seaborn and matplotlib libraries.



Here, KED plot of time is used in order to visualize the hour of days to find the possible fraud transactions.

# Hadoop

As data grow more and more sql takes more time to complete the queries. We need new way to make things easier and Hadoop comes in play here. Hadoop is a tool that helps process and store large amounts of data across multiple computers. It is great for handling big data because it splits the work among different machines. Hadoop has a special file system called HDFS that stores data across many computers, so even if one computer fails, the data is still safe. It also has a processing framework called MapReduce that breaks down tasks into smaller parts and runs them in parallel on different computers. This makes data processing faster and more efficient. Here, we have used to hadoop cluster to analyze the data in order to process and run mapreduce job for it.



Here, Hadoop cluster is introduced in order to find out the performed transaction is fraud or not using mapreduce in Hadoop cluster.

Here, Hadoop cluster is run in order to differentiate the categories using mapreduce. The results is presented in ascending order based on transaction count.

## PySpark Jobs

PySpark is a Python library that provides an interface for programming Apache Spark, a fast and distributed data processing engine. It integrates well with popular Python libraries such as NumPy, Pandas, and scikit-learn. As it has distributed computing capabilities.

```
pySparkJob.py > ...
  1    import pyspark
  2    import sys
  3    from pyspark.sql import SparkSession
  4    from pyspark.sql.functions import col, current_date, datediff, when
  5
  6    #SparkContext
  7    sc = pyspark.SparkContext()
  8    sqlContext = pyspark.sql.SQLContext(sc)
  9
 10    inputdir = sys.argv[1]
 11
 12    # Load the data into a DataFrame
 13    df = sqlContext.read.format('com.databricks.spark.csv').options(header='true', inferschema='true').load(inputdir+'transaction.csv')
 14
 15    # Step 3: Calculate the age of each cardholder and categorize them into age groups
 16    df = df.withColumn("dob", df["dob"].cast("date"))
 17    df = df.withColumn("age", datediff(current_date(), df["dob"]) / 365)
 18    df = df.withColumn(
 19        "AgeGroup",
 20        when(col("age") < 18, "Under 18")
 21        .when((col("age") >= 18) & (col("age") <= 30), "18-30")
 22        .when((col("age") >= 31) & (col("age") <= 50), "31-50")
 23        .otherwise("Over 50")
 24    )
 25
 26    # Step 4: Group by AgeGroup and calculate total transactions
 27    result_df = df.groupBy("AgeGroup").agg({"age": "count"}).withColumnRenamed("count(age)", "TotalTransactions")
 28
 29    # Display the final result
 30    print("AgeGroup with total transactions : ")
 31    result_df.show()
 32
```

```
+--------+-----------------+
|AgeGroup|TotalTransactions|
+--------+-----------------+
|   31-50|           543268|
| Over 50|           601171|
|   18-30|           152236|
+--------+-----------------+
```

This job categorizes customers into age groups, calculates total transactions per group, and displays the analysis results - demonstrating common PySpark DataFrame operations like casting, date functions, conditional logic, grouping, aggregating and renaming columns.

```
pySparkJob3.py > ...
  1    import pyspark
  2    import sys
  3    from pyspark.sql import SparkSession
  4    from pyspark.sql.window import Window
  5    from pyspark.sql import functions as F
  6    from pyspark.sql.functions import col, lag, udf
  7    from pyspark.sql.types import StringType
  8
  9    sc = pyspark.SparkContext()
 10    sqlContext = pyspark.sql.SQLContext(sc)
 11
 12    inputdir = sys.argv[1]
 13
 14    df = sqlContext.read.format('com.databricks.spark.csv').options(header='true', inferschema='true').load(inputdir+'transaction.csv')
 15
 16    windowSpec = Window.partitionBy("cc_num").orderBy("trans_date_trans_time")
 17
 18    df = df.withColumn("num_transactions", F.count("trans_num").over(windowSpec))
 19
 20    df = df.withColumn("avg_transaction_amount", F.avg("amt").over(windowSpec))
 21
 22    df = df.withColumn("std_dev_transaction_amount", F.stddev("amt").over(windowSpec))
 23
 24    df = df.withColumn("prev_transaction_amount", F.lag("amt", 1).over(windowSpec))
 25
```

```
26  def detect_fraud(amt, avg_transaction_amount, std_dev_transaction_amount, prev_transaction_amount, num_transactions):
27      if amt is None or avg_transaction_amount is None or std_dev_transaction_amount is None or prev_transaction_amount is None:
28          return 'None'
29
30      if amt > (avg_transaction_amount + (3 * std_dev_transaction_amount)):
31          return 'High Amount'
32      elif amt > (1.5 * prev_transaction_amount):
33          return 'Significant Increase'
34      elif num_transactions > 10:
35          return 'High Frequency'
36      else:
37          return 'None'
38
39  detect_fraud_udf = udf(detect_fraud, StringType())
40  df = df.withColumn("fraud_type", detect_fraud_udf("amt", "avg_transaction_amount", "std_dev_transaction_amount", "prev_transaction_amount", "num_tr
41
42  potential_fraud_df = df.filter(df['fraud_type'] != 'None')
43
44  selected_columns = ['trans_num', 'amt', 'cc_num', 'trans_date_trans_time', 'fraud_type']
45  selected_potential_fraud_df = potential_fraud_df.select(selected_columns)
46
47  selected_potential_fraud_df.show()
```

```
+--------------------+------+------------+---------------------+--------------------+
|           trans_num|   amt|      cc_num|trans_date_trans_time|          fraud_type|
+--------------------+------+------------+---------------------+--------------------+
|12554920d4920895a...| 52.67|501828204849|  2019-01-02 12:48:10|Significant Increase|
|4a52d2e61cb687153...| 31.79|501828204849|  2019-01-15 17:10:52|Significant Increase|
|e3d88c2eedd635a49...| 61.56|501828204849|  2019-01-22 11:38:05|Significant Increase|
|e75b1bb8e5ba6823e...| 50.91|501828204849|  2019-01-22 18:35:13|      High Frequency|
|1ed15a8619b0ceba8...| 64.03|501828204849|  2019-01-23 05:43:39|      High Frequency|
|367bc6d8b47753138...|118.58|501828204849|  2019-01-23 17:18:28|Significant Increase|
|66cf05aa6cb206f74 |199 27|501828204849|  2019 01 24 19:44:31|     High Frequency |
```

This PySpark code detects fraudulent transactions by calculating statistics like average, standard deviation and previous amount over each card using window functions. It then flags transactions based on rules like high amount, significant increase, or high frequency by applying a UDF. Finally, it filters and displays transactions with potential fraud.

Feature engineering is the process of transforming raw data into a format that is suitable for machine learning algorithms. It involves creating new features or modifying existing ones to improve the performance and effectiveness of machine learning models.

Here, we have implemented feature engineering and analysis in spark which give error on which hour of time and which age group is more targeted. The segmentation is done based on time where we also calculated displacement between customer and merchant.

## Conclusion

This project successfully developed a data model leveraging design patterns and integrated data analytics decision trees to build an efficient fraud detection system for the banking sector. The SQL queries provided valuable insights into transaction patterns. The implementation of this fraud detection system can significantly reduce financial losses due to fraudulent activities, while enhancing customer trust and security within the banking domain.

**Work division**

| Group Member | Task Performed |
|---|---|
| | |

| | |
|---|---|
| Robert Thapa | - Develop Hadoop cluster and run mapreduce job<br><br>- Develop database design and implementation<br><br>- Feature engineering and segmentation in pyspark<br><br>- Review the visualization performed by **Koshish** |
| Koshish Aryal | - Create visualizations using BI tool along with **Punam**<br>- Design conceptual and logical data model diagrams<br>- Review and help **Robert** feature engineering and segmentation |
| Saaz neupane | - Import and preprocess dataset in Python<br><br>- Write 2 complex SQL queries to analyze data<br><br>- Help and review **Shweta** to visualization and run Spark Jobs |
| Shweta Laljibhai Thummar | - Statistical analysis of data in SQL<br><br>- Create visualizations using Python matplotlib and seaborn<br><br>- Review and help **Saaz** for SQL queries and preprocessing |
| Punam Bhattarai | - Create Power BI reports and visualizations along with **Koshish**<br><br>- Research on pyspark jobs<br><br>- Review and help **Dipti** to prepare reports |
| Dipti Baral | - Statistical analysis of data in SQL along with **Shweta**<br><br>- Create visualizations using the on seaborn and matplotlib after spark job is run<br><br>- Review the task performed by **Punam**<br><br>- Prepare comprehensive project report and slides |
| Bipin Pandey | - Create presentation highlighting project architecture and key insights |

| | |
|---|---|
| | - Develop Spark job for business insights<br><br>- Review the final reports and help **Manan** with spark jobs |
| Manan Tushar Kapadia | - Verify the final presentations<br><br>- Develop Spark job for business insights along with **Bipin**<br><br>- Review the final reports done by **Dipti** |