

# Capstone Project

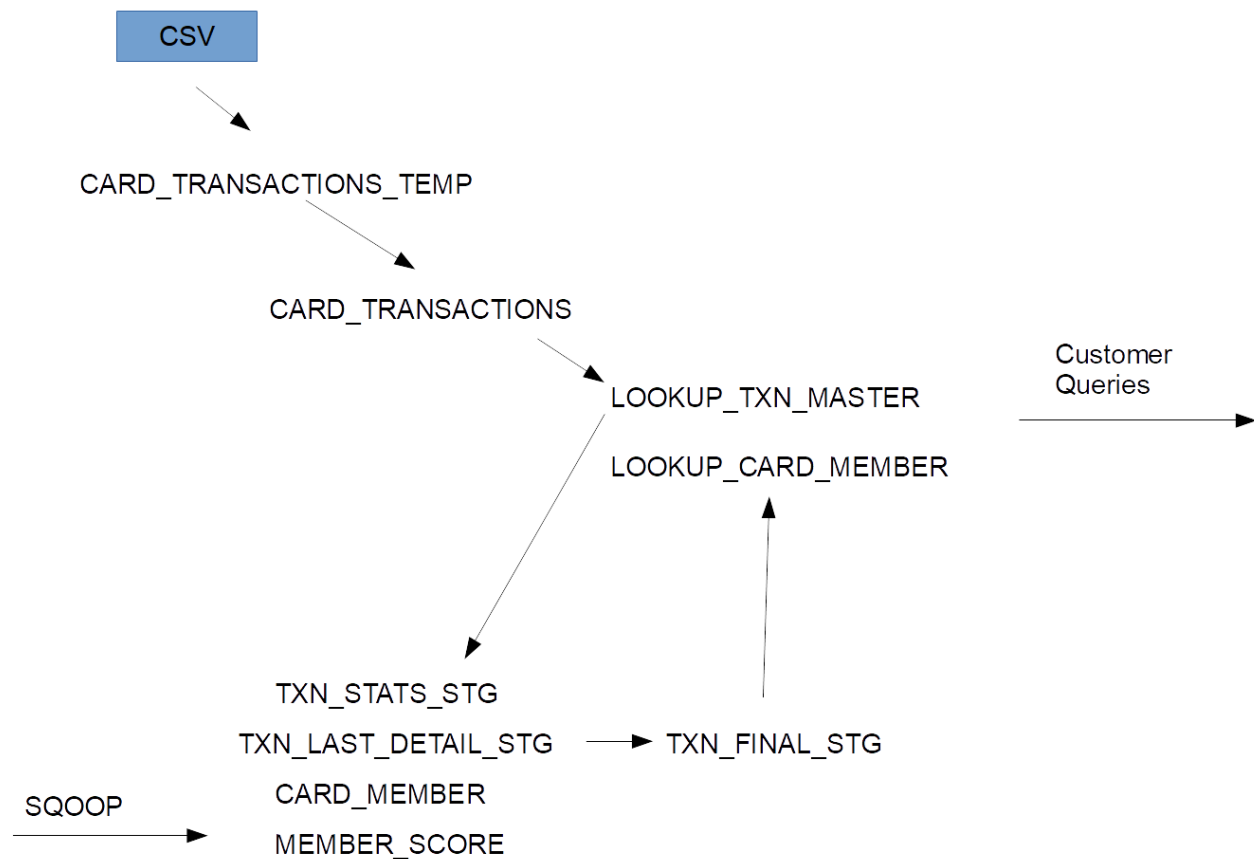
## Mid Submission

### Logic – Batch Layer Problem

By Venkatesh Jagannathan

# Table Design

Following are the tables needed with their purpose indicated: -



## LookUp

### LOOKUP\_CARD\_MEMBER(HBASE+HIVE)

card\_details:card\_id (key)  
last\_transaction\_details:ucl  
last\_transaction\_details:postcode  
last\_transaction\_details:transaction\_dt  
last\_transaction\_details:score  
member\_details:member\_id  
member\_details:member\_joining\_dt  
member\_details:card\_purchase\_dt  
member\_details:country  
member\_details:city

This is the primary lookup table to hold all card wise details like ucl of last 10 transactions, last transaction details & card member details. Keeping card and member details in same table makes it efficient for customer queries and they are 1:1 relation with the card id.

### **LOOKUP\_TXN\_MASTER**(HBASE+HIVE)

key(card\_id~timestamp)  
card\_details:card\_id  
member\_details:member\_id  
transaction\_details:amount  
transaction\_details:postcode  
transaction\_details:pos\_id  
transaction\_details:transaction\_dt  
transaction\_details:status  
transaction\_details:txrank (Ranking of transaction in descending order of time of occurrence)

This table holds last 10 “genuine” transactions along with ranking of each transaction to obtain the last transaction easily. This table along with LOOKUP\_CARD\_MEMBER becomes our lookup table for all the customer queries, efficiently without need to query transaction table.

### INTERMEDIATE TABLES (HIVE)

#### **TXN\_STATS\_STG**

card\_id  
moving\_average  
std\_dev

Stores statistics of each card id

#### **TXN\_LAST\_DETAIL\_STG**

card\_id  
postcode  
transaction\_dt

Stores last transaction detail

#### **TXN\_FINAL\_STG**

card\_id  
moving\_average  
std\_dev  
postcode  
transaction\_dt  
score  
member\_id  
member\_joining\_dt

card\_purchase\_dt  
country string  
city string

### **CARD\_TRANSACTIONS\_TEMP**

card\_id  
member\_id  
amount  
postcode  
pos\_id  
transaction\_d  
status

Intermediate table for CSV load

Stores consolidated information of statistics, last transaction details, score, card details & member score to be used as source to update lookup table back

The other tables are as per mentioned in assignment: -

CARD\_TRANSACTIONS to hold incoming transactions & initial csv upload.

CARD\_MEMBER to receive card member details ingested from RDS.

MEMBER\_SCORE to receive card member score data ingested from RDS.

## **Initial Load**

In this stage, the csv is loaded into CARD\_TRANSACTIONS with help of intermediate table to convert the transaction\_dt to timestamp(bigint). Also at this stage, top ten genuine transactions are loaded to LOOKUP\_TXN\_MASTER to facilitate processing. The top ten transactions are “**ranked**” using a rank column in descending order of occurrence (latest one on top) to help us identify the order of transaction later.

## **Data Import**

In this stage a stored scoop job is run to import data from amazon RDS database. Creating this as a sqoop job facilitates incremental imports without need to pass last value of check column manually to support increments. The import is implemented as “hive import” to directly import it as hive table. CARD\_MEMBER & MEMBER\_SCORE tables get created/updated by ingestion process. The password in assignment has been given on command line when job runs but an encrypted password file and sqoop’s **CryptoFileLoader** class will be used in production for ensuring security.

## **Calculation of UCL, Passing on last transaction details & score**

The moving average and standard deviation are calculated on last ten ranked transaction in LOOKUP\_TXN\_MASTER table which holds last ranked transactions. If the transactions for a card is less than 10, the average and standard deviation consideration the available number of records. Last transaction post code & transaction date details are extracted. These along with card member details and member score is combined to load LOOKUP\_CARD\_MEMBER hbase look up table.

## **Scheduled Job (Every 4 hours)**

The assignment solution contains output from a simple batch program that runs the hbase, hive and sqoop commands in a loop in the interval of every 4 hours. Oozie Scheduler is perfect choice for managing this workflow of different types of actions but could not be made possible in time due to infrastructure issues. The provided solution will work perfectly fine in a system if no infrastructure issues.