# Question – Project 1:

Stock Exchange Data Analysis
DESCRIPTION

**Objective**: To use hive features for data engineering or analysis and sharing the actionable insights

**Problem Statement:**
NewYork stock exchange data of seven years, between 2010 to 2016, is captured for 500+ listed companies. The data set comprises of intra-day prices and volume traded for each listed company. The data serves both for machine learning and exploratory analysis projects, to automate the trading process and to predict the next trading-day winners or losers.. The scope of this project is limited to exploratory data analysis.

**Domain**: BFSI

**Analysis to be done:** Exploratory analysis to understand how MoM or YoY companies from different sectors or industries and states have progressed in a period of 7 years

**Content:** This data set contains prices.csv and securities.csv files having the following features:

Prices.csv:

1. Date: Trading date
2. Symbol: Ticker code or listed company code on NY stock exchange
3. Open: Intra-day opening price for each listed company
4. Close: Intra-day closing price for each listed company
5. Low: Intra-day lowest price for each listed company
6. High: Intra-day highest price for each listed company
7. Volume: Number of shares traded per day per company

Securities.csv:

1. Ticker_Symbol: Country to which the customer belongs
2. Security: Legal name of the listed company
3. Sector: Business vertical of the listed company
4. Sub_Industry: Business domain of the listed company within a Sector.
5. Headquarter: Headquarters address

**Steps to perform:**

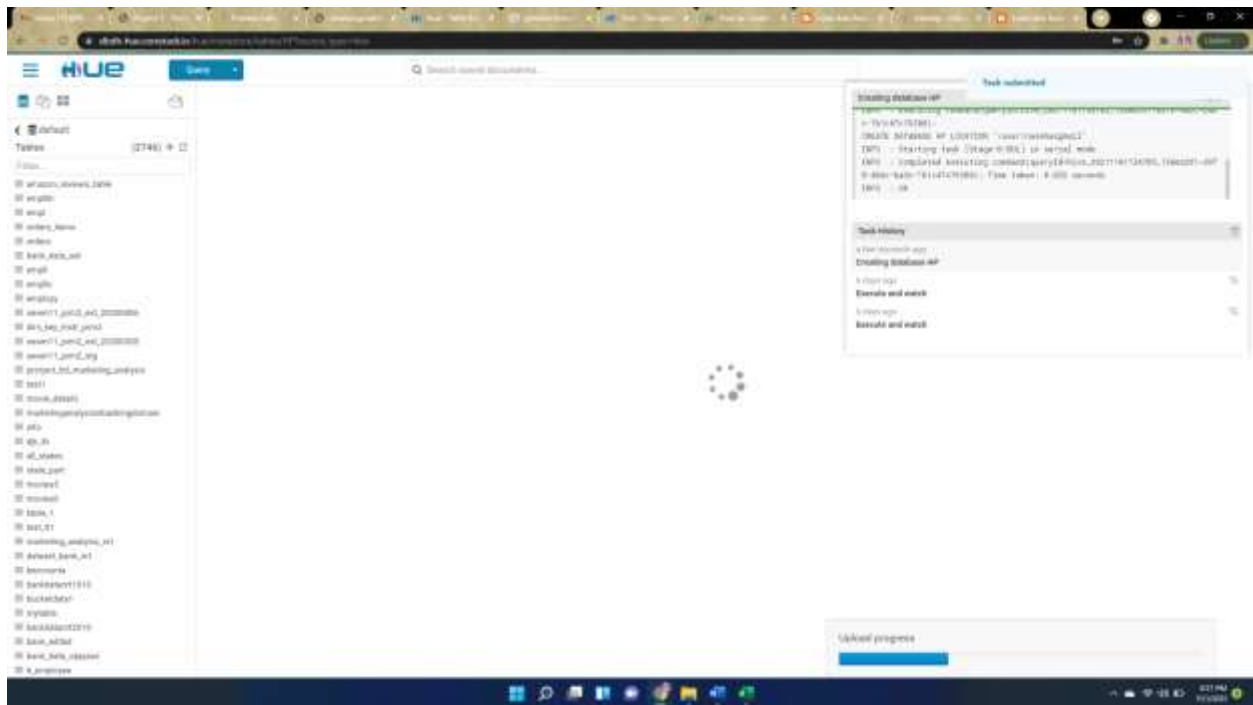   1) Create a data pipeline using sqoop to pull the data from the table below from MYSQL server into Hive.

      a. MYSQL DATABASE NAME: BDHS_PROJECT

            i. Stock_prices
            ii. Stock_companies

Check the TABLE description: STOCK_PRICES

| Column Name |
| --- |
| Trading_date |
| Symbol |

| Open |
| --- |

| Close |
| --- |

| Low |
| --- |

| High |
| --- |

| Volume |
| --- |

TABLE: STOCK_COMPANIES

| Column Name |
| --- |

| Symbol |
| --- |

| Company_name |
| --- |

| Sector |
| --- |

| Sub_industry |
| --- |

| Headquarter |
| --- |

2) Create a new hive table with the following fields by joining the above two hive tables. Please use appropriate Hive built-in functions for columns (a,b,e and h to l).

- Trading_year: Should contain YYYY for each record
- Trading_month: Should contain MM or MMM for each record
- Symbol: Ticker code
- CompanyName: Legal name of the listed company
- State: State to be extracted from headquarters value.
- Sector: Business vertical of the listed company
- Sub_Industry: Business domain of the listed company within a sector
- Open: Average of intra-day opening price by month and year for each listed company
- Close: Average of intra-day closing price by month and year for each listed company
- Low: Average of intra-day lowest price by month and year for each listed company
- High: Average of intra-day highest price by month and year for each listed company

- Volume: Average of number of shares traded by month and year for each listed company

**DATA ANALYSIS USING HIVE**

3) Find the top five companies that are good for investment

4) Show the best-growing industry by each state, having at least two or more industries mapped.

5) For each sector find the following.

- Worst year
- b. Best year
- c. Stable year

# Write Up & Screenshots

Log in using Web Console:

**MySQL Shell Command Structure**

Username: nasehasgmail

Password: nasehasgmailonz2c

At prompt:

1) Create a data pipeline using sqoop to pull the data from the table below from MYSQL server into Hive.

mysql -h sqoopdb.slbdh.cloudlabs.com -u nasehasgmail -pnasehasgmailonz2c

Show databases;

Use nasehasgmail;

Describe STOCK_COMPANIES;

Describe STOCK_PRICES;

MYSQL> Exit

```
[nasehasgmail@ip-10-0-41-79 ~]$ hadoop fs -ls
/user/hive/warehouse/hivedb_naseha.db/STOCK_* /
```



HIVE

Used Hue Interface to do the rest of the questions

Creating a Database HP

File Section -> Tables -> Create Database

Click on New -> Name the Database



Here a database is created under my user id: Home/user/nasehasgmail

Then we upload the raw files (.csv) into the database

Upload both the .csv

2) Create a new hive table with the following fields by joining the above two hive tables. Please use appropriate Hive built-in functions for columns (a,b,e and h to l).

Create table HP as SELECT

Date_format(stock_prices.trade_date, 'yyyy') as TradeYear,

date_format(stock_prices.trade_date, 'MM') as TradeMonth,

stock_prices.symbol as Symbol,

stock_companies.security as CompanyName,

split(stock_companies.HQ, '\u0059')[1] as State, stock_companies.sector as Sector,

stock_companies.sub_industry as sub_Industry,

avg(stock_prices.open) as Open,

avg(stock_prices.close) as Close,

avg(stock_prices.high) as Hi,

avg(stock_prices.low) as Lo,

avg(stock_prices.volume) as Vol

from stock_prices INNER JOIN stock_companies

ON (stock_prices.symbol = stock_companies.symbol)

GROUP BY date_format(stock_prices.trade_date, 'yyyy'),  date_format(stock_prices.trade_date, 'MM'),

stock_companies.sub_Industry, stock_companies.security, stock_prices.symbol,
stock_companies.sector, split(stock_companies.HQ, '\u0059')[1];



## 3) Find the top five companies that are good for investment

Select companyname, avg((close-open)/(open*100)) as returns FROM hp

GROUP BY companyname

ORDER BY returns DESC limit 5;

Describe formatted STOCK_PRICES;

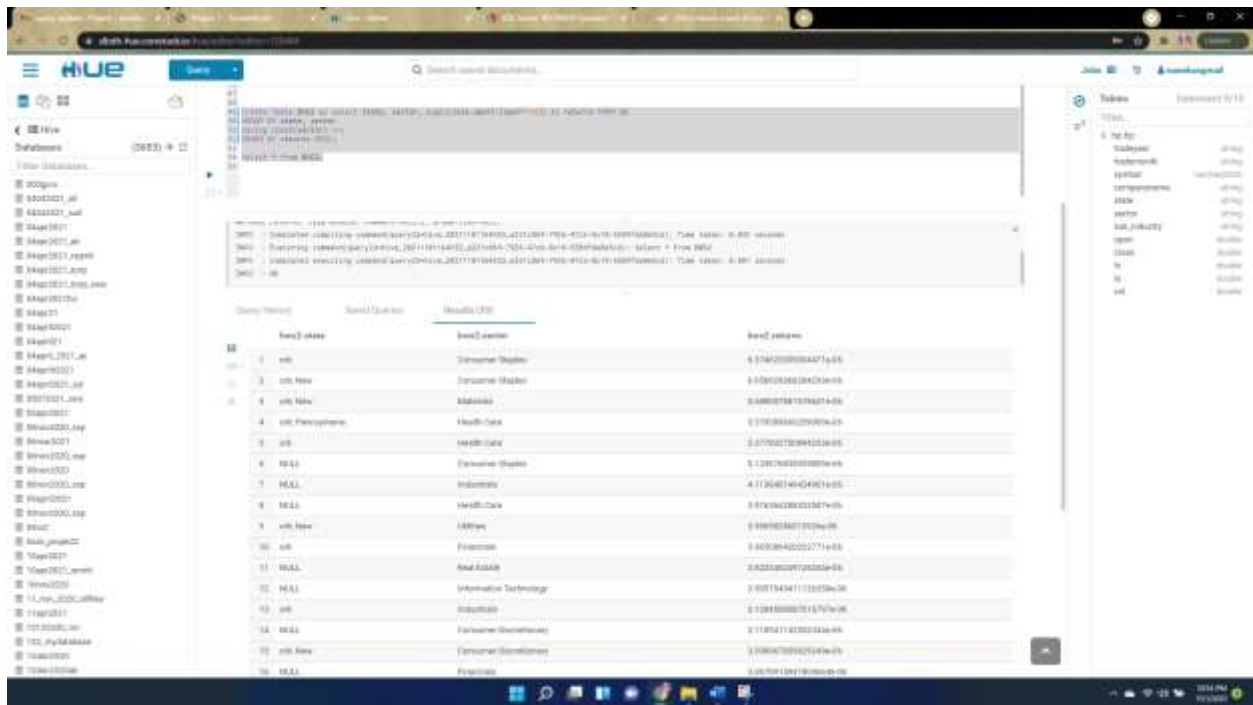## 4) Show the best-growing industry by each state, having at least two or more industries mapped

create table BWS2 as select state, sector, avg((close-open)/(open*100)) as returns FROM hp

GROUP BY state, sector

having count(sector) >=2

ORDER BY returns DESC;


Select * from BWS2;

## 5) For each sector find the following – worst, best and stable year;

To find the best and the worst years, I am creating another table with only sector, tradeyear and avg of the returns as columns.

create table BWS as select sector, tradeyear, avg((close-open)/(open*100)) as returns FROM hp

GROUP BY sector, tradeyear

ORDER BY sector, returns DESC;

select * from BWS;



Once the table is ready, using the rank function to find the best and the worst year, and cume_dist() = .50 for the median. Median is assumed to be the most stable with min variance.

## a. Worst Year

select * from (

        select sector, tradeyear, returns,

        rank() over ( partition by sector order by returns asc) as rank

        from bws ) t where rank = 1;

## b. Best Year

```
select * from (

        select sector, tradeyear, returns,

        rank() over ( partition by sector order by returns desc) as rank

        from bws ) t where rank = 1;
```

## c. Stable year

create table M as select * from (

      select sector, tradeyear, returns,

      cume_dist() over ( partition by sector order by returns desc) as CD
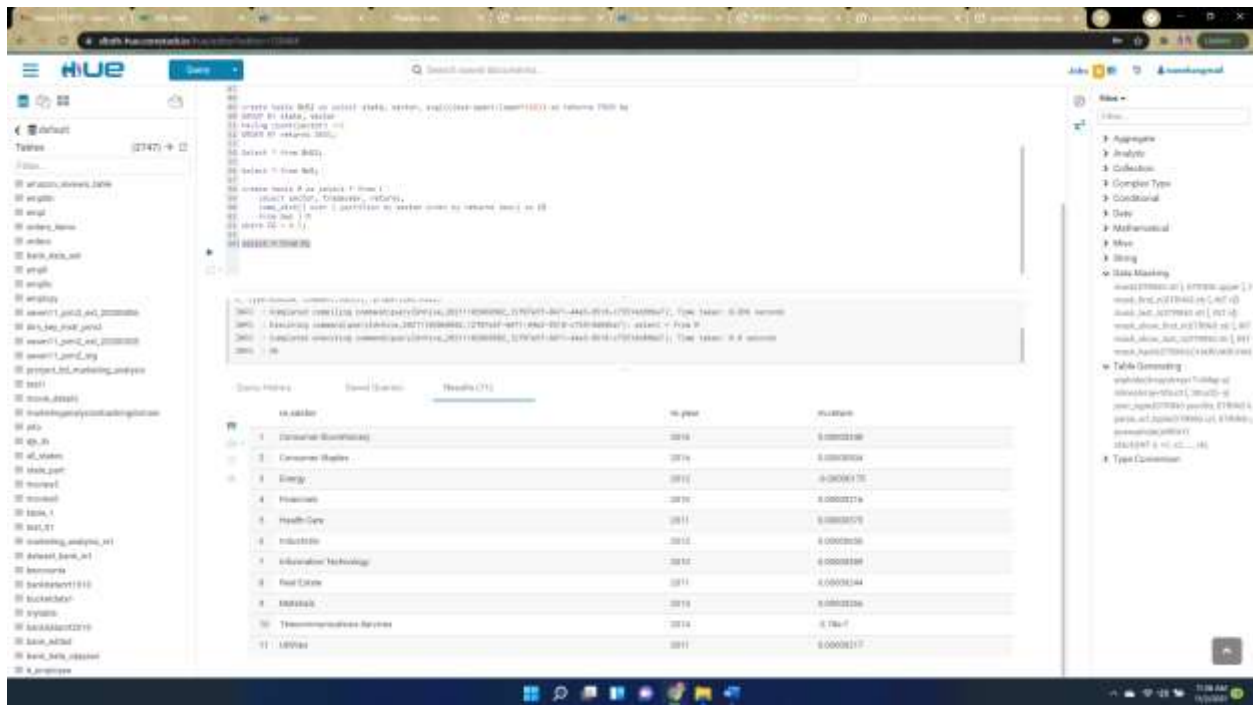
      from bws ) M

where CD = 0.5;


select * from M;

## Code:

Create table HP as SELECT

Date_format(stock_prices.trade_date, 'yyyy') as TradeYear,

date_format(stock_prices.trade_date, 'MM') as TradeMonth,

stock_prices.symbol as Symbol,

stock_companies.security as CompanyName,

split(stock_companies.HQ, '\u0059')[1] as State, stock_companies.sector as Sector,

stock_companies.sub_industry as sub_Industry,

avg(stock_prices.open) as Open,

avg(stock_prices.close) as Close,

avg(stock_prices.high) as Hi,

avg(stock_prices.low) as Lo,

avg(stock_prices.volume) as Vol

from stock_prices INNER JOIN stock_companies

```sql
ON (stock_prices.symbol = stock_companies.symbol)

GROUP BY date_format(stock_prices.trade_date, 'yyyy'),  date_format(stock_prices.trade_date, 'MM'),

stock_companies.sub_Industry, stock_companies.security, stock_prices.symbol,
stock_companies.sector, split(stock_companies.HQ, '\u0059')[1];




SElect * from HP limit 10;


Select companyname, avg((close-open)/(open*100)) as returns FROM hp

GROUP BY companyname

ORDER BY returns DESC limit 5;


select state, sector, avg((close-open)/(open*100)) as returns FROM hp

GROUP BY state, sector

having count(sector >2)

ORDER BY returns DESC;


create table BWS as select sector, tradeyear, avg((close-open)/(open*100)) as returns FROM hp

GROUP BY sector, tradeyear

ORDER BY sector, returns DESC;


select * from BWS;


select * from (

        select sector, tradeyear, returns,

        rank() over ( partition by sector order by returns desc) as rank

        from bws ) t where rank = 1;




select * from (
```

```
        select sector, tradeyear, returns,

        rank() over ( partition by sector order by returns asc) as rank

        from bws ) t where rank = 1;




create table BWS2 as select state, sector, avg((close-open)/(open*100)) as returns FROM hp

GROUP BY state, sector

having count(sector) >=2

ORDER BY returns DESC;


Select * from BWS2;


Select * from BWS;


create table M as select * from (

        select sector, tradeyear, returns,

        cume_dist() over ( partition by sector order by returns desc) as CD

        from bws ) M

where CD = 0.5;


select * from M;
```