# CS750 DISTRIBUTED DATA MANAGEMENT
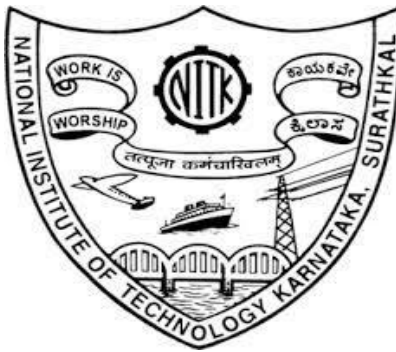
**A REPORT ON THE PROJECT ENTITLED**

**"Analysis of Stock Market Using Hadoop and Map Reduce"**



**SUBMITTED BY**

**Milind Chauhan - 212CS013**

**Sandeep Talapally - 212CS033**

**II SEMESTER M-TECH CSE**

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**

**NATIONAL INSTITUTE OF TECHNOLOGY KARNATAKA, SURATHKAL**

**2021-2022**

**INDEX**

# 1.    ABSTRACT

Stock Market has huge high risk and high profit. That is why its analysis is very important. The Stock Market contains huge amounts of data measuring in GBs or TBs. These are very complex and should be analyzed by data mining methods. Large computations are more costly and time consuming. Distributed computing provides cost effective and scalable solutions to such problems.

# 2.    INTRODUCTION

Volatility is a statistical measure of the dispersion of returns for a given security or market index. In most cases, the higher the volatility, the riskier the security. Volatility is often measured as either the standard deviation or variance between returns from that same security or market index.

To compute volatility over large data using normal computing takes a lot of time. In order to provide trade off between speed and cost of computing one can use distributed computing and solve the issue using map reduce and hadoop framework.

In the securities markets, volatility is often associated with big swings in either direction. For example, when the stock market rises and falls more than one percent over a sustained period of time, it is called a "volatile" market. An asset's volatility is a key factor when pricing options contracts.  This way one can be cautious about a certain stock using this volatility.

Following are the terms that are used in relations of the dataset of an example stock AAPL.

Date represents the date of the stock AAPL
Open represents the open price in that day of stock AAPL
High represents the highest price in that day of stock AAPL
Low represents the lowest price in that day of stock AAPL
Adj Close represents the close price in that day of stock AAPL
Volume represents the volume in that day of stock AAPL

# 3. EXISTING SYSTEM

RDBMS stores data in a structural way.  Scaling up (upgrading) is very expensive. Basic data unit is relational tables. Sequential or traditional computing takes a lot of time. It adds more cost and time.

Solution of scaling out and trade off cannot be made easily using traditional computing. One cheap processing batch unit costs around Rs. 30,000 whereas a standard machine with good power costs more than 5 times of cheap processing unit. These cheap processing units can be repaired and even if lost it won't add much expense to management.

One can even easily set up a hadoop batch with a set of 3-4 laptops and perform distributed computing which was very tough with traditional approaches.

# 3. PROPOSED SYSTEM

This project uses Mapreduce on a Hadoop environment to compute the monthly volatility of stocks. The database is of 2970 stocks on NASDAQ market for 3 years from 01/01/2012 to 12/31/2014 (except holidays, otherwise called trading days). Our job is to analyze stock price data, and find out which stocks in a certain period have higher earnings potential, etc. One characteristic that is widely used by traders is the volatility index. Our data is 2970 CSV format files (Comma Separated). Each file contains the data for one stock using its symbol as the file name.

# 4. IMPLEMENTATION DETAILS

This project uses Mapreduce paradigm of Hadoop to serialize the calculation of volatility for each month and computes the top 10 most volatile and bottom 10 least volatile stock values.

- Number of Mapper Implementation : 3
- Number of Reducer Implementation : 3

Roles of each Mapper and Reducer:

## **Mapper1**

• splits the input data and options the date and close adjacent value.

• key - stock_name + month + year

• value - date + adjacent close value

## **Reducer1**

• Since after the map step the values which have the same key are grouped together and passed to the reducer as iterable, values that correspond to specific month and year of the particular stock are grouped together.

• Beginning adjacent close value and end adjacent close value are obtained by integrating through the iterable and the value of xi for the corresponding month is computed.

• Key - Company Name

- Value - Computed $X_i$

**Mapper2**

- Now we have to consolidate all the values obtained from the reducer with respect to the company name.

- Key - Company Name

- Value - $X_i$

**Reducer2**

- All the xi corresponding to the the respective companies are grouped together.

- Volatility for the particular company is obtained from these values.

- Key - Company Name

- Value - Volatility

**Mapper3**:

- All the companies are grouped together with a common key.

- Key - Common

- Value = Company Name + Volatility

## Reducer3:

• Obtained iterable contains all the company name with values and they are sorted by a custom comparator.

• Top 10 and bottom 10 values are obtained from the List

# 5. ADVANTAGES

a) Scaling Out:   In Traditional RDBMS it is quite difficult to add more hardware, software resources i.e. scale up. In Hadoop this can be easily done i.e. scale down.

b) Transfer code to data In RDBMS generally data is moved to code and results are stored back. As data  is moving there  is always a  security threat.  In Hadoop small code is moved  to data  and it  is executed there itself. Thus data is local. Thus Hadoop correlates preprocessors and storage.

 c) Fault Tolerance: Hadoop is designed to cope up with node failures. As a large number of machines are there, a node failure is a very common problem.

d) Abstraction of complexities  Hadoop provides proper interfaces between components for proper working. e) Data protection and consistency  Hadoop handles system level challenges as it supports data consistency.

 a) Low cost  As Hadoop  is an  open-source framework,  it is  free. It uses  commodity hardware  to store and process huge data. Hence not much costly.

b) High Computing power: Hadoop uses a distributed computing model. Due to this, tasks can be distributed amongst different nodes and can be processed quickly. Cluster has thousands of nodes which gives high computing capability to Hadoop.

# 6. LIMITATIONS

This approach may lead to the issue of HDFS error due to clustering. This project formulated the overall problem as an interaction problem between namenode and data node and solved it by creating a temp directory for data node.

Apache Hadoop is for batch processing, which means it takes a huge amount of data in input, processes it and produces the result. Although batch processing is very efficient for processing a high volume of data, depending on the size of the data that processes and the computational power of the system, an output can delay significantly. Hadoop is not suitable for Real-time data processing.

This can be solved by Apache Spark and Apache Flink.

Hadoop is not so efficient for iterative processing, as Hadoop does not support cyclic data flow.

Hadoop is challenging in managing complex applications. If the user doesn't know how to enable a platform who is managing the platform, your data can be a huge risk. At storage and network levels, Hadoop is missing encryption, which is a major point of concern.

This can be solved by, third-party vendors have enabled an organization to leverage Active Directory Kerberos and LDAP for authentication.

In Hadoop, MapReduce developers need to hand code for each and every operation which makes it very difficult to work

# 7. FUTURE WORK

This project can be scaled to work in conjunction with Machine Learning Algorithms that can use these kinds of Map Reduce applications to handle large data sets with multiple Master nodes.

# 8. CONCLUSION

This project proposes a Triple Map Reduce Model that selects data from different sources with the help of Hadoop Distributed File System (HDFS) such that the stock price of volatility can be used for either minimum 10 or maximum 10.

Results:



fig: commands to put our input on hadoop distributed file system and run the program

fig: Map Reduce jobs running in parallel manner facilitated by the distributed file system

fig: Completion of All tasks



fig: Commands to get the output back from the distributed file system onto our own file system

```
copyFromHadoop2 >  ≡ part-r-00000
  1    Top 10 stocks with Minimum volatility    0.0
  2    AGZD    0.003938593878697365
  3    AXPWW   0.0044388372955839524
  4    AUMAU   0.006017144863314729
  5    AGND    0.010751963436794309
  6    AGNCP   0.01669670030720715
  7    AGNCB   0.016781408595782567
  8    ALLB    0.021866756279518028
  9    AGIIL   0.022955571847192706
 10    ASRVP   0.028529779716934052
 11    ACNB    0.028565410375761102
 12    Top 10 stocks with Maximum volatility    1.0
 13    APDN    0.3773818041663614
 14    ALDR    0.39064070974779724
 15    ANY 0.4118627840513947
 16    AFMD    0.41919685205354573
 17    AMCF    0.4279202890844407
 18    ATRA    0.42898449574226183
 19    ASPX    0.43506357182854893
 20    ADXS    0.4411702287863926
 21    APDNW   0.6975880360551902
 22    ACST    9.271589761859984
 23
```

fig: final result