# Big Data Analytics Innovative Assignment Report



**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING**

Presented to: Prof. Aparna Kumari
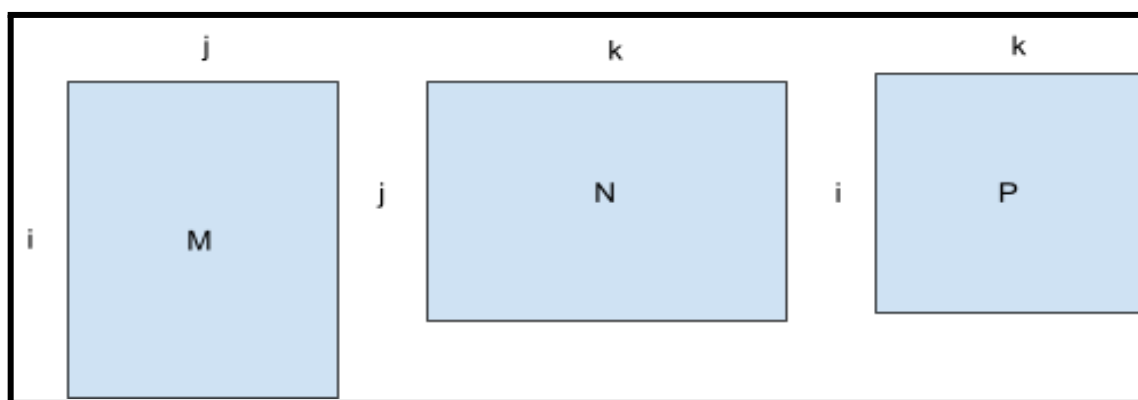
Presented by:     Dip Patel (19BCE166)
                  Jalpan Patel(19BCE177)
                  Parv Patel(19BCE190)

# Sparse Matrix Multiplication with Hadoop

## INTRODUCTION

Matrix multiplication is a fundamental operation in linear algebra with related real-life applications, such as matrix factorization, chemical system formulation, and graph analysis [14]. In addition to its naturally related applications, several problems are reducible by matrix multiplication. Thus, these problems should be investigated thoroughly to enhance the efficiency of implemented algorithms for matrix multiplication. Given the inputs of two matrices A and B, where the number of columns in A equals the number of rows in B, matrix multiplication produces matrix C with the number of rows equal to that in A and number of columns equal to that in B. The Brute-Force matrix multiplication algorithm for square matrices is given in Algorithm 1. The Brute-Force algorithm has a high processing complexity of O(n3), but suffers from the massive memory lookup process required to locate each array element for multiplication. Over the years, several matrix multiplication algorithms have been proposed to reduce the cost and time of the matrix multiplication process [2, 15].

Matrix Multiplication is when we have a **"p x q"** matrix **M**, whose element in **row i** and **column j** will be denoted $m_{ij}$ and a **"q x r"** matrix **N** whose element in **row j** and **column k** is denoted by $n_{jk}$ then the product **P = MN** will be **"p x r"** matrix **P** whose element in **row i** and **column k** will be donated by $P_{ik}$, where $P(i,k) = m_{ij}*n_{jk}$



Most matrices are sparse so large numbers of cells have value zero. When we represent matrices in this form, we do not need to keep entries for the cells that have values of zero to

save a large amount of disk space. As input data files, we store matrix M and N on HDFS in the following format:

## MAP-REDUCE

The map and reduce functions have been built using the following Algorithms:

**Algorithm 1:** The Map Function

1 **for** *each element $m_{ij}$ of M* **do**
2     produce $(key, value)$ pairs as $((i, k), (M, j, m_{ij}))$, for $k = 1, 2, 3, ..$ up to the number of columns of $N$
3 **for** *each element $n_{jk}$ of N* **do**
4     produce $(key, value)$ pairs as $((i, k), (N, j, n_{jk}))$, for $i = 1, 2, 3, ...$ up to the number of rows of $M$
5 **return** *Set of (key, value) pairs that each key, $(i, k)$, has a list with values $(M, j, m_{ij})$ and $(N, j, n_{jk})$ for all possible values of $j$*

---

**Algorithm 2:** The Reduce Function

1 **for** *each key (i,k)* **do**
2     sort values begin with $M$ by $j$ in $list_M$
3     sort values begin with $N$ by $j$ in $list_N$
4     multiply $m_{ij}$ and $n_{jk}$ for $j_{th}$ value of each list
5     sum up $m_{ij} * n_{jk}$
6 **return** $(i, k), \sum_{j=1} m_{ij} * n_{jk}$

Map function will produce key,value pairs from the input data as it is described in Algorithm 1. Reduce function uses the output of the Map function and performs the calculations and produces key,value pairs as described in Algorithm 2. All outputs are written to HDFS.

1. **MAP Task:**
   a. **For matrix M**: Map task (Algorithm 1) will produce key,value pairs as

   follows: (i,k), (M, j, $m_{ij}$)

   $m_{11}$ = 1
   (1,1), (M, 1, 1) for k = 1
   (1,2), (M, 1, 1) for k = 2

   $m_{12}$ = 1
   (1,1), (M, 2, 1) for k = 1
   (1,2), (M, 2, 1) for k = 2

$m_{13} = 1$
(1,1), (M, 3, 1) for k = 1
(1,2), (M, 3, 1) for k = 2

$m_{21} = 1$
(2,1), (M, 1, 1) for k = 1
(2,2), (M, 1, 1) for k = 2

$m_{22} = 1$
(2,1), (M, 2, 1) for k = 1
(2,2), (M, 2, 1) for k = 2

$m_{23} = 1$
(2,1), (M, 3, 1) for k = 1
(2,2), (M, 3, 1) for k = 2

b. **For Matrix N:** Map task (Algorithm 1) will produce key, value pairs as

follows: (i,k), (N, j, $n_{jk}$)

$n_{11} = 1$
(1,1), (N, 1, 1) for i = 1
(2,1), (N, 1, 1) for i = 2

$n_{21} = 1$
(1,1), (N, 2, 1) for i = 1
(2,1), (N, 2, 1) for i = 2

$n_{31} = 1$
(1,1), (N, 3, 1) for i = 1
(2,1), (N, 3, 1) for i = 2

$n_{12} = 2$
(1,2), (N, 1, 1) for i = 1
(2,2), (N, 1, 1) for i = 2

$n_{22} = 1$
(1,2), (N, 2, 1) for i = 1
(2,2), (N, 2, 1) for i = 2

$n_{32} = 1$
(1,2), (N, 3, 1) for i = 1
(2,2), (N, 3, 1) for i = 2

   **c. After combine operation the map task will return key, value pairs will look like as follows:**

( (i,k), [ (M, j, mij ), (M, j, mij ),..., (N, j, njk), (N, j, njk),   ] )

(
(1,1), [ (M, 1, 1), (M, 2, 1), (M, 3, 1), (N, 1, 1), (N, 2, 1), (N, 3, 1) ]
(1,2), [ (M, 1, 1), (M, 2, 1), (M, 3, 1), (N, 1, 1), (N, 2, 1), (N, 3, 1) ]
(2,1), [ (M, 1, 1), (M, 2, 1), (M, 3, 1), (N, 1, 1), (N, 2, 1), (N, 3, 1) ]
(2,1), [ (M, 1, 1), (M, 2, 1), (M, 3, 1), (N, 1, 1), (N, 2, 1), (N, 3, 1) ]
)

Note that the entries for the same key are grouped in the same list, which is performed by the framework. This output will be stored in HDFS and is fed to the reducer task as input.

2. **REDUCER Task:**

Reduce task takes the key,value pairs as the input and processes one key at a time. For each key it divides the values in two separate lists for M and N. As an example, it will create the following list for key (1,1), [ (M, 1, 1), (M, 2, 1), (M, 3, 1), (N, 1, 1), (N, 2, 1), (N, 3, 1) ]
$I_m$ = [ (M, 1, 1), (M, 2, 1), (M, 3, 1) ]
$I_n$ = [ (N, 1, 1), (N, 2, 1), (N, 3, 1) ]
then sums up the multiplication of $m_{ij}$ and $n_{jk}$ for each j as follows:
P(1,1) = 1 + 1 + 1 = 3
P(1,2) = 1 + 1 + 1 = 3
P(2,1) = 1 + 1 + 1 = 3
P(2,2) = 1 + 1 + 1 = 3

Upon completing all calculations we will get the following output:

# Hadoop

Hadoop is an Apache open-source platform for storing, processing, and analyzing massive amounts of data. Hadoop is a Java-based data warehouse that is not an OLAP system (online analytical processing). It's a batch/offline processing system. It could also be upscaled by adding extra nodes to the cluster.

Modules of Hadoop

- HDFS: Hadoop Distributed File System. Google published a document called GFS, and HDFS was developed on the basis of it. The data will be divided down into blocks and stored in nodes as part of the distributed design, according to the document.
- Yarn: Yet another Resource Negotiator is utilized to control the cluster and schedule jobs.
- Map Reduce: This is a framework that allows Java programs to perform parallel computing on data employing key-value pairs. The Map task turns input data into a data set that can be evaluated in Key-value pairs. The result of the Map activity is absorbed by the Reduce task, which subsequently outputs the required result.
- Hadoop Common: These Java libraries are utilized by several other Hadoop modules and can be used to launch Hadoop

## Specifications

| Title | Description |
|---|---|
| Data specification | Matrices of various sizes are tested with dimensions in a range up to 10^4. |
| Configuration of System | **RAM -** 8GB<br>**Processor -** Intel Core i5 8th Gen |
| Cluster setup | **Single node -** One DataNode running and setting up all the NameNode, DataNode, ResourceManager, and NodeManager on a single machine. |
| Parameters | **Time comparison -** So in our comparison when we ran the MP with 1 reducer the task of all reducers was completed in 69052 ms, but when we increased the number of reducers to 5 reducers the task of all reducers was completed in 152296 ms.<br><br>**Sparsity -** Sparsity of 1000*1000 input was .3 which means that around 30% of it was filled with non-zero values, the rest was filled with 0s. |

# Input-Output





Changed output

```
 C:\ Administrator: Command Prompt
cat: `/jdpoutput5/part-r-00000.txt': No such file or directory

C:\WINDOWS\system32>hadoop fs -cat /jdpoutput5/part-r-00000
2022-11-25 07:02:43,231 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostT
0,0,3.0
0,1,3.0
0,2,3.0
1,0,3.0
1,1,3.0
1,2,3.0
2,0,3.0
2,1,3.0
2,2,3.0

C:\WINDOWS\system32>hadoop fs -cat /jdpinput/M.txt
2022-11-25 07:03:08,389 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostT
M,0,0,1
M,0,1,1
M,0,2,1
M,1,0,1
M,1,1,1
M,1,2,1
M,2,0,1
M,2,1,1
M,2,2,1
C:\WINDOWS\system32>hadoop fs -cat /jdpinput/N.txt
2022-11-25 07:03:39,558 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostT
N,0,0,1
N,0,1,1
N,0,2,1
N,1,0,1
N,1,1,1
N,1,2,1
N,2,0,1
N,2,1,1
N,2,2,1
C:\WINDOWS\system32>
```

## Applications:

1) Perspective projections, which are the core of 3D animation, are created using matrix multiplication. A two-dimensional image depicting a position in 3-dimensional space is displayed on the computer display. You may switch across two and three-dimensional worlds using matrix multiplication.

2) It's frequently utilized in fields like network theory, linear system of equations, coordinate system transformation, and population modeling, to list some of them.

3) Matrix multiplication was first established in linear algebra to make calculations easier and more clear. Matrix multiplication and linear algebra have a close link that is vital in all of the math, along with chemistry, physics, engineering, and computer science.

4) In graph theory, matrix multiplication — especially, powers of a given matrix A — is a valuable tool when the matrix in hand is the adjacency matrix of a graph.

5) Matrix multiplication is a fundamental concept in quantum mechanics and physics in general. The moment of inertia tensor, Hamiltonians-based continuous-time representations of the development of physical systems, and the most general formulation of the Lorentz transformation from special relativity are some examples.

6) In mathematical finance, matrices are widely employed in a number of ways. A

correlation matrix, for example, is a table in which an item (i,j) indicates the degree to which price movements in item I and item j are associated during a given timeframe. Every day, a large number of computer cycles are expended calculating and analyzing such types of matrices in order to, in part, estimate the risk involved with such a basket of items.

7) Any linear system of equations, i.e. a plane in certain N-dimensional space, can be described by a matrix. All current physics and simulations require multiplying one plane with another plane.

    a) If you must determine the gravitational force across a path, the path matrix is multiplied by the gravity matrix.

    b) If the intensity of a magnetic field (Matrix a) at a spatial position described by matrix B must be determined, the magnetic field is multiplied by the geographical matrix.

    c) If you want to locate the best location in town for a firehouse, an As matrix A, you create a block-by-block fire hazard matrix, and as matrix B, you create a trip time matrix centered on a given place in town.

    d) To create a self-driving automobile employing neural networks in the computer vision system, rely on GPUs to interpret images employing real-time matrix multiplication.

## Challenges

Three issues unique to the sparse scenario are identified.
- The first is a problem with processor load imbalance. In the dense state, when W = N3, every processor produces a net of W/p work, but in the sparse state, such a proportion is difficult to attain. Additionally, dense GEMM algorithms are often executed in s stages, with every processor performing W/p s work per stage. It's significantly more difficult to achieve appropriate load balance per stage for sparse matrices than it is to achieve load balance for the entire operation. Considering matrices with complete diagonal, asynchronous techniques may be useful in alleviating the load balancing issue.
- The addition of submatrices poses the second problem. For submatrix additions, no additional operations are done in the dense case; the same amount of work would be required in the sequential case.
- The third difficulty is to keep the communication hidden. While this is simpler in the dense case because of the higher computation to communication proportion, it is more difficult in the sparse scenario. To reduce the communication to computation ratio of sparse algorithms, simply increasing the problem size is usually insufficient.

## Result/Observations

- Matrix of large dimensions

  On taking 1000*100 size of M, 100*1000 size of N we received 1000*1000 size of P, the time taken for execution was a couple of minutes.

- The number of Reducer increased by 5





So, it is usually preferred to keep the number of reducers such that it is:
  - It is a multiple of the block size
  - It does not take a long time to execute

        ○ Doesn't require a lot of file creations

➢       But we also know that the number of reducers can't exceed the number of partitions, so we also need to take that into consideration when deciding the number of reducers. In the instance, we make the number of reducers 10 but we have a single partition system then the output will be divided into 10 shards but those 10 outputs will be consolidated by a single reducer.

➢ The benefits of having more reducers are as follows:
        ○ Increases load balancing
        ○ Lowers the cost of failure
➢ But on the other hand, having too many reducers can also:
        ○ Increase the framework overhead
        ○ Slower start-up time
        ○ Increases the number of input for the next tasks

## Conclusion

The Hadoop framework's MapReduce component is a critical processing component. It's a rapid, scalable, and cost-effective tool that may assist data analysts and developers in processing large amounts of information.

This programming approach is useful for assessing website and e-commerce platform usage trends. This framework can be used by companies that provide online services to strengthen their marketing tactics.

Total vcore-milliseconds taken by all map tasks=7177
       Total vcore-milliseconds taken by all reduce tasks=14521

| Reducer | Map time | Reduce time | Total committed heap usage (bytes) |
|---------|----------|-------------|------------------------------------|
| 1 | 9785 | 4361 | 640155648 |
| 3 | 7177 | 14521 | 861929472 |
| 5 | 6581 | 23931 | 1101004800 |
| 8 | 6225 | 47883 | 1535639552 |
| 10 | 13324 | 111481 | 1664090112 |

reducer vs map time

| reducer | time |
|---|---|
| 1 | 9785 |
| 3 | 7177 |
| 5 | 6581 |
| 8 | 6225 |
| 10 | 13324 |



reducer vs map time

| reducer | time |
|---|---|
| 1 | 4361 |
| 3 | 14521 |
| 5 | 23931 |
| 8 | 47883 |
| 10 | 111481 |

## reducer vs memory



| Reducer | Memory |
| --- | --- |
| 1 | 640155648 |
| 3 | 861929472 |
| 5 | 1101004800 |
| 8 | 1535639552 |
| 10 | 1664090112 |

## map and reduce timeline