

Matrix Multiplication Using Hadoop MapReduce

Introduction

This document outlines the implementation of a Hadoop MapReduce application for matrix multiplication. The program processes matrices stored in distributed files and computes their product using Hadoop's Mapper and Reducer functionalities.

Program Workflow

Input Data Format

Each matrix is stored in a text file, with entries formatted as follows:

MatrixName, RowIndex, ColIndex, Value

Example Input

Matrix A and B:

A,0,0,25

A,0,1,9

B,0,0,44

B,1,0,13

Output Data Format

The output matrix $C = A \times B$ is stored with the following format:

RowIndex, ColIndex Value

Example Output:

0,0 841

1,0 1493

MapReduce Process

1. Mapper: Emits intermediate key-value pairs for each potential result cell in matrix C.
2. Reducer: Aggregates the contributions to each cell and computes the final value.

Code Description

Driver Code

The driver initializes the Hadoop job and configures the input/output paths and parameters like matrix dimensions.

```
ome ▸ v3n0m ▸ Downloads ▸ MatrixMultiplication ▸ MatrixMultiplicationDriver.java
1  import org.apache.hadoop.conf.Configuration;
2  import org.apache.hadoop.fs.Path;
3  import org.apache.hadoop.io.Text;
4  import org.apache.hadoop.mapreduce.Job;
5  import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
6  import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
7
8  public class MatrixMultiplicationDriver {
9
10     public static void main(String[] args) throws Exception {
11         if (args.length != 4) {
12             System.err.println("Usage: MatrixMultiplicationDriver <input path> <output path> <matrixA.cols> <matrixB.rows>");
13             System.exit(-1);
14         }
15
16         Configuration conf = new Configuration();
17         try {
18             int matrixACols = Integer.parseInt(args[2]);
19             int matrixBRows = Integer.parseInt(args[3]);
20             conf.setInt("matrixA.cols", matrixACols);
21             conf.setInt("matrixB.rows", matrixBRows);
22         } catch (NumberFormatException e) {
23             System.err.println("Error: matrixA.cols and matrixB.rows must be integers. "+args[2]);
24             System.exit(-1);
25         }
26
27
28         Job job = new Job(conf, "Matrix Multiplication");
29         job.setJarByClass(MatrixMultiplicationDriver.class);
30
31         job.setMapperClass(MatrixMultiplicationMapper.class);
32         job.setReducerClass(MatrixMultiplicationReducer.class);
33
34         job.setOutputKeyClass(Text.class);
35         job.setOutputValueClass(Text.class);
36
37         FileInputFormat.addInputPath(job, new Path(args[0]));
38         FileOutputFormat.setOutputPath(job, new Path(args[1]));
39
40         System.exit(job.waitForCompletion(true) ? 0 : 1);
41     }
42 }
43
```

The MatrixMultiplicationDriver class is the driver code for a Hadoop MapReduce application that performs matrix multiplication. It sets up and manages the MapReduce job's execution, including its configuration, input/output paths, and Mapper/Reducer assignments.

Key Points:

- **Input Arguments:**
 - <input path>: Path to the directory or file containing the input data.
 - <output path>: Directory where the output will be saved.
 - <matrixA.cols>: Number of columns in Matrix A.
 - <matrixB.rows>: Number of rows in Matrix B.
- **Error Handling:**

- The program ensures `matrixA.cols` and `matrixB.rows` are integers and exits with an error message if they are not valid.
- **Role of Configuration:**
 - The Configuration object is used to pass custom parameters (`matrixA.cols` and `matrixB.rows`) to the Mapper and Reducer classes, enabling flexible handling of matrix dimensions.

This driver orchestrates the MapReduce job for matrix multiplication by ensuring all required configurations, paths, and operations are correctly set up.

Reducer Code

The reducer calculates the final value of each cell in the product matrix.

```
ome ▸ v3n0m ▸ Downloads ▸ MatrixMultiplication ▸ MatrixMultiplicationReducer.java
1 import org.apache.hadoop.io.IntWritable;
2 import org.apache.hadoop.io.Text;
3 import org.apache.hadoop.mapreduce.Reducer;
4
5 import java.io.IOException;
6 import java.util.HashMap;
7
8 public class MatrixMultiplicationReducer extends Reducer<Text, Text, Text, Text> {
9
10     @Override
11     protected void reduce(Text key, Iterable<Text> values, Context context) throws IOException, InterruptedException {
12         HashMap<Integer, Double> matrixA = new HashMap<>();
13         HashMap<Integer, Double> matrixB = new HashMap<>();
14
15         for (Text value : values) {
16             String[] tokens = value.toString().split(",");
17             String matrixName = tokens[0];
18             int index = Integer.parseInt(tokens[1]);
19             double val = Double.parseDouble(tokens[2]);
20
21             if (matrixName.equals("A")) {
22                 matrixA.put(index, val);
23             } else if (matrixName.equals("B")) {
24                 matrixB.put(index, val);
25             }
26         }
27
28         // Compute the dot product for C[i][k]
29         double sum = 0.0;
30         for (int j : matrixA.keySet()) {
31             if (matrixB.containsKey(j)) {
32                 sum += matrixA.get(j) * matrixB.get(j);
33             }
34         }
35
36         // Emit the result
37         context.write(key, new Text(String.valueOf(sum)));
38     }
39 }
40
```

Key Points:

1. **Intermediate Data Handling:**
 - The Reducer processes data grouped by the cell identifiers of Matrix C (e.g., `C[i][k]`).

- All relevant values from Matrix A and Matrix B for a specific $C[i][k]$ are passed as input.
- 2. Computational Logic:**
 - The dot product of a row from Matrix A and a column from Matrix B is computed.
 - Only matching indices (j) are used in the computation.
 - 3. Efficiency:**
 - The use of HashMap ensures quick lookups for matching indices, improving computational efficiency

Mapper Code

The mapper generates key-value pairs for intermediate computation.

```
ome ▸ v3n0m ▸ Downloads ▸ MatrixMultiplication ▸ MatrixMultiplicationMapper.java
1 import org.apache.hadoop.io.Text;
2 import org.apache.hadoop.mapreduce.Mapper;
3
4 import java.io.IOException;
5
6 public class MatrixMultiplicationMapper extends Mapper<Object, Text, Text, Text> {
7
8     @Override
9     protected void map(Object key, Text value, Context context) throws IOException, InterruptedException {
10         String[] tokens = value.toString().split(",");
11         String matrixName = tokens[0];
12         int row = Integer.parseInt(tokens[1]);
13         int col = Integer.parseInt(tokens[2]);
14         double val = Double.parseDouble(tokens[3]);
15
16         int matrixASize = context.getConfiguration().getInt("matrixA.cols", 0);
17         int matrixBsize = context.getConfiguration().getInt("matrixB.rows", 0);
18
19         if (matrixName.equals("A")) {
20             // Emit (i, k) as key and "A,j,value" as value
21             for (int k = 0; k < matrixASize; k++) {
22                 context.write(new Text(row + "," + k), new Text("A," + col + "," + val));
23             }
24         } else if (matrixName.equals("B")) {
25             // Emit (i, k) as key and "B,j,value" as value
26             for (int i = 0; i < matrixBsize; i++) {
27                 context.write(new Text(i + "," + col), new Text("B," + row + "," + val));
28             }
29         }
30     }
31 }
32
```

Key Points:

- 1. Key Design:**
 - The output key uniquely identifies a cell $C[i][k]$ in the resulting matrix.
 - This ensures all relevant values for a specific $C[i][k]$ are grouped together in the Reducer.
- 2. Data Encoding:**
 - The output value encodes both the matrix name (A or B) and the relevant information (index and value).
 - This allows the Reducer to reconstruct the necessary values for the dot product computation.
- 3. Efficiency:**

- The Mapper precomputes potential key-value pairs, distributing the computational workload across nodes.

Execution Steps

1. Compile the Program

Navigate to the source directory and compile the code:

```
javac -classpath $(hadoop classpath) -d classes MatrixMultiplication*.java
```

```
jar -cvf MatrixMultiplication.jar -C classes/ .
```

```
(v3n0m@V3N0M)~/Downloads/MatrixMultiplication
$ javac -classpath $(./hadoop-3.4.1/bin/hadoop classpath) -d . MatrixMultiplicationDriver.java MatrixMultiplicationMapper.java MatrixMultiplicationReducer.java
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
Note: MatrixMultiplicationDriver.java uses or overrides a deprecated API.
Note: Recompile with -Xlint:deprecation for details.

(v3n0m@V3N0M)~/Downloads/MatrixMultiplication
$ jar -cvf MatrixMultiplication.jar *.class
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
added manifest
adding: MatrixMultiplicationDriver.class(in = 2371) (out= 1196)(deflated 49%)
adding: MatrixMultiplicationMapper.class(in = 2570) (out= 1115)(deflated 56%)
adding: MatrixMultiplicationReducer.class(in = 2643) (out= 1170)(deflated 55%)

(v3n0m@V3N0M)~/Downloads/MatrixMultiplication
$
```

Execution Steps

1. Compile the Program
2. Prepare Input Files
3. Prepare Input Files
4. Prepare Input Files
5. Prepare Input Files
6. Prepare Input Files
7. Prepare Input Files
8. Prepare Input Files
9. Prepare Input Files
10. Prepare Input Files

2. Prepare Input Files

Upload the matrix files to HDFS:

```
hadoop fs -mkdir /input
```

```
hadoop fs -put matrixA1.txt /input/
```

```
hadoop fs -put matrixA2.txt /input/
```

```

(v3n0m@V3N0M)-[~/Downloads/MatrixMultiplication]
$ hadoop fs -mkdir /input
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
mkdir: `/input': Input/output error

(v3n0m@V3N0M)-[~/Downloads/MatrixMultiplication]
$ hadoop fs -mkdir ./input
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
mkdir: `input': File exists

(v3n0m@V3N0M)-[~/Downloads/MatrixMultiplication]
$ hadoop fs -put matrixA1.txt ./input/
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true

(v3n0m@V3N0M)-[~/Downloads/MatrixMultiplication]
$ hadoop fs -put matrixA2.txt ./input/
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true

(v3n0m@V3N0M)-[~/Downloads/MatrixMultiplication]

```

3. Run the Program

Execute the Hadoop job:

```
hadoop jar MatrixMultiplication.jar MatrixMultiplicationDriver /input /output matrixA.cols
matrixB.rows
```

```

L$ hadoop jar MatrixMultiplication.jar MatrixMultiplicationDriver ./input ./output 4 5
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
2024-11-29 21:46:09,607 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2024-11-29 21:46:09,675 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2024-11-29 21:46:09,675 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2024-11-29 21:46:09,749 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner.
2024-11-29 21:46:09,817 INFO input.FileInputFormat: Total input files to process : 3
2024-11-29 21:46:09,956 INFO mapreduce.JobSubmitter: number of splits:3
2024-11-29 21:46:10,262 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local1614633083_0001
2024-11-29 21:46:10,262 INFO mapreduce.JobSubmitter: Executing with tokens: []
2024-11-29 21:46:10,480 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2024-11-29 21:46:10,487 INFO mapreduce.Job: Running job: job_local1614633083_0001
2024-11-29 21:46:10,490 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2024-11-29 21:46:10,560 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitterFactory
2024-11-29 21:46:10,561 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2024-11-29 21:46:10,561 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2024-11-29 21:46:10,561 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter
2024-11-29 21:46:10,597 INFO mapred.LocalJobRunner: Waiting for map tasks
2024-11-29 21:46:10,597 INFO mapred.LocalJobRunner: Starting task: attempt_local1614633083_0001_m_000000_0
2024-11-29 21:46:10,664 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitterFactory
2024-11-29 21:46:10,664 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2024-11-29 21:46:10,664 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2024-11-29 21:46:10,713 INFO mapred.Task: Using ResourceCalculatorProcessTree : [ ]
2024-11-29 21:46:10,718 INFO mapred.MapTask: Processing split: file:/home/v3n0m/Downloads/MatrixMultiplication/input/matrix.txt:0+59
2024-11-29 21:46:10,746 INFO mapred.MapTask: (EQUATOR) 0 kvi 26214396(104857584)
2024-11-29 21:46:10,746 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
2024-11-29 21:46:10,746 INFO mapred.MapTask: soft limit at 83886080
2024-11-29 21:46:10,746 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
2024-11-29 21:46:10,746 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
2024-11-29 21:46:10,749 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
2024-11-29 21:46:10,798 INFO mapred.LocalJobRunner:
2024-11-29 21:46:10,799 INFO mapred.MapTask: Starting flush of map output
2024-11-29 21:46:10,799 INFO mapred.MapTask: Spilling map output
2024-11-29 21:46:10,799 INFO mapred.MapTask: bufstart = 0; bufend = 334; bufvoid = 104857600
2024-11-29 21:46:10,799 INFO mapred.MapTask: kvstart = 26214396(104857584); kvend = 26214296(104857184); length = 101/6553600
2024-11-29 21:46:10,816 INFO mapred.MapTask: Finished spill 0
2024-11-29 21:46:10,825 INFO mapred.Task: Task:attempt_local1614633083_0001_m_000000_0 is done. And is in the process of committing
2024-11-29 21:46:10,828 INFO mapred.LocalJobRunner: map
2024-11-29 21:46:10,828 INFO mapred.Task: Task 'attempt_local1614633083_0001_m_000000_0' done.
2024-11-29 21:46:10,847 INFO mapred.Task: Final Counters for attempt_local1614633083_0001_m_000000_0: Counters: 17
File System Counters
FILE: Number of bytes read=4815
FILE: Number of bytes written=719127
FILE: Number of read operations=0
FILE: Number of large read operations=0

```

4. View Results

Check the output in HDFS:

```
hadoop fs -cat /output/part-r-00000
```

```
(v3n0m@V3N0M)-[~/Downloads/MatrixMultiplication]
$ hadoop fs -cat ./output/part-r-00000
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
0,0 1217.0
0,1 0.0
0,2 0.0
0,3 0.0
1,0 1585.0
1,1 0.0
1,2 0.0
1,3 0.0
2,0 0.0
3,0 0.0
4,0 0.0
(v3n0m@V3N0M)-[~/Downloads/MatrixMultiplication]
```