



BIGDATA PROGRAMMING

Lab1 Assignment

Jetty, Lalitha Kumari – Class ID 5
Vinay Chandra Vasamsetti – Class ID 25
Team 11

1. Hadoop MapReduce Algorithm

Implement MapReduce algorithm for finding Facebook common friends problem and run the MapReduce job on Apache Hadoop. Show your implementation through map-reduce diagram as shown in Lesson Plan 2 : (<https://umkc.box.com/s/epp04s3m6g8jw6ulnnb4meqgnbwzb5uc>)

Write a report including your algorithm and result screenshots.

Finding Facebook common friends: Facebook has a list of friends (note that friends are a bi-directional thing on Facebook. If I'm your friend, you're mine). They also have lots of disk space and they serve hundreds of millions of requests everyday. They've decided to pre-compute calculations when they can to reduce the processing time of requests. One common processing request is the "You and Joe have 230 friends in common" feature. When you visit someone's profile, you see a list of friends that you have in common. We're going to use MapReduce so that we can calculate everyone's common friends once a day and store those results. Later on it's just a quick lookup. We've got lots of disk, it's cheap.

Example (What is the Key/Value Pair?)

Assume the friends are stored as Person-

>[List of Friends], our friends list is then:

A -> B C D

B -> A C D E

C -> A B D E

D -> A B C E

E -> B C D

The result after reduction is:

(A B) -> (C D)

(A C) -> (B D)

(A D) -> (B C)

(B C) -> (A D E)

(B D) -> (A C E)

(B E) -> (C D)

(C D) -> (A B E)

(C E) -> (B D)

(D E) -> (B C)

When D visits B's profile, we can quickly look up (B D) and see that they have three friends in common, (A C E).

Input		Mapper				Reducer		Output
		Key Value Pair						
A : B C D		A : B C D map => (A, B) => B C D (A, C) => B C D (A, D) => B C D		(A B) -> (A C D E) (B C D) (A C) -> (A B D E) (B C D) (A D) -> (A B C E) (B C D)		A B : C D A C : B D A D : B C		A B : C D A C : B D A D : B C B C : A D E B D : A C E B E : C D C D : A B E C E : B D D E : B C
B : A C D E		(A B) -> A C D E (B C) -> A C D E (B D) -> A C D E (B E) -> A C D E		(B C) -> (A B D E) (A C D E) (B D) -> (A B C E) (A C D E)		
C : A B D E		(A C) -> A B D E (B C) -> A B D E (C D) -> A B D E (C E) -> A B D E		. . .				
D : A B C E		(A D) -> A B C E (B D) -> A B C E (C D) -> A B C E (D E) -> A B C E		. .				
E : B C D		(B E) -> B C D (C E) -> B C D (D E) -> B C D		. .				

Input

A : B C D
B : A C D E
C : A B D E
D : A B C E
E : B C D

Mapper

Generates key value pairs for each user to its friends.

A : B C D
map =>
(A, B) => B C D
(A, C) => B C D
(A, D) => B C D

Reducer

Finds the common friends to the each user from the keyvalue pairs.

(A, B) => [B C D, A C D E]
reduce =>
(A, B) => C D

Output

AB: CD
AC: BD
AD: BC
BC: ADE
BD: ACE
BE: CD
CD: ABE
CE: BD
DE: BC

Code Screenshots

The screenshot displays an IDE with the following components:

- Project Explorer:** Shows the project structure with folders for `src`, `main`, `java`, and `mypackage`. The `CommonFriend` class is highlighted in the `java` folder.
- Code Editor:** Displays the `CommonFriend.java` file. The code includes imports for `java.io.IOException`, `java.util.Arrays`, `java.util.Iterator`, `java.util.LinkedList`, `java.util.List`, `java.util.StringTokenizer`, `org.apache.hadoop.fs.Path`, `org.apache.hadoop.io.LongWritable`, `org.apache.hadoop.io.Text`, `org.apache.hadoop.mapred.FileInputFormat`, `org.apache.hadoop.mapred.FileOutputFormat`, `org.apache.hadoop.mapred.JobClient`, `org.apache.hadoop.mapred.JobConf`, `org.apache.hadoop.mapred.MapReduceBase`, `org.apache.hadoop.mapred.Mapper`, `org.apache.hadoop.mapred.OutputCollector`, `org.apache.hadoop.mapred.Reducer`, and `org.apache.hadoop.mapred.Reporter`. The `CommonFriend` class implements `Mapper<LongWritable, Text, Text, Text>` and overrides the `map` method.
- Run Console:** Shows the execution output for the `CommonFriend` class. The output includes statistics for the Map-Reduce framework, such as the number of records, bytes, and materialized bytes. The process finished with exit code 0.

```
package mypackage;

import java.io.IOException;
import java.util.Arrays;
import java.util.Iterator;
import java.util.LinkedList;
import java.util.List;
import java.util.StringTokenizer;

import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.FileInputFormat;
import org.apache.hadoop.mapred.FileOutputFormat;
import org.apache.hadoop.mapred.JobClient;
import org.apache.hadoop.mapred.JobConf;
import org.apache.hadoop.mapred.MapReduceBase;
import org.apache.hadoop.mapred.Mapper;
import org.apache.hadoop.mapred.OutputCollector;
import org.apache.hadoop.mapred.Reducer;
import org.apache.hadoop.mapred.Reporter;

public class CommonFriend {

    public static class Map extends MapReduceBase
        implements Mapper<LongWritable, Text, Text, Text>{
        public void map(LongWritable key, Text value, OutputCollector<Text, Text> output, Reporter reporter)
            throws IOException{
            StringTokenizer tokenizer = new StringTokenizer(value.toString(), "\n");
            String line = null;
            String[] lineArray = null;
            String[] friendArray = null;
            String[] tempArray = null;
            ...
        }
    }
}
```

Run: CommonFriend

FILE: number of write operations=0
Map-Reduce Framework
Map input records=5
Map output records=18
Map output bytes=284
Map output materialized bytes=246
Input split bytes=161
Combine input records=0
Combine output records=0
Reduce input groups=0
Reduce shuffle bytes=0
Reduce input records=18
Reduce output records=9
Spilled Records=36
Shuffled Maps=0
Failed Shuffles=0
Merged Map outputs=0
GC time elapsed (ms)=0
Total committed heap usage (bytes)=514858816
File Input Format Counters
Bytes Read=55
File Output Format Counters
Bytes Written=0
Process finished with exit code 0

Project: CommonFriend -> Documents\MS In Data Science\2018 Summer

Output: .SUCCESS.crc, part-00000.crc, .SUCCESS, part-00000

src: main, java, mypackage, CommonFriend, resources, test, java, CommonFriendTest

target: CommonFriend.iml, input.txt, pom.xml, recomm.log, results.log

Run: CommonFriend -> FILE: number of write operations=0

Map-Reduce Framework

Map input records=5

Map output records=18

Map output bytes=284

Map output materialized bytes=246

Input split bytes=161

Combine input records=0

Combine output records=0

Reduce input groups=9

Reduce shuffle bytes=0

Reduce input records=18

Reduce output records=9

Spilled Records=36

Shuffled Maps=0

Failed Shuffles=0

Merged Map outputs=0

GC time elapsed (ms)=0

Total committed heap usage (bytes)=514859816

File Input Format Counters

Bytes Read=55

File Output Format Counters

Bytes Written=90

Process finished with exit code 0

part-00000

```
1 A: B C D
2 B: A C D E
3 C: A B D E
4 D: A B C E
5 E: B C D
```

Project: CommonFriend -> Documents\MS In Data Science\2018 Summer

Output: .SUCCESS.crc, part-00000.crc, .SUCCESS, part-00000

src: main, java, mypackage, CommonFriend, resources, test, java, CommonFriendTest

target: CommonFriend.iml, input.txt, pom.xml, recomm.log, results.log

Run: CommonFriend -> FILE: number of write operations=0

Map-Reduce Framework

Map input records=5

Map output records=18

Map output bytes=284

Map output materialized bytes=246

Input split bytes=161

Combine input records=0

Combine output records=0

Reduce input groups=9

Reduce shuffle bytes=0

Reduce input records=18

Reduce output records=9

Spilled Records=36

Shuffled Maps=0

Failed Shuffles=0

Merged Map outputs=0

GC time elapsed (ms)=0

Total committed heap usage (bytes)=514859816

File Input Format Counters

Bytes Read=55

File Output Format Counters

Bytes Written=90

Process finished with exit code 0

part-00000

```
1 A: B C D
2 A: B C D
3 A: B C D
4 B: C A D E
5 B: C A D E
6 B: C A D E
7 C: D A B E
8 C: D A B E
9 D: E B C
```

Reference - <http://ernie55ernie.github.io/mapreduce/2016/06/30/map-reduce-common-friend.html>

2. Use Case Based No SQL Comparison

Consider one of the use cases from the below link:

<https://umkc.box.com/s/q64fvjm6yd454w5v3ky0he4854g6m1fq>

These use cases were discussed in Lecture 1: Cassandra.

- a) Consider one of the use case and use a simple dataset. Describe the use case considered based on your assumptions, report the dataset, its fields, datatype etc.
- b) Use HBase to implement a Solution for the use case. Report at least 3 queries, their input and output. The query's relevance towards solving the use case is important.
- c) Use **Cassandra** to implement a Solution for the use case. Report at least 3 queries, their input and output. The query's relevance towards solving the use case is important.
- d) Compare **Cassandra** and HBase for your use case. Present a table with comparison of your use case being implemented in both NO SQL Systems.

Use Case – Coursera

Description – Coursera is an Education Platform which partners with top universities and organizations worldwide, to offer courses online for anyone to take, for free.

Challenges –

1. My SQL was insufficient
2. Unstable performance
3. Unexpected downtime
4. Limitation in introducing new features

Solution – After evaluating emerging database technologies, it chooses Cassandra (data stax)

Reason – 100% application uptime needed and scalability (enabling storage of growing user data)

Creating Table in Hbase – Courses by Learner with column families as LearnerDetails, CourseDetails.

Commands Used:

```
create 'Coursera', 'LearnerDetails', 'CourseDetails'
```

```
put 'Coursera', '1', 'LearnerDetails:LearnerID', '111'
```

```

put 'Coursera', '1', 'LearnerDetails:LearnerName', 'Lalitha'

put 'Coursera', '2', 'LearnerDetails:LearnerID', '222'

put 'Coursera', '2', 'LearnerDetails:LearnerName', 'Vinay'

put 'Coursera', '1', 'CourseDetails:CourseID', 'ML001'

put 'Coursera', '1', 'CourseDetails:CourseName', 'Machine Learning'

put 'Coursera', '2', 'CourseDetails:CourseID', 'DL001'

put 'Coursera', '2', 'CourseDetails:CourseName', 'Deep Learning'

```

scan 'Coursera'

List

get 'Coursera', '1'

```

hbase(main):001:0> list
TABLE
Course_by_Learner
Coursera
Table01
Table2
Table2b
Table3
Table4
Table5
8 row(s) in 0.1670 seconds

=> ["Course_by_Learner", "Coursera", "Table01", "Table2", "Table2b", "Table3", "Table4", "Table5"]
hbase(main):002:0> put 'Coursera','1','LearnerDetails:LearnerID','111'
0 row(s) in 0.1420 seconds

hbase(main):003:0> put 'Coursera','1','LearnerDetails:LearnerName','Lalitha'
0 row(s) in 0.0860 seconds

hbase(main):004:0> put 'Coursera','2','LearnerDetails:LearnerID','222'
0 row(s) in 0.0138 seconds

hbase(main):005:0> put 'Coursera','2','LearnerDetails:LearnerName','Vinay'
0 row(s) in 0.0150 seconds

hbase(main):006:0> put 'Coursera','1','CourseDetails:CourseID','ML001'
0 row(s) in 0.0130 seconds

hbase(main):007:0> put 'Coursera','1','CourseDetails:CourseName','Machine Learning'
0 row(s) in 0.0130 seconds

hbase(main):008:0> put 'Coursera','2','CourseDetails:CourseID','DL001'
0 row(s) in 0.0140 seconds

hbase(main):009:0> put 'Coursera','2','CourseDetails:CourseName','Deep Learning'
0 row(s) in 0.0120 seconds

hbase(main):010:0> scan 'Coursera'
COLUMN+CELL
ROW
1
1 column=CourseDetails:CourseID, timestamp=1529338279039, value=ML001
1 column=CourseDetails:CourseName, timestamp=1529338303231, value=Machine Learning
1
1 column=LearnerDetails:LearnerID, timestamp=1529338169506, value=111
1 column=LearnerDetails:LearnerName, timestamp=1529338283202, value=Lalitha
2
2 column=CourseDetails:CourseID, timestamp=1529338325494, value=DL001
2 column=CourseDetails:CourseName, timestamp=1529338345631, value=Deep Learning
2
2 column=LearnerDetails:LearnerID, timestamp=1529338223445, value=222
2 column=LearnerDetails:LearnerName, timestamp=1529338245320, value=Vinay
2 row(s) in 0.0320 seconds

```

```

hbase(main):012:0> get 'Coursera','1'
COLUMN
CourseDetails:CourseID
CourseDetails:CourseName
LearnerDetails:LearnerID
LearnerDetails:LearnerName
4 row(s) in 0.0170 seconds

CELL
timestamp=1529338279039, value=ML001
timestamp=1529338303231, value=Machine Learning
timestamp=1529338169506, value=111
timestamp=1529338283202, value=Lalitha

hbase(main):013:0>

```

Creating Table in Cassandra – employee with ename, job title and hiredate

Creating Keyspace name Vinay

```
WARNING: console codepage must be set to cp65001 to support utf-8 encoding on Windows platforms.
If you experience encoding problems, change your console codepage with 'chcp 65001' before starting.

Connected to Test Cluster at 127.0.0.1:9042.
[cqlsh 5.0.1 | Cassandra 2.2.8 | CQL spec 3.3.1 | Native protocol v4]
Use HELP for help.
WARNING: pyreadline dependency missing. Install to enable tab completion.
cqlsh> create keyspace VINAY with replication={'class':'SimpleStrategy', 'replication_factor':1};
AlreadyExists: Keyspace 'vinay' already exists
cqlsh> use vinay;
```

Creating Table employee1

```
cqlsh:vinay> CREATE TABLE employee1 (
... empno int,
... ename text,
... jobtitle text,
... hiredate date,
... PRIMARY KEY (empno));
```

Inserting values in to table

```
SyntaxException: line 1:79 mismatched input '-06' expecting ')' (...) VALUES (1225,'katre',1993[-06]-12...)
cqlsh:vinay> INSERT INTO employee1 (empno,ename,hiredate,jobtitle) VALUES (1225,'katre','1993-06-12','softengin');
cqlsh:vinay> INSERT INTO employee1 (empno,ename,hiredate,jobtitle) VALUES (1245,'mourya','1980-12-17','manager');
cqlsh:vinay> INSERT INTO employee1 (empno,ename,hiredate,jobtitle) VALUES (1265,'santosh','1980-12-17','associate');
cqlsh:vinay> INSERT INTO employee1 (empno,ename,hiredate,jobtitle) VALUES (1265,'santosh','1980-12-17','associate');
cqlsh:vinay> SELECT * FROM employee1;
```

empno	ename	hiredate	jobtitle
1245	mourya	1980-12-17	manager
1225	katre	1993-06-12	softengin
1265	santosh	1980-12-17	associate

Adding Column using Alter command

```
cqlsh:vinay> ALTER TABLE employee1 ADD salary int ;
cqlsh:vinay> SELECT * FROM employee1;
```

empno	ename	hiredate	jobtitle	salary
1245	mourya	1980-12-17	manager	null
1225	katre	1993-06-12	softengin	null
1265	santosh	1980-12-17	associate	null

(3 rows)

Updating salaries if the employees

```
cqlsh:vinay> update employee1 set salary = 10000 where empno = 1245;
cqlsh:vinay> update employee1 set salary = 20000 where empno = 1225;
cqlsh:vinay> update employee1 set salary = 30000 where empno = 1265;
cqlsh:vinay> SELECT * FROM employee1;
```

empno	ename	hiredate	jobtitle	salary
1245	mourya	1980-12-17	manager	10000
1225	katre	1993-06-12	softengin	20000
1265	santosh	1980-12-17	associate	30000

List name salary of the employee who are clerks

```
cqlsh:vinay> INDEX employee1 (jobtitle);
SyntaxException: line 1:0 no viable alternative at input 'INDEX' ([INDEX]...)
cqlsh:vinay> INDEX employee1(jobtitle);
SyntaxException: line 1:0 no viable alternative at input 'INDEX' ([INDEX]...)
cqlsh:vinay> CREATE INDEX on employee1 (jobtitle);
cqlsh:vinay> SELECT ename,salary from employee1
... WHERE(jobtitle = 'clerk');

ename | salary
-----+-----
sammy | 4000
```

List the name,job,salary for every employee joined on dec 17 1980

```
(1 rows)
cqlsh:vinay> CREATE INDEX on employee1 (hiredate);
cqlsh:vinay> SELECT * from employee1
... WHERE(hiredate = '1980-12-17');

empno | ename   | hiredate   | jobtitle | salary
-----+-----+-----+-----+-----
1245  | mourya | 1980-12-17 | manager | 10000
1265  | sammy  | 1980-12-17 | clerk  | 4000

(2 rows)
```

Similarities and differences between Cassandra and Hbase:

	Cassandra	Hbase
Database	open-source NoSQL Database, They both are distributed database that can manage extremely large data sets and handle non-relational data.	open-source NoSQL Database, They both are distributed database that can manage extremely large data sets and handle non-relational data.
Scalability	high linear scalability. That means to handle more data, the user should simply increase the number of nodes in the cluster. Because of this feature, they both are an excellent choice for handling a large amount of data.	high linear scalability. That means to handle more data, the user should simply increase the number of nodes in the cluster. Because of this feature, they both are an excellent choice for handling a large amount of data.

Replication	But in both Cassandra and HBase, there is a safeguard that prevents data loss even after failure	But in both Cassandra and HBase, there is a safeguard that prevents data loss even after failure
Programming/Coding	Both are column-oriented databases that implement similar write paths	Both are column-oriented databases that implement similar write paths
Infrastructure	Cassandra, on the other hand, has different infrastructure and operation than Hadoop	This HBase-Hadoop infrastructure consists of several moving parts like Zookeeper, HBase master, Data nodes and Name Node.
Support	Cassandra, on the other hand, supports ordered partitioning. Cassandra is also limited in supporting range based row scans.	HBase, do not supports ordered partitioning. HBase offers a coprocessor capability

Reference - <https://data-flair.training/blogs/hbase-vs-cassandra/>