

Spring 2019

## E599/B649 High Performance Big Data Systems

### Lab Tasks



---

# Lab 3 of SP19-BL-ENGR-E599-30563

## Goal

---

- Walk through of the K-means example in Harp
- Memory configuration of mapper
- From single mapper to multiple mappers
- Use different communication operations
- Learn how to debug the program

## Deliverables

---

Submit the experiment results to Canvas.

Run K-means example (Points=1000, Centroids=10, Dimension=100, Iterations=100) and record the execution time in milli-seconds, excluding the data loading time.

- Thread number fixed to 24 and mapper number varying from 1 to 4 (with correct memory configuration).
- Thread number fixed to 24 and mapper number to 4 with different communication operations and record the time spent in *communication* for each communication operations (allreduce, regroupallgather, bcastreduce, pushpull).
  - The time evaluation is already done in allreduce model, use it as an example.
  - After modifying the codes, you shall re-compile the harp codes before launching test script.

## Evaluation

---

Lab participation: credit for 1 point based upon a successful completion of the lab tasks

## Download the Harp Codes for Lab Session

---

Git pull the codes from *labsession* branch (Please run *git clone* command on the login node)

```
cd ~
mkdir -p Labsession
cd Labsession
git clone -b labsession https://github.com/DSC-SPIDAL/harp.git
```

Compile the codes by Maven

```
cd ~/labsession/harp
mvn clean package -Phadoop-2.6.0
## copy compiled jars to Hadoop directory
cp core/harp-hadoop/target/harp-hadoop-0.1.0.jar $HADOOP_HOME/share/hadoop/mapreduce/
cp core/harp-collective/target/harp-collective-0.1.0.jar $HADOOP_HOME/share/hadoop/mapreduce/
cp core/harp-daal-interface/target/harp-daal-interface-0.1.0.jar $HADOOP_HOME/share/hadoop/mapreduce/
cp third_party/*.jar $HADOOP_HOME/share/hadoop/mapreduce/
```

## The Harp API

---

Harp provides communication and computation modules for distributed implementations of data analytics algorithms. We first introduce fundamental data structures, the inter-mapper communication API. In the next Lab session, we shall cover the multi-threading schedulers.

### Data Structure in Harp

Harp provides three levels of data structures: arrays and objects, partition, and table.

#### Primitive Data Type

Arrays and Serializable objects are the basic data structures, which include:

1. `ByteArray` : an array with byte-type elements
2. `ShortArray` : an array with short-type elements
3. `IntArray` : an array with int-type elements
4. `FloatArray` : an array with float-type elements
5. `LongArray` : an array with long-type elements

6. `DoubleArray` : an array with double-type elements
7. `Writable` : serializable object

## Partition

`Partition` is a wrapper of the data structures shown above. Every partition has an ID. In collective communication, partitions from different processors with the same ID will be merged. The merge operation is defined by `PartitionCombiner`.

## Table

`Table` is a container for partitions. It is a high-level data structure and the unit for collective communication.

An example of how to construct a table is:

```
95 Table<DoubleArray> table = new Table<>(0, new DoubleArrPlus());
96 for (int i = 0; i < numPartitions; i++) {
97     DoubleArray array = DoubleArray.create(size, false);
98     table.addPartition(new Partition<>(i, array));
99 }
```

## Inter-Mapper Communication Operations

All of the communication APIs are defined at

`core/harp-hadoop/src/main/java/org/apache/hadoop/mapred/CollectiveMapper.java` We select five operations used in the K-means. The rest of the APIs are introduced with details at our [website tutorial page](#).

### Allreduce

`allreduce` aims to first combine tables from other workers and then broadcast the accumulated table. All workers should run it concurrently. The definition of the method is:

```

467  /**
468   * Allreduce partitions of the tables to all the
469   * local tables.
470   *
471   * @param contextName
472   *         the name of the operation context
473   * @param operationName
474   *         the name of the operation
475   * @param table
476   *         the table to hold the partitions
477   * @return a boolean tells if the operation
478   *         succeeds
479   */
480  public <P extends Simple> boolean allreduce(
481      String contextName, String operationName,
482      Table<P> table) {
483      boolean isSuccess =
484          AllreduceCollective.allreduce(contextName,
485              operationName, table, dataMap, workers);
486      dataMap.cleanOperationData(contextName,
487          operationName);
488      return isSuccess;
489  }

```

```

//Example
allreduce("k-means", "allreduce_centroids", table);

```

## Allgather

`allgather` aims to first collect tables from other workers and then broadcast the collection. All workers should run it concurrently. The definition of the method is:

```

443  /**
444   * Allgather partitions of the tables to all the
445   * local tables.
446   *
447   * @param contextName
448   *         the name of the operation context
449   * @param operationName
450   *         the name of the operation
451   * @param table
452   *         the table to hold the partitions
453   * @return a boolean tells if the operations
454   *         succeeds
455   */
456  public <P extends Simple> boolean allgather(
457      String contextName, String operationName,
458      Table<P> table) {
459      boolean isSuccess =
460          AllgatherCollective.allgather(contextName,
461              operationName, table, dataMap, workers);
462      dataMap.cleanOperationData(contextName,
463          operationName);
464      return isSuccess;
465  }

```

```

//example
allgather("k-means", "allgather_centroids", table);

```

## Broadcast

`broadcast` aims to share a table in one worker with others. All workers should run it concurrently. The definition of the method is:

```

385  /**
386   * Broadcast the partitions of the table on a
387   * worker to other workers.
388   *
389   * @param contextName
390   *         the name of the operation context
391   * @param operationName
392   *         the name of the operation
393   * @param table
394   *         the table used to hold the
395   *         partitions
396   * @param bcastWorkerID
397   *         the worker ID of broadcasting data
398   * @param useMSTBcast
399   *         if minimum-spanning tree algorithm
400   *         is used
401   * @return a boolean tells if the operation
402   *         succeeds
403   */
404  public <P extends Simple> boolean broadcast(
405      String contextName, String operationName,
406      Table<P> table, int bcastWorkerID,
407      boolean useMSTBcast) {
408      boolean isSuccess =
409          BcastCollective.broadcast(contextName,
410              operationName, table, bcastWorkerID,
411              useMSTBcast, dataMap, workers);
412      dataMap.cleanOperationData(contextName,
413          operationName);
414      return isSuccess;
415  }

```

```

// example to bcast data from mapper 0 and to use minimum-spanning tree
broadcast("k-means", "bcast_centroids", table, 0, true);

```

## Push

`push` aims to collect all workers' partitions and store them in global. All workers should run it concurrently. The definition of the method is:

```

552  /**
553   * Push the partitions of local tables to the
554   * global table.
555   *
556   * @param contextName
557   *         the name of the operation context
558   * @param operationName
559   *         the name of operation
560   * @param localTable
561   *         the local tables
562   * @param globalTable
563   *         the global table which acts like a
564   *         distributed dataset, each partition
565   *         in this table is unique
566   * @param partitioner
567   *         when a partitioner is used, the
568   *         local partitions are sent to the
569   *         global table even without partition
570   *         ID association
571   * @return a boolean tells if the operation
572   *         succeeds
573   */
574  public <P extends Simple, PT extends Partitioner>
575    boolean push(String contextName,
576                String operationName, Table<P> localTable,
577                Table<P> globalTable, PT partitioner) {
578    boolean isSuccess =
579      LocalGlobalSyncCollective.push(contextName,
580      operationName, localTable, globalTable,
581      partitioner, dataMap, workers);
582    dataMap.cleanOperationData(contextName,
583      operationName);
584    return isSuccess;
585  }

```

```

//example
push("k-means", "push_local_centroids", table_local, table_global, new Partitioner(th
is.getNumWorkers()));

```

## Pull

`pull` aims to distribute the partitions stored in global to all workers. All workers should run it concurrently. The definition of the method is:

```

518  /**
519   * Pull partitions in the global table to the
520   * local tables. If any partition ID conflicts
521   * with the partition ID in the local table,
522   * combines them.
523   *
524   * @param contextName
525   *         the name of the operation context
526   * @param operationName
527   *         the name of the operation
528   * @param localTable
529   *         the local tables
530   * @param globalTable
531   *         the global table, acts like a
532   *         distributed dataset
533   * @param useBcast
534   *         if using broadcast in scattering the
535   *         partitions
536   * @return a boolean tells if the operation
537   *         succeeds
538   */
539  public <P extends Simple> boolean pull(
540      String contextName, String operationName,
541      Table<P> localTable, Table<P> globalTable,
542      boolean useBcast) {
543      boolean isSuccess =
544          LocalGlobalSyncCollective.pull(contextName,
545              operationName, localTable, globalTable,
546              useBcast, dataMap, workers);
547      dataMap.cleanOperationData(contextName,
548          operationName);
549      return isSuccess;
550  }

```

```

// example
pull("k-means", "global-local", table_local, table_global, true);

```

## A Walk Through of K-means in Harp

The file of K-means example is at the following path

```
cd harp/contrib/src/main/java/edu/iu/kmeans/
```



The `common` folder contains the entrance of the K-means program, where it configures the Harp mapper function and launch the Hadoop job. The other folders refer to different distributed implementations of K-means. The file *common/KmeansMapCollective.java* configures and launches a harp job upon Hadoop system.

```

122 private void runKMeans(int numOfDataPoints,
123     int numCentroids, int vectorSize,
124     int numIterations, int JobID, int numMapTasks,
125     Configuration configuration, Path workDirPath,
126     Path dataDir, Path cDir, Path outDir,
127     String operation)
128     throws IOException, URISyntaxException,
129     InterruptedException, ClassNotFoundException {

...

136 do {
137 // -----

...

144 // configure kmeans job
145 Job kmeansJob = configureKMeansJob(
146     numOfDataPoints, numCentroids, vectorSize,
147     numMapTasks, configuration, workDirPath,
148     dataDir, cDir, outDir, JobID,
149     numIterations, operation);

...

    // run kmeans job
159     jobSuccess =
160         kmeansJob.waitForCompletion(true);

...

    // check whether job is successful
174     if (!jobSuccess) {
175         System.out.println(
176             "KMeans Job failed. Job ID:" + JobID);
177         jobRetryCount++;
178         if (jobRetryCount == 3) {
179             break;
180         }
181     } else {
182         break;
183     }
184 } while (true);
185 }

```

## Job Configuration

Line 226 to 245 configures the mapper function we will use in K-means. Each mapper corresponds to a type of synchronization pattern in distributed system, namely the computation model. Find more details in our website [documentation](#)

```
226    // use different kinds of mappers
227    if (operation.equalsIgnoreCase("allreduce")) {
228        job.setMapperClass(
229            edu.iu.kmeans.allreduce.KmeansMapper.class);
230    } else if (operation
231        .equalsIgnoreCase("regroup-allgather")) {
232        job.setMapperClass(
233            edu.iu.kmeans.regrouppallgather.KmeansMapper.class);
234    } else if (operation
235        .equalsIgnoreCase("broadcast-reduce")) {
236        job.setMapperClass(
237            edu.iu.kmeans.bcastreduce.KmeansMapper.class);
238    } else if (operation
239        .equalsIgnoreCase("push-pull")) {
240        job.setMapperClass(
241            edu.iu.kmeans.pushpull.KmeansMapper.class);
242    } else { // by default, allreduce
243        job.setMapperClass(
244            edu.iu.kmeans.allreduce.KmeansMapper.class);
245    }
```

Secondly, we must setup the memory allocated to each mapper in Hadoop and the ratio of JVM heap memory from Line 265 to 274,

```
265    jobConf.setInt(
266        "mapreduce.map.collective.memory.mb", this.mem_per_mapper);
267    // mapreduce.map.collective.java.opts
268    int xmx = (int) Math.ceil((this.mem_per_mapper)*0.8);
269    int xmn = (int) Math.ceil(0.25 * xmx);
270    jobConf.set(
271        "mapreduce.map.collective.java.opts",
272        "-Xmx" + xmx + "m -Xms" + xmx + "m"
273        + " -Xmn" + xmn + "m");
274    return job;
```

where *mapreduce.map.collective.memory.mb* is the memory allocated to each mapper in MB, which overrides the Hadoop mapper configuration xml files *mapred-site.xml*. Here the default value is 2GB, and make sure the

value memory per mapper times the number of mappers does not exceed the total memory allocated to a node specified by `yarn.nodemanager.resource.memory-mb` in `yarn-site.xml`

The `xmx` and `xmn` refers to the maximum and minimum value of memory available to the JVM heap. By default, it is 80%.

## Job launch

The invocation of function `waitForCompletion` will launch the mapper function defined in the job configuration step and wait for its return.

```
159     jobSuccess =  
160         kmeansJob.waitForCompletion(true);
```

Here we do not setup the reducer function in the job configuration step at line 254, because Harp uses in-memory collective

`java 254 job.setNumReduceTasks(0);` communication operation such as *allreduce* to execute the similar task of a *Reducer* in the standard MapReduce framework.

Whenever a job is launched, the control flow is taken by the mapper function defined in

`allreduce/KmeansMapper.java`, which shall execute steps as follows

- Data Loading (All the mappers)
- Local Computation (All the mappers)
- Collective Communication (Among mappers)
- Output Result (Master mapper)

Except for data loading and output of result, the other steps are iteratively executed in a for loop.

## Data Loading

K-means has two types of data. 1) The training data (data points) is loaded from HDFS, 2) the model data (centroids) is loaded by master mapper and broadcasted to the slave mappers.

```
110     // load data  
111     ArrayList<DoubleArray> dataPoints =  
112         loadData(fileNames, vectorSize, conf);  
113     numPoints = dataPoints.size();
```

```

...

92     Table<DoubleArray> cenTable =
93         new Table<>(0, new DoubleArrPlus());
94     if (this.isMaster()) {
95         loadCentroids(cenTable, vectorSize,
96             conf.get(KMeansConstants.CFILE), conf);
97     }

...

102     // broadcast centroids
103     broadcastCentroids(cenTable);
...

```

## Local Computation

Each node holds a portion of training data and a complete copy of the centroids data. The local computation computes the MSE value (mean square error) and updates the centroids table by the partial results.

```

116     // iterations
117     for (int iter = 0; iter < iteration; iter++) {
118         previousCenTable = cenTable;
119         cenTable =
120             new Table<>(0, new DoubleArrPlus());

...

        // compute new partial centroid table using
        // previousCenTable and data points
130         MSE = computation(cenTable, previousCenTable,
131             dataPoints);
...

    }

```

## Inter-mapper Communication

The inter-mapper communication at each iteration plays the same role as the *reducer* in MapReduce.

```

116     // iterations
117     for (int iter = 0; iter < iteration; iter++) {

...

        // AllReduce;
138         allreduce("main", "allreduce_" + iter,
139                 cenTable);

...

145         calculateCentroids(cenTable);

...

    }

```

*allreduce* is a communication operation API defined by Harp, which first *reduce* the partial results to Mapper 0 and then broadcast the results to all of the other mappers. *allreduce* adopts the minimum spanning tree algorithm in reducing the data, which makes it fast when compared to standard *Reduce* operation.

After the *allreduce* communication, the mapper needs to re-calculate the centroids table, which could also be done within the *allreduce* communication operation via the user-defined *allreduce* operation.

## Output Results

When all of the iterations are completed, calculate the final MSE value and write back the centroids table to HDFS. This step is only executed at the Master mapper.

```

154    // evaluate the results
155    calcMSE(conf);
156
157    if (this.isMaster()) {
158        outputCentroids(cenTable, conf, context);
159    }
160
161    //save evaluation info
162    if (this.isMaster()){
163
164        FileSystem fs = FileSystem.get(conf);
165        Path path = new Path(conf.get(KMeansConstants.WORK_DIR) + "/evaluation");
166        FSDataOutputStream output =
167            fs.create(path, true);
168        BufferedWriter writer = new BufferedWriter(
169            new OutputStreamWriter(output));
170        writer.write("MSE : " + this.finalMSE + "\n");
171        writer.write("Compute Time : " + this.compute_time + "\n");
172        writer.write("Comm Time : " + this.comm_time + "\n");
173        writer.close();
174
175    }

```

## Run K-means on Multiple Mappers

---

We provide a bash script to run the K-means example. Please check the file

```

## to edit the script
vim contrib/test_scripts/k-means.sh
## to run the script
contrib/test_scripts/k-means.sh

```

The first part of the script file defines a variety of paths related to our execution.

```
1 #!/bin/bash
2
3 #
4 # test script for KMeans, using random generated dataset
5 #
6
7 #get the startup directory
8 startdir=$(dirname $0)
9 harproot=$(readlink -m $startdir/../../)
10 bin=$harproot/contrib/target/contrib-0.1.0.jar
11 hdfsroot=/harp-test
12 hdfsoutput=$hdfsroot/km/
13
14 if [ ! -f $bin ] ; then
15     echo "harp contrib app not found at "$bin
16     exit -1
17 fi
18 if [ -z ${HADOOP_HOME+x} ];then
19     echo "HADOOP not setup"
20     exit
21 fi
22
23 hdfs dfsadmin -safemode get | grep -q "ON"
24 if [[ "$?" = "0" ]]; then
25     hdfs dfsadmin -safemode leave
26 fi
27
28 #workdir
29 workdir=test_km
30
31 mkdir -p $workdir
32 cd $workdir
```

The second part of the script defines a function to call the program execution. Here \$1, \$2, ... represents the arguments fed to the function. At the end of the function, the output files are downloaded from HDFS to your local folder



```

37 runtest()
38 {
39     # <numOfDataPoints>: the number of data points you want to generate randomly
40     # <num of centriods>: the number of centroids you want to clustering the data to
41     # <size of vector>: the number of dimension of the data
42     # <number of map tasks>: number of map tasks
43     # <number of iteration>: the number of iterations to run
44     # <work dir>: the root directory for this running in HDFS
45     # <local dir>: the harp kmeans will firstly generate files which contain data points to local directory. Set this argument
46     # <communication operation> includes:
47     #     [allreduce]: use allreduce operation to synchronize centroids
48     #     [regroup-allgather]: use regroup and allgather operation to synchronize centroids
49     #     [broadcast-reduce]: use broadcast and reduce operation to synchronize centroids
50     #     [push-pull]: use push and pull operation to synchronize centroids
51     # <mem per mapper>
52     hadoop jar $bin edu.iu.kmeans.common.KmeansMapCollective $1 $2 $3 $4 $5 $hdfsoutput/$6 /tmp/kmeans $6 $7
53
54     if [ $? -ne 0 ]; then
55         echo "run km failure"
56         exit -1
57     fi
58
59     ## retrieve the results from HDFS
60     if [ -d ./$6 ];then
61         rm -rf ./$6/*
62     else
63         mkdir -p ./$6
64     fi
65
66     hdfs dfs -get $hdfsoutput/$6/evaluation ./$6
67     hdfs dfs -get $hdfsoutput/$6/centroids ./$6
68 }

```

The third part of the script runs the experiment with 1000 data points, 10 centroids, feature dimension of 100, 2 mappers, 100 iterations, allreduce model, and 2000MB memory per mapper.

```

#run test
runtest 1000 10 100 2 100 allreduce 2000

```

The output result is at `./test_km/allreduce/evaluation`

## Debugging of the Harp workload

Harp workload is running within the Hadoop runtime. When a program fails, we must first collect any useful log information, and secondly insert print out statements into the source codes to identify and isolate the problematic code sections.

### Collect Log Information

There are two types of log files related to the Harp workload.

The Hadoop system logs are located at `$HADOOP_HOME/logs` folder by default, where each of your node keeps a log file that traces all of the events, including any exceptions (e.g., fail to start the node). For instance, a file named `hadoop-1c37-datanode-j-070.log` refers to the log information of datanode *j-070* owned by user *1c37*. Similarly, there are log files for Namenode, NodeManager, and ResourceManager.

1. `hadoop-${user}-namenode-${node}.log`
2. `hadoop-${user}-secondarynameNode-${node}.log`
3. `hadoop-${user}-datanode-${node}.log`
4. `yarn-${user}-nodemanager-${node}.log`
5. `yarn-${user}-resourcemanager-${node}.log`

The Application specific log files are defined at file `yarn-site.xml`

```
<property>
  <name>yarn.nodemanager.log-dirs</name>
  <value>/local-path/userlogs</value>
</property>
```

The log files named by the Application ID (Job ID). For instance, a folder named `application_1524272601886_0164` contains all the log files for Application 1524272601886\_0164. It includes three subfolders

1. `container1524272601886016401000001` records the log information of the AppMaster container process.
2. `container1524272601886016401000002` records the log information of the first YarnChild container process.
3. `container1524272601886016401000003` records the log information of the second YarnChild container process.

There are two YarnChild containers because job `application_1524272601886_0164` launches two

mappers. By checking the contents of these subfolders (e.g., the syslog file), you may figure out what causes the crash of the mapper functions.

## Print out Intermediate Values

A most straightforward yet effective way of debugging parallel and distributed program is to print out any intermediate values of variables or setup check points. In Harp, you may use the `LOG` class to record variable status and check them in the *syslog* file mentioned above.

```
LOG.info("Start collective mapper.");
```

For example, setup a check point at the beginning of the mapper function execution. If the program crashes, and you are not able to find the message *Start collective mapper.* in syslog file of a YarnChild, then the error occurs even before the launch of the mapper function.

Also, you could use normal Java printing function such as the `System.out.println` to record information, and these information could be found at `stdout` and `stderr` files alongside the `syslog` within the Application specific log folder.

```
System.out.println("Errors: " + err);
```