Zayd Hammoudeh

(24) +3 / 25

# CS286 Exam I

1. What are the *3 V's* of big data?
   a. Value, velocity, volume
   b. Velocity, volume, veracity
   c. Velocity, volume, variety
   d. Veracity, value, volume

   C

2. Which of the following is considered *unstructured* data?
   a. video
   b. email
   c. log
   d. table

   A

3. Which of the following scaling strategies introduces *error*?
   a. Scale-up
   b. Scale-out
   c. parallelizing
   d. sampling

   D

4. The `reduce()` method *reads one*
   a. Split at a time
   b. Block at a time
   c. Key-value pair at a time
   d. Spill at a time

   C

5. The combiner
   a. Always improves performance
   b. Is called before the `map()` method
   c. Is called after the `reduce()` method
   d. Is an optional optimization

   D

6. Which statement is *true* about MapReduce?  — may
   a. MapReduce was invented at Google
   b. MapReduce programs must be written in Java
   c. MapReduce programs are scale-dependent
   d. MapReduce is not supported by YARN

   B

7. Which of the following *extends* `InputFormat` class?
   a. `LineRecordReader`
   b. `FileInputFormat`
   c. `TextInputFormat`
   d. `InputSplit`

   C

(1)

8. The *number of mappers* that run for a job
   a. Is equal to the number of blocks
   b. Is equal to the 1 by default
   c. Is equal to the number of reducers
   (d.) Is equal to the number of input splits

9. The *number of reducers* that run for a job
   a. Depends on the number of input splits
   b. Is equal to the number of mappers
   (c.) Is 1 by default
   d. Depends on the amount of RAM installed

10. Which of the following are **NOT** part of the Hadoop strategy?
   a. Distribute data
   (b.) Synchronize data
   c. Tolerate failure
   d. Distribute computation

11. What does a *virtual file system* provide?
   a. Read-only access
   b. Append-only access
   c. Read-write access
   (d.) POSIX compliance

12. What is included as part of the output from *every* MapReduce job?
   a. _SUCCCESS file
   (b.) _logs directory
   c. part-r-xxxxx files
   d. part-m-xxxxx files

13. Which statement is *true* of the Mapper class?
   (a.) The output of the Mapper must match the input of the Reducer.
   b. The input of the Mapper must match the input of the Reducer.
   c. The input of the Mapper must match the output of the Reducer.
   d. The output of the Mapper must match the output of the Reducer.

14. Which statement is *true* of the Reducer class?
   a. The Reducer object calls the reduce() method once per split.
   b. The Reducer object calls the run() method once per key.
   c. The Reducer object calls the method once per key
   (d.) The Reducer object calls the reduce() ~~cleanup()~~ method once per key.

Cleanup

15. Which statement is true of the *driver* class?

    **D**
    - a. The driver must check the command-line syntax.
    - b. The driver must parse the command-line arguments.
    - c. The driver must use the ToolRunner interface.
    - (d.) The driver defines the Mapper and Reducer classes.

16. Which Java statement converts a `Text` parameter (`value`) to a *list* of tokens?

    **B**
    - a. `new StringTokenizer(value, "\\s+");`
    - (b.) `new String(value, "\\s+");`
    - c. `value.toString();`
    - d. `new String(value);`

17. How many ~~total~~ copies of a block are created *by default* in HDFS and MapR-FS?

    **C**
    - a. 1
    - b. 2
    - (c.) 3
    - d. HDFS and MapR-FS do not have the same default number of replicas

18. Which Java statement launches a MapReduce job *synchronously*?

    **A**
    - (a.) `return job.waitForCompletion(true) ? 0:1;`
    - b. `return job.waitForCompletion(false) ? 0:1;`
    - c. `job.submit(false);`
    - d. `job.submit(true);`

19. Which of the following is *different* between `mapred` and `mapreduce` packages?

    **B**
    - a. YARN support
    - (b.) `map()` method signature
    - c. `reduce()` method signature
    - d. Number of slots or containers required

20. Which statement is *true* of the Hadoop MapReduce data types?

    **C**
    - a. Values must implement the `WritableComparable`
    - b. The `Text` class implements the `Serializable` interface
    - (c.) Keys must implement the `WritableComparable` interface
    - d. The `BooleanWritable` implements the `Interruptable` interface

21. Which of the following is **NOT** a feature of HDFS?

    **A** (crossed out)
    - (a.) Compression
    - b. Replication
    - c. Authorization
    - d. Encryption

(6)

22. How many *mappers* will the following configuration instantiate at the job runtime? : 4 files (10KB, 250MB, 600MB, 1000MB); block size = 256MB
    a. 8
    b. 9
    c. 4
    d. 5

23. Which of the following statements is *true* regarding how records are read by the mappers?
    a. There is one mapper instantiated per record
    b. There is one map () method called per record
    c. Records must fall on input split boundaries
    d. End users must define the record terminator

24. Which of the following is a counter that the Hadoop framework will track by default?
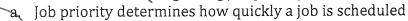    a. Number of RPC packets sent
    b. Physical memory consumption
    c. Number of RPC packets received
    d. Number of CPU pipeline stalls

25. Which statement is ~~true~~ False regarding job priority?
    a. Job priority determines how quickly a job is scheduled
    b. You can modify a job priority in the driver class
    c. You can modify a job priority when you submit the job
    d. You cannot modify a job priority in Hadoop

```
1. public class VoterDriver extends Configured implements Tool {
2.    public int run(String[] args) throws Exception {
3.       if (args.length != 2) {
4.          System.err.println("usage: hadoop jar -classpath $CLASSPATH:Voter.jar Voter.VoterDriver <inputfile>
5. <outputdir>");
6.          System.exit(1);
7.       }
8.       Job job = new Job(getConf());
9.       job.setJarByClass(VoterDriver.class);
10.       job.setMapperClass(VoterMapper.class);
11.       job.setReducerClass(VoterReducer.class);
12.       job.setInputFormatClass(TextInputFormat.class);
13.       job.setOutputKeyClass(Text.class);
14.       job.setOutputValueClass(IntWritable.class);
15.       FileInputFormat.addInputPath(job, new  Path(args[0]));
16.       FileOutputFormat.setOutputPath(job, new Path(args[1]));
17.       return job.waitForCompletion(true) ? 0 : 1;
18.    }

19.    public static void main(String[] args) throws Exception {
20.       Configuration conf = new Configuration();
21.      conf.set("mapreduce.output.key.field.separator", ",");
22.       System.exit(ToolRunner.run(conf, new VoterDriver(), args));
23.    }
24. }
```

```
1. public class VoterMapper  extends Mapper <LongWritable,Text,Text,IntWritable> {
2.    String tempString=null;
3.    private static Log log = LogFactory.getLog(VoterMapper.class);
4.    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
5.       StringTokenizer iterator = new StringTokenizer(value.toString(),",");
6.       if(iterator.countTokens() != 6) {
7.          context.getCounter("MYGROUP", "bad_num_tokens").increment(1);
8.          return;
9.       }
10.       iterator.nextToken();   // Takes off first number
11.       iterator.nextToken();   // Takes off name
12.       tempString = iterator.nextToken().toString();   // Takes age
13.       Integer ageInteger = new Integer(tempString);   // converts age to int
14.       int ageInt = ageInteger.intValue();
15.       if(ageInt < 16 || ageInt > 120) {
16.          log.error("incorrect number of tokens:" + value.toString());
17.          context.getCounter("MYGROUP", "bad_age").increment(1);
18.          return;
19.       }
20.       IntWritable age = new IntWritable(ageInt);
21.       String party = iterator.nextToken().toString();
22.       context.write(new Text(party), age);
23.    }

24. }
```

```
1.        public void reduce(Text key, Iterable<IntWritable> values, Context context) throws IOException,
2.InterruptedException {
3.        long sum = 0L;
4.        int count = 0;
5.        int tempValue = 0;
6.        int max=Integer.MIN_VALUE;
7.        int min=Integer.MAX_VALUE;
8.        for (IntWritable value: values) {
9.            tempAge = value;
10.           tempValue = (new Integer(value.toString())).intValue();
11.        if(tempValue < min) {
12.            min=tempValue;
13.        }
14.        if(tempValue > max) {
15.            max=tempValue;
16.        }
17.        sum+=tempValue;
18.        count++;
19.    }
20.    float mean=sum/count;
21.    context.write(key, new FloatWritable(mean));
22.    }
23.}
```

missing final field

```
 999991 mike falkner 38 socialist 703.09
 999992 luke falkner 61 libertarian 965.43 15683
 999993 calvin xylophone 20 democrat 794.84 21569
 999994 wendy underhill 18 democrat 343.46 4758
 999995 nick  young 57 democrat 324.15 3956
 999996 ethan brown 37 democrat 117.57 19228
 999997 ulysses nixon 55 independent 819.34 27477
 999998 calvin laertes 22 democrat 542.47 16730
 999999 irene thompson 70 green 158.22 10006

1000000 priscilla zipper 13 libertarian 862.71 7137
```

Bad

Socialist
libertarian
democrat
Independent
green

Based on the sample data, `VoterDriver`, `VoterMapper`, and `VoterReducer` classes provided above, write a short 1-2 sentence answer for each question below.

1. What is the record delimiter for the mapper?

"\n"   (I do not see you changing the default).

2. What is the field delimiter for the mapper?

Comma ","   (from StringTokenizer declaration – line #5)

3. What is the output key type and output value type in the mapper?

Key: Text       Value: IntWritable

4. What is the input key type and input value type in the reducer?

Key: Text       Value: IntWritable for each value Collection is Iterable<IntWritable>

5. Assuming there is only one reducer and based strictly on the data provided, how many calls to the `reduce()` method would be made?   since using MapReduce package

+1

Four (socialist is excluded as a bad record)

6. What would be the output from the `map()` method for the following data point?

+1 999991,mike falkner,38,socialist,703.09,20560

Key value pair. Key is "socialist" (type Text) and value is 38 (type int writable)

7. What would be the output from the `reduce()` method for the key democrat?

+1

$2 + 37 + 57 + 18 + 20 = 47 + 57 = 154/5 = 30$ (integer division is what I think Java uses)

Key value pair:
Key: "democrat" (Text)
Value: 30   (FloatWritable)

"democrat, 30.0" to file

8. In the `VoterReducer` code line 20
   a. What is the precision of the mean?

   0 – long divided by int

   b. Rewrite line 20 so that you will see significant digits after the decimal point.  float mean = Sum; mean = mean/count;   (Two lines but this is guaranteed to work. Not most elegant).

9. What is the impact of defining the `sum` and `count` variables as instance variables rather than local variables in the `map()` method?

10. What is the value of the counters `MYGROUP.bad_num_tokens` and `MYGROUP.bad_age` based on the data provided?

One bad age (Priscilla zipper since less than 16). No one is older than 120.

One bad num token (Mike falkner) as he only has five fields. Did not see any others missing a field.