

# **Analysis of reviews on amazon of electronic products.**

Pooja Yendhe

# Proposal

This project proposal is to analyze the reviews of electronics products on amazon to get better understanding of products available under electronics category of amazon. By analyzing the reviews various insights could be found like –

1. Number of Verified and non-verified purchases of product.
2. Average rating of each product.
3. Find user who reviewed the product.
4. Reviews per year.
5. Most reviewed products.
6. Records in each ratings.
7. Find maximum and minimum rating verified product.
8. Total count of reviews.
9. Reviews per productID.
10. Count of reviews on daily basis for all products.
11. Total number of products per rating.

# Proposal

To get this insights, mapreduce in mongodb, hadoop, hive, pig, tableau is implemented on amazon review data of electronics products.

The data set for this analysis is taken from amazon product reviews data(s3)

[https://s3.amazonaws.com/amazon-reviews-pds/tsv/amazon\\_reviews\\_us\\_Electronics\\_v1\\_00.tsv.gz](https://s3.amazonaws.com/amazon-reviews-pds/tsv/amazon_reviews_us_Electronics_v1_00.tsv.gz)

# Features of Dataset

Data format:

Tab ('\t') is the separator in the file, without quote or escape characters. First line in file is header and each line after that corresponds to 1 record.

Columns :

1. marketplace - 2 letter country code; here US as this data is from United States.
2. customer\_id - arbitrary identifier that can be used to group reviews by a single author
3. review\_id - ID of review
4. product\_id. - The unique Product ID the review pertains to.
5. product\_parent - Random identifier that can be used to aggregate reviews for same product.
6. product\_title. - Title of the product.
7. Product\_category. - Broad product category. ie. Electronics.
8. star\_rating. - Rating of the review (1-5 star).
9. helpful\_votes. - Number of helpful reviews.
10. total\_votes. - Total number of reviews received by the product.
11. Vine - Review was written as part of the Vine program.
12. verified\_purchase. - The review is on a verified purchase.
13. review\_headline. - The title of the review.
14. review\_body. - The review text.
15. review\_date. - The date the review was written.

# MapReduce in mongoDB

## 1. Number of Verified and non-verified purchases of product.

```
poojayendhe@Poojas-MacBook-Pro:~/Desktop$ mongoimport --db amazonreviewdata --collection reviews --type tsv --headerline --file '/Users/poojayendhe/Downloads/amazon_reviews_us_Electronics_v1_00.tsv';
2022-08-19T02:37:03.823-0700    connected to: mongodb://localhost/
2022-08-19T02:37:06.824-0700    [.....] amazonreviewdata.reviews    23.5MB/1.61GB (1.4%)
2022-08-19T02:37:09.823-0700    [.....] amazonreviewdata.reviews    47.9MB/1.61GB (2.9%)
2022-08-19T02:37:12.823-0700    [#.....] amazonreviewdata.reviews    73.8MB/1.61GB (4.5%)
2022-08-19T02:37:15.823-0700    [#.....] amazonreviewdata.reviews    100MB/1.61GB (6.1%)
2022-08-19T02:37:18.823-0700    [#.....] amazonreviewdata.reviews    125MB/1.61GB (7.6%)
2022-08-19T02:37:21.822-0700    [##.....] amazonreviewdata.reviews    149MB/1.61GB (9.0%)
2022-08-19T02:37:24.822-0700    [##.....] amazonreviewdata.reviews    174MB/1.61GB (10.6%)
2022-08-19T02:37:27.822-0700    [##.....] amazonreviewdata.reviews    199MB/1.61GB (12.1%)
2022-08-19T02:37:30.822-0700    [##.....] amazonreviewdata.reviews    224MB/1.61GB (13.6%)
2022-08-19T02:37:33.822-0700    [##.....] amazonreviewdata.reviews    248MB/1.61GB (15.1%)
2022-08-19T02:37:36.822-0700    [##..] amazonreviewdata.reviews    270MB/1.61GB (16.4%)
2022-08-19T02:37:39.823-0700    [##..] amazonreviewdata.reviews    296MB/1.61GB (18.0%)
2022-08-19T02:37:42.823-0700    [##..] amazonreviewdata.reviews    319MB/1.61GB (19.4%)
2022-08-19T02:37:45.822-0700    [####..] amazonreviewdata.reviews    343MB/1.61GB (20.9%)
2022-08-19T02:37:48.823-0700    [####..] amazonreviewdata.reviews    369MB/1.61GB (22.4%)
2022-08-19T02:37:51.822-0700    [####..] amazonreviewdata.reviews    393MB/1.61GB (23.9%)
2022-08-19T02:37:54.822-0700    [####..] amazonreviewdata.reviews    416MB/1.61GB (25.3%)
2022-08-19T02:37:57.821-0700    [####..] amazonreviewdata.reviews    440MB/1.61GB (26.8%)
2022-08-19T02:38:00.821-0700    [####..] amazonreviewdata.reviews    464MB/1.61GB (28.2%)
2022-08-19T02:38:03.822-0700    [####..] amazonreviewdata.reviews    489MB/1.61GB (29.7%)
2022-08-19T02:38:06.822-0700    [#####..] amazonreviewdata.reviews    515MB/1.61GB (31.3%)
2022-08-19T02:38:09.821-0700    [#####..] amazonreviewdata.reviews    547MB/1.61GB (33.3%)
2022-08-19T02:38:12.821-0700    [#####..] amazonreviewdata.reviews    584MB/1.61GB (35.5%)
2022-08-19T02:38:15.821-0700    [#####..] amazonreviewdata.reviews    620MB/1.61GB (37.7%)
2022-08-19T02:38:18.821-0700    [#####..] amazonreviewdata.reviews    653MB/1.61GB (39.7%)
2022-08-19T02:38:21.822-0700    [#####..] amazonreviewdata.reviews    686MB/1.61GB (41.7%)
2022-08-19T02:38:24.821-0700    [#####..] amazonreviewdata.reviews    720MB/1.61GB (43.7%)
2022-08-19T02:38:27.821-0700    [#####..] amazonreviewdata.reviews    754MB/1.61GB (45.8%)
2022-08-19T02:38:31.821-0700    [#####..] amazonreviewdata.reviews    785MB/1.61GB (47.7%)
2022-08-19T02:38:33.822-0700    [#####..] amazonreviewdata.reviews    811MB/1.61GB (49.3%)
2022-08-19T02:38:36.821-0700    [#####..] amazonreviewdata.reviews    850MB/1.61GB (51.6%)
2022-08-19T02:38:39.821-0700    [#####..] amazonreviewdata.reviews    885MB/1.61GB (53.7%)
2022-08-19T02:38:42.822-0700    [#####..] amazonreviewdata.reviews    918MB/1.61GB (55.8%)
2022-08-19T02:38:45.822-0700    [#####..] amazonreviewdata.reviews    953MB/1.61GB (57.9%)
2022-08-19T02:38:48.821-0700    [#####..] amazonreviewdata.reviews    989MB/1.61GB (60.1%)
2022-08-19T02:38:51.821-0700    [#####..] amazonreviewdata.reviews    1.000GB/1.61GB (62.3%)
2022-08-19T02:38:54.821-0700    [#####..] amazonreviewdata.reviews    1.040GB/1.61GB (64.9%)
2022-08-19T02:38:57.821-0700    [#####..] amazonreviewdata.reviews    1.098GB/1.61GB (67.7%)
2022-08-19T02:39:00.821-0700    [#####..] amazonreviewdata.reviews    1.136GB/1.61GB (70.0%)
2022-08-19T02:39:03.821-0700    [#####..] amazonreviewdata.reviews    1.178GB/1.61GB (72.5%)
2022-08-19T02:39:06.821-0700    [#####..] amazonreviewdata.reviews    1.210GB/1.61GB (75.0%)
2022-08-19T02:39:09.820-0700    [#####..] amazonreviewdata.reviews    1.258GB/1.61GB (77.7%)
2022-08-19T02:39:12.821-0700    [#####..] amazonreviewdata.reviews    1.298GB/1.61GB (80.2%)
2022-08-19T02:39:15.820-0700    [#####..] amazonreviewdata.reviews    1.338GB/1.61GB (82.5%)
2022-08-19T02:39:18.821-0700    [#####..] amazonreviewdata.reviews    1.380GB/1.61GB (85.7%)
2022-08-19T02:39:21.820-0700    [#####..] amazonreviewdata.reviews    1.430GB/1.61GB (88.7%)
2022-08-19T02:39:24.820-0700    [#####..] amazonreviewdata.reviews    1.468GB/1.61GB (90.9%)
2022-08-19T02:39:27.820-0700    [#####..] amazonreviewdata.reviews    1.516GB/1.61GB (93.9%)
2022-08-19T02:39:30.821-0700    [#####..] amazonreviewdata.reviews    1.568GB/1.61GB (97.1%)
2022-08-19T02:39:33.266-0700    [#####..] amazonreviewdata.reviews    1.610GB/1.61GB (100.0%)
2022-08-19T02:39:33.267-0700    3093869 document(s) imported successfully. 0 document(s) failed to import.
```

# MapReduce in mongoDB

1. Number of Verified and non-verified purchases of product.

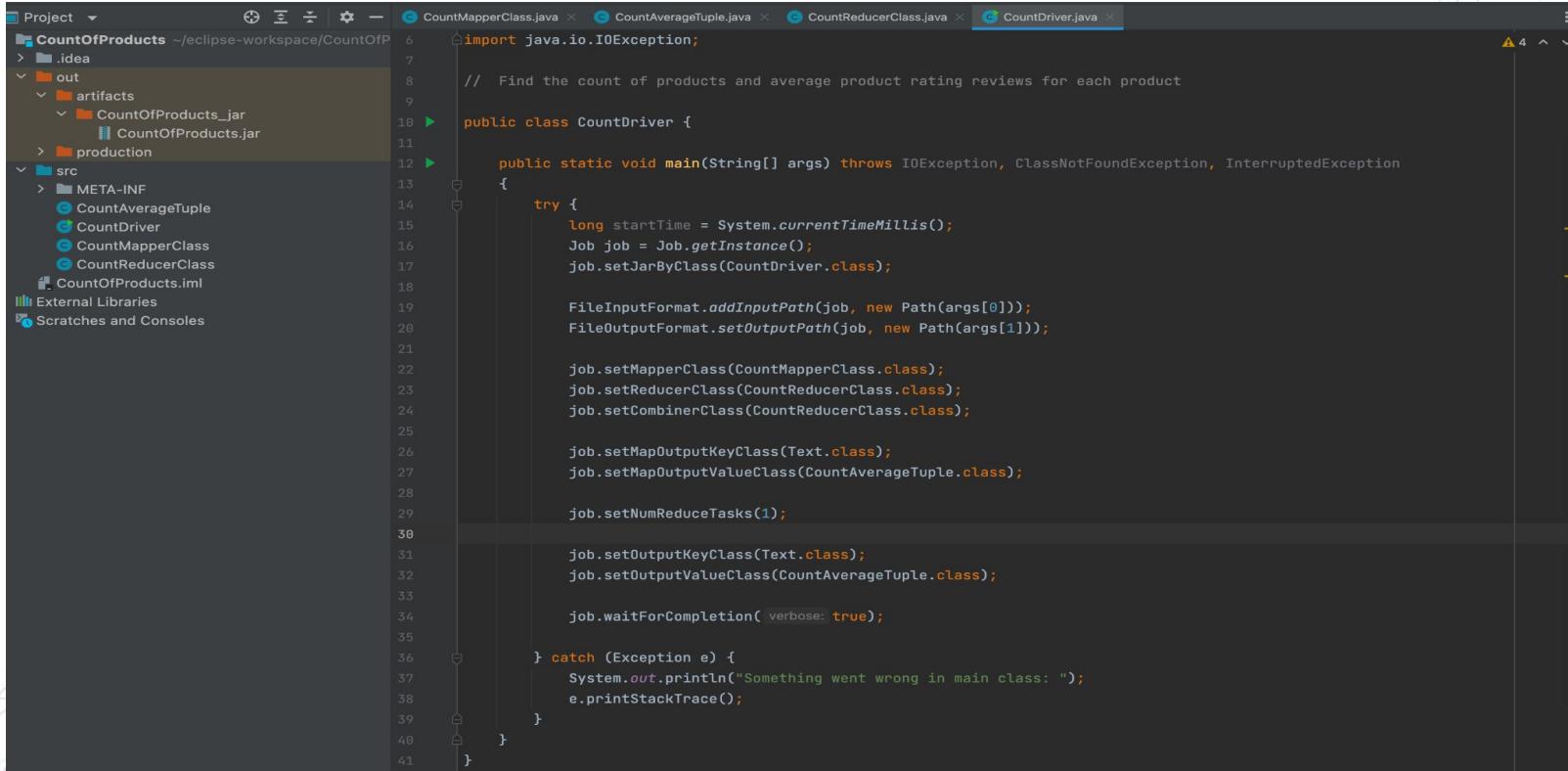
```
[amazonreviewdata> map1 = function () { emit(this.verified_purchase, this.product_id); }
[Function: map1]
[amazonreviewdata> reduce1 = function (key, value) { var cnt = 0; for (var i = 0; i < value.length; i++) { cnt++; } return cnt; }
[Function: reduce1]
[amazonreviewdata> db.reviews.mapReduce(map1, reduce1, {out : "CountOfVerifiedPurchases"});
DeprecationWarning: Collection.mapReduce() is deprecated. Use an aggregation instead.
See https://docs.mongodb.com/manual/core/map-reduce for details.
{ result: 'CountOfVerifiedPurchases', ok: 1 }
amazonreviewdata>

amazonreviewdata>

[amazonreviewdata> db.CountOfVerifiedPurchases.find();
[ { _id: 'Y', value: 396610 }, { _id: 'N', value: 493639 } ]
amazonreviewdata> ]
```

# MapReduce in Hadoop

## 1. Average rating of each product. (Numerical Summarization)



The screenshot shows an IDE interface with a project named "CountOfProducts". The project structure includes a ".idea" folder, an "out" directory containing "artifacts" and "CountOfProducts\_jar/CountOfProducts.jar", a "production" directory, a "src" directory containing "META-INF/CountAverageTuple", "CountDriver", "CountMapperClass", "CountReducerClass", and "CountOfProducts.java", and an "External Libraries" section. The "CountDriver.java" file is open in the editor, displaying the following code:

```
import java.io.IOException;

// Find the count of products and average product rating reviews for each product

public class CountDriver {

    public static void main(String[] args) throws IOException, ClassNotFoundException, InterruptedException {
        try {
            long startTime = System.currentTimeMillis();
            Job job = Job.getInstance();
            job.setJarByClass(CountDriver.class);

            FileInputFormat.addInputPath(job, new Path(args[0]));
            FileOutputFormat.setOutputPath(job, new Path(args[1]));

            job.setMapperClass(CountMapperClass.class);
            job.setReducerClass(CountReducerClass.class);
            job.setCombinerClass(CountReducerClass.class);

            job.setMapOutputKeyClass(Text.class);
            job.setMapOutputValueClass(CountAverageTuple.class);

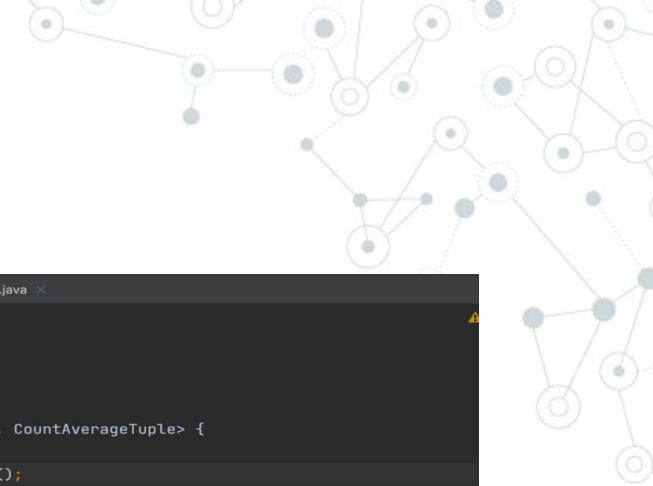
            job.setNumReduceTasks(1);

            job.setOutputKeyClass(Text.class);
            job.setOutputValueClass(CountAverageTuple.class);

            job.waitForCompletion(true);
        } catch (Exception e) {
            System.out.println("Something went wrong in main class: ");
            e.printStackTrace();
        }
    }
}
```

# MapReduce in Hadoop

## 1. Average rating of each product.



The screenshot shows an Eclipse IDE interface with a Java project named "CountOfProducts". The project structure is as follows:

- Project: CountOfProducts (~/eclipse-workspace/CountOfP)
- Artifacts: CountOfProducts.jar, CountOfProducts.iml
- src: META-INF (CountAverageTuple, CountDriver, CountMapperClass, CountReducerClass), CountOfProducts.java
- External Libraries
- Scratches and Consoles

The code editor displays CountMapperClass.java:

```
import java.io.*;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.mapreduce.Mapper;

public class CountMapperClass extends Mapper<LongWritable, Text, Text, CountAverageTuple> {

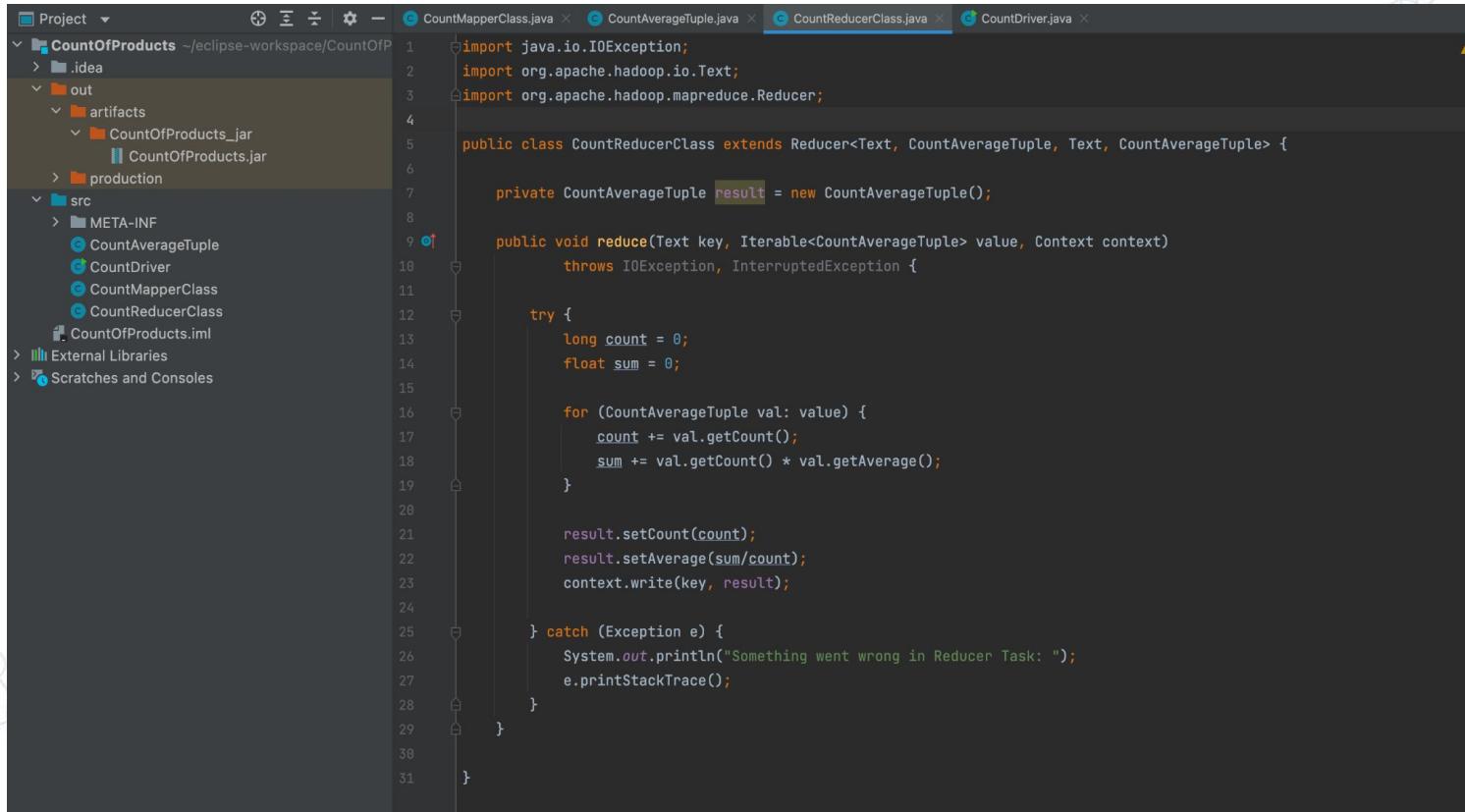
    private CountAverageTuple outCountAverage = new CountAverageTuple();
    private Text id = new Text();

    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
        try {
            String input[] = value.toString().split( regex: "\\\n");
            String productId = input[3].trim();

            if (!productId.isEmpty()) {
                id.set(productId);
                outCountAverage.setCount(Long.valueOf(1));
                outCountAverage.setAverage(Float.valueOf(input[7].trim()));
                context.write(id, outCountAverage);
            }
        } catch (Exception e) {
            System.out.println("Something went wrong in Mapper Task: ");
            e.printStackTrace();
        }
    }
}
```

# MapReduce in Hadoop

## 1. Average rating of each product.



```
import java.io.IOException;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

public class CountReducerClass extends Reducer<Text, CountAverageTuple, Text, CountAverageTuple> {

    private CountAverageTuple result = new CountAverageTuple();

    public void reduce(Text key, Iterable<CountAverageTuple> value, Context context)
        throws IOException, InterruptedException {

        try {
            long count = 0;
            float sum = 0;

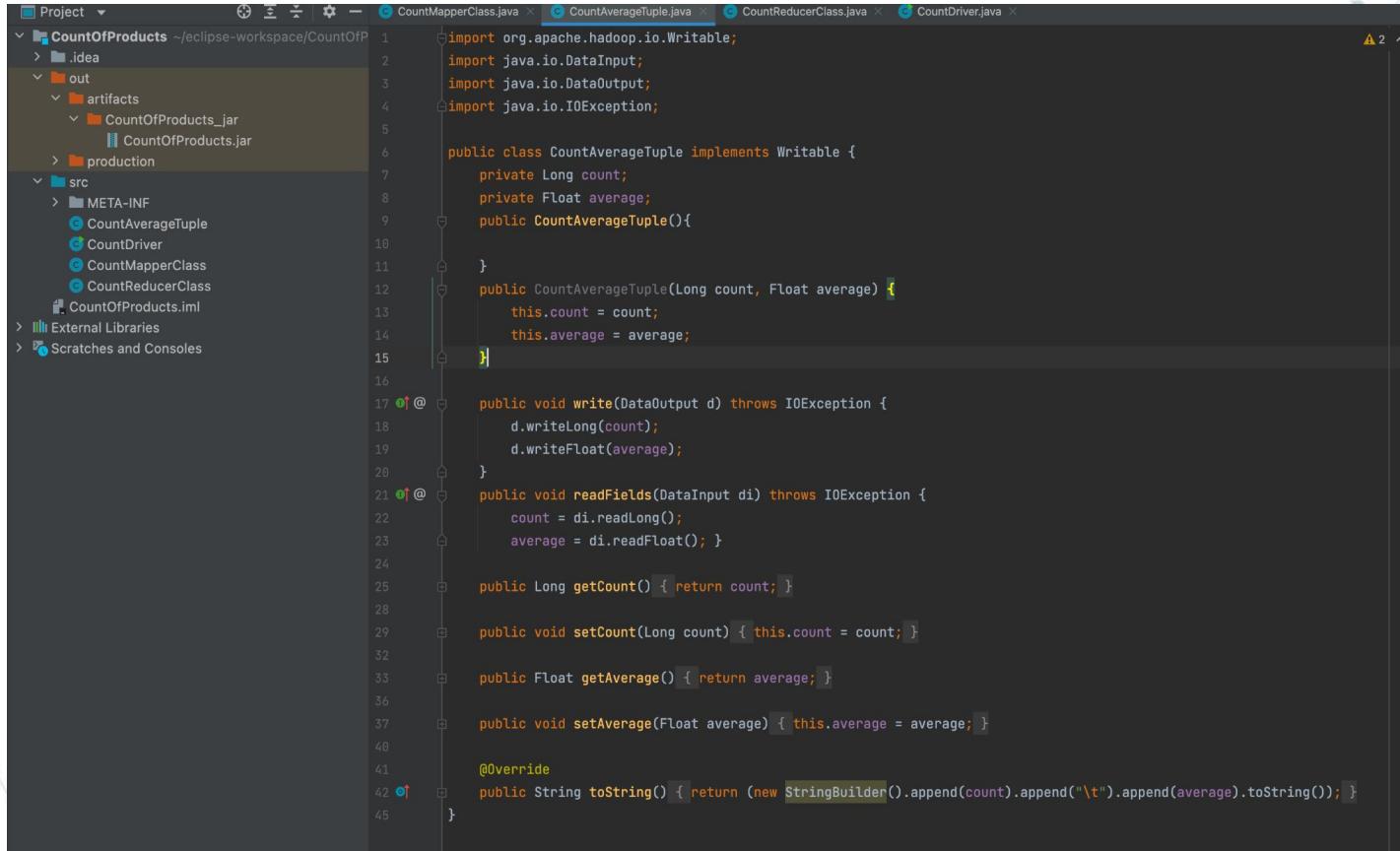
            for (CountAverageTuple val: value) {
                count += val.getCount();
                sum += val.getCount() * val.getAverage();
            }

            result.setCount(count);
            result.setAverage(sum/count);
            context.write(key, result);

        } catch (Exception e) {
            System.out.println("Something went wrong in Reducer Task: ");
            e.printStackTrace();
        }
    }
}
```

# MapReduce in Hadoop

## 1. Average rating of each product.



The screenshot shows the Eclipse IDE interface with a Java project named "CountOfProducts". The project structure is visible on the left, showing files like CountMapperClass.java, CountDriver.java, etc. The main code editor displays CountAverageTuple.java, which defines a class that implements the Writable interface. The class has private fields for count (Long) and average (Float). It includes methods for writing and reading fields from DataInput and DataOutput streams, and methods for getting and setting these values. The code uses standard Java imports and annotations like @Override and @ToString.

```
import org.apache.hadoop.io.Writable;
import java.io.DataInput;
import java.io.DataOutput;
import java.io.IOException;

public class CountAverageTuple implements Writable {
    private Long count;
    private Float average;

    public CountAverageTuple() {
    }

    public CountAverageTuple(Long count, Float average) {
        this.count = count;
        this.average = average;
    }

    @Override
    public void write(DataOutput d) throws IOException {
        d.writeLong(count);
        d.writeFloat(average);
    }

    @Override
    public void readFields(DataInput di) throws IOException {
        count = di.readLong();
        average = di.readFloat();
    }

    public Long getCount() { return count; }

    public void setCount(Long count) { this.count = count; }

    public Float getAverage() { return average; }

    public void setAverage(Float average) { this.average = average; }

    @Override
    public String toString() { return new StringBuilder().append(count).append("\t").append(average).toString(); }
}
```

# MapReduce in Hadoop

## 1. Average rating of each product.

```
poojayendhe@Poojas-MacBook-Pro ~ % 
poojayendhe@Poojas-MacBook-Pro ~ % hadoop jar /Users/poojayendhe/eclipse-workspace/CountOfProducts/out/artifacts/CountOfProducts_jar/CountOfProducts.jar hdfs://datadump/amazon_reviews_us_Electronics_v1_00.tsv hdfs://insights/AverageRating
2022-08-18 01:58:03,999 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
2022-08-18 01:58:04,581 INFO client.DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at /0.0.0:8032
2022-08-18 01:58:05,405 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2022-08-18 01:58:05,459 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/poojayendhe/.staging/job_1660809193636_0005
2022-08-18 01:58:06,503 INFO input.FileInputFormat: Total input files to process : 1
2022-08-18 01:58:07,564 INFO mapreduce.JobSubmitter: number of splits:13
2022-08-18 01:58:08,214 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1660809193636_0005
2022-08-18 01:58:08,214 INFO mapreduce.JobSubmitter: Executing with tokens: []
2022-08-18 01:58:08,450 INFO conf.Configuration: resource-types.xml not found
2022-08-18 01:58:08,451 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2022-08-18 01:58:08,532 INFO impl.YarnClientImpl: Submitted application application_1660809193636_0005
2022-08-18 01:58:08,574 INFO mapreduce.Job: The url to track the job: http://Poojas-MacBook-Pro.local:8088/proxy/application_1660809193636_0005/
2022-08-18 01:58:08,575 INFO mapreduce.Job: Running job: job_1660809193636_0005
2022-08-18 01:58:17,849 INFO mapreduce.Job: Job job_1660809193636_0005 running in uber mode : false
2022-08-18 01:58:17,854 INFO mapreduce.Job: map 0% reduce 0%
2022-08-18 01:58:31,358 INFO mapreduce.Job: map 15% reduce 0%
2022-08-18 01:58:32,382 INFO mapreduce.Job: map 46% reduce 0%
2022-08-18 01:58:42,727 INFO mapreduce.Job: map 54% reduce 0%
2022-08-18 01:58:43,745 INFO mapreduce.Job: map 85% reduce 0%
2022-08-18 01:58:44,760 INFO mapreduce.Job: map 92% reduce 0%
2022-08-18 01:58:48,848 INFO mapreduce.Job: map 100% reduce 0%
2022-08-18 01:58:51,937 INFO mapreduce.Job: map 100% reduce 100%
2022-08-18 01:58:54,031 INFO mapreduce.Job: Job job_1660809193636_0005 completed successfully
2022-08-18 01:58:54,187 INFO mapreduce.Job: Counters: 50
  File System Counters
    FILE: Number of bytes read=12034631
    FILE: Number of bytes written=27932274
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=1726043507
    HDFS: Number of bytes written=3471362
    HDFS: Number of read operations=44
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
    HDFS: Number of bytes read erasure-coded=0
  Job Counters
    Launched map tasks=13
    Launched reduce tasks=1
    Data-local map tasks=13
    Total time spent by all maps in occupied slots (ms)=136038
    Total time spent by all reduces in occupied slots (ms)=7880
    Total time spent by all map tasks (ms)=136038
    Total time spent by all reduce tasks (ms)=7880
    Total vcore-milliseconds taken by all map tasks=136038
    Total vcore-milliseconds taken by all reduce tasks=7880
    Total megabyte-milliseconds taken by all map tasks=139302912
    Total megabyte-milliseconds taken by all reduce tasks=8069120
Map-Reduce Framework
```

# MapReduce in Hadoop

1. Average rating of each product.

```
Total megabyte-milliseconds taken by all map tasks=107802712
Total megabyte-milliseconds taken by all reduce tasks=8069120
Map-Reduce Framework
  Map input records=3093870
  Map output records=3093869
  Map output bytes=71158987
  Map output materialized bytes=12034703
  Input split bytes=1755
  Combine input records=3093869
  Combine output records=481385
  Reduce input groups=185852
  Reduce shuffle bytes=12034703
  Reduce input records=481385
  Reduce output records=185852
  Spilled Records=962770
  Shuffled Maps =13
  Failed Shuffles=0
  Merged Map outputs=13
  GC time elapsed (ms)=2158
  CPU time spent (ms)=0
  Physical memory (bytes) snapshot=0
  Virtual memory (bytes) snapshot=0
  Total committed heap usage (bytes)=5812781056
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0
File Input Format Counters
  Bytes Read=1726041752
File Output Format Counters
  Bytes Written=3471362
poojayendhe@Poojas-MacBook-Pro ~bin %
```

# MapReduce in Hadoop

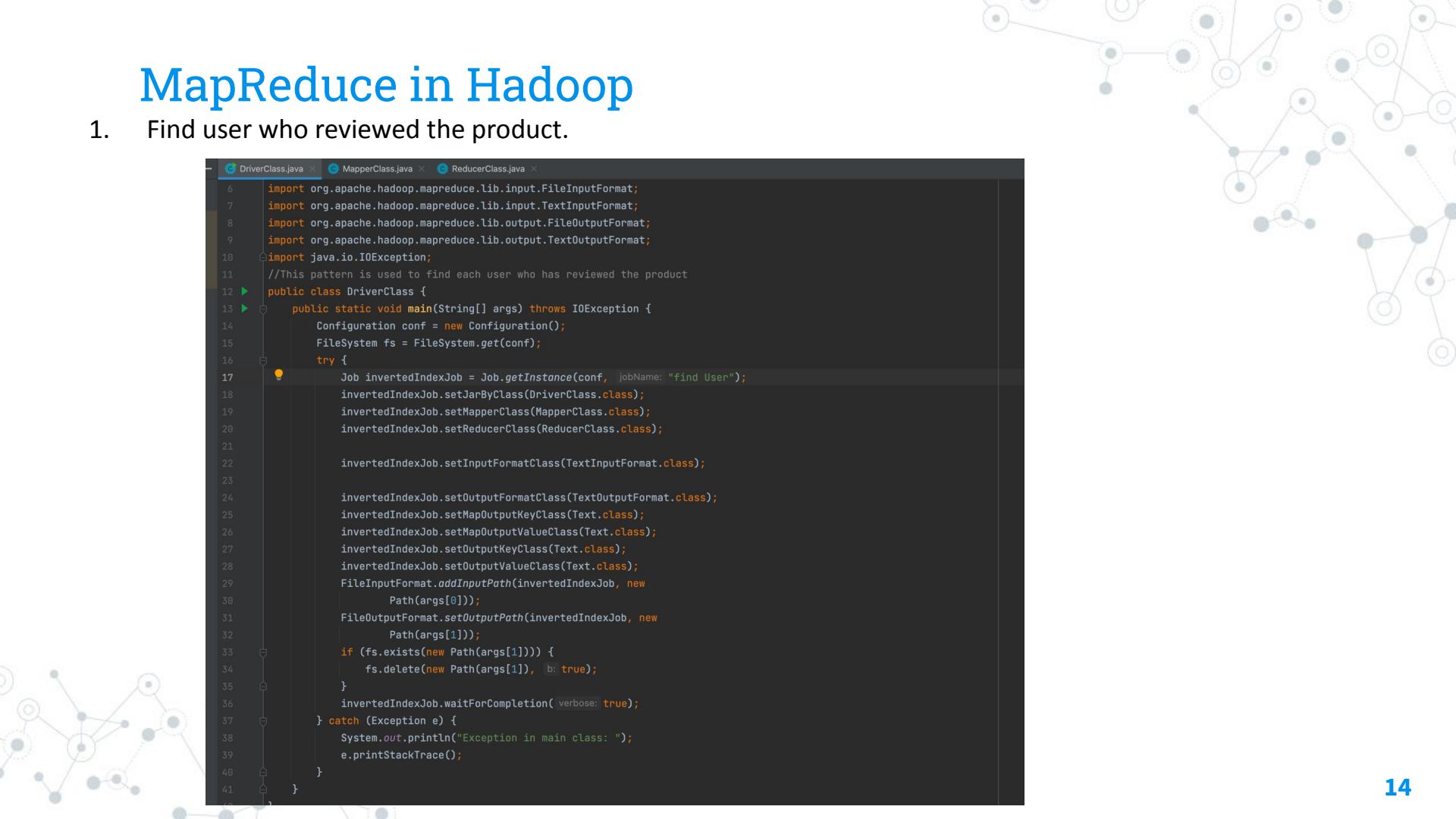
1. Average rating of each product.

The screenshot shows the Hadoop Web UI interface. At the top, there's a navigation bar with links for Overview, Datanodes, Datanode Volume Failures, Snapshot, Startup Progress, and Utilities. Below the navigation bar, the main area displays a "Browse Directory" for the path `/insights/AverageRating`. The directory listing shows two entries: `_SUCCESS` and `part-r-00000`. A modal window titled "File information - part-r-00000" is open over the directory list. This modal contains sections for "Block information" (specifically Block 0), "File contents", and links to "Download", "Head the file (first 32K)", and "Tail the file (last 32K)". The "Block information" section provides details such as Block ID (1073741903), Block Pool ID (BP-499674093-10.0.0.18-1660809147878), Generation Stamp (1079), Size (3471362), and Availability (10.0.0.18). The "File contents" section shows the following data:

ID	Rating
0141186178	5.0
0303532572	4.5
043964383X	5.0
0511189877	4.387755
0528881469	3.1470587
0558835155	3.0
0594296420	4.6842103
0594451647	4.137931

# MapReduce in Hadoop

1. Find user who reviewed the product.



```
DriverClass.java MapperClass.java ReducerClass.java
6 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
7 import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
8 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
9 import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
10 import java.io.IOException;
11 //This pattern is used to find each user who has reviewed the product
12 public class DriverClass {
13     public static void main(String[] args) throws IOException {
14         Configuration conf = new Configuration();
15         FileSystem fs = FileSystem.get(conf);
16         try {
17             Job invertedIndexJob = Job.getInstance(conf, jobName: "find User");
18             invertedIndexJob.setJarByClass(DriverClass.class);
19             invertedIndexJob.setMapperClass(MapperClass.class);
20             invertedIndexJob.setReducerClass(ReducerClass.class);
21
22             invertedIndexJob.setInputFormatClass(TextInputFormat.class);
23
24             invertedIndexJob.setOutputFormatClass(TextOutputFormat.class);
25             invertedIndexJob.setMapOutputKeyClass(Text.class);
26             invertedIndexJob.setMapOutputValueClass(Text.class);
27             invertedIndexJob.setOutputKeyClass(Text.class);
28             invertedIndexJob.setOutputValueClass(Text.class);
29             FileInputFormat.addInputPath(invertedIndexJob, new
30                 Path(args[0]));
31             FileOutputFormat.setOutputPath(invertedIndexJob, new
32                 Path(args[1]));
33             if (fs.exists(new Path(args[1]))) {
34                 fs.delete(new Path(args[1]), true);
35             }
36             invertedIndexJob.waitForCompletion( verbose: true);
37         } catch (Exception e) {
38             System.out.println("Exception in main class: ");
39             e.printStackTrace();
40         }
41     }
}
```

# MapReduce in Hadoop

1. Find user who reviewed the product.

```
DriverClass.java ✘ MapperClass.java ✘ ReducerClass.java ✘
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;
import java.io.IOException;
public class MapperClass extends Mapper<LongWritable, Text, Text, Text> {
    Text prod_cat = new Text();
    private Text productId = new Text();
    private Text userId = new Text();
    @Override
    protected void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException
    {
        if(key.get()==0){
            return;
        }
        try{
            String[] tokens = value.toString().split( regex: "\\\t");
            userId.set(tokens[1]);
            productId.set(tokens[3]);
            context.write(productId, userId);
        } catch(Exception e){
            System.out.println("Exception in Mapper Task:");
            e.printStackTrace();
        }
    }
}
```

# MapReduce in Hadoop

1. Find user who reviewed the product.



```
DriverClass.java × MapperClass.java × ReducerClass.java ×
1 import org.apache.hadoop.io.Text;
2 import org.apache.hadoop.mapreduce.Reducer;
3 import java.io.IOException;
4 public class ReducerClass extends Reducer<Text, Text, Text, Text> {
5     private Text result = new Text();
6     @Override
7     public void reduce(Text key, Iterable<Text> values, Context
8         context)
9         throws IOException, InterruptedException {
10        try {
11            StringBuilder sb = new StringBuilder();
12            boolean first = true;
13            for (Text id : values) {
14                if (first) {
15                    first = false;
16                } else {
17                    sb.append(" , ");
18                }
19                sb.append(id.toString());
20            }
21            result.set(sb.toString());
22            context.write(new Text(string: "product id : "+key), new Text(string: "users : "+result));
23        } catch (Exception e) {
24            System.out.println("Exception in Reducer Task: ");
25            e.printStackTrace();
26        }
27    }
28}
```

# MapReduce in Hadoop

## 1. Find user who reviewed the product.

```
poojayendhe@Poojas-MacBook-Pro ~bin % hadoop jar /Users/poojayendhe/eclipse-workspace/FindUser/out/artifacts/FindUser_jar/FindUser.jar hdfs:/datadump/amazon_reviews_us_Electronics_v1_00.tsv hdfs:/insights:/FindUser
2022-08-18 01:37:51,836 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
2022-08-18 01:37:52,515 INFO client.DefaultNoharmFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8832
2022-08-18 01:37:53,363 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2022-08-18 01:37:53,414 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/poojayendhe/.staging/job_1660809193636_0003
2022-08-18 01:37:54,439 INFO input.FileInputFormat: Total input files to process : 1
2022-08-18 01:37:55,053 INFO mapreduce.JobSubmitter: number of splits:13
2022-08-18 01:37:55,722 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1660809193636_0003
2022-08-18 01:37:55,723 INFO mapreduce.JobSubmitter: Executing with tokens: []
2022-08-18 01:37:55,944 INFO conf.Configuration: resource-types.xml not found
2022-08-18 01:37:55,945 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2022-08-18 01:37:56,023 INFO impl.YarnClientImpl: Submitted application application_1660809193636_0003
2022-08-18 01:37:56,059 INFO mapreduce.Job: The url to track the job: http://Poojas-MacBook-Pro.local:8088/proxy/application_1660809193636_0003/
2022-08-18 01:37:56,869 INFO mapreduce.Job: Running job: job_1660809193636_0003
2022-08-18 01:38:06,355 INFO mapreduce.Job: Job job_1660809193636_0003 running in uber mode : false
2022-08-18 01:38:06,362 INFO mapreduce.Job: map 0% reduce 0%
2022-08-18 01:38:19,944 INFO mapreduce.Job: map 46% reduce 0%
2022-08-18 01:38:31,239 INFO mapreduce.Job: map 92% reduce 0%
2022-08-18 01:38:36,299 INFO mapreduce.Job: map 100% reduce 0%
2022-08-18 01:38:39,338 INFO mapreduce.Job: map 100% reduce 100%
2022-08-18 01:38:42,469 INFO mapreduce.Job: Job job_1660809193636_0003 completed successfully
2022-08-18 01:38:42,633 INFO mapreduce.Job: Counters: 51
  File System Counters
    FILE: Number of bytes read=67725883
    FILE: Number of bytes written=139318628
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=1726943507
    HDFS: Number of bytes written=39268878
    HDFS: Number of read operations=44
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
    HDFS: Number of bytes read erasure-coded=0
  Job Counters
    Killed map tasks=1
    Launched map tasks=13
    Launched reduce tasks=1
    Data-local map tasks=13
    Total time spent by all maps in occupied slots (ms)=129452
    Total time spent by all reduces in occupied slots (ms)=7776
    Total time spent by all map tasks (ms)=129452
    Total time spent by all reduce tasks (ms)=7776
    Total vcore-milliseconds taken by all map tasks=129452
    Total vcore-milliseconds taken by all reduce tasks=7776
    Total megabyte-milliseconds taken by all map tasks=132558848
    Total megabyte-milliseconds taken by all reduce tasks=796264
Map-Reduce Framework
```

# MapReduce in Hadoop

1. Find user who reviewed the product.

```
Total megabyte-milliseconds taken by all reduce tasks=7962624
Map-Reduce Framework
    Map input records=3093870
    Map output records=3093869
    Map output bytes=61538139
    Map output materialized bytes=67725955
    Input split bytes=1755
    Combine input records=0
    Combine output records=0
    Reduce input groups=185852
    Reduce shuffle bytes=67725955
    Reduce input records=3093869
    Reduce output records=185852
    Spilled Records=6187738
    Shuffled Maps =13
    Failed Shuffles=0
    Merged Map outputs=13
    GC time elapsed (ms)=2120
    CPU time spent (ms)=0
    Physical memory (bytes) snapshot=0
    Virtual memory (bytes) snapshot=0
    Total committed heap usage (bytes)=5967970304
Shuffle Errors
    BAD_ID=0
    CONNECTION=0
    IO_ERROR=0
    WRONG_LENGTH=0
    WRONG_MAP=0
    WRONG_REDUCE=0
File Input Format Counters
    Bytes Read=1726041752
File Output Format Counters
    Bytes Written=39268878
poojayendhe@Poojas-MacBook-Pro ~ %
```

# MapReduce in Hadoop

- Find user who reviewed the product.

The screenshot shows the Hadoop File Explorer interface at `localhost:9870/explorer.html#/insights/FindUser`. A modal window titled "File information - part-r-00000" is open, displaying details about the file. The file has a single block, Block ID: 1073741864, with a Block Pool ID of BP-499674093-10.0.0.18-1660809147878. The generation stamp is 1040, size is 39268878 bytes, and it is available on 10.0.0.18. The file contents are listed below, showing a list of product IDs and user counts.

File information - part-r-00000

Download Head the file (first 32K) Tail the file (last 32K)

Block information -- Block 0

Block ID: 1073741864  
Block Pool ID: BP-499674093-10.0.0.18-1660809147878  
Generation Stamp: 1040  
Size: 39268878  
Availability:  
• 10.0.0.18

File contents

```
product id : 0141186178 users : 19408333
product id : 0303532572 users : 2623484 , 26540092
product id : 0439643837 users : 51686608
product id : 0511189877 users : 22381419 , 5044849 , 2623484 , 26540092 , 3321647 , 37833702 , 46495327 , 16957335 , 31268649 , 13867422 , 10047506 , 31618543 , 2096025 , 48373663 , 18415860 , 34598399 , 10095583 , 6281129 , 14808085 , 23849584 , 50157819 , 51078201 , 16947487 , 24038851 , 9889330 , 39474848 , 49447505 , 43882211 , 19772951 , 13258290 , 36108855 , 35686500 , 19222719 , 51766470 , 43297860 , 42388387 ,
```

# MapReduce in Hadoop

## 1. Reviews per year. ( Partitioning)

```
public class DriverClass {  
  
    @ Public static void main(String[] args) throws IOException, InterruptedException, ClassNotFoundException {  
  
        Configuration conf = new Configuration();  
        FileSystem fs = FileSystem.get(conf);  
        Job job = Job.getInstance(conf, "Partitioning method by year");  
  
        job.setJarByClass(DriverClass.class);  
  
        job.setMapperClass(MapperClass.class);  
        job.setMapOutputKeyClass(Text.class);  
        job.setMapOutputValueClass(Text.class);  
  
        //Custom Partitioner:  
        job.setPartitionerClass(YearPartitionerTuple.class);  
  
        job.setReducerClass(ReducerClass.class);  
        job.setOutputKeyClass(Text.class);  
        job.setOutputValueClass(NullWritable.class);  
        job.setNumReduceTasks(14);  
  
        FileInputFormat.addInputPath(job, new Path(args[0]));  
        FileOutputFormat.setOutputPath(job, new Path(args[1]));  
        if (fs.exists(new Path(args[1]))) {  
            fs.delete(new Path(args[1]), b: true);  
        }  
  
        System.exit(job.waitForCompletion(verbose: true) ? 0 : 1);  
    }  
}
```

# MapReduce in Hadoop

## 1. Reviews per year.

```
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;
import java.io.IOException;

public class MapperClass extends Mapper<LongWritable, Text, Text, Text> {

    private Text inputRecord = new Text();
    private Text year = new Text();

    @Override
    protected void map(LongWritable key, Text value, Mapper.Context context) throws IOException, InterruptedException{

        if(key.get()==0){
            return;
        }

        String[] line = value.toString().split( regex: "\\\t");
        String[] yearPart = line[14].split( regex: "-");
        String yearVal = yearPart[2].trim();

        year.set(yearVal);
        inputRecord.set(value);

        context.write(year, inputRecord);
    }
}
```

# MapReduce in Hadoop

## 1. Reviews per year.

```
import org.apache.hadoop.io.NullWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;
import java.io.IOException;

public class ReducerClass extends Reducer<Text, Text, Text, NullWritable> {

    protected void reduce(Text key, Iterable<Text> values, Reducer.Context context) throws IOException, InterruptedException{
        for(Text t: values){
            context.write(t, NullWritable.get());
        }
    }
}
```

# MapReduce in Hadoop

## 1. Reviews per year.

```
public class YearPartitionerTuple extends Partitioner<Text, Text> {  
    @Override  
    public int getPartition(Text key, Text value, int numPartitions){  
        int n=1;  
        if(numPartitions==0){  
            return 0;}  
        else if(key.equals(("99"))){  
            return n % numPartitions;}  
        else if(key.equals(new Text( string: "00"))){  
            return 2 % numPartitions;}  
        else if(key.equals(new Text( string: "01"))){  
            return 3 % numPartitions ; }  
        else if(key.equals(new Text( string: "02"))){  
            return 4 % numPartitions; }  
        else if(key.equals(new Text( string: "03"))){  
            return 5 % numPartitions; }  
        else if(key.equals(new Text( string: "04"))){  
            return 6 % numPartitions; }  
        else if(key.equals(new Text( string: "05"))){  
            return 7 % numPartitions; }  
        else if(key.equals(new Text( string: "06"))){  
            return 8 % numPartitions; }  
        else if(key.equals(new Text( string: "07"))){  
            return 9 % numPartitions; }  
        else if(key.equals(new Text( string: "08"))){  
            return 10 % numPartitions; }  
        else if (key.equals(new Text( string: "09"))){  
            return 11 % numPartitions; }  
        else if (key.equals(new Text( string: "10"))){  
            return 12 % numPartitions; }  
        else if (key.equals(new Text( string: "11"))){  
            return 13 % numPartitions; }  
        else  
        {return 14 % numPartitions; }  
    }
```

# MapReduce in Hadoop

## 1. Reviews per year.

```
poojayendhe@Poojas-MacBook-Pro ~bin %  
poojayendhe@Poojas-MacBook-Pro ~bin % hadoop jar /Users/poojayendhe/eclipse-workspace/YearPartitioner/out/artifacts/YearPartitioner_jar/YearPartitioner.jar hdfs://datadump/amazon_reviews_us_Electronics_v1_00.tsv hdfs://insights/YearPartitionOutput  
2022-08-18 01:50:48.076 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable  
2022-08-18 01:50:48.731 INFO client.DefaultClientProvider: Connecting to ResourceManager at /0.0.0.0:8032  
2022-08-18 01:50:49.551 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.  
2022-08-18 01:50:49.601 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/poojayendhe/.staging/job_1660809193636_0004  
2022-08-18 01:50:50.637 INFO input.FileInputFormat: Total input files to process : 1  
2022-08-18 01:50:51.778 INFO mapreduce.JobSubmission: number of splits:13  
2022-08-18 01:50:52.445 INFO mapreduce.JobSubmitter: Submitting task for job: job_1660809193636_0004  
2022-08-18 01:50:52.448 INFO mapreduce.JobSubmitter: Executing with tokens: []  
2022-08-18 01:50:52.664 INFO conf.Configuration: resource-types.xml not found  
2022-08-18 01:50:52.664 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.  
2022-08-18 01:50:52.743 INFO impl.YarnClientImpl: Submitted application application_1660809193636_0004  
2022-08-18 01:50:52.774 INFO mapreduce.Job: The url to track the job: http://Poojas-MacBook-Pro.local:8088/proxy/application_1660809193636_0004/  
2022-08-18 01:50:52.774 INFO mapreduce.Job: Running job: job_1660809193636_0004  
2022-08-18 01:51:02.032 INFO mapreduce.Job: Job job_1660809193636_0004 running in uber mode : false  
2022-08-18 01:51:02.035 INFO mapreduce.Job: map 0% reduce 0%  
2022-08-18 01:51:22.087 INFO mapreduce.Job: map 46% reduce 0%  
2022-08-18 01:51:36.565 INFO mapreduce.Job: map 77% reduce 0%  
2022-08-18 01:51:37.584 INFO mapreduce.Job: map 92% reduce 0%  
2022-08-18 01:51:45.769 INFO mapreduce.Job: map 100% reduce 0%  
2022-08-18 01:51:46.899 INFO mapreduce.Job: map 100% reduce 7%  
2022-08-18 01:51:47.866 INFO mapreduce.Job: map 100% reduce 14%  
2022-08-18 01:51:48.881 INFO mapreduce.Job: map 100% reduce 29%  
2022-08-18 01:51:54.818 INFO mapreduce.Job: map 100% reduce 35%  
2022-08-18 01:51:56.075 INFO mapreduce.Job: map 100% reduce 42%  
2022-08-18 01:51:57.187 INFO mapreduce.Job: map 100% reduce 58%  
2022-08-18 01:51:58.142 INFO mapreduce.Job: map 100% reduce 64%  
2022-08-18 01:51:59.156 INFO mapreduce.Job: map 100% reduce 71%  
2022-08-18 01:52:03.256 INFO mapreduce.Job: map 100% reduce 79%  
2022-08-18 01:52:04.278 INFO mapreduce.Job: map 100% reduce 100%  
2022-08-18 01:52:06.406 INFO mapreduce.Job: Job job_1660809193636_0004 completed successfully  
2022-08-18 01:52:06.547 INFO mapreduce.Job: Counters: 52  
File System Counters  
FILE: Number of bytes read=3505836171  
FILE: Number of bytes written=5265026081  
FILE: Number of read operations=0  
FILE: Number of large read operations=0  
FILE: Number of write operations=0  
HDFS: Number of bytes read=1726043507  
HDFS: Number of bytes written=1735269919  
HDFS: Number of read operations=109  
HDFS: Number of large read operations=0  
HDFS: Number of write operations=28  
HDFS: Number of bytes read erasure-coded=0  
Job Counters  
Killed map tasks=1  
Killed reduce tasks=1  
Launched map tasks=13  
Launched reduce tasks=14  
Data-local map tasks=13  
Total time spent by all maps in occupied slots (ms)=192926  
Total time spent by all reduces in occupied slots (ms)=124420  
Total time spent by all map tasks (ms)=192926  
Total time spent by all reduce tasks (ms)=124420  
Total vcore-milliseconds taken by all map tasks=192926  
Total vcore-milliseconds taken by all reduce tasks=124420  
Total megabyte-milliseconds taken by all map tasks=19756224  
Total megabyte-milliseconds taken by all reduce tasks=127406080  
Map-Reduce Framework
```

# MapReduce in Hadoop

## 1. Reviews per year.

```
Total megabyte-milliseconds taken by all map tasks=197556224
Total megabyte-milliseconds taken by all reduce tasks=127406080

Map-Reduce Framework
  Map input records=3093870
  Map output records=3093869
  Map output bytes=1740794819
  Map output materialized bytes=1752523888
  Input split bytes=1755
  Combine input records=0
  Combine output records=0
  Reduce input groups=31
  Reduce shuffle bytes=1752523888
  Reduce input records=3093869
  Reduce output records=3093869
  Spilled Records=9281607
  Shuffled Maps =182
  Failed Shuffles=0
  Merged Map outputs=182
  GC time elapsed (ms)=4458
  CPU time spent (ms)=0
  Physical memory (bytes) snapshot=0
  Virtual memory (bytes) snapshot=0
  Total committed heap usage (bytes)=9079095296

Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0

File Input Format Counters
  Bytes Read=1726041752
File Output Format Counters
  Bytes Written=1735269919
poojayendhe@Poojas-MacBook-Pro ~bin %
```

# MapReduce in Hadoop

1. Reviews per year.

localhost:9870/explorer.html#/insights/YearPartitionOutput

## Browse Directory

/insights/YearPartitionOutput

Show 25 entries

Search:

	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
<input type="checkbox"/>	-rw-r--r--	poojayendhe	supergroup	0 B	Aug 18 01:52	1	128 MB	_SUCCESS	
<input type="checkbox"/>	-rw-r--r--	poojayendhe	supergroup	1.03 GB	Aug 18 01:51	1	128 MB	part-r-00000	
<input type="checkbox"/>	-rw-r--r--	poojayendhe	supergroup	0 B	Aug 18 01:51	1	128 MB	part-r-00001	
<input type="checkbox"/>	-rw-r--r--	poojayendhe	supergroup	0 B	Aug 18 01:51	1	128 MB	part-r-00002	
<input type="checkbox"/>	-rw-r--r--	poojayendhe	supergroup	50.81 MB	Aug 18 01:51	1	128 MB	part-r-00003	
<input type="checkbox"/>	-rw-r--r--	poojayendhe	supergroup	54.02 MB	Aug 18 01:51	1	128 MB	part-r-00004	
<input type="checkbox"/>	-rw-r--r--	poojayendhe	supergroup	55.69 MB	Aug 18 01:51	1	128 MB	part-r-00005	
<input type="checkbox"/>	-rw-r--r--	poojayendhe	supergroup	54.82 MB	Aug 18 01:51	1	128 MB	part-r-00006	
<input type="checkbox"/>	-rw-r--r--	poojayendhe	supergroup	53.93 MB	Aug 18 01:51	1	128 MB	part-r-00007	
<input type="checkbox"/>	-rw-r--r--	poojayendhe	supergroup	54.74 MB	Aug 18 01:51	1	128 MB	part-r-00008	
<input type="checkbox"/>	-rw-r--r--	poojayendhe	supergroup	55.36 MB	Aug 18 01:51	1	128 MB	part-r-00009	
<input type="checkbox"/>	-rw-r--r--	poojayendhe	supergroup	54.59 MB	Aug 18 01:52	1	128 MB	part-r-00010	

# MapReduce in Hadoop

## 1. Most reviewed products.

```
public class MostReviewDriver {
    public static void main(String[] args) throws IOException {

        Configuration conf = new Configuration();
        FileSystem fs = FileSystem.get(conf);
        try {
            Job topNProductsJob = Job.getInstance(conf, jobName: "Top N products with most reviews.");
            topNProductsJob.setJarByClass(MostReviewDriver.class);

            int N = 10;
            topNProductsJob.getConfiguration().setInt( name: "N", N);
            topNProductsJob.setInputFormatClass(TextInputFormat.class);
            topNProductsJob.setOutputFormatClass(TextOutputFormat.class);

            topNProductsJob.setMapperClass(MostReviewMapper.class);
            topNProductsJob.setSortComparatorClass(CountComparator.class);
            topNProductsJob.setReducerClass(MostReviewReducer.class);
            topNProductsJob.setNumReduceTasks(1);

            topNProductsJob.setMapOutputKeyClass(IntWritable.class);
            topNProductsJob.setMapOutputValueClass(Text.class);
            topNProductsJob.setOutputKeyClass(IntWritable.class);
            topNProductsJob.setOutputValueClass(Text.class);

            FileInputFormat.setInputPaths(topNProductsJob, new Path(args[0]));
            FileOutputFormat.setOutputPath(topNProductsJob, new Path(args[1]));
            if (fs.exists(new Path(args[1]))) {
                fs.delete(new Path(args[1]), b: true);
            }
            topNProductsJob.waitForCompletion( verbose: true);
        } catch (Exception e) {
            System.out.println("Exception in main class: ");
            e.printStackTrace();
        }
    }
}
```

# MapReduce in Hadoop

## 1. Most reviewed products.

```
public class MostReviewMapper extends Mapper<LongWritable, Text, IntWritable, Text> {  
  
    public void map(LongWritable key, Text value, Context context){  
  
        String[] row = value.toString().split( regex: "\\t");  
        String product_Id;  
        int count;  
        try {  
            product_Id = row[3].trim();  
            count = Integer.parseInt(row[9].trim());  
        }  
        catch (NumberFormatException e){  
            return;  
        }  
        try{  
            Text id = new Text(product_Id);  
            IntWritable productRating = new IntWritable(count);  
            context.write(productRating, id);  
  
        }catch(Exception e){  
            System.out.println("Exception in Mapper Task: ");  
            e.printStackTrace();  
        }  
    }  
}
```

# MapReduce in Hadoop

## 1. Most reviewed products.

```
public class MostReviewReducer extends Reducer<IntWritable, Text, IntWritable, Text> {  
  
    int count = 0;  
    private int N = 10;  
  
    @Override  
    protected void setup(Context context) throws IOException, InterruptedException {  
        // default = 10  
        this.N = context.getConfiguration().getInt("N", defaultValue: 10);  
    }  
  
    @Override  
    public void reduce(IntWritable key, Iterable<Text> value, Context context)  
        throws IOException, InterruptedException{  
        try {  
            for(Text val: value){  
                if(count < N)  
                {  
                    context.write(key, new Text(string: "ProductID: "+val));  
                }  
                count++;  
            }  
        } catch (Exception e) {  
            System.out.println("Exception in Reducer Task: ");  
            e.printStackTrace();  
        }  
    }  
}
```

# MapReduce in Hadoop

1. Most reviewed products.

```
public class CountComparator extends WritableComparator {  
  
    protected CountComparator() {  
  
        super(IntWritable.class, createInstances: true);  
    }  
  
    public int compare(WritableComparable w1, WritableComparable w2) {  
        IntWritable cmp1 = (IntWritable) w1;  
        IntWritable cmp2 = (IntWritable) w2;  
  
        int result = cmp1.get() < cmp2.get() ? 1 : cmp1.get() == cmp2.get() ? 0 : -1;  
        return result;  
    }  
}
```

# MapReduce in Hadoop

## 1. Most reviewed products.

```
2022-08-18 00:56:08,782 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
poojayendhe@Poojas-MacBook-Pro ~bin % hadoop jar /Users/poojayendhe/eclipse-workspace/MostReviewedProducts/out/artifacts/MostReviewedProducts_jar/MostReviewedProducts.jar hdfs://datadump/amazon_reviews_us_Electronics_v1_00.tsv hdfs://insights/mostReviewedProductList
2022-08-18 01:22:28,646 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
2022-08-18 01:22:29,382 INFO client.DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
2022-08-18 01:22:30,264 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2022-08-18 01:22:30,321 INFO mapreduce.ResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/poojayendhe.staging/job_1660809193636_0001
2022-08-18 01:22:30,987 INFO input.FileInputFormat: Total input files to process : 1
2022-08-18 01:22:31,568 INFO mapreduce.JobSubmitter: number of splits:13
2022-08-18 01:22:32,241 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1660809193636_0001
2022-08-18 01:22:32,242 INFO mapreduce.JobSubmitter: Executing with tokens: []
2022-08-18 01:22:32,482 INFO conf.Configuration: resource-types.xml not found
2022-08-18 01:22:32,483 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2022-08-18 01:22:32,938 INFO impl.YarnClientImpl: Submitted application application_1660809193636_0001
2022-08-18 01:22:32,981 INFO mapreduce.Job: The url to track the job: http://Poojas-MacBook-Pro.local:8088/proxy/application_1660809193636_0001/
2022-08-18 01:22:32,982 INFO mapreduce.Job: Running job: job_1660809193636_0001
2022-08-18 01:22:45,418 INFO mapreduce.Job: Job job_1660809193636_0001 running in uber mode : false
2022-08-18 01:22:45,427 INFO mapreduce.Job: map 0% reduce 0%
2022-08-18 01:22:58,104 INFO mapreduce.Job: map 15% reduce 0%
2022-08-18 01:22:59,169 INFO mapreduce.Job: map 46% reduce 0%
2022-08-18 01:23:07,347 INFO mapreduce.Job: map 54% reduce 0%
2022-08-18 01:23:09,434 INFO mapreduce.Job: map 69% reduce 0%
2022-08-18 01:23:10,454 INFO mapreduce.Job: map 92% reduce 0%
2022-08-18 01:23:14,518 INFO mapreduce.Job: map 100% reduce 0%
2022-08-18 01:23:18,599 INFO mapreduce.Job: map 100% reduce 100%
2022-08-18 01:23:20,703 INFO mapreduce.Job: Job job_1660809193636_0001 completed successfully
2022-08-18 01:23:20,858 INFO mapreduce.Job: Counters: 51
  File System Counters
    FILE: Number of bytes read=52595779
    FILE: Number of bytes written=109063488
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=1726043507
    HDFS: Number of bytes written=161
    HDFS: Number of read operations=44
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
    HDFS: Number of bytes read erasure-coded=0
  Job Counters
    Killed map tasks=1
    Launched map tasks=13
    Launched reduce tasks=1
    Data-local map tasks=13
    Total time spent by all maps in occupied slots (ms)=124811
    Total time spent by all reduces in occupied slots (ms)=10107
    Total time spent by all map tasks (ms)=124811
    Total time spent by all reduce tasks (ms)=10107
    Total vcore-milliseconds taken by all map tasks=124811
    Total vcore-milliseconds taken by all reduce tasks=10107
    Total megabyte-milliseconds taken by all map tasks=127886464
    Total megabyte-milliseconds taken by all reduce tasks=10349568
  Map-Reduce Framework
```

# MapReduce in Hadoop

## 1. Most reviewed products.

```
Total megabyte-milliseconds taken by all reduce tasks=10349568
Map-Reduce Framework
  Map input records=3093870
  Map output records=3093869
  Map output bytes=46408035
  Map output materialized bytes=52595851
  Input split bytes=1755
  Combine input records=0
  Combine output records=0
  Reduce input groups=934
  Reduce shuffle bytes=52595851
  Reduce input records=3093869
  Reduce output records=10
  Spilled Records=6187738
  Shuffled Maps =13
  Failed Shuffles=0
  Merged Map outputs=13
  GC time elapsed (ms)=1919
  CPU time spent (ms)=0
  Physical memory (bytes) snapshot=0
  Virtual memory (bytes) snapshot=0
  Total committed heap usage (bytes)=5808586752
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0
File Input Format Counters
  Bytes Read=1726041752
File Output Format Counters
  Bytes Written=161
poojayendhe@Poojas-MacBook-Pro ~ %
```

# MapReduce in Hadoop

## 1. Most reviewed products.

The screenshot shows a web-based Hadoop file explorer interface at `localhost:9870/explorer.html#/insights/mostReviewedProductList`. The main page displays a "Browse Directory" view with a list of two entries under the path `/insights/mostReviewedProductList`. The entries are "part-r-00000" and "SUCCESS". The "part-r-00000" entry has a "Download" link and a preview of its contents. A modal window titled "File information - part-r-00000" is open, showing details for Block 0: Block ID: 1073741854, Block Pool ID: BP-499674093-10.0.0.18-1660809147878, Generation Stamp: 1030, Size: 271, and Availability: 10.0.0.18. Below this, the "File contents" section lists product IDs and their review counts:

ProductID	Count
B0001iX6PM	12944
B000J36XR2	9072
B003EM6AOG	8680
B001FA1O18	6353
B0001iX6PM	5546
B002M3SOBU	4595
B0054JJ0QW	4556
B003ELYQGG	4341

# MapReduce in Hadoop

## 1. Records in each ratings.

```
public class BinningByRatingDriverClass {  
  
    public static void main(String[] args) throws IOException, InterruptedException, ClassNotFoundException {  
  
        Configuration conf = new Configuration();  
        FileSystem fs = FileSystem.get(conf);  
  
        try {  
  
            Job binningJob = Job.getInstance(conf, jobName: "Binning Pattern");  
            binningJob.setJarByClass(BinningByRatingDriverClass.class);  
  
            binningJob.setMapperClass(BinningRatingOfProductMapperClass.class);  
            binningJob.setMapOutputKeyClass(Text.class);  
            binningJob.setMapOutputValueClass(NullWritable.class);  
            binningJob.setNumReduceTasks(1);  
  
            FileInputFormat.setInputPaths(binningJob, new Path(args[0]));  
            FileOutputFormat.setOutputPath(binningJob, new Path(args[1]));  
            if (fs.exists(new Path(args[1]))) {  
                fs.delete(new Path(args[1]), b: true);  
            }  
  
            MultipleOutputs.addNamedOutput(binningJob, namedOutput: "bins", TextOutputFormat.class, Text.class, NullWritable.class);  
            MultipleOutputs.setCountersEnabled(binningJob, enabled: true);  
  
            System.exit(binningJob.waitForCompletion( verbose: true) ? 0 : 1);  
  
        } catch (Exception e ) {  
            System.out.println("Exception in driver class: ");  
            e.printStackTrace();  
        }  
    }  
}
```

# MapReduce in Hadoop

## 1. Records in each ratings.

```
public class BinningRatingOfProductMapperClass extends Mapper<LongWritable, Text, Text, NullWritable> {  
    private MultipleOutputs<Text, NullWritable> output = null;  
  
    @Override  
    protected void setup(Context context){  
        output = new MultipleOutputs(context);  
    }  
  
    @Override  
    protected void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException{  
        try {  
            if(key.get()==0) { return; }  
            String[] token = value.toString().split( regex: "\\Wt");  
            String rating = token[7].trim();  
            if(rating.equals("1")){  
                output.write( namedOutput: "bins", value, NullWritable.get(), baseOutputPath: "RatingStar1"); }  
            if(rating.equals("2")){  
                output.write( namedOutput: "bins", value, NullWritable.get(), baseOutputPath: "RatingStar2"); }  
            if(rating.equals("3")){  
                output.write( namedOutput: "bins", value, NullWritable.get(), baseOutputPath: "RatingStar3"); }  
            if(rating.equals("4")){  
                output.write( namedOutput: "bins", value, NullWritable.get(), baseOutputPath: "RatingStar4"); }  
            if(rating.equals("5")){  
                output.write( namedOutput: "bins", value, NullWritable.get(), baseOutputPath: "RatingStar5"); }  
        } catch (Exception e) {  
            System.out.println("Exception in Mapper Task: ");  
            e.printStackTrace(); }  
    }  
  
    @Override  
    protected void cleanup(Context context) throws IOException, InterruptedException{  
        output.close();  
    }  
}
```

# MapReduce in Hadoop

## 1. Records in each ratings.

```
poojayendhe@Poojas-MacBook-Pro ~bin %  
poojayendhe@Poojas-MacBook-Pro ~bin % hadoop jar /Users/poojayendhe/eclipse-workspace/BinningRatingOfProduct/out/artifacts/BinningRatingOfProduct_jar/BinningRatingOfProduct.jar hdfs://datadump/amazon_reviews_us_Electronics_v1_00.tsv  
WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable  
2022-08-18 02:14:23,669 INFO client.DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8832  
2022-08-18 02:14:24,356 INFO client.DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8832  
2022-08-18 02:14:25,180 INFO mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.  
2022-08-18 02:14:25,760 INFO input.FileInputFormat: Total input files to process : 1  
2022-08-18 02:14:26,884 INFO mapreduce.JobSubmitter: number of splits:13  
2022-08-18 02:14:27,538 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1660809193636_0006  
2022-08-18 02:14:27,540 INFO mapreduce.JobSubmitter: Executing with tokens: []  
2022-08-18 02:14:27,779 INFO conf.Configuration: resource-types.xml not found  
2022-08-18 02:14:27,780 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.  
2022-08-18 02:14:28,891 INFO impl.YarnClientImpl: Submitted application application_1660809193636_0006  
2022-08-18 02:14:28,144 INFO mapreduce.Job: The url to track the job: http://Poojas-MacBook-Pro.local:8088/proxy/application_1660809193636_0006  
2022-08-18 02:14:28,145 INFO mapreduce.Job: Running job: job_1660809193636_0006  
2022-08-18 02:14:38,510 INFO mapreduce.Job: Job job_1660809193636_0006 running in uber mode : false  
2022-08-18 02:14:38,516 INFO mapreduce.Job: map 0% reduce 0%  
2022-08-18 02:14:55,156 INFO mapreduce.Job: map 23% reduce 0%  
2022-08-18 02:14:56,181 INFO mapreduce.Job: map 46% reduce 0%  
2022-08-18 02:15:10,716 INFO mapreduce.Job: map 69% reduce 0%  
2022-08-18 02:15:11,735 INFO mapreduce.Job: map 92% reduce 0%  
2022-08-18 02:15:17,822 INFO mapreduce.Job: map 100% reduce 100%  
2022-08-18 02:15:19,916 INFO mapreduce.Job: Job job_1660809193636_0006 completed successfully  
2022-08-18 02:15:20,066 INFO mapreduce.Job: Counters: 55  
File System Counters  
FILE: Number of bytes read=6  
FILE: Number of bytes written=3867168  
FILE: Number of read operations=0  
FILE: Number of large read operations=0  
FILE: Number of write operations=0  
HDFS: Number of bytes read=1726643507  
HDFS: Number of bytes written=1725988312  
HDFS: Number of read operations=148  
HDFS: Number of large read operations=0  
HDFS: Number of write operations=132  
HDFS: Number of bytes read erasure-coded=0  
Job Counters  
Launched map tasks=13  
Launched reduce tasks=1  
Data-local map tasks=13  
Total time spent by all maps in occupied slots (ms)=182201  
Total time spent by all reduces in occupied slots (ms)=6428  
Total time spent by all map tasks (ms)=182201  
Total time spent by all reduce tasks (ms)=6428  
Total vcore-milliseconds taken by all map tasks=182201  
Total vcore-milliseconds taken by all reduce tasks=6428  
Total megabyte-milliseconds taken by all map tasks=186573824  
Total megabyte-milliseconds taken by all reduce tasks=6582272  
Map-Reduce Framework  
MapReduce MapSpills=0  
MapReduce ReduceSpills=0
```

# MapReduce in Hadoop

## 1. Records in each ratings.

```
Total megabyte-milliseconds taken by all reduce tasks=6582272
Map-Reduce Framework
  Map input records=3093870
  Map output records=0
  Map output bytes=0
  Map output materialized bytes=78
  Input split bytes=1755
  Combine input records=0
  Combine output records=0
  Reduce input groups=0
  Reduce shuffle bytes=78
  Reduce input records=0
  Reduce output records=0
  Spilled Records=0
  Shuffled Maps =13
  Failed Shuffles=0
  Merged Map outputs=13
  GC time elapsed (ms)=3386
  CPU time spent (ms)=0
  Physical memory (bytes) snapshot=0
  Virtual memory (bytes) snapshot=0
  Total committed heap usage (bytes)=5865734144
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0
File Input Format Counters
  Bytes Read=1726041752
File Output Format Counters
  Bytes Written=0
org.apache.hadoop.mapreduce.lib.output.MultipleOutputs
  RatingStar1=358120
  RatingStar2=179180
  RatingStar3=238587
  RatingStar4=536821
  RatingStar5=1781161
poojayendhe@Poojas-MacBook-Pro ~ %
```

# MapReduce in Hadoop

## 1. Records in each ratings.

The screenshot shows a Hadoop File Explorer interface on a web browser at `localhost:9870/explorer.html#/insights/binningByRatings`. The main view displays a directory listing for the path `/insights/binningByRatings`, showing 25 entries. A modal dialog is open, providing detailed information about a specific file:

**Block information -- Block 0**

Block ID: 1073741918  
Block Pool ID: BP-499674093-10.0.0.18-1660809147878  
Generation Stamp: 1094  
Size: 15629187  
Availability: • 10.0.0.18

**File contents**

```
US 38487968 R1EBPM82ENI67M B000NU4OTA 72265257 LIMTECH Wall
charger + USB Hotsync & Charging Dock Cradle desktop Charger for Apple IPOD Shuffle 2nd
Generation MP3 Player Electronics 1 0 0 N Y One Star Did not work at all.
2015-08-31
US 31463514 R3T9GZS2TMXZGM B0035PBHX6 249533961 Cobay 8 GB 1.8-Inch
Video MP3 Player with FM Radio Electronics 1 0 0 N Y One Star Breaks
very easily, and takes a while to load music 2015-08-31
US 47537250 R20TC495KA8WVA B00990Z4W6 57455227 Crosley CR8005D-
```

A "Close" button is located at the bottom right of the modal.

The right side of the interface shows a sidebar with icons for folder, file, search, and other navigation options. A search bar is also present.

# Apache Pig

## 1. Count of reviews on daily basis for all products.

```
2022-08-18 09:54:33,549 [main] WARN org.apache.pig.PigServer - ATS is disabled since yarn.timeline-service.enabled set to false
grunt> data = LOAD 'hdfs://datadump/amazon_reviews_us_Electronics_v1_00.tsv' AS (marketplace, customer_id,review_id, product_id, product_parent, product_title, product_category,star_rating, helpful_votes, total_votes, vine, verified_purchase, review_headline,review_body, review_date);
grunt> grouped = GROUP data by review_date;
grunt> daily_reviews = FOREACH grouped GENERATE group as review_date, COUNT(data.review_id) as count;
grunt> order_by_data = ORDER daily_reviews BY count DESC;
grunt> store order_by_data into 'hdfs://pig/output/dailyReviews';
2022-08-18 09:55:16,854 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.textoutputformat.separator is deprecated. Instead, use mapreduce.output.textoutputformat.separator
2022-08-18 09:55:16,875 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: GROUP_BY, ORDER_BY
2022-08-18 09:55:16,917 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set.. will not generate code.
2022-08-18 09:55:16,946 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, ConstantCalculator, GroupByConstParallelSetter, LimitOptimizer, LoadtypeCastInserter, MergeFilter, MergeForEach, NestedLimitOptimizer, PartitionFilterOptimizer, PredicatePushdownOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInsertor]}
2022-08-18 09:55:16,998 [main] INFO org.apache.pig.impl.util.SpillableMemoryManager - Selected heap (PS Old Gen) of size 699400192 to monitor. collectionUsageThreshold = 489580128, usageThreshold = 489580128
2022-08-18 09:55:17,066 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - File concatenation threshold: 100 optimistic? false
2022-08-18 09:55:17,128 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.CombineOptimizerUtil - Choosing to move algebraic foreach to combiner
2022-08-18 09:55:17,145 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.SecondaryKeyOptimizerMR - Using Secondary Key Optimization for MapReduce node scope=49
2022-08-18 09:55:17,156 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 3
2022-08-18 09:55:17,156 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 3
2022-08-18 09:55:17,234 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8032
2022-08-18 09:55:17,421 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
2022-08-18 09:55:17,432 [main] INFO org.apache.pig.tools.pigstats.MRScriptState - Pig script settings are added to the job
2022-08-18 09:55:17,442 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.reduce.markreset.buffer.percent is deprecated. Instead, use mapreduce.reduce.markreset.buffer.percent
2022-08-18 09:55:17,443 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3
2022-08-18 09:55:17,446 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.output.compress is deprecated. Instead, use mapreduce.output.fileoutputformat.compress
2022-08-18 09:55:17,449 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Reduce phase detected, estimating # of required reducers.
2022-08-18 09:55:17,450 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Using reducer estimator: org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputSizeReducerEstimator
2022-08-18 09:55:17,464 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputSizeReducerEstimator - BytesPerReducer=1000000000 maxReducers=999 totalInputFileSize=1725988504
2022-08-18 09:55:17,464 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting Parallelism to 2
2022-08-18 09:55:17,465 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.reduce.tasks is deprecated. Instead, use mapreduce.job.reduces
2022-08-18 09:55:17,465 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - This job cannot be converted run in-process
2022-08-18 09:55:17,472 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.submit.replication is deprecated. Instead, use mapreduce.client.submit.file.replication
2022-08-18 09:55:17,856 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Added jar file:/opt/homebrew/Cellar/pig/0.17.0_1/libexec/pig-0.17.0-core-h2.jar to DistributedCache through /tmp/tmp891994362/tmp-1112323073/pig-0.17.0-core-h2.jar
2022-08-18 09:55:18,325 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Added jar file:/opt/homebrew/Cellar/pig/0.17.0_1/libexec/lib/automaton-1.11-8.jar to DistributedCache through /tmp/tmp891994362/tmp-882908173/automaton-1.11-8.jar
2022-08-18 09:55:18,815 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Added jar file:/opt/homebrew/Cellar/pig/0.17.0_1/libexec/lib/antlr-runtime-3.4.jar to DistributedCache through /tmp/tmp891994362/tmp1366334706/antlr-runtime-3.4.jar
2022-08-18 09:55:19,327 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Added jar file:/opt/homebrew/Cellar/pig/0.17.0_1/libexec/lib/joda-time-2.9.3.jar to DistributedCache through /tmp/tmp891994362/tmp1800291863/joda-time-2.9.3.jar
2022-08-18 09:55:19,382 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Setting up single store job
2022-08-18 09:55:19,387 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.
2022-08-18 09:55:19,388 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache
2022-08-18 09:55:19,388 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Setting key [pig.schematuple.classes] with classes to deserilize []
2022-08-18 09:55:19,488 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.
2022-08-18 09:55:19,491 [JobControl] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8032
2022-08-18 09:55:19,519 [JobControl] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.task.id is deprecated. Instead, use mapreduce.task.attempt.id
2022-08-18 09:55:19,980 [JobControl] INFO org.apache.hadoop.mapreduce.JobResourceUploader - Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/poojayendhe/.staging/job_1660836347592_0009
2022-08-18 09:55:20,000 [JobControl] WARN org.apache.hadoop.mapreduce.JobResourceUploader - No job jar file set. User classes may not be found. See Job or Job#setJar(String).
2022-08-18 09:55:20,041 [JobControl] INFO org.apache.pig.builtin.PigStorage - Using PigTextInputFormat
2022-08-18 09:55:20,048 [JobControl] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input files to process : 1
2022-08-18 09:55:20,048 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2022-08-18 09:55:20,092 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths (combined) to process : 13
2022-08-18 09:55:21,054 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - number of splits:13
2022-08-18 09:55:21,138 [JobControl] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
```

# Apache Pig

## 1. Count of reviews on daily basis for all products.

```
2022-08-18 09:55:23,712 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - Submitting tokens for job: job_1660836347592_0009
2022-08-18 09:55:23,713 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - Executing with tokens: []
2022-08-18 09:55:23,861 [JobControl] INFO org.apache.hadoop.mapred.YARNRunner - Job jar is not present. Not adding any jar to the list of resources.
2022-08-18 09:55:23,932 [JobControl] INFO org.apache.hadoop.conf.Configuration - resource-types.xml not found
2022-08-18 09:55:23,932 [JobControl] INFO org.apache.hadoop.yarn.util.resource.ResourceUtils - Unable to find 'resource-types.xml'.
2022-08-18 09:55:22,630 [JobControl] INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl - Submitted application application_1660836347592_0009
2022-08-18 09:55:22,664 [JobControl] INFO org.apache.hadoop.mapreduce.Job - The url to track the job: http://Poojas-MacBook-Pro.local:8088/proxy/application_1660836347592_0009/
2022-08-18 09:55:22,665 [main] INFO org.apache.pig.backend.hadoop.executionengine.MapReduceLauncher - HadoopJobId: job_1660836347592_0009
2022-08-18 09:55:22,665 [main] INFO org.apache.pig.backend.hadoop.executionengine.MapReduceLauncher - Processing aliases daily_reviews,data,grouped
2022-08-18 09:55:22,665 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - detailed locations: M: data[1,7],data[-1,-1],daily_reviews[3,16],grouped[2,10] C: daily_reviews[3,16]
2022-08-18 09:55:22,684 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - % complete
2022-08-18 09:55:22,684 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Running jobs are [job_1660836347592_0009]
2022-08-18 09:55:49,449 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 5% complete
2022-08-18 09:55:49,455 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Running jobs are [job_1660836347592_0009]
2022-08-18 09:56:04,785 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 15% complete
2022-08-18 09:56:04,787 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Running jobs are [job_1660836347592_0009]
2022-08-18 09:56:14,333 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 33% complete
2022-08-18 09:56:14,333 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Running jobs are [job_1660836347592_0009]
2022-08-18 09:56:17,418 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8032
2022-08-18 09:56:17,467 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 09:56:19,280 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8032
2022-08-18 09:56:19,301 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 09:56:19,301 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8032
2022-08-18 09:56:19,360 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 09:56:19,450 [main] INFO org.apache.pig.tools.pigstats.mapreduce.MRScriptState - Pig script settings are added to the job
2022-08-18 09:56:19,453 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3
2022-08-18 09:56:19,458 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Reduce phase detected, estimating # of required reducers.
2022-08-18 09:56:19,458 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Using reducer estimator: org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputSizeReducerEstimator
2022-08-18 09:56:19,497 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputSizeReducerEstimator - BytesPerReducer=1000000000 maxReducers=999 totalInputFileSize=109271
2022-08-18 09:56:19,498 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting Parallelism to 1
2022-08-18 09:56:19,498 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - This job cannot be converted run in-process
2022-08-18 09:56:28,029 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Added jar file:/opt/homebrew/Cellar/pig/0.17.0.1/libexec/pig-0.17.0-core-h2.jar to DistributedCache through /tmp/temp891994362/tmp-1082209908/pig-0.17.0-core-h2.jar
2022-08-18 09:56:28,086 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Added jar file:/opt/homebrew/Cellar/pig/0.17.0.1/libexec/lib/automaton-1.11-8.jar to DistributedCache through /tmp/temp891994362/tmp2054683026/automaton-1.11-8.jar
2022-08-18 09:56:28,124 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Added jar file:/opt/homebrew/Cellar/pig/0.17.0.1/libexec/lib/antlr-runtime-3.4.jar to DistributedCache through /tmp/temp891994362/tmp1231263656/antlr-runtime-3.4.jar
2022-08-18 09:56:28,154 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Added jar file:/opt/homebrew/Cellar/pig/0.17.0.1/libexec/lib/joda-time-2.9.3.jar to DistributedCache through /tmp/temp891994362/tmp-1198089679/joda-time-2.9.3.jar
2022-08-18 09:56:28,174 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting up single store job
2022-08-18 09:56:28,177 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.
2022-08-18 09:56:28,178 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache
2022-08-18 09:56:28,178 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Setting key [pig.schematuple.classes] with classes to deserialize []
2022-08-18 09:56:28,233 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.
2022-08-18 09:56:28,244 [JobControl] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8032
2022-08-18 09:56:20,258 [JobControl] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
2022-08-18 09:56:20,300 [JobControl] INFO org.apache.hadoop.mapreduce.JobResourceUploader - Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/poojayendhe/.staging/job_1660836347592_0010
2022-08-18 09:56:20,316 [JobControl] WARN org.apache.hadoop.mapreduce.JobResourceUploader - No job jar file set. User classes may not be found. See Job or Job#setJar(String).
2022-08-18 09:56:20,350 [JobControl] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input files to process : 2
2022-08-18 09:56:20,351 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 2
2022-08-18 09:56:20,351 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths (combined) to process : 1
2022-08-18 09:56:20,426 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - number of splits:1
2022-08-18 09:56:20,926 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - Submitting tokens for job: job_1660836347592_0010
2022-08-18 09:56:20,927 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - Executing with tokens: []
2022-08-18 09:56:20,939 [JobControl] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Job jar is not present. Not adding any jar to the list of resources.
```

# Apache Pig

## 1. Count of reviews on daily basis for all products.

```
2022-08-18 09:56:50.391 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - detailed locations: M: order_by_data[4,16] C: R: 83% complete
2022-08-18 09:57:04.784 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Running jobs are [job_1660836347592_0011]
2022-08-18 09:57:09.737 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Running jobs are [job_1660836347592_0011]
2022-08-18 09:57:15.897 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8832
2022-08-18 09:57:15.946 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 09:57:16.271 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8832
2022-08-18 09:57:16.284 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 09:57:16.377 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8832
2022-08-18 09:57:16.388 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 09:57:16.424 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
2022-08-18 09:57:16.494 [main] INFO org.apache.pig.tools.pigstats.mapreduce.SimplePigStats - Script Statistics:

HadoopVersion PigVersion UserId StartedAt FinishedAt Features
3.3.3 0.17.0 poojayendhe 2022-08-18 09:55:17 2022-08-18 09:57:16 GROUP_BY,ORDER_BY

Success!

Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReducetime Alias Feature Outputs
job_1660836347592_0009 13 2 17 5 14 13 8 8 8 daily_reviews,data_grouped GROUP_BY,COMBINER
job_1660836347592_0010 1 1 3 3 3 4 4 4 4 order_by_data SAMPLER
job_1660836347592_0011 1 1 3 3 3 3 3 3 3 order_by_data ORDER_BY hdfs:/pig_output/dailyReviews,
Input(s):
Successfully read 3093870 records (1726046913 bytes) from: "hdfs:/datadump/amazon_reviews_us_Electronics_v1_00.tsv"

Output(s):
Successfully stored 5905 records (86483 bytes) in: "hdfs:/pig_output/dailyReviews"

Counters:
Total records written : 5905
Total bytes written : 86483
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1660836347592_0009 -> job_1660836347592_0010,
job_1660836347592_0010 -> job_1660836347592_0011,
job_1660836347592_0011

2022-08-18 09:57:16.497 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8832
2022-08-18 09:57:16.588 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 09:57:16.549 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8832
2022-08-18 09:57:16.556 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 09:57:16.587 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8832
2022-08-18 09:57:16.592 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 09:57:16.622 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8832
2022-08-18 09:57:16.638 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 09:57:16.658 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8832
2022-08-18 09:57:16.653 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 09:57:16.673 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8832
2022-08-18 09:57:16.697 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 09:57:16.728 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8832
2022-08-18 09:57:16.733 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 09:57:16.752 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8832
2022-08-18 09:57:16.769 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 09:57:16.794 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8832
2022-08-18 09:57:16.802 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 09:57:16.832 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
```

# Apache Pig

1. Count of reviews on daily basis for all products.

localhost:9870/explorer.html#/pig\_output/dailyReviews

Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities

## Browse Directory

/pig\_output/dailyReviews

Show 25 entries

Permission	Owner
-rW-r--r--	poojayer
-rW-r--r--	poojayer

Showing 1 to 2 of 2 entries

Hadoop, 2022.

File information - part-r-00000

Download Head the file (first 32K) Tail the file (last 32K)

Block information -- Block 0

Block ID: 1073742172  
Block Pool ID: BP-499674093-10.0.0.18-1660809147878  
Generation Stamp: 1348  
Size: 86483  
Availability:

- 10.0.0.18

File contents

Date	Count
2015-01-03	5498
2015-01-05	5403
2014-12-29	5292
2015-01-07	5096
2015-01-04	5009
2014-12-31	4805
2015-01-09	4765
2015-01-01	4705

# Apache Pig

## 1. Total number of products per rating.

```
grunt> data = LOAD 'hdfs:/datadump/amazon_reviews_us_Electronics_v1_00.tsv' AS (marketplace, customer_id,review_id, product_id, product_parent, product_title, product_category,star_rating, helpful_votes, total_votes, vine, verified_purchase, review_headline,review_body, review_date);
2022-08-18 10:02:28,326 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
2022-08-18 10:02:28,491 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
2022-08-18 10:04:48,517 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
2022-08-18 10:04:48,593 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: GROUP_BY
2022-08-18 10:04:48,619 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
2022-08-18 10:04:48,622 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2022-08-18 10:04:48,626 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, ConstantCalculator, GroupByConstParallelSetter, LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForEach, NestedLimitOptimizer, PartitionFilterOptimizer, PredicatePushdownOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInsertor]}
2022-08-18 10:04:48,641 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.RMCompiler - File concatenation threshold: 100 optimistic? false
2022-08-18 10:04:48,656 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.CombinerOptimizerUtil - Choosing to move algebraic foreach to combiner
2022-08-18 10:04:48,655 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1
2022-08-18 10:04:48,655 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 1
2022-08-18 10:04:48,685 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
2022-08-18 10:04:48,694 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8032
2022-08-18 10:04:48,707 [main] INFO org.apache.pig.tools.pigstats.mapreduce.MRScriptState - Pig script settings are added to the job
2022-08-18 10:04:48,709 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3
2022-08-18 10:04:48,711 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Reduce phase detected, estimating # of required reducers.
2022-08-18 10:04:48,711 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Using reducer estimator: org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputSizeReducerEstimator
2022-08-18 10:04:48,714 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputSizeReducerEstimator - BytesPerReducer=1000000000 maxReducers=999 totalInputFileSize=1725988504
2022-08-18 10:04:48,714 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting Parallelism to 2
2022-08-18 10:04:48,714 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - This job cannot be converted run in-process
2022-08-18 10:04:48,928 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Added jar file:/opt/homebrew/Cellar/pig/0.17.0_1/libexec/pig-0.17.0-core-h2.jar to DistributedCache through /tmp/temp891994362/tmp229389327/pig-0.17.0-core-h2.jar
2022-08-18 10:04:49,388 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Added jar file:/opt/homebrew/Cellar/pig/0.17.0_1/libexec/lib/automaton-1.11-8.jar to DistributedCache through /tmp/temp891994362/tmp583072744/automaton-1.11-8.jar
2022-08-18 10:04:49,877 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Added jar file:/opt/homebrew/Cellar/pig/0.17.0_1/libexec/lib/antlr-runtime-3.4.jar to DistributedCache through /tmp/temp891994362/tmp-1132402459/antlr-runtime-3.4.jar
2022-08-18 10:04:50,354 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Added jar file:/opt/homebrew/Cellar/pig/0.17.0_1/libexec/lib/joda-time-2.9.3.jar to DistributedCache through /tmp/temp891994362/tmp1528044819/joda-time-2.9.3.jar
2022-08-18 10:04:50,377 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting up single store job
2022-08-18 10:04:50,384 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.
2022-08-18 10:04:50,384 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache
2022-08-18 10:04:50,384 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Setting key [pig.schematuple.classes] with classes to deserialize []
2022-08-18 10:04:50,418 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.
2022-08-18 10:04:50,427 [JobControl] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8032
2022-08-18 10:04:50,508 [JobControl] INFO org.apache.hadoop.mapreduce.JobResourceUploader - Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/poojayandhe/.staging/job_1660836347592_0012
2022-08-18 10:04:50,518 [JobControl] WARN org.apache.hadoop.mapreduce.JobResourceUploader - No job jar file set. User classes may not be found. See Job or Job#setJar(String).
2022-08-18 10:04:50,536 [JobControl] INFO org.apache.pig.builtin.PigStorage - Using PigTextInputFormat
2022-08-18 10:04:50,542 [JobControl] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input files to process : 1
2022-08-18 10:04:50,542 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
```

# Apache Pig

## 1. Total number of products per rating.

```
2022-08-18 10:04:50.542 [JobControl] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input files to process : 1
2022-08-18 10:04:50.542 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2022-08-18 10:04:51.547 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - number of splits:13
2022-08-18 10:04:52.052 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - Submitting tokens for job: job_1660836347592_0012
2022-08-18 10:04:52.054 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - Executing with tokens: []
2022-08-18 10:04:52.081 [JobControl] INFO org.apache.hadoop.mapred.YARNRunner - Job jar is not present. Not adding any jar to the list of resources.
2022-08-18 10:04:52.379 [JobControl] INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl - Submitted application application_1660836347592_0012
2022-08-18 10:04:52.385 [JobControl] INFO org.apache.hadoop.mapreduce.Job - The url to track the job: http://Poojas-MacBook-Pro.local:8088/proxy/application_1660836347592_0012/
2022-08-18 10:04:52.385 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - HadoopJobId: job_1660836347592_0012
2022-08-18 10:04:52.386 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Processing aliases data,groupeddata,productcount
2022-08-18 10:04:52.386 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - detailed locations: M: data[6,7],data[-1,-1],productcount[8,15],groupeddata[7,14] C: p
productcount[14,15],groupeddata[7,14] R: productcount[8,15]
2022-08-18 10:04:52.407 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 0% complete
2022-08-18 10:04:52.407 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Running jobs are [job_1660836347592_0012]
2022-08-18 10:05:19.614 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 23% complete
2022-08-18 10:05:19.620 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Running jobs are [job_1660836347592_0012]
2022-08-18 10:05:34.902 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 46% complete
2022-08-18 10:05:34.903 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Running jobs are [job_1660836347592_0012]
2022-08-18 10:05:41.988 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Running jobs are [job_1660836347592_0012]
2022-08-18 10:05:48.129 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8032
2022-08-18 10:05:48.173 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 10:05:48.487 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8032
2022-08-18 10:05:48.502 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 10:05:48.585 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8032
2022-08-18 10:05:48.609 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 10:05:48.652 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
2022-08-18 10:05:48.655 [main] INFO org.apache.tools.pigstats.mapreduce.SimplePigStats - Script Statistics:

HadoopVersion PigVersion UserId StartedAt FinishedAt Features
3.3.3 0.17.0 poojayendhe 2022-08-18 10:04:48 2022-08-18 10:05:48 GROUP_BY

Success!

Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime      MinMapTime    AvgMapTime   MedianMapTime  MaxReduceTime  MinReduceTime  AvgReduceTime  MedianReducetime   Alias  Feature Outputs
job_1660836347592_0012 13      2          16           5            14          14           7             7            7           data,groupeddata,productcount GROUP_BY,COMBINER          hdfs:/pig_output/ProductsEachRating.

Input(s):
Successfully read 3093870 records (1726046913 bytes) from: "hdfs:/datadump/amazon_reviews_us_Electronics_v1_00.tsv"

Output(s):
Successfully stored 6 records (60 bytes) in: "hdfs:/pig_output/ProductsEachRating"

Counters:
Total records written : 6
Total bytes written : 60
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1660836347592_0012

2022-08-18 10:05:48.660 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8032
2022-08-18 10:05:48.664 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 10:05:48.718 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8032
2022-08-18 10:05:48.725 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 10:05:48.773 [main] INFO org.apache.hadoop.yarn.client.DefaultNoHARMFailoverProxyProvider - Connecting to ResourceManager at /0.0.0.0:8032
2022-08-18 10:05:48.780 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2022-08-18 10:05:48.821 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
grunt>
```

# Apache Pig

1. Total number of products per rating.

The screenshot shows the Apache Hadoop Web UI interface. At the top, the URL is `localhost:9870/explorer.html#/pig_output/ProductsEachRating`. The main page displays a "Browse Directory" view for the path `/pig_output/ProductsEachRating`, showing three entries: `_SUCCESS`, `part-r-00000`, and `part-r-00001`. A modal window titled "File information - part-r-00001" is open, providing detailed information about Block 0 of the file. The modal includes fields for "Block ID" (1073742186), "Block Pool ID" (BP-499674093-10.0.0.18-1660809147878), "Generation Stamp" (1362), "Size" (28), and "Availability" (10.0.0.18). Below the modal, the "File contents" section shows the following data:

Rating	Count
1	358120
3	238587
5	1781161

At the bottom left, it says "Showing 1 to 3 of 3 entries". The footer of the page reads "Hadoop, 2022."

**File information - part-r-00001**

Download Head the file (first 32K) Tail the file (last 32K)

Block information -- Block 0

Block ID: 1073742186  
Block Pool ID: BP-499674093-10.0.0.18-1660809147878  
Generation Stamp: 1362  
Size: 28  
Availability:  
• 10.0.0.18

File contents

Rating	Count
1	358120
3	238587
5	1781161

Showing 1 to 3 of 3 entries

Hadoop, 2022.

# Apache Pig

1. Total number of products per rating.

localhost:9870/explorer.html#/pig\_output/ProductsEachRating

Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities ▾

## Browse Directory

/pig\_output/ProductsEachRating

Show 25 entries

Permission	Owner
-rw-r--r--	poojayer
-rw-r--r--	poojayer
-rw-r--r--	poojayer

Showing 1 to 3 of 3 entries

Hadoop, 2022.

File information - part-r-00000

Download Head the file (first 32K) Tail the file (last 32K)

Block information -- Block 0

Block ID: 1073742185  
Block Pool ID: BP-499674093-10.0.0.18-1660809147878  
Generation Stamp: 1361  
Size: 32  
Availability:

- 192.168.7.204

File contents

2 179180
4 536821

Block Size Name

MB	_SUCCESS
MB	part-r-00000
MB	part-r-00001

Previous 1 Next

# Apache Hive

- Find maximum and minimum rating verified product.

```
hive> show tables;
OK
Time taken: 0.596 seconds
hive> CREATE TABLE IF NOT EXISTS amazonreviewdata (marketplace String,
> customer_id String, review_id String, product_id String, product_parent String,
> product_title String, product_category String, star_rating String, helpful_votes
> String, total_votes String, vine String, verified_purchase String, review_headline
> String, review_body String, review_date String) ROW FORMAT DELIMITED
> FIELDS TERMINATED BY '\t'
> LINES TERMINATED BY '\n'
| > STORED AS TEXTFILE tblproperties("skip.header.line.count" = "1");
OK
Time taken: 1.147 seconds
hive> show tables;
OK
amazonreviewdata
Time taken: 0.117 seconds, Fetched: 1 row(s)
hive> Load data inpath 'hdfs://datadump/amazon_reviews_us_Electronics_v1_00.tsv' into table amazonreviewdata;
Loading data to table default.amazonreviewdata
OK
Time taken: 1.424 seconds
hive> INSERT OVERWRITE DIRECTORY 'ElectronicsHiveout.tsv' ROW FORMAT
> DELIMITED FIELDS TERMINATED BY ',' SELECT product_id, Max(star_rating), Min
> (star_rating), SUM(helpful_votes) from amazonreviewdata where verified_purchase = 'Y'
| > GROUP BY product_id ;
Query ID = poojayendhe_20220818085153_7472c277-fa8c-40d7-9879-e5c32ac276f9
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 7
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1660836347592_0001, Tracking URL = http://Poojas-MacBook-Pro.local:8088/proxy/application_1660836347592_0001/
Kill Command = /opt/homebrew/Cellar/hadoop/3.3.3/libexec/bin/mapred job -kill job_1660836347592_0001
Hadoop job information for Stage-1: number of mappers: 7; number of reducers: 7
2022-08-18 08:52:10,825 Stage-1 map = 0%, reduce = 0%
2022-08-18 08:52:28,984 Stage-1 map = 43%, reduce = 0%
2022-08-18 08:52:30,037 Stage-1 map = 86%, reduce = 0%
2022-08-18 08:52:37,494 Stage-1 map = 100%, reduce = 0%
2022-08-18 08:52:41,813 Stage-1 map = 100%, reduce = 14%
2022-08-18 08:52:42,915 Stage-1 map = 100%, reduce = 29%
2022-08-18 08:52:44,043 Stage-1 map = 100%, reduce = 43%
2022-08-18 08:52:46,162 Stage-1 map = 100%, reduce = 71%
2022-08-18 08:52:47,229 Stage-1 map = 100%, reduce = 86%
2022-08-18 08:52:49,337 Stage-1 map = 100%, reduce = 100%
Ended Job = job_1660836347592_0001
Moving data to directory ElectronicsHiveout.tsv
MapReduce Jobs Launched:
Stage-Stage-1: Map: 7 Reduce: 7 HDFS Read: 1726134863 HDFS Write: 3105855 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
Time taken: 58.255 seconds
hive>
```

# Apache Hive

- Find maximum and minimum rating verified product.

The screenshot shows the Apache Hive Web UI interface. A modal dialog box is open, titled "File information - 000001\_0". Inside the dialog, there are three buttons: "Download", "Head the file (first 32K)", and "Tail the file (last 32K)". Below these buttons, a green header bar says "Block information -- Block 0". Underneath, the following details are listed:

- Block ID: 1073742003
- Block Pool ID: BP-499674093-10.0.0.18-1660809147878
- Generation Stamp: 1179
- Size: 446504
- Availability:
  - 192.168.7.204

In the bottom left corner of the dialog, there is a "File contents" section containing a list of numerical values. In the bottom right corner of the dialog, there is a "Close" button. The background of the page shows a list of files in a directory, with the file "000001\_0" selected. The list includes columns for Name, Size, and Availability. The "Name" column contains entries like "000001\_0", "000002\_0", etc. The "Size" column shows 28 MB for each entry. The "Availability" column shows a trash bin icon for each entry. At the bottom of the list, there are buttons for "Previous", "1" (highlighted in blue), and "Next".

Name	Size	Availability
000001_0	28 MB	
000002_0	28 MB	
000003_0	28 MB	
000004_0	28 MB	
000005_0	28 MB	
000006_0	28 MB	

# Apache Hive

## 1. Total count of reviews.

```
Time taken: 39.384 seconds
hive> select count(*) from amazonreviewdata;
Query ID = poojayendhe_20220818085918_2f664fe4-6af1-41d1-930f-de4d266b7a6b
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1660836347592_0002, Tracking URL = http://Poojas-MacBook-Pro.local:8088/proxy/application_1660836347592_0002/
Kill Command = /opt/homebrew/Cellar/hadoop/3.3.3/libexec/bin/mapred job -kill job_1660836347592_0002
Hadoop job information for Stage-1: number of mappers: 7; number of reducers: 1
2022-08-18 08:59:32,822 Stage-1 map = 0%,  reduce = 0%
2022-08-18 08:59:47,285 Stage-1 map = 71%,  reduce = 0%
2022-08-18 08:59:48,346 Stage-1 map = 86%,  reduce = 0%
2022-08-18 08:59:52,591 Stage-1 map = 100%,  reduce = 0%
2022-08-18 08:59:54,740 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_1660836347592_0002
MapReduce Jobs Launched:
Stage-Stage-1: Map: 7  Reduce: 1  HDFS Read: 1726107730  HDFS Write: 107 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
3093869
Time taken: 39.384 seconds, Fetched: 1 row(s)
hive>
```

# Apache Hive

## 1. Total count of product id in category

```
Time taken: 47.344 seconds, Fetched: 1 row(s)
hive> select product_category, count(distinct product_id) from amazonreviewdata group by product_category;
Query ID = poojayendhe_20220818090616_f06ad72b-e077-4d9d-ab9b-6cc7961a901d
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 7
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1660836347592_0004, Tracking URL = http://Poojas-MacBook-Pro.local:8088/proxy/application_1660836347592_0004/
Kill Command = /opt/homebrew/Cellar/hadoop/3.3.3/libexec/bin/mapred job -kill job_1660836347592_0004
Hadoop job information for Stage-1: number of mappers: 7; number of reducers: 7
2022-08-18 09:06:28,794 Stage-1 map = 0%,  reduce = 0%
2022-08-18 09:06:44,070 Stage-1 map = 86%,  reduce = 0%
2022-08-18 09:06:51,363 Stage-1 map = 100%,  reduce = 0%
2022-08-18 09:06:54,571 Stage-1 map = 100%,  reduce = 29%
2022-08-18 09:06:57,791 Stage-1 map = 100%,  reduce = 57%
2022-08-18 09:06:58,849 Stage-1 map = 100%,  reduce = 71%
2022-08-18 09:07:00,930 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_1660836347592_0004
MapReduce Jobs Launched:
Stage-Stage-1: Map: 7  Reduce: 7  HDFS Read: 1726123334 HDFS Write: 640 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
Electronics      185852
Time taken: 47.344 seconds, Fetched: 1 row(s)
```

# Apache Hive

## 1. Count of reviews per product\_id

```
total MapReduce CPU time Spent: 0 msec
hive> INSERT OVERWRITE DIRECTORY 'productreviewcountHiveout.tsv' ROW FORMAT
    > DELIMITED FIELDS TERMINATED BY ','
[  > select product_id, count(review_id) from amazonreviewdata group by product_id;
Query ID = poojayendhe_20220819050853_8cd668fc-8976-4f72-b9f4-b4a18fbf52eb
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 7
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1660910676561_0002, Tracking URL = http://Poojas-MacBook-Pro.local:8088/proxy/application_1660910676561_0002/
Kill Command = /opt/homebrew/Cellar/hadoop/3.3.3/libexec/bin/mapred job -kill job_1660910676561_0002
Hadoop job information for Stage-1: number of mappers: 7; number of reducers: 7
2022-08-19 05:09:03,142 Stage-1 map = 0%,  reduce = 0%
2022-08-19 05:09:21,580 Stage-1 map = 71%,  reduce = 0%
2022-08-19 05:09:22,705 Stage-1 map = 86%,  reduce = 0%
2022-08-19 05:09:31,470 Stage-1 map = 100%,  reduce = 0%
2022-08-19 05:09:34,888 Stage-1 map = 100%,  reduce = 14%
2022-08-19 05:09:35,968 Stage-1 map = 100%,  reduce = 43%
2022-08-19 05:09:37,112 Stage-1 map = 100%,  reduce = 71%
2022-08-19 05:09:40,287 Stage-1 map = 100%,  reduce = 86%
2022-08-19 05:09:41,374 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_1660910676561_0002
Moving data to directory productreviewcountHiveout.tsv
MapReduce Jobs Launched:
Stage-Stage-1: Map: 7  Reduce: 7  HDFS Read: 1726155576 HDFS Write: 2458318 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
Time taken: 51.738 seconds
hive> █
```

# Apache Hive

1. Count of reviews per product\_id

The screenshot shows the Apache Hadoop File Explorer interface at `localhost:9870/explorer.html#/user/poojayendhe/productreviewcountHiveout.tsv`. A modal window titled "File information - 000001\_0" is open, displaying details about a specific data block. The modal includes tabs for "Download", "Head the file (first 32K)", and "Tail the file (last 32K)". The "Block information" tab is selected, showing:

- Block ID: 1073742217
- Block Pool ID: BP-499674093-10.0.0.18-1660809147878
- Generation Stamp: 1393
- Size: 351611
- Availability:
  - 192.168.7.204

The main browser view shows a "Browse Directory" listing for the path `/user/poojayendhe/productreviewcountHiveout`. It lists 7 entries, each with a checkbox, permission (e.g., `-rw-r--r--`), owner (`poojayendhe`), and a timestamp. The entries are numbered 000001\_0 through 000006\_0, all with a size of 28 MB.

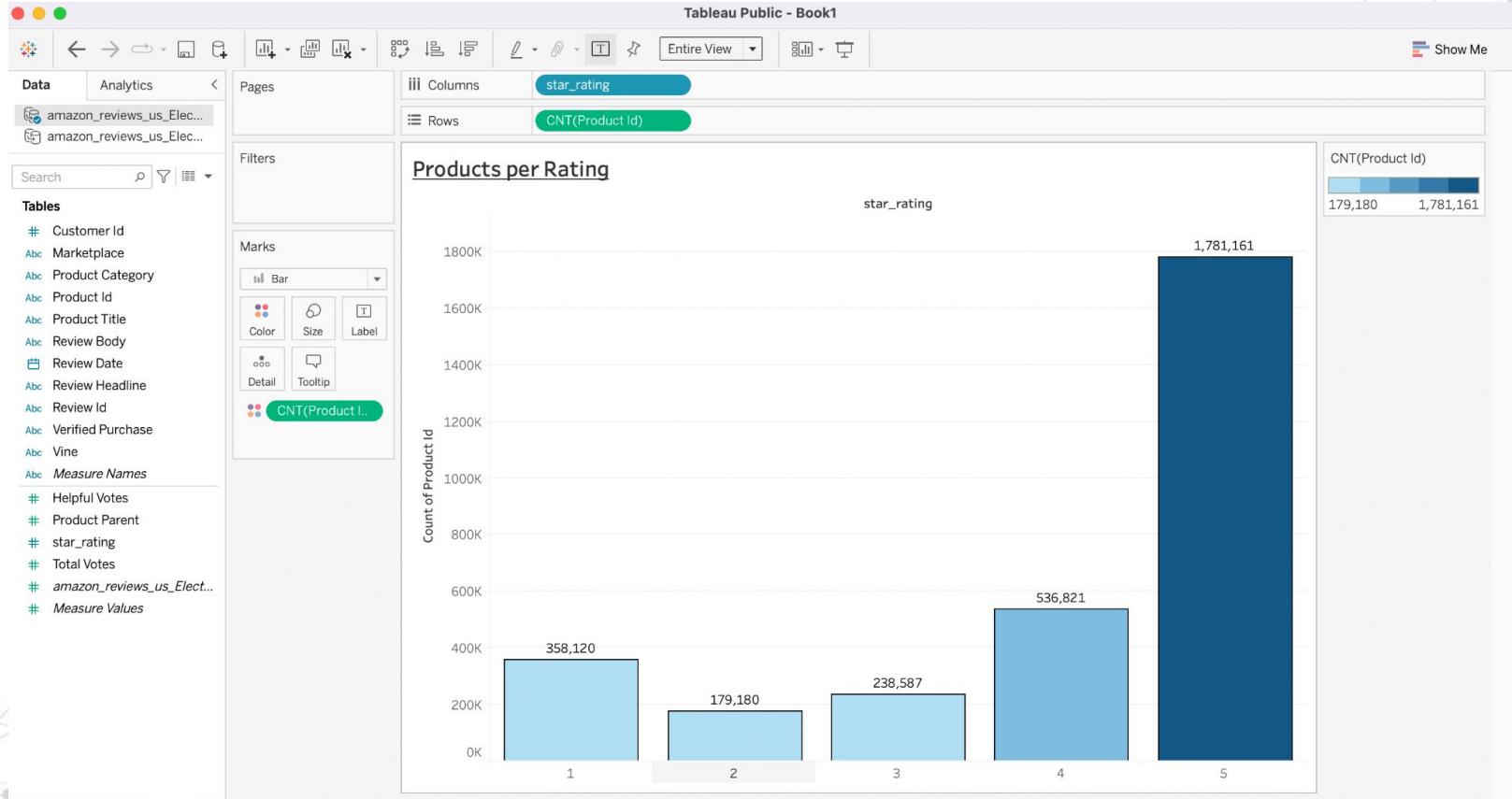
Block Size	Name	Action
28 MB	000000_0	trash
28 MB	000001_0	trash
28 MB	000002_0	trash
28 MB	000003_0	trash
28 MB	000004_0	trash
28 MB	000005_0	trash
28 MB	000006_0	trash

At the bottom of the page, there are navigation links for "Previous", "1", and "Next".

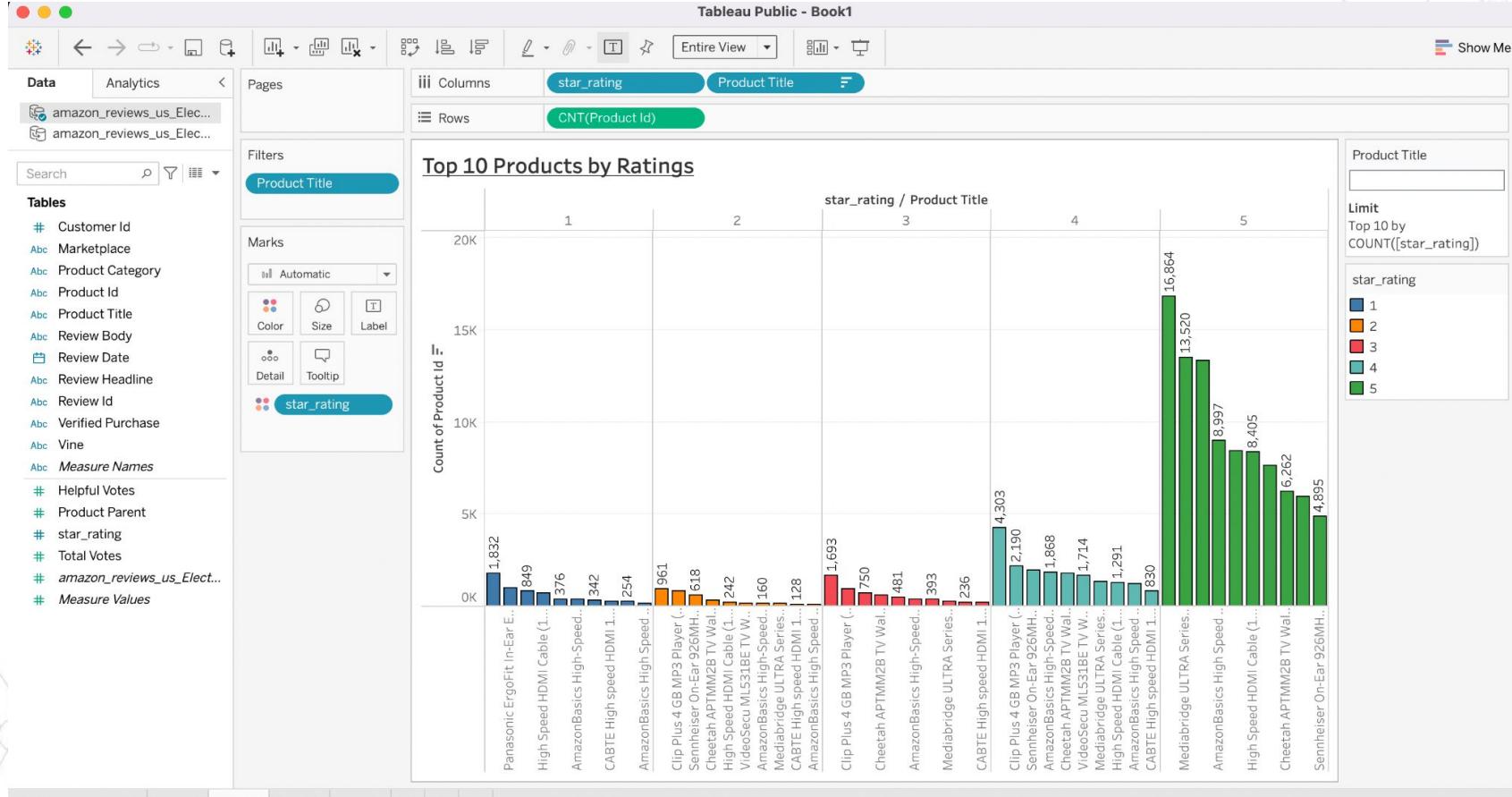
File contents (partial):

```
0558835155,1
0594451647,29
1060078031,4
1585747416,1
1590650670,1
1610130804,1
1934805912,1
1938401387,1
```

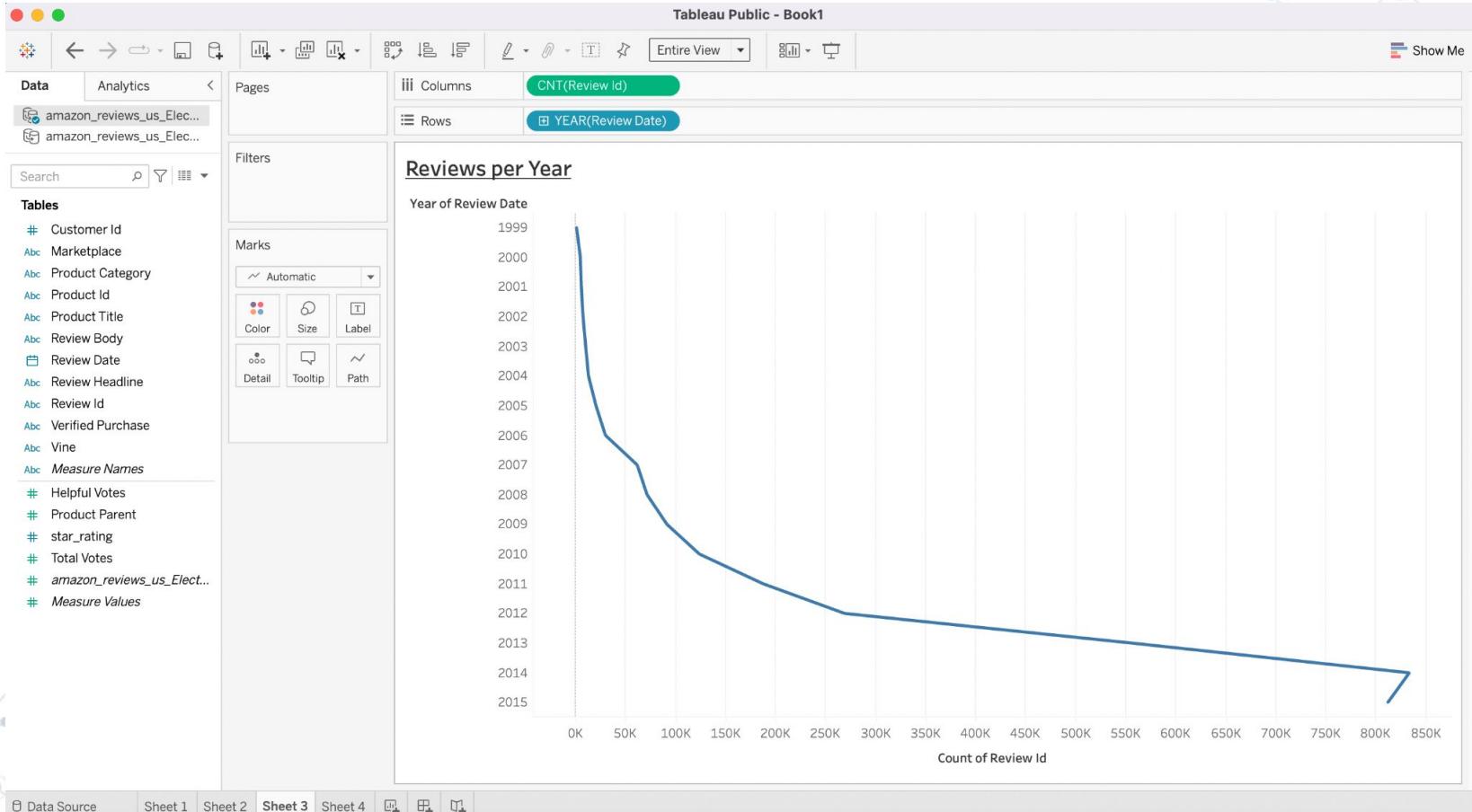
# Tableau



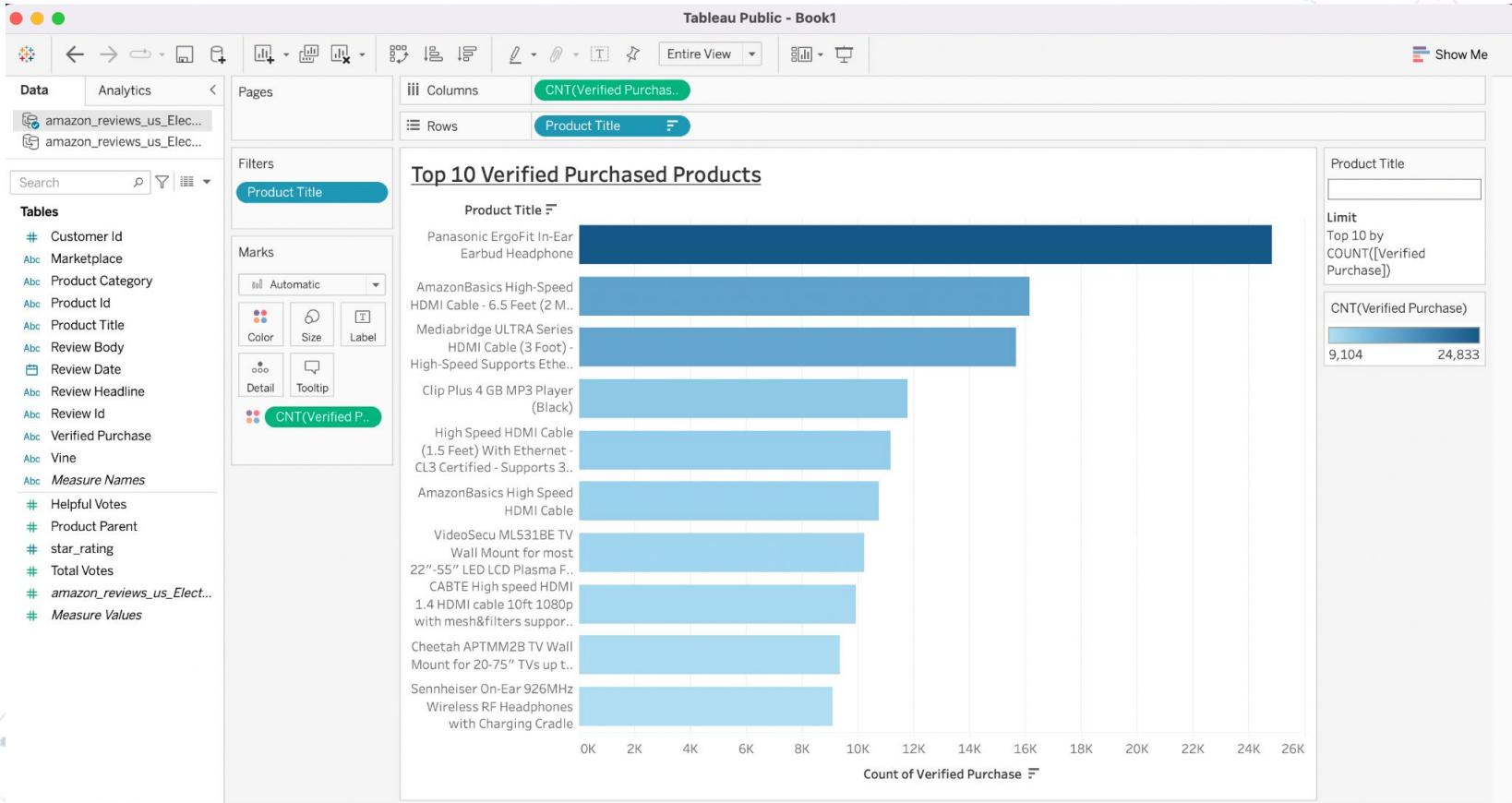
# Tableau



# Tableau



# Tableau





# Thank You !