

INFO 7275 38035

Advanced Database Management System

Section 01

A Final Project Report

on

Analysis of ‘Amazon Fine Food Review’

Dataset

Tanisha Jain

Northeastern University

MS-IS

(NU Id: 001617268)

Summary:

This project is based on the Analysis of the '**Amazon Fine Food Review**' dataset. The raw dataset consists of reviews of fine foods from Amazon. The data span a period of more than 10 years, including all ~500,000 reviews up to October 2012. Reviews include product id and user information, ratings, and a plaintext review.

In the project, first the analysis of the raw dataset using Tableau is done. Then the cleaning of the data files is done using RScript on RStudio. The MapReduce analysis is performed using several MapReduce patterns such as Partitioning, Distinct etc. Also, in this project, Apache Hive and Apache Mahout is implemented. Amazon Elastic Map Reduce is also used to implement MapReduce algorithms on cloud provided by Amazon Web Services.

The analysis on the dataset can help Amazon products with its sales, scope of improvement and honest user reviews for many other users buying similar products.

Dataset Link: <https://www.kaggle.com/snap/amazon-fine-food-reviews> (originally this dataset was published on <http://snap.stanford.edu/data/web-FineFoods.html>)

Attributes of dataset:

Number of reviews	568,454
Number of users	256,059
Number of products	74,258
Users with > 50 reviews	260
Median no. of words per review	56
Timespan	Oct 1999 - Oct 2012

Amazon Fine Food Review Dataset:

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	
1	Id	ProductId	UserId	Helpful	HelpfulN	Score	Time	Summary	Review											
2	1	B001E4KFG0	A35GXH7AUH8GW	1	1	5	1303862400	Good Quality	I have bought several of the Vitality canned dog food products and have found them all to be of good quality. The product looks more like a stew than a processed r											
3	2	B00813GGR4	A1D87F6ZCV5NK	0	0	1	1346976000	Not as Adver	Product arrived labeled as Jumbo Salted Peanuts...the peanuts were actually small sized unsalted. Not sure if this was an error or if the vendor intended to represent											
4	3	B000L0OCHO	ABKLWMIJXXAIN	1	1	4	1219017600	Delight says	I This is a conf pillowy citrus and very flava. The Witch and The Wardrobe - this is the treat that seduces Edmund into selling out his Brother and Sisters to the Witch.											
5	4	B000UA00IQ	A395B0RC6FGVX	3	3	2	1307923200	Cough Medic	If you are looking for the secret ingredient in Robitussin I believe I have found it. I got this in addition to the Root Beer Extract I ordered (which was good) and made											
6	5	B006K2ZZ7K	A1UQRSCLF8GW1T	0	0	5	1350777600	Great taffy	Great taffy a this is a deal.											
7	6	B006K2ZZ7K	ADT05RK11MGOEU	0	0	4	1340251200	Nice Taffy	I got a wild h root beer melon peppermint grape etc. My only my kids and my hust this lasted only two weeks! I would recommend this br											
8	7	B006K2ZZ7K	A152VKFXXR01	0	0	5	1340150400	Great!	Just a This saltwater which did not Fralinger's. Would highly recommend this candy! I served it at a beach-themed party and everyone loved it!											
9	8	B006K2ZZ7K	A3JRQVEQN31B	0	0	5	1336003200	Wonderful	tasty taffy	This taffy is so good. It is very soft and chewy. The flavors are amazing. I would definitely recommend you buying it. Very satisfying!!										
10	9	B000E7L2R4	A1MYO79ZK0BB1	1	1	5	1322006400	Yay Barley	Right now I'm mostly just sprouting this so my cats can eat the grass. They love it. I rotate it around with Wheatgrass and Rye too											
11	10	B00171APVA	A21BT40VZCCYT4	0	0	5	1351209600	Healthy Dog	This is a very healthy dog food. Good for their digestion. Also good for small puppies. My dog eats her required amount at every feeding.											
12	11	B0001P9FFE	A3HDKO70WQWNK4	1	1	5	1107820800	The Best Hot I don't know	but the flavor because of t we have a ci but don't grab a bottle you will nevi incredible service!											
13	12	B0009XLVGO	A2725IBAY99JEB	4	4	5	1282867200	My cats LOVI	One of my bx and the proct no-by-product food up higher where only my skinny boy can jump. The higher food sits going stale. They both really go for this food. Ar											
14	13	B0009XLVGO	A327PCT23YH90	1	1	5	1339545600	My Cats Are	My cats have I now need to find a new food that my cats will eat.											
15	14	B001GVISJM	A18CWX2RJTHUE	2	2	4	1288915200	fresh and gre good flavor!	These came securely packed... they were fresh and delicious! I love these Twizzlers!											
16	15	B001GVISJM	A2MUGFV2TDQ47K	4	5	5	1268352000	Strawberry	The Strawberry Twizzlers are my guilty pleasure - yummy. Six pounds will be around for a while with my son and I.											
17	16	B001GVISJM	A1CZXC3CPBKQU	4	5	5	1262044800	Lots of wizzi	Just what my daughter loves twizzlers and this shipment of six pounds really hit the spot. It's exactly what you would expect...six packages of strawberry twizzler											
18	17	B001GVISJM	A3KLWF6WQSBNYO	0	0	2	1348099200	poor taste	I love eating them and they are good for watching TV and looking at movies! It is not too sweet. I like to transfer them to a zip lock baggie so they stay fresh so I can											
19	18	B001GVISJM	A3KHF14U97ZQ60	0	0	5	1282672000	Love it!	I am very satisfied with my Twizzler purchase. I shared these with others and we have all enjoyed them. I will definitely be ordering more.											
20	19	B001GVISJM	A2A9XS8GZGTBLP	0	0	5	1324598400	GREAT SWEET Twizzlers	Strawberry i made in Lan Inc. one of ti now a Subsi the Compan they also m Green Color I like them a the longest I Inc. This Rec 1998. This Produ											
21	20	B001GVISJM	A3IV7CL2C13K2U	0	0	5	1318032000	Home delive	Candy was delivered very fast and was purchased at a reasonable price. I was home bound and unable to get to a store so this was perfect for me.											
22	21	B001GVISJM	A1W0OKGLPR5PV6	0	0	5	1313452800	Always fresh	My husband packed well and arrive in a timely manner.											
23	22	B001GVISJM	AZOF9E17RGZHB	0	0	5	1308960000	TWIZZLERS	I bought thes and appetan and this was Strawberry 16-Ounce Bags (Pack of 6)											
24	23	B001GVISJM	ARYVQLAN73TA1	0	0	5	1304899200	Delicious prc	I can remember buying this candy as a kid and the quality hasn't dropped in all these years. Still a superb product you won't be disappointed with.											
25	24	B001GVISJM	A1613OLZZUG7V	0	0	5	1304467200	Twizzlers	I love this candy. After weight watchers I had to cut back but still have a craving for it.											
26	25	B001GVISJM	A22P2J09N9HKE	0	0	5	1295481600	Please sell th i have lived c and I so mis	I always stop more often than I'm able to buy them right now.											
27	26	B001GVISJM	A3FPONPR03HP3I	0	0	5	1288310400	Twizzlers - St Produc	reco Strawberry 16-Ounce Bags (Pack of 6)											
28	27	B001GVISJM	A3RXAU2NBKV45G	0	1	5	1332633600	Nasty No fla	The candy is No flavor. Just plan and chewy . I would never buy them again											
29	28	B001GVISJM	AAAS38898HNMK	0	1	4	1331856000	Great Bargai	I was so glad Amazon carried these batteries. I have a hard time finding them elsewhere because they are such a unique size. I need them for my garage door open											
30	29	B0014AC10S	A2FL4ZVGFLD10B	0	0	5	1338854400	YUMMY!	I got this for and my fath you would never guess that they're sugar-free and it's so great that you can eat them pretty much guilt free! I was so impressed that I've											
31	30	B0001P99FY	A3HDKO70WQWNK4	1	1	5	1107820800	The Best Hot I don't know	but the flav because of t we have a ci but don't grab a bottle you will nevi incredible service!											
32	31	B003F6U07K	A1MF09480F04W	0	0	5	1297641600	Great machi	I have never my mother (too as a usally non-coffee drinker). The little Dolce Gusto Machine is super easy to use and prepares a really good Coffee/Latt											
33	32	B003F6U07K	A13Q0Q709M207	0	1	5	1288310400	THIS IS MY T!	This offer is e thanks Amazon for selling this product. Staral											
34	33	B001EO5QW8	AOVR0BZBNTP7	19	19	4	1163376000	Best of the Ir McCann's Ins	however that even th the McCann all-natural b though is McCann's tasty oatmeal after sitting the instant McCann's becomes too thick an											
35	34	B001EO5QW8	A3PMMPNFEV1GK9	13	13	4	1166313600	Good Instant	This is a goor so not only c but some doctors now say that this form of sugar is better for you. Great on a cold morning when you don't have time to make McCann'											

Amazon User Profiles:

	A	B	C	D	E	F	G
1	UserId	FirstName	LastName	Age	Street	City	State
2	A3SGXH7AUHU8GW	Isabelle	Colon	38	Emazur Way	Hupobad	OK
3	A1D87F6ZCVE5NK	Katie	Hamilton	36	Pungut Plaza	Hizigu	NY
4	ABXLMWJIXXAIN	Bill	Jackson	62	Dutaf Highway	Humcobpof	MI
5	A395BORC6FGVXV	Russell	Phelps	33	Pezok Pass	Silikomi	SC
6	A1UQRSCLF8GW1T	Sam	Rodriquez	30	Ezoid Boulev	Koofvo	NM
7	ADT0SRK1MGOEU	Duane	Ferguson	39	Vifimo Place	Jemoke	NJ
8	A1SP2KVFKXXRU1	Trevor	McCarthy	60	Usoc Parkwa	Zesevroc	IA
9	A3JRGQVEQN3IQ	Hilda	Hardy	40	Kearu Boulev	Novotukor	MT
10	A1MZY09TZK0BBI	Jeff	Diaz	27	Duli Square	Mufojbu	MA
11	A21BT40VZCCYT4	Zachary	Johnston	46	Pada Loop	Esfokaj	AR
12	A3HDKO7OW0QNK4	Juan	Boone	31	Vigo Circle	Zijnagken	IA
13	A2725IB4YY9JEB	Terry	Daniels	18	Nese Terrace	Sehelri	WI
14	A327PCT23YH90	Josephine	Estrada	30	Dotza Key	Zuhatuca	OK
15	A18ECVX2RJ7HUE	Norman	Lamb	41	Solma Glen	Cuhdakbe	AR
16	A2MUGFV2TDQ47K	Benjamin	Harvey	64	Rophu Drive	Nejugan	UT
17	A1CZX3CP8IKQIJ	Leo	Daniels	52	Hake Ridge	Ugelehwus	UT
18	A3KLWF6WQ5BNYO	Louise	Clark	35	Ebuuj Boulev	Egizevzo	NV
19	AFKW14U97Z6QO	Rosie	Roy	35	Fusef Street	Abomaoz	TX
20	A2A9X58G2GTBLP	Luella	Kim	38	Figzo Avenue	Ursadpas	WY
21	A3IV7CL2C13K2U	Phoebe	Williamson	44	Muiz Circle	Kafjuwu	DE
22	A1WOO0KGLPR5PV6	Mitchell	Silva	54	Zapok View	Riavamas	FL
23	AZOF9E17RGZH8	Polly	Watkins	48	Porap Junction	Tuwojug	KS
24	ARYVQL4N737A1	Mae	Mendez	37	Arepoo Loop	Zapodiso	FL
25	AJ613OLZZUG7V	Zachary	Howell	23	Rozi Way	Ebejipvi	UT
26	A22P2J09NJ9HKE	Bill	Butler	46	Sefin Way	Lohtaelo	SC
27	A3FONPR03H3PJS	Mayme	Carpenter	49	Erugu Extens	Semojid	ID
28	A3RXAU2N8KV45G	Logan	Long	25	Vurbel View	Vavenza	MI
29	AAAS38B98HMIK	Herbert	Maxwell	50	Kinju Boulev	Ezewiidi	MT
30	A2F4LZVGFLD1OB	Isaac	Harrington	55	Gefibi Highw	Faukte	NC
31	A3HDKO7OW0QNK4	Dennis	McGee	22	Fici Street	Mimhibam	NM

Data Cleaning using R script in RStudio:

- 1) The raw data file with columns - Id, ProductId, UserId, ProfileName, HelpfulnessNumerator, HelpfulnessDenominator, Score, Time, Summary, Text and rows-568454. File is first read into RStudio and stored in a variable.
- 2) The time column which showed a Unix time field is replaced with its corresponding date value. (Formula Used to calculate time from unix to date format =(A1/86400)+25569+(-5/24))
- 3) The data is binded in a dataframe and written to the csv file.
- 4) Another file with users' profile is used, binded into a dataframe and is written to csv.

Find the RScript below-

```
#New R Script for cleaning the dataset and storing as a cleansed csv data file (Reviews.csv)
File1<-read.csv(file.choose(),header=TRUE)
View(File1)

#Data with Date to be appended
File2<-read.table(file.choose(),header=TRUE)
View(File2)

#Binding date with File1 and dropping columns 'Profile Name' and 'Time'
DataBind1<-cbind.data.frame(Id = File1$Id, ProductId = File1$ProductId,
                             UserId = File1$UserId, Date = File2>Date,
                             HelpfulnessNumerator = File1$HelpfulnessNumerator,
                             HelpfulnessDenominator = File1$HelpfulnessDenominator,
                             Rating = File1$Score, Summary = File1$Summary, Text = File1$Text)
View(DataBind1)

write.table(DataBind1, file = "AmazonFineFoodReview500k.csv", sep = ",", col.names = NA,
           qmethod = "double")
AmazonFFR<-read.table("AmazonFineFoodReview500k.csv", header = TRUE, sep = ",",
                       row.names = 1)
View(AmazonFFR)

#Trim the dataset
DataBind1<-head(AmazonFFR, 90000)
DataTable1<-data.frame(DataBind1)
View(DataTable1)

#Writing the first 90000 rows in a new file
write.table(DataTable1, file = "FineFoodReview90k.csv", sep = ",", col.names = NA,
            qmethod = "double")

#Editing new generated file for Users Profile
File3<-read.csv(file.choose(), header=TRUE)
View(File3)
```

```
file3Modified<-cbind.data.frame(Seq=File3$Seq, UserId=DataTable1$UserId, First=File3$First,
Last=File3$Last, Age=File3$Age, Street=File3$Street, City=File3$City, State=File3$State)
```

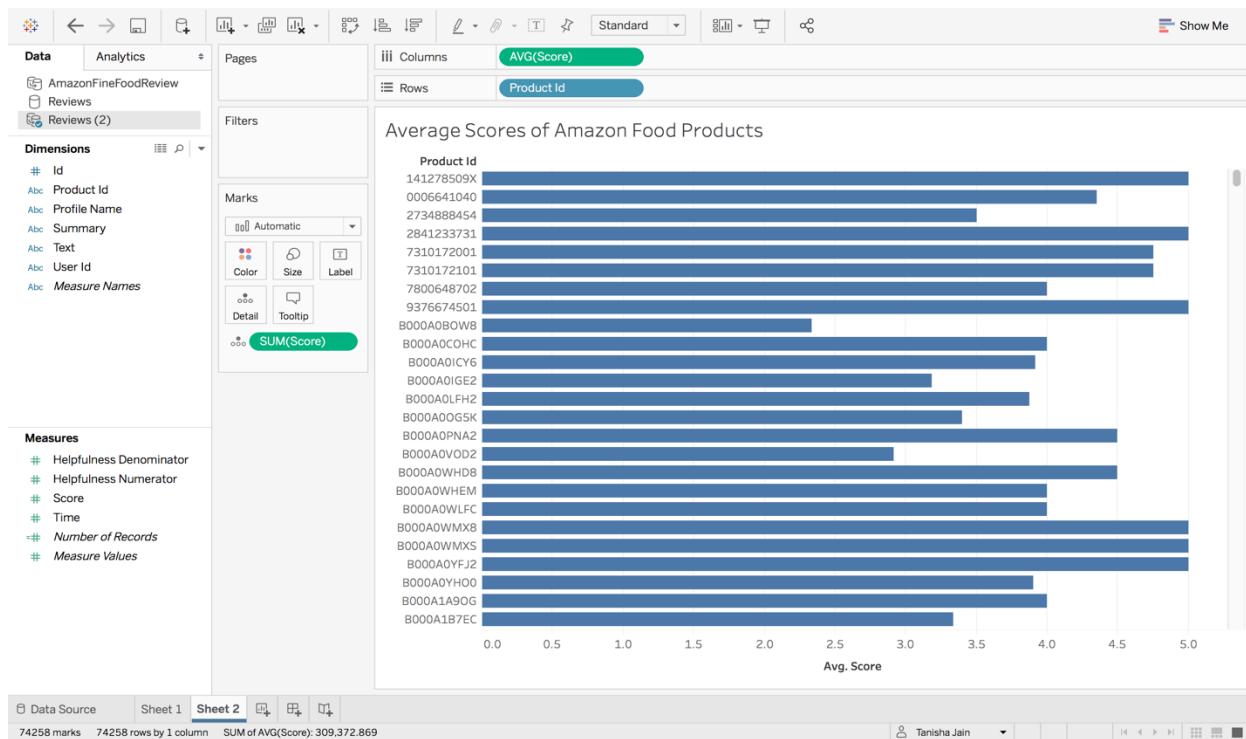
```
View(file3Modified)
```

```
#Writing the user modified file
```

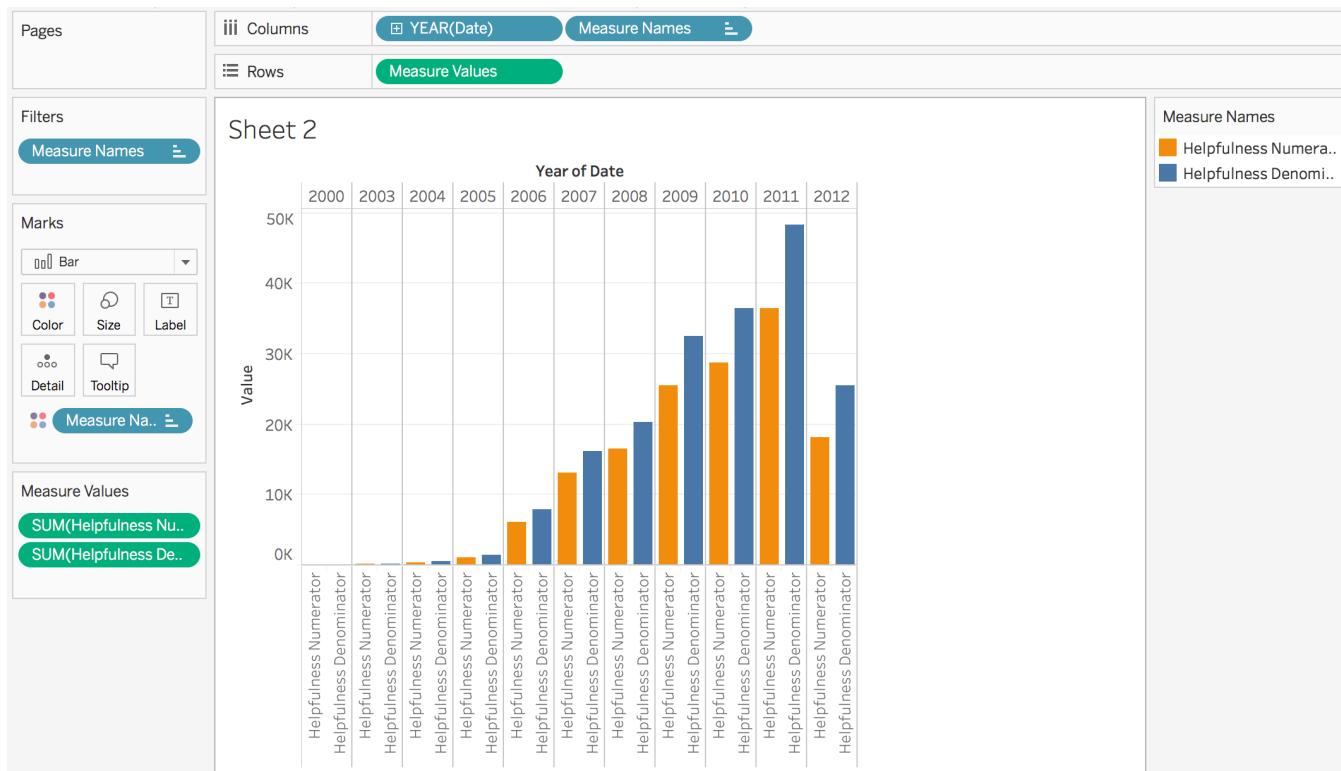
```
write.csv(file3Modified, "UserProfiles90k.csv", row.names = FALSE, quote = FALSE)
```

Data Analysis Using Tableau:

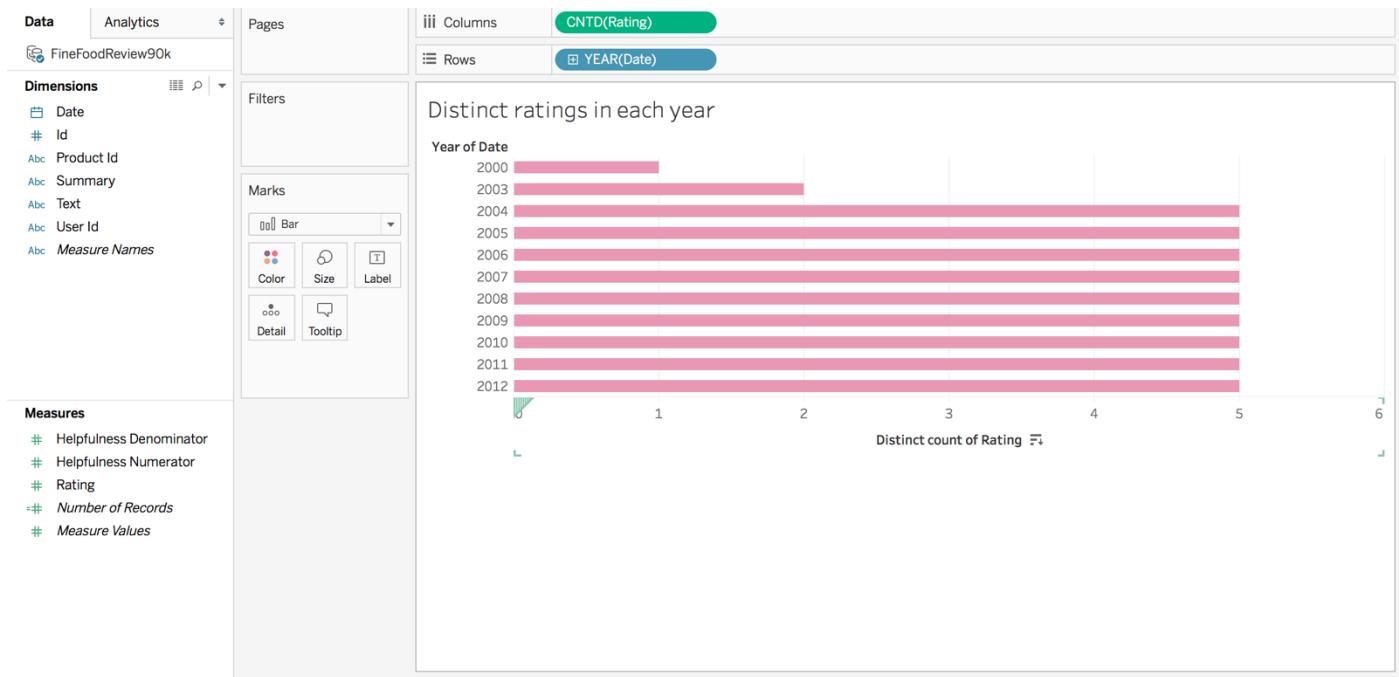
- 1) This is the average rating scores for the Amazon Food Products. By this analysis we can see which product has what average rating.



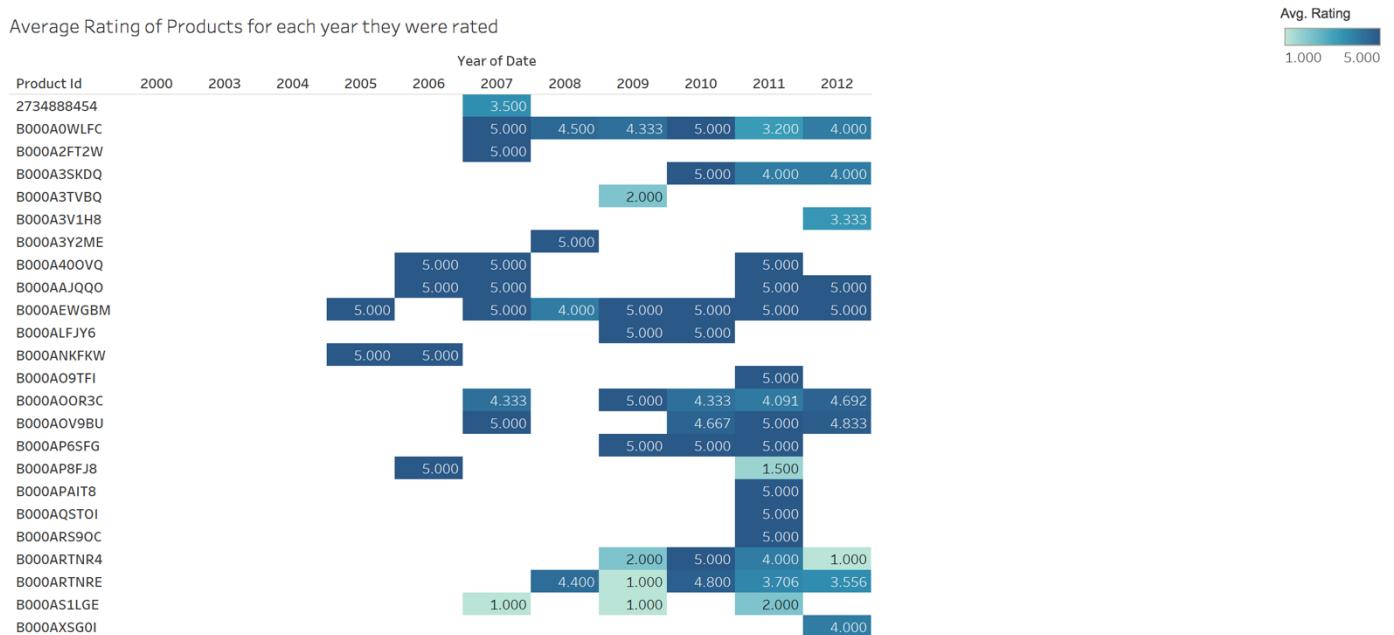
- 2) The columns in the dataset – Helpfulness Numerator and the Helpfulness Denominator indicates how many people liked the review out of what respectively. So here, the analysis is the overall helpfulness numerator and denominator value in a particular year.



- 3) The analysis graph below shows how many distinct ratings users gave for the products of that year. For example - Products received only 2 different ratings (can be either 2 from) in the year 2000 and all 5 different ratings in the year 2004 through 2012.



- 4) This graph below shows the average rating of a particular product in the particular year. The color strength shows the score of the rating 1 being low score and 5 being the best.

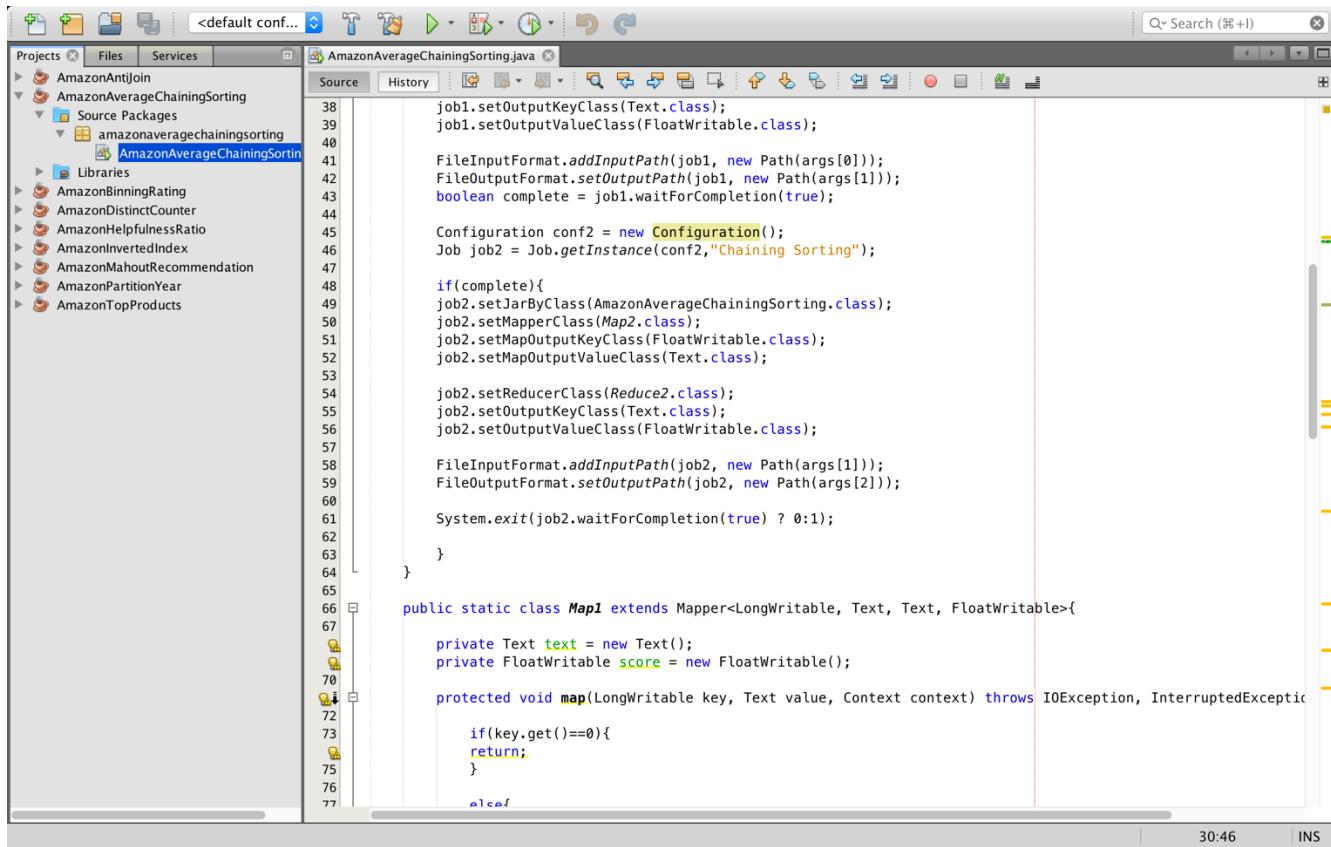


All the above analysis are created in Tableau and are also exported to Tableau Public.
(URL: <https://public.tableau.com/profile/tanishajain>)

MapReduce Analysis on 'Amazon Fine Food Review' Dataset:

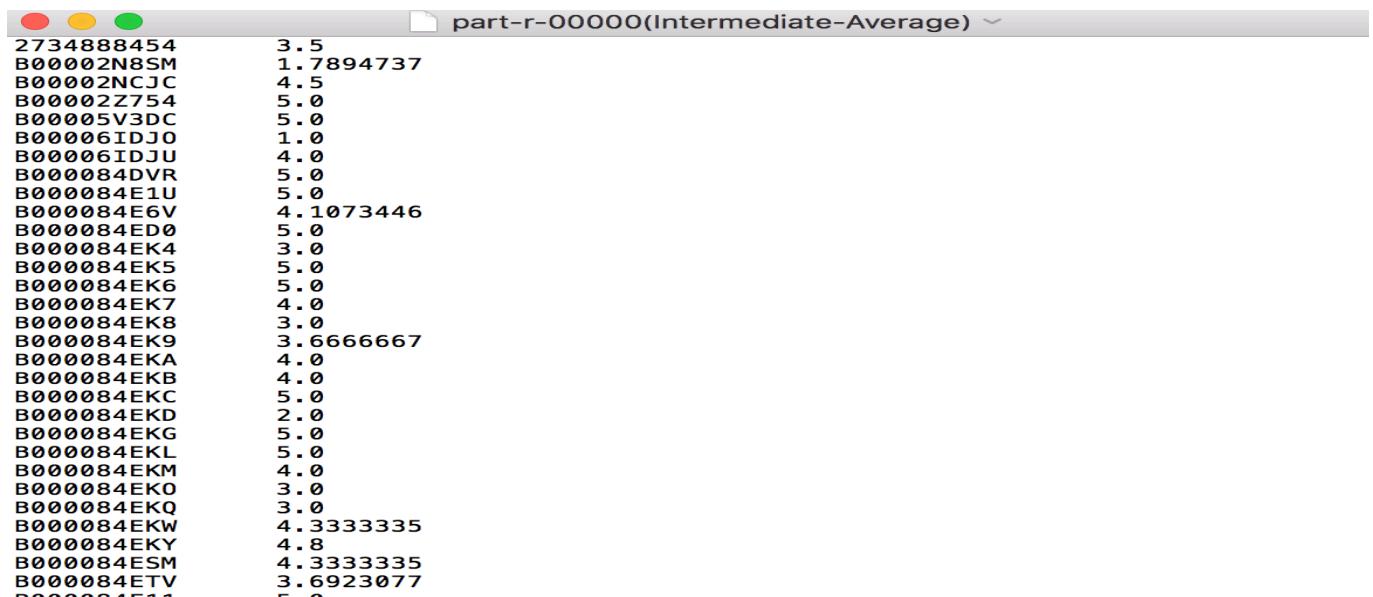
1) Amazon Fine Food Review Data: Average-Chaining-Sorting: (Numerical Summarization Pattern & Job Chaining)

This MR job uses the chaining where the output of one MapReduce job goes to the other before finally outputting it to the user. There are 2 mappers and 2 reducers.



```
38     job1.setOutputKeyClass(Text.class);
39     job1.setOutputValueClass(FloatWritable.class);
40
41     FileInputFormat.addInputPath(job1, new Path(args[0]));
42     FileOutputFormat.setOutputPath(job1, new Path(args[1]));
43     boolean complete = job1.waitForCompletion(true);
44
45     Configuration conf2 = new Configuration();
46     Job job2 = Job.getInstance(conf2,"Chaining Sorting");
47
48     if(complete){
49         job2.setJarByClass(AmazonAverageChainingSorting.class);
50         job2.setMapperClass(Map2.class);
51         job2.setMapOutputKeyClass(FloatWritable.class);
52         job2.setMapOutputValueClass(Text.class);
53
54         job2.setReducerClass(Reduce2.class);
55         job2.setOutputKeyClass(Text.class);
56         job2.setOutputValueClass(FloatWritable.class);
57
58         FileInputFormat.addInputPath(job2, new Path(args[1]));
59         FileOutputFormat.setOutputPath(job2, new Path(args[2]));
60
61         System.exit(job2.waitForCompletion(true) ? 0:1);
62     }
63 }
64
65 public static class Map1 extends Mapper<LongWritable, Text, Text, FloatWritable>{
66
67     private Text text = new Text();
68     private FloatWritable score = new FloatWritable();
69
70     protected void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException{
71
72         if(key.get()==0){
73             return;
74         }
75
76         else{
77
78             String[] rating = value.toString().split(",");
79             float sum = 0.0f;
80             int count = 0;
81
82             for(int i=1; i<rating.length; i++){
83                 sum += Float.parseFloat(rating[i]);
84                 count++;
85             }
86
87             score.set(sum/count);
88             text.set(key.toString());
89             context.write(text, score);
90
91         }
92     }
93 }
```

The first MR output averages the ratings of a particular product with productId as key and average rating as value.



productId	average rating
2734888454	3.5
B00002N8SM	1.7894737
B00002NCJC	4.5
B00002Z754	5.0
B00005V3DC	5.0
B00006IDJO	1.0
B00006IDJU	4.0
B000084DVR	5.0
B000084E1U	5.0
B000084E6V	4.1073446
B000084ED0	5.0
B000084EK4	3.0
B000084EK5	5.0
B000084EK6	5.0
B000084EK7	4.0
B000084EK8	3.0
B000084EK9	3.6666667
B000084EKA	4.0
B000084EKB	4.0
B000084EKC	5.0
B000084EKD	2.0
B000084EKG	5.0
B000084EKL	5.0
B000084EKM	4.0
B000084EKO	3.0
B000084EKQ	3.0
B000084EKW	4.3333335
B000084EKY	4.8
B000084ESM	4.3333335
B000084ETV	3.6923077
B000084E11	5.0

The second MR output shows the averages of the products ascendingly.

part-r-00000 (Final-sortedAverageAscending) ▾	
B001EQ4VXG	4.909091
B000COMPIY	4.9130435
B000F0FZCI	4.9166665
B000Q0IMOK	4.9166665
B000F4DKB2	4.9166665
B0050HZN4	4.9166665
B001QEFA4XC	4.9166665
B003KLSZGW	4.923077
B001964970	4.923077
B0002YGSA0	4.923077
B001EQ582E	4.9259257
B001EQ4G6I	4.928571
B0002BKIRW	4.928571
B002AUBJFI	4.9333334
B000YV9QMI	4.9333334
B000E682LY	4.9333334
B0030MNHEG	4.9375
B003DNL9U6	4.9411764
B008NDSNAU	4.9411764
B0056AL6G6	4.9473686
B002LMD8P4	5.0
B002LG5TWK	5.0
B002LMJ44S	5.0
B002LML5IQ	5.0
B002LDA2TS	5.0
B002LMNYG2	5.0
B002LN282M	5.0
B002LN3AOM	5.0
B002LN3DKI	5.0
B002LN4M8A	5.0
B002LN4YOC	5.0
B002L6KDQM	5.0
B002L68KDK	5.0
B002L3GKX0	5.0
B002LNZ7J8	5.0
B002KXDK48	5.0
B002L05SFA	5.0
B002LQAZGA	5.0
B002K08114	5.0
B002KGN4LE	5.0
B002KGE08G	5.0
B002KE92UI	5.0

2) Amazon Fine Food Review Data: Helpfulness Numerator/Helpfulness Denominator (Numerical Summarization Pattern)

In this analysis, I have taken ratio of Helpfulness numerator and Helpfulness Denominator for each product and averaged the value. The ratio is then multiplied with 100 and the output only shows upto 2 decimal places for the float values.

The screenshot shows an IDE interface with the following details:

- Project Explorer:** Shows various Amazon-related projects like "AmazonAntiJoin", "AmazonAverageChainingSorting", etc., under "Source Packages".
- Code Editor:** Displays the Java code for "AmazonHelpfulnessRatio.java".
- Code Content:** The code defines two classes: `FoodMapper` and `FoodReducer`.
 - `FoodMapper` extends `Mapper<LongWritable, Text, Text, FloatWritable>`. It has a `text` field and a `score` field. The `map` method processes a line from input, splits it by commas, and calculates the ratio of `helpnum` to `helpden`. If `helpden` is zero, the ratio is set to 0.0. The `text` and `score` fields are then written to context.
 - `FoodReducer` extends `Reducer<Text, FloatWritable, Text, Text>`. It has a `result` field. The `reduce` method is overridden but not fully shown in the screenshot.

The output file is shown below:

		part-r-00000(Percentage) ▾
100.0	%	B00008DFRS
32.65	%	B00008JOL0
12.5	%	B00008MOJ2
0.0	%	B00008036H
0.0	%	B00008RCMI
0.0	%	B000093HOV
0.0	%	B00009608Y
0.0	%	B0000960IV
33.33	%	B000090LE2
0.0	%	B000090LEW
0.0	%	B000090LFC
100.0	%	B0000A0BS8
0.0	%	B0000AH3OX
100.0	%	B0000AH3QT
28.57	%	B0000AH3QW
50.0	%	B0000AH3RR
66.67	%	B0000CDBRH
71.43	%	B0000CDBRJ
22.22	%	B0000CDBRP
28.57	%	B0000CDBRV
62.5	%	B0000CDBRY
62.5	%	B0000CDBRZ
100.0	%	B0000CE29E
100.0	%	B0000CEOKH
73.91	%	B0000CEPDP
33.33	%	B0000CERHH
100.0	%	B0000CERPU
0.0	%	B0000CERXB
25.0	%	B0000CERXU
20.0	%	B0000CERZK
0.0	%	B0000CGE3G
66.67	%	B0000CGE3R
0.0	%	B0000CGFCD
100.0	%	B0000CGFSC
0.0	%	B0000CGFV4
0.0	%	B0000CNU18
0.0	%	B0000CNU1A

3) Amazon Fine Food Review Data: Inverted Index (Numerical Summarization Pattern)

In this analysis, Inverted Indexing is used. It is very helpful as the output can show what Product are reviewed by which users. The key value will have multiple values.

```
public static class InvertedMapper extends Mapper<LongWritable, Text, Text, Text>{
    private Text product = new Text();
    private Text userId = new Text();

    public void map(LongWritable key, Text values, Context context){
        if(key.get()==0){
            return;
        }

        try{
            String[] tokens = values.toString().split(",");
            userId.set(tokens[2]);
            product.set(tokens[1]);

            context.write(product, userId);
        }
        catch(IOException | InterruptedException ex){
            System.out.println("Error in Mapper" + ex.getMessage());
        }
    }
}

public static class InvertedReducer extends Reducer<Text,Text,Text,Text>{
    private Text result = new Text();

    @Override
    public void reduce(Text key, Iterable<Text> values, Context context) throws IOException, InterruptedException{
        StringBuilder sb=new StringBuilder();
        boolean first = true;

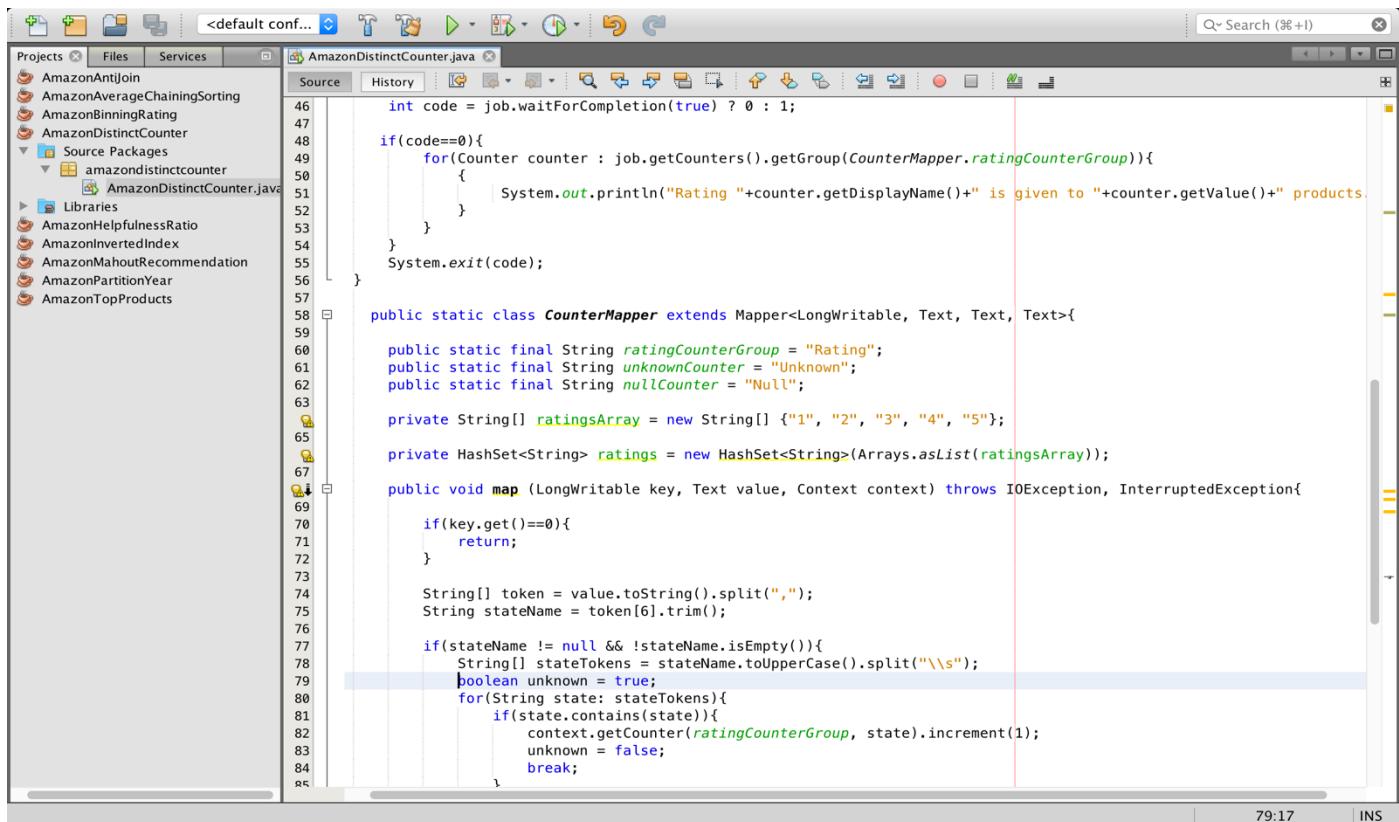
        for(Text id:values){
            if(first){
                first = false;
            }
            else{
                sb.append(",");
            }
            sb.append(id);
        }
        result.set(sb.toString());
        context.write(key, result);
    }
}
```

The file shows the ProductId and all the UserIds of the users who reviewed that particular product.

part-r-00000(InvertedIndex) ~														
B00002N8SM	A1GQUBNRZK3RH	A23K00L75SE03K	A2WCKCP2TY3SE	A17UIV0KM8JSV	A5W8ZRE07C1WK	A2579IQYG8Y9YF	A1S64U84E0YFTK	A1LU93B4X1PZZ	A32D342WBj6BX					
A1J67DKGFR7Z4	ASF1MKA0N60	A8BY0FGJ6IP0Z	A3EAHSVSDAY02K	A1G98WP2G0CMC	A29GWJL72GXXZ	A2B7IZJ243KWF	A1W48SP4GC100	AZKXR73X1CVIM	A3B28AM0CW7B3					
A11Q0UC2ULFN7R	AXV0SCCA43HE99	A8XKSGA2EADGJ	A2281X5GBCKA6	A7E1DY01QXM12	A2Z0052PK2P04H	A1L227LG53GTLB	A392XPJUTJDHSJ	A1NHUGER6N9EL	A19Q006CSFT011					
A1FYH4502BW7FN	AUE8TBSVHS6ZV	A1VEBVZS81UCI	A3R3T5T1K8IN6P	A11DUBT1ZXPRT7	A71C9WY0242QU	A1CFZ65HHRPRUY	AGF5R0NX8BXPH	A1J4110YC5IFAB						
B00002NCJ3	A196A7J9UEASJN	A13RRPE79XKFH												
B000027754	A29ZP19BWZPU3	A3B8RCIEI0FXF16												
B00005V3DC	A8K7Y548EWTHLW	A2ZYCEEYBUQZND	APASCXWTM041											
B000061D10	A1NKR5R76B11AF	A2CA4TJ6RL332H												
B000061D10	A1NKR5R76B11AF	A2CA4TJ6RL332H												
B000084DVR	A1UGDJP1ZJ3WVPF	A3DKGXWUEP1A2												
B000084E1U	A3DH85EYHM4AOH													
B000084E1U	A27VY0152W7SP7	A53ML2E3JY0DX	A2X09ABXBN2S6TZ	A1UW10D4G8752Q	A1IMWSB875ZK18	A07LLVTUYYMMW7	A4PH5BJ45GF62	AHB5QUSGHA8F2	A1302H6WKBA20H					
A1ZZ57TVB393K6	A1HJ3LZCFJXG14	A520HHSLSLP08	A2NPQH672LHELI	A1STJ270SC7HP6	A10AOY26U1RYLL	A2DWWQ1SF25USC5	A4VXMT07TDQHBC	A0Q0MVBBP17N	AYYWKPJM2XZ2H					
A12M1L2XMX75LL	A3G5T9AHNN7E5	A2DF47YCA4EM75	A1LB1ZLE79CW11	A1D3X1EHNZHY0	A2CT65CRHNS977	A31HD62RV56U80	AHXFOHC9YAK9M	A2CF6NAQ8G7B58	A2HBKE72530UZR					
A2EY1MR55IE4NG	A30UMX34RFPW098	A1GUW6NT1Z9N7	AFG2J3X5R69285	A1LV3GTRK0W3EC	A20TPV1KXXZP0E	A3H1K0HJ6BGXRL	A21IM576KY1VPC	A2Q0UTY4V48HZ55	A1X2GN2X0DGP2					
A1TK45KUTN1	AD5G5PZK7M8U9Y9	A1SKFCY29LC1T	A2K3C825MH5PDR	A3VK221B2M505B	AFX2WEKN91ANR	A285880PREK07	A2K1SNOHEK0T0Y	A25AOKSSQS678	A36JD4W1T7K57					
A5ANSLZBVSRHOF	A1AGU05PPZQLO6S	A52P0PSAT3B1F	A1RFEP3D8UZJ7	A1LE761RVN2Y7C	A1MFXL0D6TN0PM	A2604M802K8XT	A1PYK03XUR6WER	A1V3Q73621HHI0	A689HJRJVEK00					
A230K5MMCMT6W	A1Z2KHMMEHBBZ0	A1SG5JRXTBK66U	A1B5C93Q2BHK3R	AIDZU3LW6YIII	A9B0W1HS884L3	A2J1KVQYDTEKZZ	A23M122U0685D	A20H5NAANEBCYF	A1KG1G7PA9U0YEH					
A3UHYTRDPHMR5D	A5FRQDX2DEWXG	ABGHK3Q2BHK37	A2WRC10S9JW0X	A2B885100QIF61H	AEBMHARILGYOF	A161I0K7CAWRUW	A8BGU3MWHYZ16	ALAWLB1P5FV8	A1L1L0CERMZDQP					
A195YFLWLTZC9D	A2LZMFLWLTZC9D	A3A6MEPNP0VYKK	A1RN9FGWEUHDP5	A3RHPGU1HN16G	A1VEMH5514J3VZL	A2PBNUNZRBC6Y7	A30CH7A99FLRN	A162K7K008G0AD	A2LN4MDWXZM7CG					
A1A54IS5LK7V7B	A1ARV0539P37NI	A1TQJBUS26GA17S	A3UDDPF74X7E1U5	A36HVLNN9HWZ4X	A7B48AJT61C0A	A2HRGWAAMXYU51	A395523BF5JGAN	AY3KER02MFGPS	APX698XCS9UYS					
A1P023KZDR7W7	A2B5M06ZDR7W7	A2K5AFTU4YOPPW	ADE3J0DAUY2YHG	AFLUPRLM26N9	A2MV3HRSRUBWA2B	A1J50NM36XH7P	A3Z145L35GCY83	A1GL3V23LMM8KD	A50DQ3WXS9C7					
A30DC3EGHNCNCB	A23006PBBH78AF	A2KAMN8A0JZKH1	A8891HVRDJAM6	A1C5JH384RHVRD	A1GSB2MRA6L6VB6	A19X700NHL5ZNE	A2N5DQEVMG29KB	A3573T0S0BENAB	A2Y1LXP64VF0NF					
A1ITZK68SG0B8	A3IMFWY02E2K8	AKHFT3B0T2ENA	A1Z4F1J01WKB7B	A2P11H7383JX8Y	A1F950N7HHSU50	A841ZE5BVA1S	A3ETJXG12FL000	AP57AXG5HXZEO	A442AGZLAGPU					
A37FHDU0H1G5J	A1Z5V8SHPX3Y9L9	A8Y2FCFN5YSW	A2DN0D3F3E4K4M	A221XUA0701D7	A2Z18EXTM0GZK2	A1S29177DRDWLY	A3GKX52MS52BW1	A34AMF05J2083J	A2HMB0PTMV0ZGP					
A3GR3PAPFBNZX4	A1V6ENRZBNFR0K	A21VA2CMIP2UHI	A8891HVRDJAM6	A1PUM736GQ441P	A141T3TRKVRP60	A3ECY09BLZCQ28	A1E0D5IRU597GG	A2PKN7PKPJKE05	A24PNF4Q2TKB3					
A3L2PTMSUJZ2F3T	AU80C0U2W4059N	A2WTNYGSV4XDBP	A1XHF8JBTBTP0V0	A0HADDGUZ24Z	A25B1G19V71Y2L2	A2XR1H3Z16HAAV	A34TNH97309X70	A0ES3U50KACKHP	AVX34RCA3URC					
A1NKMGJGISMU70	AH1G606B1LF1B	A3D64VZH70E2E5	A11JWHLCL50Q40	A2RB8CF9F3MNZ0	A3GMN66FGTDLQY	A3T8X4P4CHAYVQL	A1V6ENRZBNFR0K	A3647014B4UBD0	A1K3V7QK565NR					
APF20D6YURK3G	A10FMK3MG4002	A2URXZ0PM5BBAU	AMPUN4RSR543V	A19MVNNI6ZEV0B	A1KPDK7VYRE3B	A3FSGHS5ZYDRK	A2F6HCTFKKJ9							
B000084ED0	AC21DFP21CJX													
B000084EK4	A1Z54EM24Y40LL													
B000084EK5	A1Z54EM24Y40LL													
B000084EK6	A1Z54EM24Y40LL													
B000084EK7	A1Z54EM24Y40LL													
B000084EK8	A1Z54EM24Y40LL													
B000084EK9	A1Z54EM24Y40LL													
B000084EKA	A1Z54EM24Y40LL		A20VA909VD90P6	AUQIKXJAWM05										
B000084EKB	A1Z54EM24Y40LL													
B000084EKC	A1Z54EM24Y40LL													
B000084EKD	A1Z54EM24Y40LL													

4) Amazon Fine Food Review Data: Amazon Counting with Counters (Numerical Summarization Pattern)

The counting with counters is a numerical summarization pattern where each attribute's number of occurrence can be counted. In this analysis, the counters are used to count the number of ratings appeared in the code from 1 – 5.



```
int code = job.waitForCompletion(true) ? 0 : 1;
if(code==0){
    for(Counter counter : job.getCounters().getGroup(CounterMapper.ratingCounterGroup)){
        System.out.println("Rating "+counter.getDisplayName()+" is given to "+counter.getValue()+" products.");
    }
}
System.exit(code);

public static class CounterMapper extends Mapper<LongWritable, Text, Text, Text>{
    public static final String ratingCounterGroup = "Rating";
    public static final String unknownCounter = "Unknown";
    public static final String nullCounter = "Null";

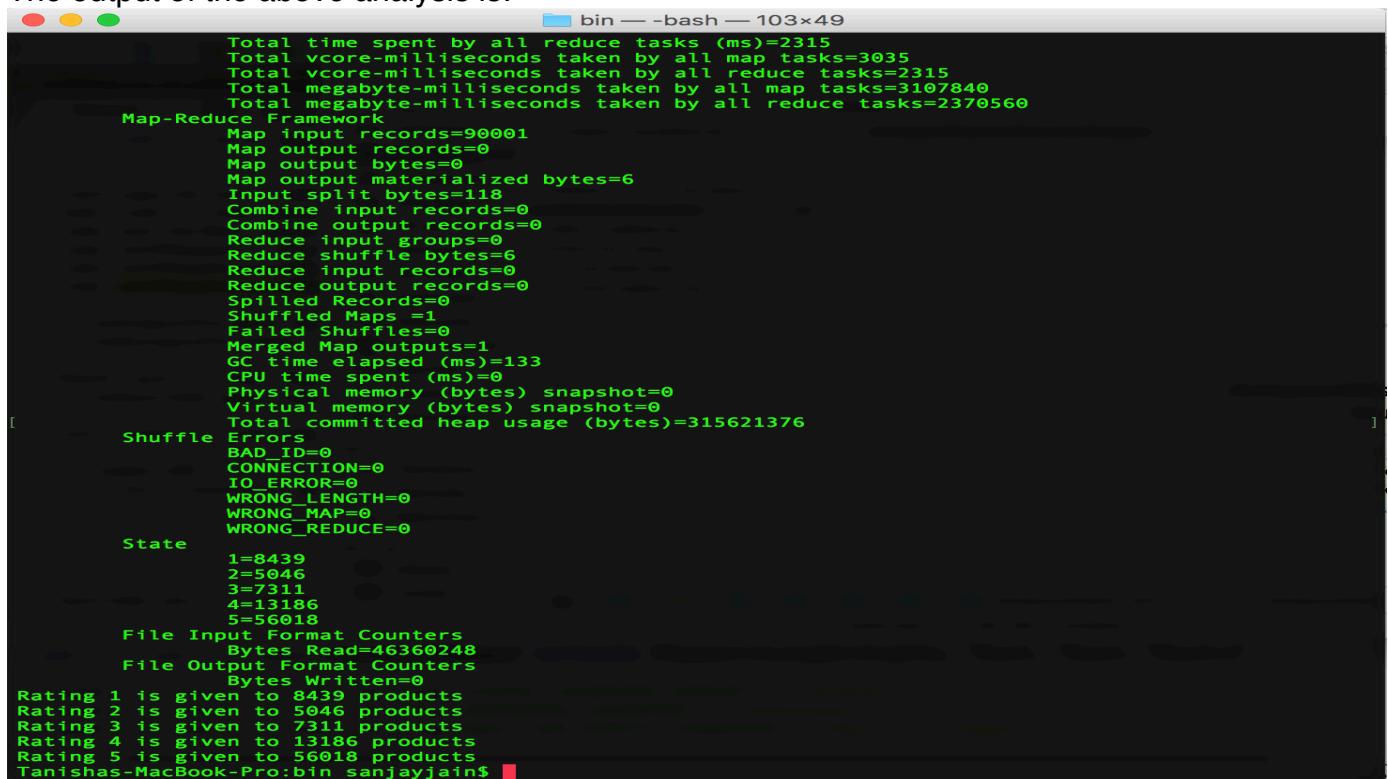
    private String[] ratingsArray = new String[] {"1", "2", "3", "4", "5"};
    private HashSet<String> ratings = new HashSet<String>(Arrays.asList(ratingsArray));

    public void map (LongWritable key, Text value, Context context) throws IOException, InterruptedException{
        if(key.get()==0){
            return;
        }

        String[] token = value.toString().split(",");
        String stateName = token[6].trim();

        if(stateName != null && !stateName.isEmpty()){
            String[] stateTokens = stateName.toUpperCase().split("\\s");
            boolean unknown = true;
            for(String state: stateTokens){
                if(state.contains(state)){
                    context.getCounter(ratingCounterGroup, state).increment(1);
                    unknown = false;
                    break;
                }
            }
        }
    }
}
```

The output of the above analysis is:



```
Total time spent by all reduce tasks (ms)=2315
Total vcore-milliseconds taken by all map tasks=3035
Total vcore-milliseconds taken by all reduce tasks=2315
Total megabyte-milliseconds taken by all map tasks=3107840
Total megabyte-milliseconds taken by all reduce tasks=2370560
Map-Reduce Framework
  Map input records=90001
  Map output records=0
  Map output bytes=0
  Map output materialized bytes=6
  Input split bytes=118
  Combine input records=0
  Combine output records=0
  Reduce input groups=0
  Reduce shuffle bytes=6
  Reduce input records=0
  Reduce output records=0
  Spilled Records=0
  Shuffled Maps =1
  Failed Shuffles=0
  Merged Map outputs=1
  GC time elapsed (ms)=133
  CPU time spent (ms)=0
  Physical memory (bytes) snapshot=0
  Virtual memory (bytes) snapshots=0
  Total committed heap usage (bytes)=315621376
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0
State
  1=8439
  2=5046
  3=7311
  4=13186
  5=56018
File Input Format Counters
  Bytes Read=46360248
File Output Format Counters
  Bytes Written=0
Rating 1 is given to 8439 products
Rating 2 is given to 5046 products
Rating 3 is given to 7311 products
Rating 4 is given to 13186 products
Rating 5 is given to 56018 products
Tanishas-MacBook-Pro:bin sanjayjain$
```

5) Amazon Fine Food Review Data: Reduce Side Join (Anti Join & Inner Join) (Join Pattern)

This pattern removes the common attributes of 2 files from a whole joining of 2 files. Anti-join is basically the full outer join minus the inner join.

Inner Join shows the output of the 2 files as one joining them on the common attributes only.

```
AmazonAntijoin.java
115     public void reduce(Text key, Iterable<Text> values, Context context) throws IOException {
116         //Clearing the lists
117         listA.clear();
118         listB.clear();
119
120         while (values.iterator().hasNext()) {
121             tmp = values.iterator().next();
122
123             if (Character.toString((char) tmp.charAt(0)).equals("A")) {
124                 listA.add(new Text(tmp.toString().substring(1)));
125             }
126             if (Character.toString((char) tmp.charAt(0)).equals("B")) {
127                 listB.add(new Text(tmp.toString().substring(1)));
128             }
129         }
130         executeJoinLogic(context);
131     }
132
133     private void executeJoinLogic(Context context) throws IOException {
134         //This is removal of common attributes of 2 files from the whole join
135         //anti join = full outer join - inner join;
136         if (joinType.equalsIgnoreCase("antijoin")) {
137             if (listA.isEmpty() ^ listB.isEmpty()) {
138
139                 for (Text A : listA) {
140                     context.write(A, EMPTY_TEXT);
141                 }
142                 for (Text B : listB) {
143                     context.write(EMPTY_TEXT, B);
144                 }
145             }
146         }
147     }
148
149
150
151
152
153 }
```

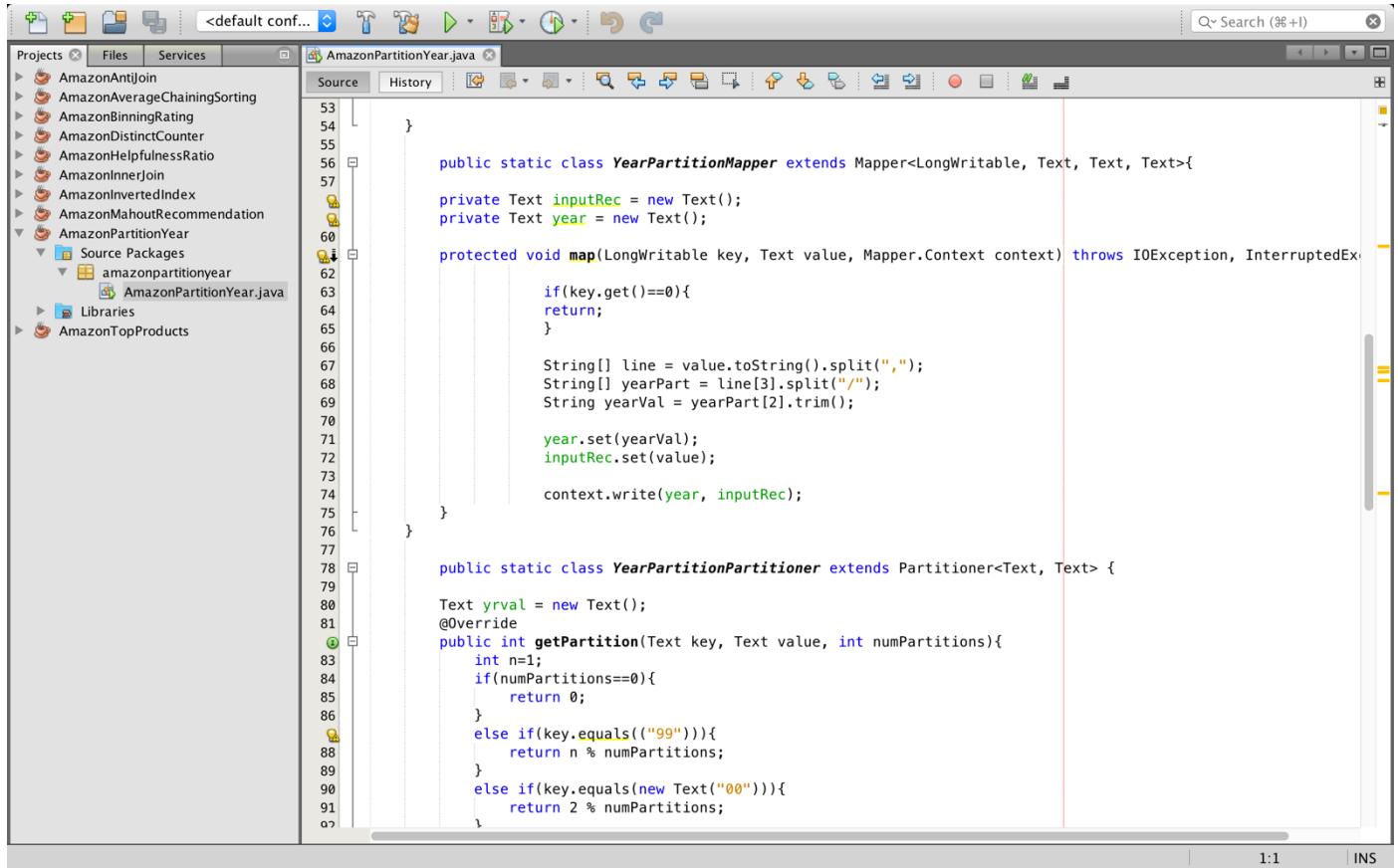
```
AmazonInnerjoin.java
119     //Clearing the lists
120     listA.clear();
121     listB.clear();
122
123     while (values.iterator().hasNext()) {
124         tmp = values.iterator().next();
125
126         if (Character.toString((char) tmp.charAt(0)).equals("A")) {
127             listA.add(new Text(tmp.toString().substring(1)));
128         }
129         if (Character.toString((char) tmp.charAt(0)).equals("B")) {
130             listB.add(new Text(tmp.toString().substring(1)));
131         }
132     }
133     executeJoinLogic(context);
134
135     private void executeJoinLogic(Context context) throws IOException {
136         if (joinType.equalsIgnoreCase("inner")) {
137             if (!listA.isEmpty() && !listB.isEmpty()) {
138                 for (Text A : listA) {
139                     for (Text B : listB) {
140                         context.write(A, B);
141                     }
142                 }
143             }
144         }
145     }
146
147
148
149
150
151
152
153
154
155
156
157 }
```

Output of Inner Join:

Output of Anti-Join:

6) Amazon Fine Food Review Data: Partitioning (Organization Pattern)

Partitioning is a pattern where data can be partitioned based on the number of partitioners mentioned in the Custom partitioner.

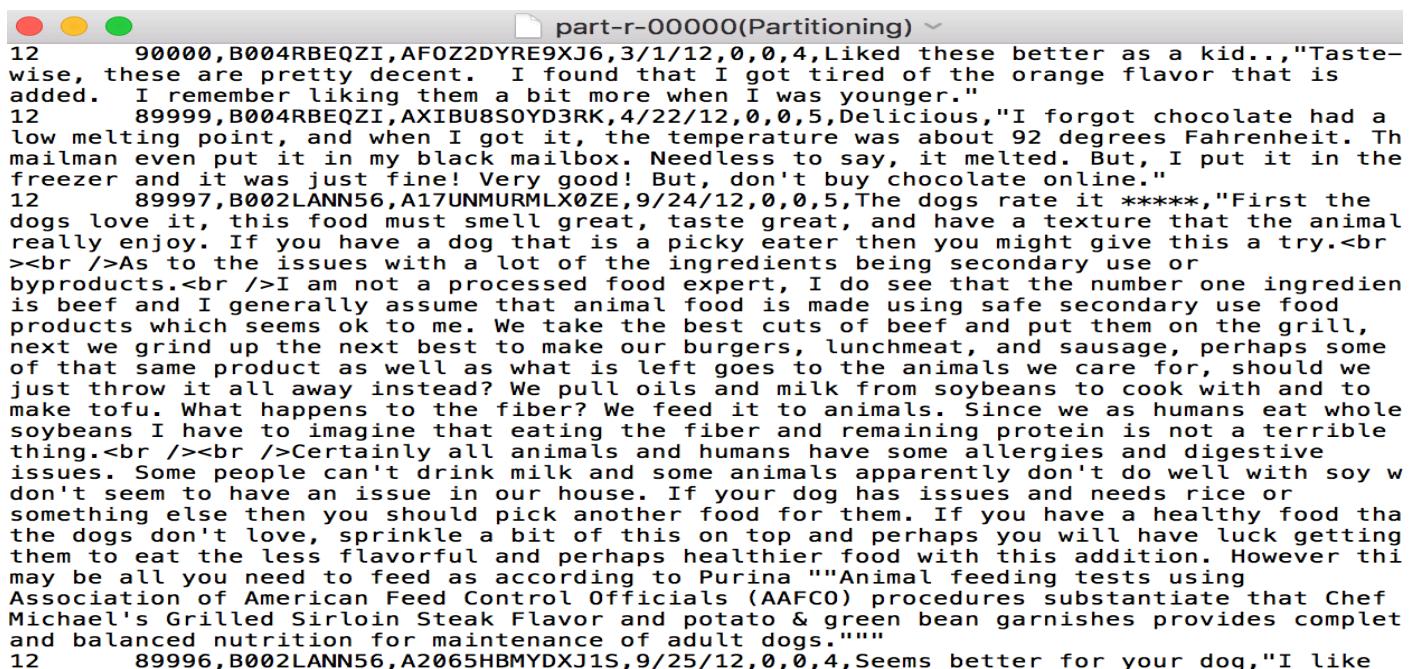


The screenshot shows an IDE interface with the following details:

- Project Explorer:** Shows various Amazon-related Java classes and packages.
- Code Editor:** Displays the `AmazonPartitionYear.java` file.
- Code Content:** The code defines a custom partitioner for handling reviews by year. It includes a `YearPartitionMapper` class for mapping key-value pairs and a `YearPartitionPartitioner` class for determining the partition index based on the review year.

```
53     }
54 }
55
56     public static class YearPartitionMapper extends Mapper<LongWritable, Text, Text, Text>{
57
58         private Text inputRec = new Text();
59         private Text year = new Text();
60
61         protected void map(LongWritable key, Text value, Mapper.Context context) throws IOException, InterruptedException {
62             if(key.get()==0){
63                 return;
64             }
65
66             String[] line = value.toString().split(",");
67             String[] yearPart = line[3].split("/");
68             String yearVal = yearPart[2].trim();
69
70             year.set(yearVal);
71             inputRec.set(value);
72
73             context.write(year, inputRec);
74         }
75     }
76
77     public static class YearPartitionPartitioner extends Partitioner<Text, Text> {
78
79         Text yrval = new Text();
80         @Override
81         public int getPartition(Text key, Text value, int numPartitions){
82             int n=1;
83             if(numPartitions==0){
84                 return 0;
85             }
86             else if(key.equals(("99"))){
87                 return n % numPartitions;
88             }
89             else if(key.equals(new Text("00"))){
90                 return 2 % numPartitions;
91             }
92         }
93     }
94 }
```

The number of partitioners is equal to the number of reducers and so there are 13 output files. Two of them are shown below:



The screenshot shows two lines of Amazon Fine Food Review Data:

```
part-r-00000(Partitioning) ~
12 90000,B004RBEQZI,AF0Z2DYLE9XJ6,3/1/12,0,0,4,Liked these better as a kid...,"Taste-
wise, these are pretty decent. I found that I got tired of the orange flavor that is
added. I remember liking them a bit more when I was younger."
12 89999,B004RBEQZI,AXIBU8S0YD3RK,4/22/12,0,0,5,Delicious,"I forgot chocolate had a
low melting point, and when I got it, the temperature was about 92 degrees Fahrenheit. Th
mailman even put it in my black mailbox. Needless to say, it melted. But, I put it in the
freezer and it was just fine! Very good! But, don't buy chocolate online."
```

part-r-00013(Partitioning) ▾

11 58693,B0013P3KC6,A3NJI70WKI0VIT,11/17/11,2,3,5,use this for nasal irrigatiin,I use xylitol for nasal irrigation. Xylitol kills bacteria and is very soothing to the nasal passages. Clear makes some xylitol packets for this purposes but it is far cheaper to buy this product in bulk. I like the now foods brand as it is pure xylitol and of good consistently.

11 19366,B003ZNRDPE,A14CEAW83UN2B3,12/15/11,0,0,4,My Pup Loves it!!,My pup loves it!! The minute I open my amazon package he smelled the bones! He instantly sat and waited patiently for me to open the pack of Smartbones.

I am so glad I finally found a better alternative to rawhide. I feel more ease at heart leaving the house with this treat in his paws.

The only reason I am giving it four stars it's because my pup was able to finish quite quickly. Otherwise I am very happy with this product and it actually smells yummy. The smell reminds me of sweet milk.

11 88,B0019CW0HE,A1T4L0Q470Z9N,11/9/11,1,1,5,Great for stomach problems!,My shepherd/collie mix has IBS. Our vet recommended a limited ingredient food. This has really helped her symptoms and she likes it. I will always buy it from Amazon...it's \$10 cheaper and free shipping!

11 58686,B0013P3KC6,A20303DEQ2NIF4,7/18/11,1,1,5,Xylitol,Amazon made it so easy to compare prices and order the Xylitol. Quick ship as well. Very pleased with entire transaction.

11 58685,B0013P3KC6,A3QVG8XW5BEEQH,10/22/11,1,1,5,sugar,My husband uses this sugar since he is a diabetic. He likes it in his coffee and tea. It dissolves easy. It is better than using anything that isn't aspartame.

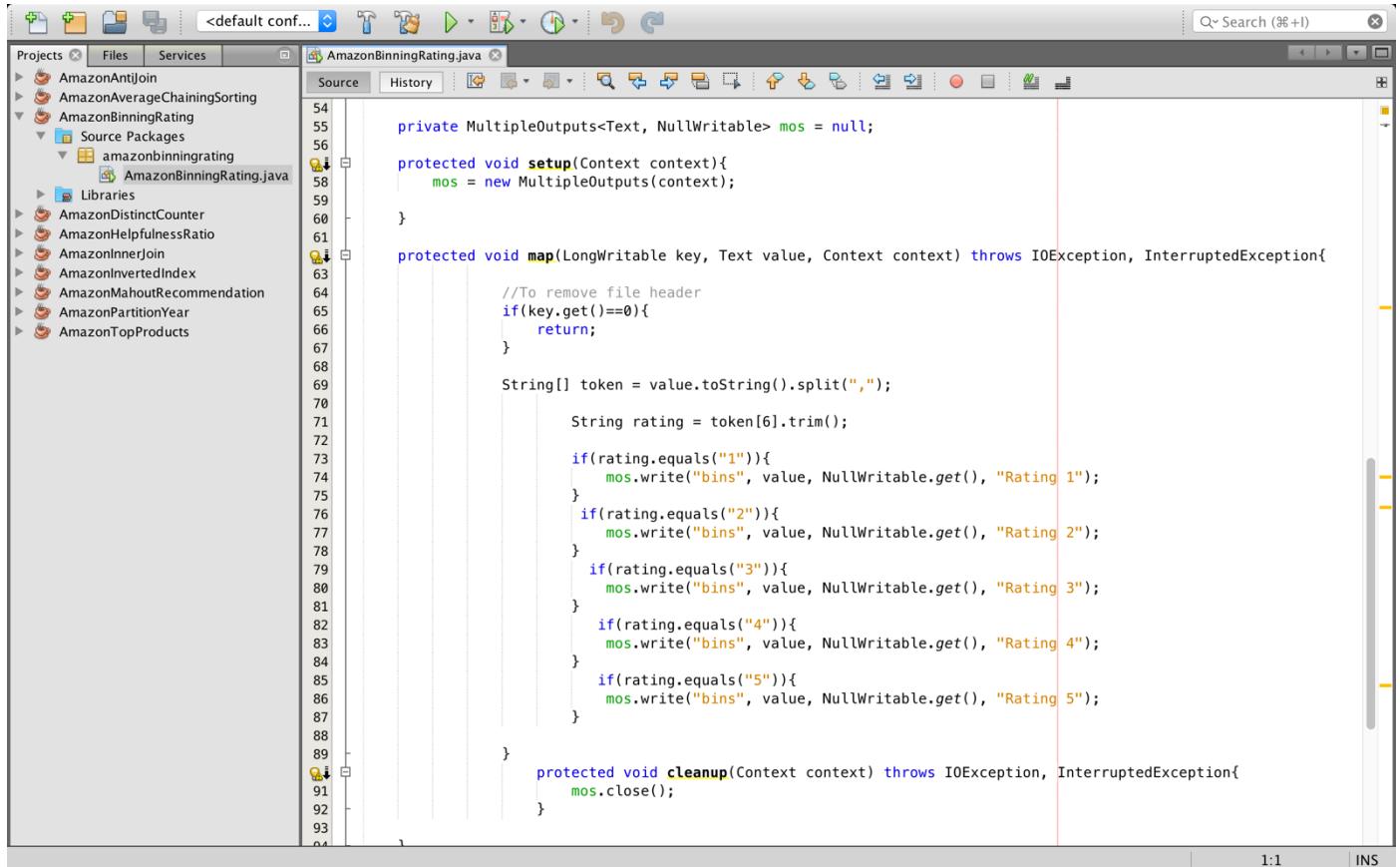
11 58684,B0013P3KC6,A2C9T12UA659SU,12/9/11,1,1,4,it tastes good,"I'm not sure if it helps teeth, but it tastes good.
I used it to sweeten tea, yogurt and other things."

11 217,B002TDK0VK,A3GRP8QRFGDH8T,1/26/11,1,4,1,Price cannot be correct,"Hey, the description says 360 grams - that is roughly 13 ounces at under \$4.00 per can. No way - that is the approximate price for a 100 gram can."

11 6249,B000E650F6,AHV9XJHC0DRTC,3/22/11,1,1,5,"Soothing, calming, allows you to relax and rest.", "Tea has a great flavor - almost what I'd call 'spearmint' flavor which is great with honey or sugar with or without cream. I drink it hot and if it does cool.

7) Amazon Fine Food Review Data: Binning (Organization Pattern)

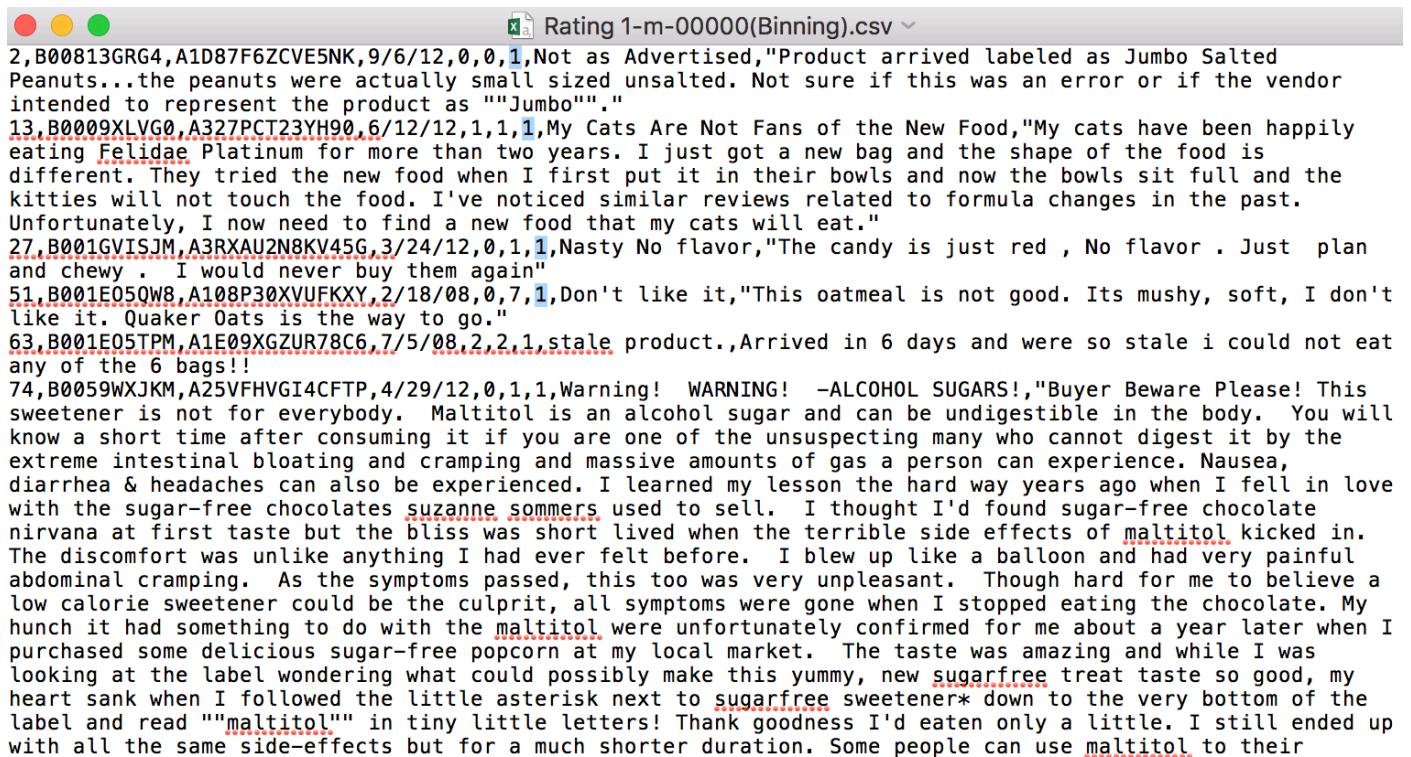
This is similar to the partitioning pattern described above but the only difference is that this is a mapper-only job and the bins have to be specified in the mapper setup method.



The screenshot shows an IDE interface with the file 'AmazonBinningRating.java' open. The code implements a mapper for binning food reviews based on their rating. It uses a 'MultipleOutputs' class to write different parts of the review to different output files based on the rating value.

```
private MultipleOutputs<Text, NullWritable> mos = null;
protected void setup(Context context){
    mos = new MultipleOutputs(context);
}
protected void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException{
    String[] token = value.toString().split(",");
    String rating = token[6].trim();
    if(rating.equals("1")){
        mos.write("bins", value, NullWritable.get(), "Rating 1");
    }
    if(rating.equals("2")){
        mos.write("bins", value, NullWritable.get(), "Rating 2");
    }
    if(rating.equals("3")){
        mos.write("bins", value, NullWritable.get(), "Rating 3");
    }
    if(rating.equals("4")){
        mos.write("bins", value, NullWritable.get(), "Rating 4");
    }
    if(rating.equals("5")){
        mos.write("bins", value, NullWritable.get(), "Rating 5");
    }
}
protected void cleanup(Context context) throws IOException, InterruptedException{
    mos.close();
}
```

The outputs of the bins are:



The screenshot shows a CSV file named 'Rating 1-m-00000(Binning).csv'. The file contains several rows of data, each representing a food review. The reviews are categorized into five bins based on their rating (1, 2, 3, 4, or 5), which corresponds to the number of stars given. The reviews include various comments from consumers about the taste, texture, and quality of the products.

Review ID	Product ID	Rating	Review Text
2	B00813GRG4	1	Not as Advertised,"Product arrived labeled as Jumbo Salted Peanuts...the peanuts were actually small sized unsalted. Not sure if this was an error or if the vendor intended to represent the product as ""Jumbo""."
13	B0009XLVG0	1	My Cats Are Not Fans of the New Food,"My cats have been happily eating Felidae Platinum for more than two years. I just got a new bag and the shape of the food is different. They tried the new food when I first put it in their bowls and now the bowls sit full and the kitties will not touch the food. I've noticed similar reviews related to formula changes in the past. Unfortunately, I now need to find a new food that my cats will eat."
27	B001GVTSJM	1	Nasty No flavor,"The candy is just red , No flavor . Just plan and chewy . I would never buy them again"
51	B001E05QW8	1	Don't like it,"This oatmeal is not good. Its mushy, soft, I don't like it. Quaker Oats is the way to go."
63	B001E05TPM	1	Arrived in 6 days and were so stale i could not eat any of the 6 bags!"
74	B0059WXJKM	1	Warning! WARNING! -ALCOHOL SUGARS!,"Buyer Beware Please! This sweetener is not for everybody. Maltitol is an alcohol sugar and can be undigestible in the body. You will know a short time after consuming it if you are one of the unsuspecting many who cannot digest it by the extreme intestinal bloating and cramping and massive amounts of gas a person can experience. Nausea, diarrhea & headaches can also be experienced. I learned my lesson the hard way years ago when I fell in love with the sugar-free chocolates suzanne sommers used to sell. I thought I'd found sugar-free chocolate nirvana at first taste but the bliss was short lived when the terrible side effects of maltitol kicked in. The discomfort was unlike anything I had ever felt before. I blew up like a balloon and had very painful abdominal cramping. As the symptoms passed, this too was very unpleasant. Though hard for me to believe a low calorie sweetener could be the culprit, all symptoms were gone when I stopped eating the chocolate. My hunch it had something to do with the maltitol were unfortunately confirmed for me about a year later when I purchased some delicious sugar-free popcorn at my local market. The taste was amazing and while I was looking at the label wondering what could possibly make this yummy, new sugarfree treat taste so good, my heart sank when I followed the little asterisk next to sugarfree sweetener* down to the very bottom of the label and read ""maltitol"" in tiny little letters! Thank goodness I'd eaten only a little. I still ended up with all the same side-effects but for a much shorter duration. Some people can use maltitol to their

 Rating 5-m-00000(Binning).csv

1,B001E4KFG0,A3SGXH7AUHU8GW,4/26/11,1,1,5,Good Quality Dog Food,I have bought several of the Vitality canned dog food products and have found them all to be of good quality. The product looks more like a stew than a processed meat and it smells better. My Labrador is finicky and she appreciates this product better than most.

5,B006K2ZZ7K,A1U0RSCLF8GW1T,10/20/12,0,0,5,Great taffy,"Great taffy at a great price. There was a wide assortment of yummy taffy. Delivery was very quick. If your a taffy lover, this is a deal."

7,B006K2ZZ7K,A1SP2KVKFXXRU1,6/19/12,0,0,5,Great! Just as good as the expensive brands!,"This saltwater taffy had great flavors and was very soft and chewy. Each candy was individually wrapped well. None of the candies were stuck together, which did happen in the expensive version, Fralinger's. Would highly recommend this candy! I served it at a beach-themed party and everyone loved it!"

8,B006K2ZZ7K,A3JRGQVEON31I0,5/2/12,0,0,5,"Wonderful, tasty taffy",This taffy is so good. It is very soft and chewy. The flavors are amazing. I would definitely recommend you buying it. Very satisfying!!

9,B000E7L2R4,A1MZY09TZK0BBI,11/22/11,1,1,5,Yay Barley,Right now I'm mostly just sprouting this so my cats can eat the grass. They love it. I rotate it around with Wheatgrass and Rye too

10,B00171APVA,A21BT40VZCCYT4,10/25/12,0,0,5,Healthy Dog Food,This is a very healthy dog food. Good for their digestion. Also good for small puppies. My dog eats her required amount at every feeding.

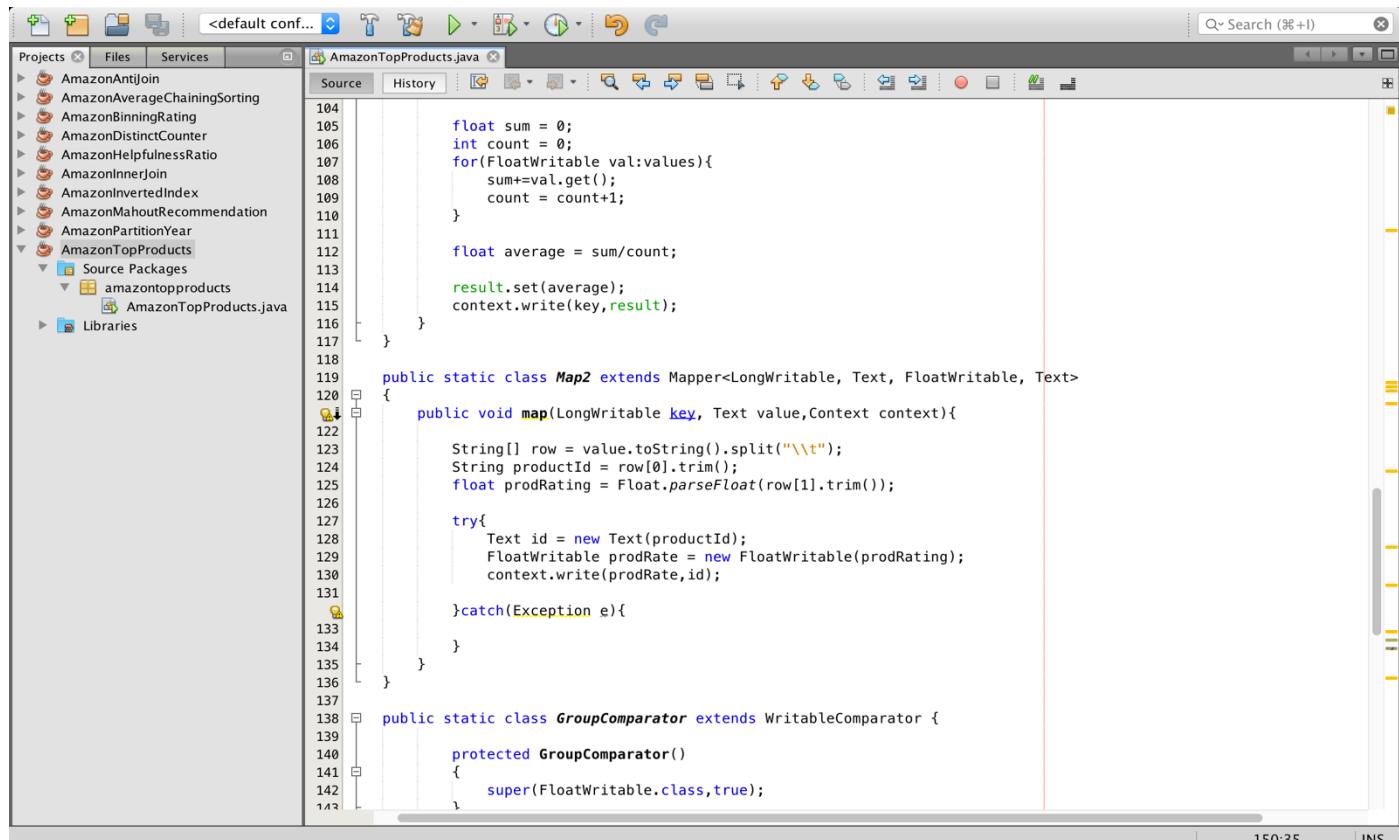
11,B0001PB9FE,A3HDK070W0QNk4,2/7/05,1,1,5,The Best Hot Sauce in the World,"I don't know if it's the cactus or the tequila or just the unique combination of ingredients, but the flavour of this hot sauce makes it one of a kind! We picked up a bottle once on a trip we were on and brought it back home with us and were totally blown away! When we realized that we simply couldn't find it anywhere in our city we were bummed.

Now, because of the magic of the internet, we have a case of the sauce and are ecstatic because of it.

If you love hot sauce..I mean really love hot sauce, but don't want a sauce that tastelessly burns your throat... grab a bottle of Tequila Picante Gourmet do

8) Amazon Fine Food Review Data: Top 100 Products (Filtering Pattern)

This filtering pattern outputs the top products and the output file shows the ProductIds and their corresponding rating.



The screenshot shows the Apache Studio IDE interface. The left pane displays a project tree with several Amazon-related components like AntiJoin, AverageChainingSorting, BinningRating, DistinctCounter, HelpfulnessRatio, InInnerJoin, InvertedIndex, MahoutRecommendation, PartitionYear, and the specific component AmazonTopProducts. The right pane shows the code for `AmazonTopProducts.java`. The code implements a MapReduce job to calculate the average rating for each product. It includes a Mapper class that splits each input line into productId and rating, converts them to `FloatWritable` objects, and writes them to context. A Reducer class is also partially defined. The code uses Java's `Map` and `Reduce` interfaces.

```
104     float sum = 0;
105     int count = 0;
106     for(FloatWritable val:values){
107         sum+=val.get();
108         count = count+1;
109     }
110
111     float average = sum/count;
112
113     result.set(average);
114     context.write(key,result);
115 }
116
117 }
118
119 public static class Map2 extends Mapper<LongWritable, Text, FloatWritable, Text>
120 {
121     public void map(LongWritable key, Text value,Context context){
122
123         String[] row = value.toString().split("\n");
124         String productId = row[0].trim();
125         float prodRating = Float.parseFloat(row[1].trim());
126
127         try{
128             Text id = new Text(productId);
129             FloatWritable prodRate = new FloatWritable(prodRating);
130             context.write(prodRate,id);
131         }catch(Exception e){}
132     }
133 }
134
135 }
136
137
138 public static class GroupComparator extends WritableComparator {
139
140     protected GroupComparator()
141     {
142         super(FloatWritable.class,true);
143     }
144 }
```

Output:



The screenshot shows a terminal window displaying the contents of the file `part-r-00000`. The file is titled `(Top100productswithrating)`. The output lists 100 product IDs, each preceded by a rating of 5.0. The products are listed in descending order of rating.

```
5.0 B000084EKC
5.0 B009C16ZP2
5.0 B000084EKG
5.0 B000084EKL
5.0 B00995MODI
5.0 B00993CW3M
5.0 B0098MEKQQ
5.0 B0097KJ6ZY
5.0 B000084F11
5.0 B000087BJ4
5.0 B00008DF6C
5.0 B0096RP84Q
5.0 B00008DFR5
5.0 B0096E5196
5.0 B00961CUXO
5.0 B00959DMWK
5.0 B0094K4VNE
5.0 B0093NIWVO
5.0 B00009608Y
5.0 B0093ASRU6
5.0 B0093A5PDK
5.0 B00934WBNI
5.0 B0092XAMDQ
5.0 B0000960IV
5.0 B000090LEW
5.0 B0092X7B5S
5.0 B0092VQCM8
5.0 B0092SQT2E
5.0 B000090LFC
5.0 B0091J1QRC
5.0 B0091ITGW0
5.0 B0000AH3OX
5.0 B00910LUQS
5.0 B0000CDBRH
5.0 B008ZL3XS2
5.0 B0000CDBRJ
5.0 B008Z1EREG
5.0 B008YKMCNG
5.0 B008YGWIJM
5.0 B008YFFNJ6
5.0 B008YC9SMW
5.0 B0000CDBRP
5.0 B008YAITNA
5.0 B008Y5NE10
5.0 B0000CE29E
5.0 B0000CEOKE
5.0 B0000CERHH
5.0 B008VRDF42
5.0 B0000CERXB
5.0 B008V20UV8
5.0 B008UF5JQC
5.0 B0000CGFCD
5.0 B0000CGFSC
5.0 B0000HDPW6
```

9) Apache Mahout:

Apache Mahout is mainly used for the recommendations and in this analysis, Mahout libraries and Java code helps find the recommendations for the users. The usual alphanumeric Ids cannot be used and the input is an edited file with three columns – userId (long), productId (long) and product ratings (int). The recommendation here is user-based.

The screenshot shows the Eclipse IDE interface. The left pane displays the project structure under 'AmazonMahoutRecommendation'. The 'Source' tab is selected in the center editor area, showing the Java code for 'AmazonMahoutRecommendation.java'. The code reads a CSV file ('InputForMahout.csv') containing user preferences and uses Mahout's GenericUserBasedRecommender to generate recommendations for each user. The right pane shows the 'Output' window with the generated recommendations for various user IDs.

```
28 /**
29  * @param args the command line arguments
30  */
31 public static void main(String[] args) throws IOException, TasteException{
32
33     File userPreferencesFile = new File("/Users/sanjayjain/Desktop/ADBMS_PROJECT/Trials/InputForMahout.csv");
34     DataModel dataModel = new FileDataModel(userPreferencesFile);
35
36     UserSimilarity userSimilarity = new PearsonCorrelationSimilarity(dataModel);
37     UserNeighborhood userNeighborhood = new ThresholdUserNeighborhood(0.1, userSimilarity, dataModel);
38
39     // Create a generic user based recommender with the dataModel, the userNeighborhood and the userSimilarity
40     Recommender genericRecommender = new GenericUserBasedRecommender(dataModel, userNeighborhood, userSimilarity);
41
42     for (LongPrimitiveIterator iterator = dataModel.getUserIDs(); iterator.hasNext();)
43     {
44         long userId = iterator.nextLong();
45
46         //5 recommendations for each user
47         List<RecommendedItem> itemRecommendations = genericRecommender.recommend(userId, 5);
48
49         System.out.format("User Id: %d%n", userId);
50
51         if (itemRecommendations.isEmpty())
52         {
53             System.out.println("There is no recommended item for this user.");
54         }
55     }
56 }
```

User Id: 1272
Recommended Item Id: 196. Strength of preference: 4.481771
Recommended Item Id: 158. Strength of preference: 4.392509
Recommended Item Id: 201. Strength of preference: 4.381534
Recommended Item Id: 167. Strength of preference: 4.340056
Recommended Item Id: 207. Strength of preference: 4.320653

User Id: 1273
There is no recommended item for this user.

User Id: 1274
Recommended Item Id: 152. Strength of preference: 5.000000
Recommended Item Id: 154. Strength of preference: 4.400000

The screenshot shows a terminal window titled 'MahoutRecommendationsOutput(Terminal)'. It displays the command-line output of the Java application, showing recommendations for various user IDs. The output includes user IDs, recommended item IDs, and their corresponding strength of preference.

```
User Id: 730
No recommendations for this user.
User Id: 737
No recommendations for this user.
User Id: 739
Recommended Item Id 210. Strength of the preference: 5.000000
Recommended Item Id 207. Strength of the preference: 4.555555
Recommended Item Id 196. Strength of the preference: 4.509434
Recommended Item Id 203. Strength of the preference: 4.423077
Recommended Item Id 158. Strength of the preference: 4.414286
User Id: 742
No recommendations for this user.
User Id: 743
Recommended Item Id 152. Strength of the preference: 4.881478
Recommended Item Id 165. Strength of the preference: 4.555309
Recommended Item Id 167. Strength of the preference: 4.492374
Recommended Item Id 193. Strength of the preference: 4.347790
Recommended Item Id 196. Strength of the preference: 4.344056
User Id: 744
No recommendations for this user.
User Id: 745
No recommendations for this user.
User Id: 749
No recommendations for this user.
User Id: 752
No recommendations for this user.
User Id: 753
Recommended Item Id 194. Strength of the preference: 4.453411
Recommended Item Id 201. Strength of the preference: 4.448694
Recommended Item Id 158. Strength of the preference: 4.425107
```

10) Apache Hive

Apache Hive provides a SQL-like interface to data stored in HDP(Hortonworks Data Platform) whose component is Hive.

Hive can run on the localhost:50070 and the data is stored inside the user folder created automatically upon installation and creation of tables.

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
drwxr-xr-x	sanjayjain	supergroup	0 B	4/26/2017, 1:52:39 AM	0	0 B	AmazonDir
drwxr-xr-x	sanjayjain	supergroup	0 B	4/26/2017, 1:54:11 AM	0	0 B	MRoutputsDir
drwxr-xr-x	sanjayjain	supergroup	0 B	4/24/2017, 2:35:38 AM	0	0 B	hive
drwx-----	sanjayjain	supergroup	0 B	4/24/2017, 11:34:52 AM	0	0 B	tmp
drwxr-xr-x	sanjayjain	supergroup	0 B	4/24/2017, 11:36:23 AM	0	0 B	user

Hive queries are used to create schema, table, load bulk data and join patterns. Here the code of hive inner join and full outer join is shown:

```
localhost: starting nodemanager, logging to /usr/local/bin/hadoop-2.7.3/logs/yarn-sanjayjain-nodemanager...
Tanishas-MacBook-Pro:local.out Tanishas-MacBook-Pro:sbin sanjayjain$ jps
31396 SecondaryNameNode
31509 ResourceManager
31637 Jps
31209 NameNode
31291 DataNode
31596 NodeManager
Tanishas-MacBook-Pro:sbin sanjayjain$ cd ..
Tanishas-MacBook-Pro:hadoop-2.7.3 sanjayjain$ cd ..
Tanishas-MacBook-Pro:bin sanjayjain$ cd apache-hive-2.1.0-bin/
Tanishas-MacBook-Pro:apache-hive-2.1.0-bin sanjayjain$ cd bin
Tanishas-MacBook-Pro:bin sanjayjain$ ./hive

Logging initialized using configuration in jar:file:/usr/local/bin/apache-hive-2.1.0-bin/lib/hive-common-2.1.0.jar!hive-log4j2.properties Async: true
Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
hive> use review;
OK
Time taken: 0.971 seconds
hive> select * from amazonfoodreview a
   > join profiles b on a.userid = b.userid;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
Query ID = sanjayjain_20170424174246_1897b6f1-5406-453e-bb9d-91014e613869
Total jobs = 1
2017-04-24 17:42:55      Starting to launch local task to process map join;          maximum memory = 477626368
2017-04-24 17:42:57      Dump the side-table for tag: 1 with group count: 64693 into file: file:/var/folders/w9/4ccpbff0x3zl_dnpvyv9p3l3c0000gn/T/sanjayjain/88aa9c3b-d3fa-41cf-9918-24598370885e/
hive_2017-04-24_17-42-46_579_3191866684908100166-1/-local-10004/HashTable-Stage-3/MapJoin-mapfile01--.hashtable
2017-04-24 17:42:57      Uploaded 1 File to: file:/var/folders/w9/4ccpbff0x3zl_dnpvyv9p3l3c0000gn/T/sanjayjain/88aa9c3b-d3fa-41cf-9918-24598370885e/hive_2017-04-24_17-42-46_579_3191866684908100166-1/-local-10004/HashTable-Stage-3/MapJoin-mapfile01--.hashtable (6037966 bytes)
2017-04-24 17:42:57      End of local task; Time Taken: 2.152 sec.
Execution completed successfully
MapredLocal task succeeded
Launching Job 1 out of 1
Number of reduce tasks is set to 0 since there's no reduce operator
Starting Job = job_1493070008739_0001, Tracking URL = http://Tanishas-MacBook-Pro.local:8088/proxy/application_1493070008739_0001/
Kill Command = /usr/local/bin/hadoop-2.7.3/bin/hadoop job -kill job_1493070008739_0001
```

```

InnerJoinOutput(TimeTaken30secs) ▾
Tanishas-MacBook-Pro:hadoop_2.7.3 sanjayjain$ cd ..
Tanishas-MacBook-Pro:bin sanjayjain$ cd apache-hive-2.1.0-bin/
Tanishas-MacBook-Pro:apache-hive-2.1.0-bin sanjayjain$ cd bin
Tanishas-MacBook-Pro:bin sanjayjain$ ./hive

Logging initialized using configuration in jar:file:/usr/local/bin/apache-hive-2.1.0-bin/lib/hive-
common-2.1.0.jar!hive-log4j2.properties Async: true
Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different
execution engine (i.e. spark, tez) or using Hive 1.X releases.
hive> use review;
OK
Time taken: 0.971 seconds
hive> select * from amazonfoodreview a
   > join profiles b on a.userid = b.userid;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a
different execution engine (i.e. spark, tez) or using Hive 1.X releases.
Query ID = sanjayjain_20170424174246_1897b6f1-5406-453e-bb9d-91014e613869
Total jobs = 1
2017-04-24 17:42:55      Starting to launch local task to process map join;      maximum memory = 477626368
2017-04-24 17:42:57      Dump the side-table for tag: 1 with group count: 64693 into file: file:/var/folders/
w9/4ccpbfox3zl_dnpvyv9p3l3c0000gn/T/sanjayjain/88aa9c3b-d3fa-41cf-9918-24598370885e/
hive_2017-04-24_17-42-46_579_3191866684908100166-1-local-10004/HashTable-Stage-3/MapJoin-mapfile01--.hashtable
2017-04-24 17:42:57      Uploaded 1 File to: file:/var/folders/w9/4ccpbfox3zl_dnpvyv9p3l3c0000gn/T/sanjayjain/
88aa9c3b-d3fa-41cf-9918-24598370885e/hive_2017-04-24_17-42-46_579_3191866684908100166-1-local-10004/HashTable-
Stage-3/MapJoin-mapfile01--.hashtable (6037966 bytes)
2017-04-24 17:42:57      End of local task; Time Taken: 2.152 sec.
Execution completed successfully
MapredLocal task succeeded
Launching Job 1 out of 1
Number of reduce tasks is set to 0 since there's no reduce operator
Starting Job = job_1493070008739_0001, Tracking URL = http://Tanishas-MacBook-Pro.local:8088/proxy/
application_1493070008739_0001/
Kill Command = /usr/local/bin/hadoop-2.7.3/bin/hadoop job -kill job_1493070008739_0001
Hadoop job information for Stage-3: number of mappers: 1; number of reducers: 0
2017-04-24 17:43:08,636 Stage-3 map = 0%, reduce = 0%
2017-04-24 17:43:17,097 Stage-3 map = 100%, reduce = 0%
Ended Job = job_1493070008739_0001
MapReduce Jobs Launched:
Stage-Stage-3: Map: 1  HDFS Read: 46369769 HDFS Write: 25115070 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
```

```

## 11) Amazon Elastic MapReduce:

Amazon Web Services provides a platform for executing MapReduce programs using its services. S3 () is used for storing the data into buckets and EMR platform is used for creating clusters and executing programming scripts. In this analysis, the execution of few java jar files and with input and output path specified on S3 is done. The cluster is created with EMR.

S3 bucket:

The screenshot shows the Amazon S3 console interface. At the top, there's a navigation bar with 'Services', 'Resource Groups', and user information. A banner below it says 'S3 Object Tagging makes it easier to manage your data. Read the blog'. Below the banner is the 'Amazon S3' logo. To the right are links for 'Switch to the old console', 'Discover the new console', and 'Quick tips'. A search bar labeled 'Search for buckets' is present. Below the search bar are buttons for '+ Create bucket', 'Delete bucket', and 'Empty bucket'. On the right, it shows '1 Buckets' and '1 Regions'. The main table lists the bucket 'hadoopmapreducebucket' with details: Region 'US East (N. Virginia)', Date created 'Apr 25, 2017 1:01:54 AM'.

Files inside S3 bucket:

The screenshot shows the 'hadoopmapreducebucket' overview page. The top navigation bar is identical to the previous one. The main area shows the bucket name 'hadoopmapreducebucket'. Below it is a navigation menu with tabs for 'Objects' (which is selected), 'Properties', 'Permissions', and 'Management'. A search bar at the top of the list area contains the placeholder 'Type a prefix and press Enter to search. Press ESC to clear.' Below the search bar are buttons for 'Upload', '+ Create folder', 'More', 'All' (selected), and 'Deleted objects'. To the right, it shows the region 'US East (N. Virginia)' and a refresh icon. The main table lists 8 objects, each with a checkbox, name, last modified date, size, and storage class. The objects listed are: 'Review', 'logs', 'output', 'AmazonFineFoodReview5lacs.csv', 'AmazonInvertedIndex.jar', 'AmazonPartitionYear.jar', 'BasicMRFoodReview.jar', and 'FineFoodReview90k.csv'.

## Creating Cluster with Configurations:

The screenshot shows the 'Create cluster' configuration page in the AWS EMR console. Key settings include:

- Cluster name:** My cluster
- Logging:** Enabled (checked), S3 folder: s3://hadoopmapreducebucket/logs/
- Launch mode:** Step execution (selected)
- Add steps:** Step type dropdown set to 'Select a step'.
- Software configuration:** Vendor: Amazon, Release: emr-5.4.0, Applications: Hadoop 2.7.3.
- Hardware configuration:** Instance type: m3.xlarge, Number of instances: 3 (1 master and 2 core nodes).
- Security and access:** Not explicitly shown in the screenshot.

## Running Clusters: (The output is right when it says Terminated All Steps completed)

The screenshot shows the 'Clusters' page in the AWS EMR console, displaying a list of clusters:

|                          | Name                            | ID              | Status                                     | Creation time (UTC-4)    | Elapsed time        |
|--------------------------|---------------------------------|-----------------|--------------------------------------------|--------------------------|---------------------|
| <input type="checkbox"/> | Partitioning-5 lacs dataset     | j-5L1379F0PCUS  | Terminated<br>All steps completed          | 2017-04-26 02:32 (UTC-4) | 8 minutes           |
| <input type="checkbox"/> | Partitioning Large Dataset      | j-166F9NFVM9ODR | Terminated with errors<br>Validation error | 2017-04-26 02:03 (UTC-4) | 1 minute            |
| <input type="checkbox"/> | Partitioning                    | j-VMFEZX7BCY7Z  | Terminated<br>All steps completed          | 2017-04-26 01:48 (UTC-4) | 9 minutes           |
| <input type="checkbox"/> | Inverted Index - 5 lacs         | j-15FL7LNBJXJ8E | Terminated<br>All steps completed          | 2017-04-26 01:34 (UTC-4) | 10 minutes          |
| <input type="checkbox"/> | Inverted Indexing - Larger Data | j-CL9WFQYUUXEN  | Terminated                                 | 2017-04-26 00:51 (UTC-4) | 9 minutes           |
| <input type="checkbox"/> | Inverted Indexing               | j-3JXG78MZRTQ7L | Terminated<br>All steps completed          | 2017-04-26 00:29 (UTC-4) | 9 minutes           |
| <input type="checkbox"/> | Inverted Indexing - MR          | j-LK25W7GF0N17  | Terminated                                 | 2017-04-25 23:56 (UTC-4) | 8 minutes           |
| <input type="checkbox"/> | My cluster                      | j-3TRW08W2DHMRX | Terminated<br>All steps completed          | 2017-04-25 10:49 (UTC-4) | 8 minutes           |
| <input type="checkbox"/> | Partitioning                    | j-2BS3YVHN6CLHM | Terminated                                 | 2017-04-25 10:26 (UTC-4) | 7 minutes           |
| <input type="checkbox"/> | Partitioning                    | j-2FF3UMOFJ3FI3 | Terminated with errors<br>Step failure     | 2017-04-25 09:59 (UTC-4) | 8 minutes           |
| <input type="checkbox"/> | My cluster                      | j-36RHVIYPTZOO8 | Terminated<br>User request                 | 2017-04-25 01:42 (UTC-4) | 8 hours, 16 minutes |
| <input type="checkbox"/> | ffr-cluster                     | j-W3CRS39LPAHV  | Terminated with errors<br>Validation error | 2017-04-25 01:37 (UTC-4) | 2 minutes           |

## **Conclusions:**

The MapReduce algorithms performed in the project are MapReduce java codes with implementation on Hadoop and filesystem HDFS. The programs are run on local machine with single node cluster and the output files are accessible on the localhost:50070.

Apache Mahout is a tool used for recommendations and is a simple java file. It can be run without Hadoop with file taken from the local machine as the input. It can display user-based recommendations.

Apache hive is a SQL-based query tool and is very fast when performing queries such as creating databases, creating tables and inserting bulk loads.

AWS EMR(Elastic Map Reduce) is very useful for storing data and creating clusters with no additional installations. This is very useful for handling large datasets.

## **Source Codes:**

### **Program 1 – Average-Chaining-Sorting**

---

```
/*
 * To change this license header, choose License Headers in Project Properties.
 * To change this template file, choose Tools | Templates
 * and open the template in the editor.
 */
package amazonaveragechainingsorting;

import java.io.IOException;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.FloatWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

/**
 *
 * @author Tanisha_Jain
 */
public class AmazonAverageChainingSorting {

 /**
 * @param args the command line arguments
 */
 public static void main(String[] args) throws IOException, InterruptedException,
 ClassNotFoundException{
 Configuration conf1 = new Configuration();
 Job job1 = Job.getInstance(conf1,"Amazon Average");
 job1.setJarByClass(AmazonAverageChainingSorting.class);
 job1.setMapperClass(Map1.class);
 job1.setMapOutputKeyClass(Text.class);
 job1.setMapOutputValueClass(FloatWritable.class);

 job1.setReducerClass(Reduce1.class);
 job1.setOutputKeyClass(Text.class);
 job1.setOutputValueClass(FloatWritable.class);

 FileInputFormat.addInputPath(job1, new Path(args[0]));
 FileOutputFormat.setOutputPath(job1, new Path(args[1]));
 boolean complete = job1.waitForCompletion(true);
 }
}
```

```

Configuration conf2 = new Configuration();
Job job2 = Job.getInstance(conf2,"Chaining Sorting");

if(complete){
job2.setJarByClass(AmazonAverageChainingSorting.class);
job2.setMapperClass(Map2.class);
job2.setMapOutputKeyClass(FloatWritable.class);
job2.setMapOutputValueClass(Text.class);

job2.setReducerClass(Reduce2.class);
job2.setOutputKeyClass(Text.class);
job2.setOutputValueClass(FloatWritable.class);

FileInputFormat.addInputPath(job2, new Path(args[1]));
FileOutputFormat.setOutputPath(job2, new Path(args[2]));

System.exit(job2.waitForCompletion(true) ? 0:1);

}

}

public static class Map1 extends Mapper<LongWritable, Text, Text, FloatWritable>{

private Text text = new Text();
private FloatWritable score = new FloatWritable();

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

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

else{

String[] line = value.toString().split(",");
String productId = line[1].trim();
float ratingVal = Float.valueOf(line[6].trim());

text.set(productId);
score.set(ratingVal);

context.write(text, score);

}
}
}

public static class Reduce1 extends Reducer<Text, FloatWritable, Text, FloatWritable>{

private FloatWritable result = new FloatWritable();


```

```

@Override
protected void reduce(Text key, Iterable<FloatWritable> values, Context context) throws
IOException, InterruptedException{
 float sum = 0;
 int count = 0;
 for(FloatWritable val:values){
 sum+=val.get();
 count = count+1;
 }
 float average = sum/count;
 result.set(average);
 context.write(key,result);
}
}

public static class Map2 extends Mapper<LongWritable, Text, FloatWritable, Text>{
 public void map(LongWritable key, Text value, Context context){
 String[] row =value.toString().split("\t");
 Text movieId = new Text(row[0]);
 float movieRatings = Float.valueOf(row[1].trim());
 try{
 FloatWritable count = new FloatWritable(movieRatings);
 context.write(count, movieId);
 }
 catch(Exception e){
 }
 }
}

public static class Reduce2 extends Reducer<FloatWritable, Text, Text, FloatWritable>{
 public void reduce(FloatWritable key, Iterable<Text> value, Context context) throws
IOException, InterruptedException{
 for(Text val : value){
 context.write(val,key);
 }
 }
}

```

}

---

## Program 2 - Inverted Index(Numerical Summarization Pattern)

---

```
/*
 * To change this license header, choose License Headers in Project Properties.
 * To change this template file, choose Tools | Templates
 * and open the template in the editor.
 */
package amazoninvertedindex;

import java.io.IOException;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

/**
 *
 * @author Tanisha_Jain
 */
public class AmazonInvertedIndex {

 /**
 * @param args the command line arguments
 */
 public static void main(String[] args) throws IOException, InterruptedException,
 ClassNotFoundException {
 // TODO code application logic here
 Configuration conf = new Configuration();
 Job job = Job.getInstance(conf, "Inverted Index");
 job.setJarByClass(AmazonInvertedIndex.class);
 job.setMapperClass(InvertedMapper.class);
 job.setReducerClass(InvertedReducer.class);

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

 job.setOutputKeyClass(Text.class);
 job.setOutputValueClass(Text.class);
 }
}
```

```

FileInputFormat.addInputPath(job, new Path(args[0]));
FileOutputFormat.setOutputPath(job, new Path(args[1]));
System.exit(job.waitForCompletion(true) ? 0 : 1);
}

public static class InvertedMapper extends Mapper<LongWritable, Text, Text, Text>{

 private Text product = new Text();
 private Text userId = new Text();

 public void map(LongWritable key, Text values, Context context){

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

 try{
 String[] tokens = values.toString().split(",");
 userId.set(tokens[2]);
 product.set(tokens[1]);

 context.write(product, userId);
 }
 catch(IOException | InterruptedException ex){
 System.out.println("Error in Mapper" + ex.getMessage());
 }
 }
}

public static class InvertedReducer extends Reducer<Text,Text,Text,Text>{

 private Text result = new Text();

 @Override
 public void reduce(Text key, Iterable<Text> values, Context context) throws IOException,
 InterruptedException{

 StringBuilder sb=new StringBuilder();
 boolean first = true;

 for(Text id:values){
 if(first){
 first = false;
 }
 else{
 sb.append(" ");
 }
 sb.append(id.toString());
 }
 }
}

```

```
 result.set(sb.toString());
 context.write(key, result);
 }
}
}
```

---

## Program 3 - Amazon Counting with Counters (Numerical Summarization Pattern)

---

```
/*
 * To change this license header, choose License Headers in Project Properties.
 * To change this template file, choose Tools | Templates
 * and open the template in the editor.
 */
package amazondistinctcounter;

import java.io.IOException;
import java.util.Arrays;
import java.util.HashSet;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.FileSystem;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Counter;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

/**
 *
 * @author Tanisha_Jain
 */
public class AmazonDistinctCounter {

 /**
 * @param args the command line arguments
 */
 public static void main(String[] args) throws IOException, InterruptedException,
 ClassNotFoundException {
 Configuration conf = new Configuration();
 Job job = Job.getInstance(conf, "Counting Product Ratings With Counters");
 job.setJarByClass(AmazonDistinctCounter.class);
 job.setMapperClass(CounterMapper.class);
 }
}
```

```

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

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

Path out = new Path(args[1]);
FileSystem.get(conf).delete(out, true);

int code = job.waitForCompletion(true) ? 0 : 1;

if(code==0){
 for(Counter counter : job.getCounters().getGroup(CounterMapper.ratingCounterGroup)){
 {
 System.out.println("Rating "+counter.getDisplayName()+" is given to
"+counter.getValue()+" products.");
 }
 }
}
System.exit(code);
}

public static class CounterMapper extends Mapper<LongWritable, Text, Text, Text>{

 public static final String ratingCounterGroup = "Rating";
 public static final String unknownCounter = "Unknown";
 public static final String nullCounter = "Null";

 private String[] ratingsArray = new String[] {"1", "2", "3", "4", "5"};
 private HashSet<String> ratings = new HashSet<String>(Arrays.asList(ratingsArray));

 public void map (LongWritable key, Text value, Context context) throws IOException,
InterruptedException{

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

 String[] token = value.toString().split(",");
 String stateName = token[6].trim();

 if(stateName != null && !stateName.isEmpty()){
 String[] stateTokens = stateName.toUpperCase().split("\\s");
 boolean unknown = true;
 for(String state: stateTokens){
 if(state.contains(state)){
 context.getCounter(ratingCounterGroup, state).increment(1);
 unknown = false;
 break;
 }
 }
 }
 }
}

```

```
 }

 if(unknown){
 context.getCounter(ratingCounterGroup, unknownCounter).increment(1);
 }
 } else{
 context.getCounter(ratingCounterGroup, nullCounter).increment(1);
 }
}

}
```

---

## Program 4 - Reduce Side Join (Anti Join & Inner Join) (Join Pattern)

---

### AntiJoin:

```
/*
 * To change this license header, choose License Headers in Project Properties.
 * To change this template file, choose Tools | Templates
 * and open the template in the editor.
 */
package amazonantijoin;

import java.io.IOException;
import java.util.ArrayList;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.MultipleInputs;
import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;

/**
 *
 * @author Tanisha_Jain
 */
public class AmazonAntiJoin {

 /**
 * @param args the command line arguments
 */
 public static void main(String[] args) throws IOException, InterruptedException,
 ClassNotFoundException {
```

---

```
Configuration conf = new Configuration();
Job job = Job.getInstance(conf, "Anti Join");
job.setJarByClass(AmazonAntiJoin.class);

MultipleInputs.addInputPath(job, new Path(args[0]), TextInputFormat.class,
AntiJoinMapper1.class);
MultipleInputs.addInputPath(job, new Path(args[1]), TextInputFormat.class,
AntiJoinMapper2.class);
job.getConfiguration().set("join.type", "antijoin");

job.setReducerClass(UserRatingIdReducer.class);

job.setOutputFormatClass(TextOutputFormat.class);
TextOutputFormat.setOutputPath(job, new Path(args[2]));

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

System.exit(job.waitForCompletion(true)? 0: 2);
}

public static class AntiJoinMapper1 extends Mapper<LongWritable, Text, Text, Text>{

private Text outkey = new Text();
private Text outvalue = new Text();

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

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

//UserIds in Amazon Fine Food Review
String[] separatedInput = value.toString().split(",");
String userId = separatedInput[2].trim();

if(userId== null) {
 return;
}

outkey.set(userId);
// Flag this record for the reducer and then output
outvalue.set("A" + value.toString());
context.write(outkey, outvalue);
}
}

public static class AntiJoinMapper2 extends Mapper<LongWritable, Text, Text, Text> {
```

---

```

private Text outkey = new Text();
private Text outvalue = new Text();

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

 if(key.get()==0){
 return;
 }
 //UserId in User Profiles
 String[] separatedInput = value.toString().split(",");
 String ratingId = separatedInput[1];

 if (ratingId == null) {
 return;
 }
 // The foreign join key is the user ID
 outkey.set(ratingId);
 // Flag this record for the reducer and then output
 outvalue.set("B" + value.toString());
 context.write(outkey, outvalue);
}
}

public static class UserRatingIdReducer extends Reducer<Text, Text, Text, Text>{

private static final Text EMPTY_TEXT = new Text("");
private Text tmp = new Text();
private ArrayList<Text> listA = new ArrayList<Text>();
private ArrayList<Text> listB = new ArrayList<Text>();

private String joinType = null;

public void setup(Reducer.Context context) {

 // Get the type of join from our configuration
 joinType = context.getConfiguration().get("join.type");
}

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

 //Clearing the lists
 listA.clear();
 listB.clear();

 while (values.iterator().hasNext()) {
 tmp = values.iterator().next();

 if (Character.toString((char) tmp.charAt(0)).equals("A")) {

```

---

```

 listA.add(new Text(tmp.toString().substring(1)));
 }
 if (Character.toString((char) tmp.charAt(0)).equals("B")) {

 listB.add(new Text(tmp.toString().substring(1)));
 }

}
executeJoinLogic(context);
}

private void executeJoinLogic(Context context) throws IOException, InterruptedException{

//This is removal of common attributes of 2 files from a whole joining of 2 files
//anti join = full outer join - inner join;
if (joinType.equalsIgnoreCase("antijoin")) {
 if (listA.isEmpty() ^ listB.isEmpty()) {

 for (Text A : listA) {
 context.write(A, EMPTY_TEXT);
 }
 for (Text B : listB) {
 context.write(EMPTY_TEXT, B);
 }
 }
}

}

}

}

```

---

### **Inner Join:**

```

/*
 * To change this license header, choose License Headers in Project Properties.
 * To change this template file, choose Tools | Templates
 * and open the template in the editor.
 */
package amazoninnerjoin;

import java.io.IOException;
import java.util.ArrayList;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;

```

---

```

import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.MultipleInputs;
import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;

/**
 *
 * @author Tanisha_Jain
 */
public class AmazonInnerJoin {

 /**
 * @param args the command line arguments
 */
 public static void main(String[] args) throws IOException, InterruptedException,
 ClassNotFoundException {
 Configuration conf = new Configuration();
 Job job = Job.getInstance(conf, "Inner Join");
 job.setJarByClass(AmazonInnerJoin.class);

 MultipleInputs.addInputPath(job, new Path(args[0]), TextInputFormat.class,
 InnerJoinMapper1.class);
 MultipleInputs.addInputPath(job, new Path(args[1]), TextInputFormat.class,
 InnerJoinMapper2.class);
 job.getConfiguration().set("join.type", "inner");
 //job.setNumReduceTasks(0);
 job.setReducerClass(UserRatingIdReducer.class);

 job.setOutputFormatClass(TextOutputFormat.class);
 TextOutputFormat.setOutputPath(job, new Path(args[2]));

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

 System.exit(job.waitForCompletion(true)? 0: 2);
 }

 public static class InnerJoinMapper1 extends Mapper<LongWritable, Text, Text, Text>{

 private Text outkey = new Text();
 private Text outvalue = new Text();

 public void map(LongWritable key, Text value, Mapper.Context context) throws IOException,
 InterruptedException{
 if(key.get()==0){
 return;
 }
 }
 }
}

```

```

//UserIds in Amazon Fine Food Review
String[] separatedInput = value.toString().split(",");
String userId = separatedInput[2].trim();

if(userId== null) {
 return;
}

outkey.set(userId);
// Flag this record for the reducer and then output
outvalue.set("A" + value.toString());
context.write(outkey, outvalue);
}
}

public static class InnerJoinMapper2 extends Mapper<LongWritable, Text, Text, Text> {

 private Text outkey = new Text();
 private Text outvalue = new Text();

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

 if(key.get()==0){
 return;
 }
 //UserId in User Profiles
 String[] separatedInput = value.toString().split(",");

 String ratingId = separatedInput[1];
 if (ratingId == null) {
 return;
 }
 // The foreign join key is the user ID
 outkey.set(ratingId);
 // Flag this record for the reducer and then output
 outvalue.set("B" + value.toString());
 context.write(outkey, outvalue);
 }
}

public static class UserRatingIdReducer extends Reducer<Text, Text, Text, Text>{

 private static final Text EMPTY_TEXT = new Text("");
 private Text tmp = new Text();
 private ArrayList<Text> listA = new ArrayList<Text>();
 private ArrayList<Text> listB = new ArrayList<Text>();

 private String joinType = null;

```

```

public void setup(Reducer.Context context) {
 // Get the type of join from our configuration
 joinType = context.getConfiguration().get("join.type");
}

public void reduce(Text key, Iterable<Text> values, Context context) throws IOException, InterruptedException {
 //Clearing the lists
 listA.clear();
 listB.clear();

 while (values.iterator().hasNext()) {
 tmp = values.iterator().next();

 if (Character.toString((char) tmp.charAt(0)).equals("A")) {
 listA.add(new Text(tmp.toString().substring(1)));
 }
 if (Character.toString((char) tmp.charAt(0)).equals("B")) {
 listB.add(new Text(tmp.toString().substring(1)));
 }
 }
 executeJoinLogic(context);
}

private void executeJoinLogic(Context context) throws IOException, InterruptedException{
 if (joinType.equalsIgnoreCase("inner")) {
 if (!listA.isEmpty() && !listB.isEmpty()) {
 for (Text A : listA) {
 for (Text B : listB) {
 context.write(A, B);
 }
 }
 }
 }
}
}
}

```

## Program 5 - Partitioning (Organization Pattern)

---

```
/*
 * To change this license header, choose License Headers in Project Properties.
 * To change this template file, choose Tools | Templates
 * and open the template in the editor.
 */
package amazonpartitionyear;

import java.io.IOException;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.NullWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Partitioner;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

/**
 *
 * @author Tanisha_Jain
 */
public class AmazonPartitionYear {

 /**
 * @param args the command line arguments
 */
 public static void main(String[] args) throws IOException, InterruptedException,
 ClassNotFoundException {
 Configuration conf = new Configuration();
 Job job = Job.getInstance(conf, "Partitioning");

 job.setJarByClass(AmazonPartitionYear.class);

 job.setMapperClass(YearPartitionMapper.class);
 job.setMapOutputKeyClass(Text.class);
 job.setMapOutputValueClass(Text.class);

 //Custom Partitioner::
 job.setPartitionerClass(YearPartitionPartitioner.class);

 job.setReducerClass(YearPartitionReducer.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]));

System.exit(job.waitForCompletion(true) ? 0 : 1);

}

public static class YearPartitionMapper extends Mapper<LongWritable, Text, Text, Text>{

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

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

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

 String[] line = value.toString().split(",");
 String[] yearPart = line[3].split("/");
 String yearVal = yearPart[2].trim();

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

 context.write(year, inputRec);
}
}

public static class YearPartitionPartitioner extends Partitioner<Text, Text> {

Text yrval = new 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("00"))){
 return 2 % numPartitions;
 }
 else if(key.equals(new Text("01"))){
 return 3 % numPartitions ;
 }
 else if(key.equals(new Text("02"))){

```

---

```
 return 4 % numPartitions;
 }
 else if(key.equals(new Text("03"))){
 return 5 % numPartitions;
 }
 else if(key.equals(new Text("04"))){
 return 6 % numPartitions;
 }
 else if(key.equals(new Text("05"))){
 return 7 % numPartitions;
 }
 else if(key.equals(new Text("06"))){
 return 8 % numPartitions;
 }
 else if(key.equals(new Text("07"))){
 return 9 % numPartitions;
 }
 else if(key.equals(new Text("08"))){
 return 10 % numPartitions;
 }
 else if (key.equals(new Text("09"))){
 return 11 % numPartitions;
 }
 else if (key.equals(new Text("10"))){
 return 12 % numPartitions;
 }
 else if (key.equals(new Text("11"))){
 return 13 % numPartitions;
 }
 else
 {
 return 14 % numPartitions;
 }
}
}
```

```
public static class YearPartitionReducer 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());
 }
 }
}
```

---

## Program 6 – Amazon Binning Rating

---

```
/*
 * To change this license header, choose License Headers in Project Properties.
 * To change this template file, choose Tools | Templates
 * and open the template in the editor.
 */
package amazonbinningrating;

import java.io.IOException;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.NullWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.MultipleOutputs;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;

/**
 *
 * @author Tanisha_Jain
 */
public class AmazonBinningRating {

 /**
 * @param args the command line arguments
 */
 public static void main(String[] args) throws IOException, InterruptedException,
 ClassNotFoundException {
 // TODO code application logic here
 Configuration conf = new Configuration();
 Job job = Job.getInstance(conf, "Binning Hour");
 job.setJarByClass(AmazonBinningRating.class);

 //Setting Mapper Class and the output key and value
 job.setMapperClass(BinningMapper.class);
 job.setMapOutputKeyClass(Text.class);
 job.setMapOutputValueClass(NullWritable.class);

 //No combiner, partitioner or reducer is used in this pattern!
 job.setNumReduceTasks(0);

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

```

 MultipleOutputs.addNamedOutput(job, "bins", TextOutputFormat.class, Text.class,
NullWritable.class);
 MultipleOutputs.setCountersEnabled(job, true);

 System.exit(job.waitForCompletion(true) ? 0 : 1);
 }

public static class BinningMapper extends Mapper<LongWritable, Text, Text, NullWritable>{
 private MultipleOutputs<Text, NullWritable> mos = null;

 protected void setup(Context context){
 mos = new MultipleOutputs(context);

 }

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

 //To remove file header
 if(key.get()==0){
 return;
 }

 String[] token = value.toString().split(",");
 String rating = token[6].trim();

 if(rating.equals("1")){
 mos.write("bins", value, NullWritable.get(), "Rating 1");
 }
 if(rating.equals("2")){
 mos.write("bins", value, NullWritable.get(), "Rating 2");
 }
 if(rating.equals("3")){
 mos.write("bins", value, NullWritable.get(), "Rating 3");
 }
 if(rating.equals("4")){
 mos.write("bins", value, NullWritable.get(), "Rating 4");
 }
 if(rating.equals("5")){
 mos.write("bins", value, NullWritable.get(), "Rating 5");
 }

 }

 protected void cleanup(Context context) throws IOException, InterruptedException{
 mos.close();
 }

}

```

}

---

## Program 7 - Top 100 Products (Filtering Pattern)

---

```
/*
 * To change this license header, choose License Headers in Project Properties.
 * To change this template file, choose Tools | Templates
 * and open the template in the editor.
 */
package amazontopproducts;

import java.io.IOException;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.FloatWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.io.WritableComparable;
import org.apache.hadoop.io.WritableComparator;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

/**
 *
 * @author Tanisha_Jain
 */
public class AmazonTopProducts {

 /**
 * @param args the command line arguments
 */
 public static void main(String[] args) throws IOException, InterruptedException,
 ClassNotFoundException {
 Configuration conf1 = new Configuration();
 Job job1 = Job.getInstance(conf1,"Top 100 Products");
 job1.setJarByClass(AmazonTopProducts.class);
 job1.setMapperClass(Map1.class);
 job1.setMapOutputKeyClass(Text.class);
 job1.setMapOutputValueClass(FloatWritable.class);

 job1.setReducerClass(Reduce1.class);
 job1.setOutputKeyClass(Text.class);
 job1.setOutputValueClass(FloatWritable.class);
```

```

FileInputFormat.addInputPath(job1, new Path(args[0]));
FileOutputFormat.setOutputPath(job1, new Path(args[1]));

boolean complete = job1.waitForCompletion(true);

Configuration conf2 = new Configuration();

if(complete){
 Job job2 = Job.getInstance(conf2,"Top 100 Products");
 job2.setJarByClass(AmazonTopProducts.class);
 job2.setMapperClass(Map2.class);
 job2.setMapOutputKeyClass(FloatWritable.class);
 job2.setMapOutputValueClass(Text.class);

 job2.setSortComparatorClass(GroupComparator.class);

 job2.setReducerClass(Reduce2.class);
 job2.setOutputKeyClass(FloatWritable.class);
 job2.setOutputValueClass(Text.class);

 FileInputFormat.addInputPath(job2, new Path(args[1]));
 FileOutputFormat.setOutputPath(job2, new Path(args[2]));

 System.exit(job2.waitForCompletion(true) ? 0 : 1);
}
}

public static class Map1 extends Mapper<LongWritable, Text, Text, FloatWritable>{

 private Text text = new Text();
 private FloatWritable score = new FloatWritable();

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

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

 else{

 String[] line = value.toString().split(",");
 String productId = line[1].trim();
 float ratingVal = Float.valueOf(line[6].trim());

 text.set(productId);
 score.set(ratingVal);

 context.write(text, score);
 }
 }
}

```

```

 }
 }

public static class Reduce1 extends Reducer<Text, FloatWritable, Text, FloatWritable>{
 private FloatWritable result = new FloatWritable();

 @Override
 protected void reduce(Text key, Iterable<FloatWritable> values, Context context) throws
 IOException, InterruptedException{
 float sum = 0;
 int count = 0;
 for(FloatWritable val:values){
 sum+=val.get();
 count = count+1;
 }

 float average = sum/count;

 result.set(average);
 context.write(key,result);
 }
}

public static class Map2 extends Mapper<LongWritable, Text, FloatWritable, Text>
{
 public void map(LongWritable key, Text value,Context context){

 String[] row = value.toString().split("\t");
 String productId = row[0].trim();
 float prodRating = Float.parseFloat(row[1].trim());

 try{
 Text id = new Text(productId);
 FloatWritable prodRate = new FloatWritable(prodRating);
 context.write(prodRate,id);

 }catch(Exception e){

 }
 }
}

public static class GroupComparator extends WritableComparator {
 protected GroupComparator()
 {
 super(FloatWritable.class,true);
 }
}

```

```

 }

 public int compare(WritableComparable w1, WritableComparable w2)
 {
 FloatWritable cw1 = (FloatWritable) w1;
 FloatWritable cw2 = (FloatWritable) w2;
 int result = cw1.get() < cw2.get() ? 1 : cw1.get() == cw2.get() ? 0 : -1;
 return result;
 }
}

public static class Reduce2 extends Reducer<FloatWritable, Text, FloatWritable, Text>{
 int count=0;
 public void reduce(FloatWritable key, Iterable<Text> value,Context context)
 throws IOException, InterruptedException{

 for(Text val: value){

 if(count<100)
 {
 context.write(key,val);
 }
 count++;
 }
 }
}

```

---

## Program 9 - Apache Mahout

---

```

/*
 * To change this license header, choose License Headers in Project Properties.
 * To change this template file, choose Tools | Templates
 * and open the template in the editor.
 */
package amazonmahoutrecommendation;

import java.io.File;
import java.io.IOException;
import java.util.List;
import org.apache.mahout.cf.taste.common.TasteException;
import org.apache.mahout.cf.taste.impl.common.LongPrimitiveIterator;

```

---

```
import org.apache.mahout.cf.taste.impl.model.file.FileDataModel;
import org.apache.mahout.cf.taste.impl.neighborhood.ThresholdUserNeighborhood;
import org.apache.mahout.cf.taste.impl.recommender.GenericUserBasedRecommender;
import org.apache.mahout.cf.taste.impl.similarity.PearsonCorrelationSimilarity;
import org.apache.mahout.cf.taste.model.DataModel;
import org.apache.mahout.cf.taste.neighborhood.UserNeighborhood;
import org.apache.mahout.cf.taste.recommender.RecommendedItem;
import org.apache.mahout.cf.taste.recommender.Recommender;
import org.apache.mahout.cf.taste.similarity.UserSimilarity;

/**
 *
 * @author Tanisha_Jain
 */
public class AmazonMahoutRecommendation {

 /**
 * @param args the command line arguments
 */
 public static void main(String[] args) throws IOException, TasteException{

 File userPreferencesFile = new
File("/Users/sanjayjain/Desktop/ADBMS_PROJECT/Trials/InputForMahout.csv");
 DataModel dataModel = new FileDataModel(userPreferencesFile);

 UserSimilarity userSimilarity = new PearsonCorrelationSimilarity(dataModel);
 UserNeighborhood userNeighborhood = new ThresholdUserNeighborhood(0.1, userSimilarity,
dataModel);

 // Create a generic user based recommender with the dataModel, the userNeighborhood and
the userSimilarity
 Recommender genericRecommender = new GenericUserBasedRecommender(dataModel,
userNeighborhood, userSimilarity);

 for (LongPrimitiveIterator iterator = dataModel.getUserIDs(); iterator.hasNext();)
{
 long userId = iterator.nextLong();

 //5 recommendations for each user
 List<RecommendedItem> itemRecommendations =
genericRecommender.recommend(userId, 5);

 System.out.format("User Id: %d%n", userId);

 if (itemRecommendations.isEmpty())
{
 System.out.println(" There is no recommended item for this user.");
}
 else
{
 // Display the list of recommendations

```

---

```

 for (RecommendedItem recommendedItem : itemRecommendations)
 {
 System.out.format(" Recommened Item Id: %d. Strength of preference: %f%n",
recommendedItem.getItemId(), recommendedItem.getValue());
 }
 }
}
}
}
}

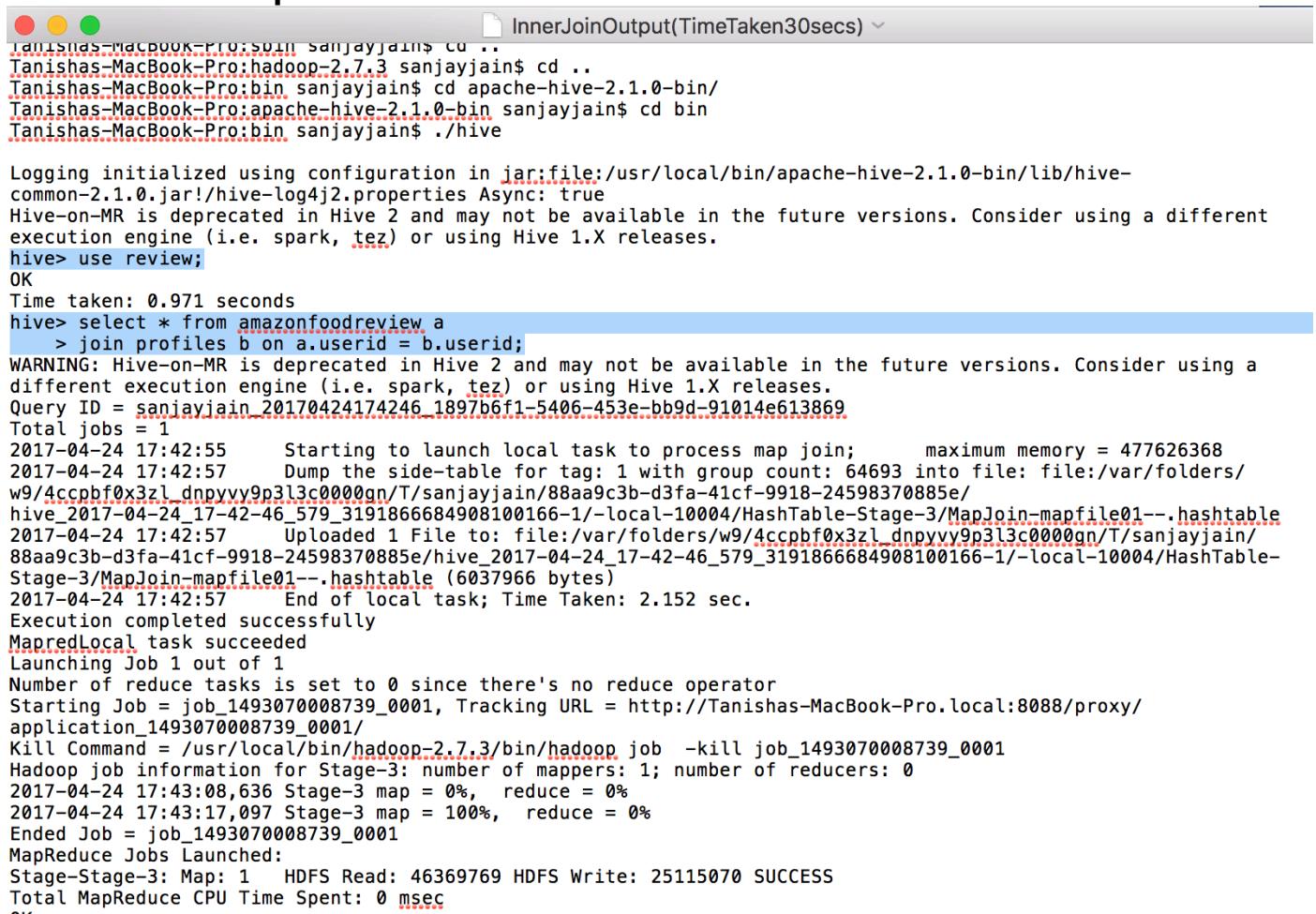
```

---

## Program 10 – Apache Hive

---

### Queries and output:



```

InnerJoinOutput(TimeTaken30secs) ~
Tanishas-MacBook-Pro:~ sanjayjain$ cd ..
Tanishas-MacBook-Pro:hadoop-2.7.3 sanjayjain$ cd ..
Tanishas-MacBook-Pro:bin sanjayjain$ cd apache-hive-2.1.0-bin/
Tanishas-MacBook-Pro:apache-hive-2.1.0-bin sanjayjain$ cd bin
Tanishas-MacBook-Pro:bin sanjayjain$./hive

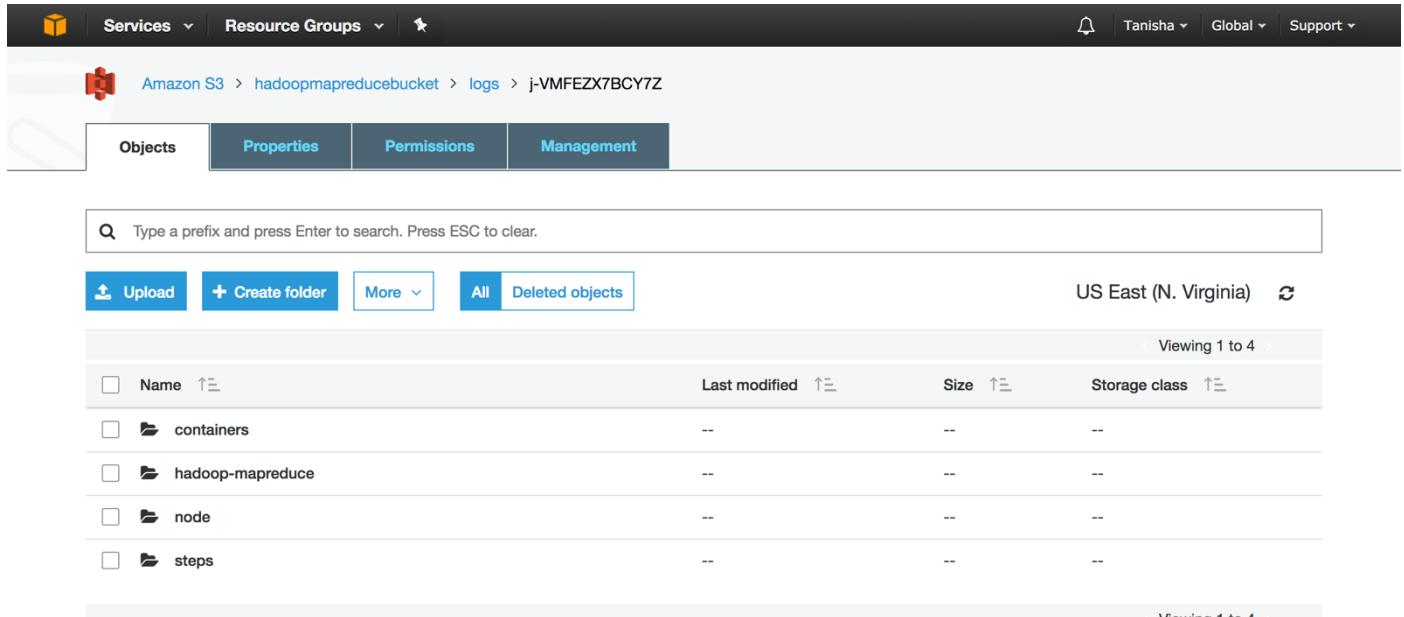
Logging initialized using configuration in jar:file:/usr/local/bin/apache-hive-2.1.0-bin/lib/hive-
common-2.1.0.jar!/hive-log4j2.properties Async: true
Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different
execution engine (i.e. spark, tez) or using Hive 1.X releases.
hive> use review;
OK
Time taken: 0.971 seconds
hive> select * from amazonfoodreview a
 > join profiles b on a.userid = b.userid;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a
different execution engine (i.e. spark, tez) or using Hive 1.X releases.
Query ID = sanjayjain_20170424174246_1897b6f1-5406-453e-bb9d-91014e613869
Total jobs = 1
2017-04-24 17:42:55 Starting to launch local task to process map join; maximum memory = 477626368
2017-04-24 17:42:57 Dump the side-table for tag: 1 with group count: 64693 into file: file:/var/folders/
w9/4ccpbfox3zl_dnpvyv9p3l3c0000gn/T/sanjayjain/88aa9c3b-d3fa-41cf-9918-24598370885e/
hive_2017-04-24_17-42-46_579_3191866684908100166-1-local-10004/HashTable-Stage-3/MapJoin-mapfile01--.hashtable
2017-04-24 17:42:57 Uploaded 1 File to: file:/var/folders/w9/4ccpbfox3zl_dnpvyv9p3l3c0000gn/T/sanjayjain/
88aa9c3b-d3fa-41cf-9918-24598370885e/hive_2017-04-24_17-42-46_579_3191866684908100166-1-local-10004/HashTable-
Stage-3/MapJoin-mapfile01--.hashtable (6037966 bytes)
2017-04-24 17:42:57 End of local task; Time Taken: 2.152 sec.
Execution completed successfully
MapredLocal task succeeded
Launching Job 1 out of 1
Number of reduce tasks is set to 0 since there's no reduce operator
Starting Job = job_1493070008739_0001, Tracking URL = http://Tanishas-MacBook-Pro.local:8088/proxy/
application_1493070008739_0001/
Kill Command = /usr/local/bin/hadoop-2.7.3/bin/hadoop job -kill job_1493070008739_0001
Hadoop job information for Stage-3: number of mappers: 1; number of reducers: 0
2017-04-24 17:43:08,636 Stage-3 map = 0%, reduce = 0%
2017-04-24 17:43:17,097 Stage-3 map = 100%, reduce = 0%
Ended Job = job_1493070008739_0001
MapReduce Jobs Launched:
Stage-Stage-3: Map: 1 HDFS Read: 46369769 HDFS Write: 25115070 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
...

```

---

# Program 11 – AWS Elastic Map Reduce

## Program logs and output files screenshots:



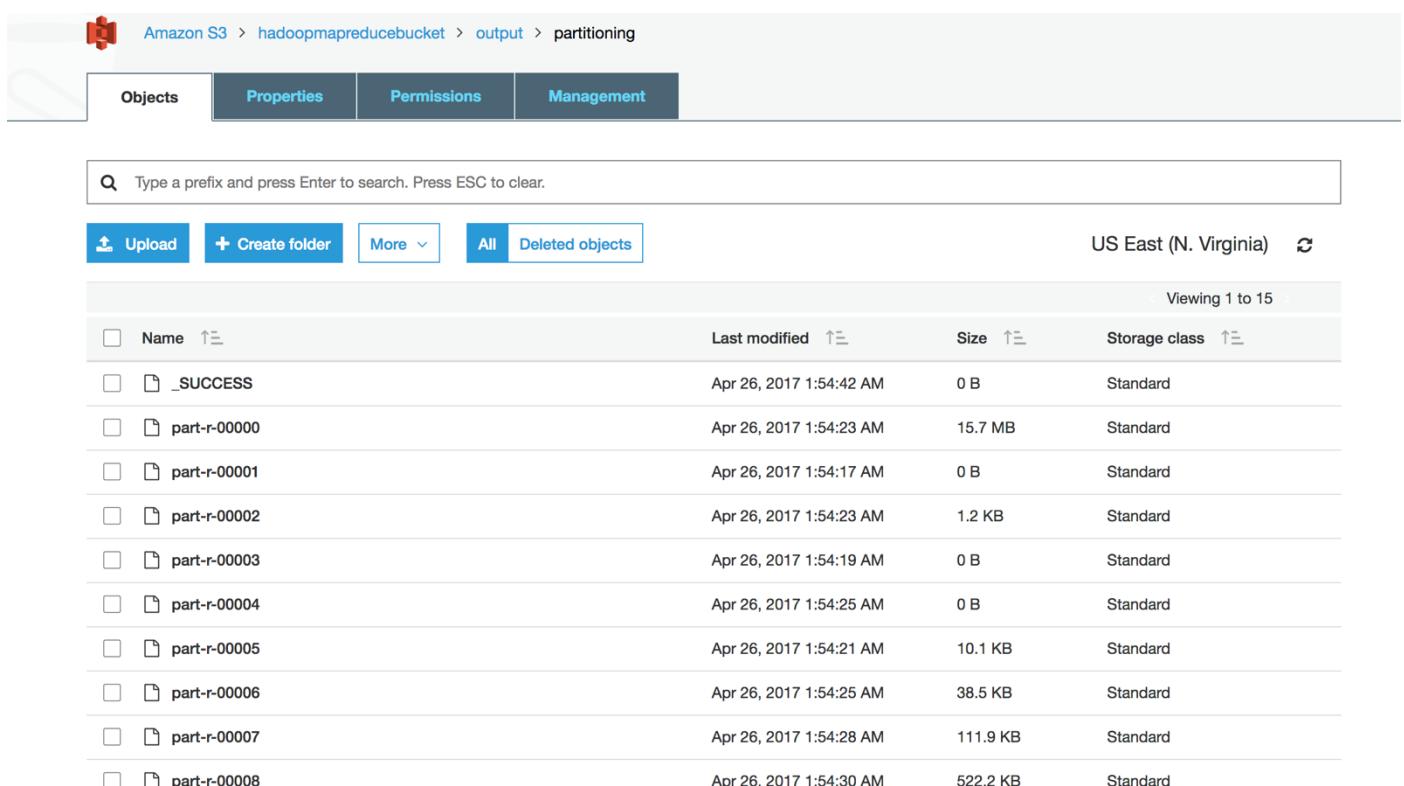
Amazon S3 > hadoopmapreducebucket > logs > j-VMFEZX7BCY7Z

Objects Properties Permissions Management

Upload Create folder More All Deleted objects US East (N. Virginia) 

| Name             | Last modified | Size | Storage class |
|------------------|---------------|------|---------------|
| containers       | --            | --   | --            |
| hadoop-mapreduce | --            | --   | --            |
| node             | --            | --   | --            |
| steps            | --            | --   | --            |

Viewing 1 to 4

Amazon S3 > hadoopmapreducebucket > output > partitioning

Objects Properties Permissions Management

Upload Create folder More All Deleted objects US East (N. Virginia) 

| Name         | Last modified           | Size     | Storage class |
|--------------|-------------------------|----------|---------------|
| _SUCCESS     | Apr 26, 2017 1:54:42 AM | 0 B      | Standard      |
| part-r-00000 | Apr 26, 2017 1:54:23 AM | 15.7 MB  | Standard      |
| part-r-00001 | Apr 26, 2017 1:54:17 AM | 0 B      | Standard      |
| part-r-00002 | Apr 26, 2017 1:54:23 AM | 1.2 KB   | Standard      |
| part-r-00003 | Apr 26, 2017 1:54:19 AM | 0 B      | Standard      |
| part-r-00004 | Apr 26, 2017 1:54:25 AM | 0 B      | Standard      |
| part-r-00005 | Apr 26, 2017 1:54:21 AM | 10.1 KB  | Standard      |
| part-r-00006 | Apr 26, 2017 1:54:25 AM | 38.5 KB  | Standard      |
| part-r-00007 | Apr 26, 2017 1:54:28 AM | 111.9 KB | Standard      |
| part-r-00008 | Apr 26, 2017 1:54:30 AM | 522.2 KB | Standard      |

Viewing 1 to 15

Amazon S3 > hadoopmapreducebucket > output > invertedindexoutput

Objects Properties Permissions Management

Type a prefix and press Enter to search. Press ESC to clear.

Upload  + Create folder  More  All  Deleted objects US East (N. Virginia)

Viewing 1 to 8

| <input type="checkbox"/> | Name         | Last modified           | Size   | Storage class |
|--------------------------|--------------|-------------------------|--------|---------------|
| <input type="checkbox"/> | _SUCCESS     | Apr 26, 2017 1:41:58 AM | 0 B    | Standard      |
| <input type="checkbox"/> | part-r-00000 | Apr 26, 2017 1:41:52 AM | 1.4 MB | Standard      |
| <input type="checkbox"/> | part-r-00001 | Apr 26, 2017 1:41:52 AM | 1.5 MB | Standard      |
| <input type="checkbox"/> | part-r-00002 | Apr 26, 2017 1:41:53 AM | 1.5 MB | Standard      |
| <input type="checkbox"/> | part-r-00003 | Apr 26, 2017 1:41:55 AM | 1.5 MB | Standard      |
| <input type="checkbox"/> | part-r-00004 | Apr 26, 2017 1:41:57 AM | 1.5 MB | Standard      |
| <input type="checkbox"/> | part-r-00005 | Apr 26, 2017 1:41:53 AM | 1.4 MB | Standard      |
| <input type="checkbox"/> | part-r-00006 | Apr 26, 2017 1:41:58 AM | 1.5 MB | Standard      |