

# Big data analytics

# Assignment 2

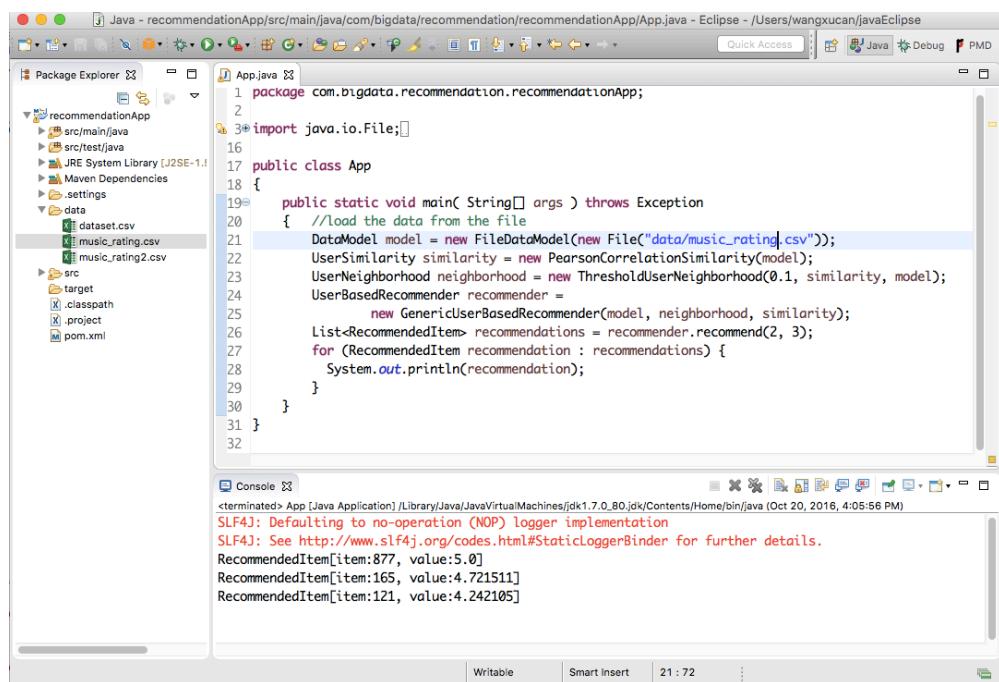
Name: Xucan Wang

Uni: xw2401

## Part 1: Recommendation

There are three parts for the recommendation. First is the DataModel which is created from .csv file. The Csv file needs to contain the userId, itemId and the ratings for a specific item for the user. The UserSimilarity defines the similarity between users, which is used to get the neighborhoods of a specific user. The UserNeighborhood computes the neighborhood of a user given a user. With DataModel, UserSimiliarity and UserNeighborhood we can create the UserBasedRecommender object, which is then used to make recommendations to the user.

### 1) PearsonCorrelationSimilarity & ThresholdUserNeignborhood



The screenshot shows the Eclipse IDE interface with the following details:

- Package Explorer:** Shows the project structure for "recommendationApp". It includes a "src" folder containing "main/java/com/bigdata/recommendation/recommendationApp/App.java", "test/java", and "JRE System Library [J2SE-1.8]". Inside "src/main/java", there is a "data" folder containing "dataset.csv", "music\_rating.csv", and "music\_rating2.csv".
- App.java:** The code implements a main method that loads data from "dataset.csv", creates a DataModel, sets up a UserSimilarity (PearsonCorrelationSimilarity), a UserNeighborhood (ThresholdUserNeighborhood), and a UserBasedRecommender (GenericUserBasedRecommender). It then recommends items for user 2 and prints them to the console.
- Console:** Displays the output of the application execution. It shows SLF4J logger information, followed by three recommended items: RecommendedItem[item:877, value:5.0], RecommendedItem[item:165, value:4.721511], and RecommendedItem[item:121, value:4.242105].

```
1 package com.bigdata.recommendation.recommendationApp;
2
3 import java.io.File;
4
5 public class App
6 {
7     public static void main( String[] args ) throws Exception
8     {
9         //Load the data from the file
10        DataModel model = new FileDataModel(new File("data/music_rating.csv"));
11        UserSimilarity similarity = new PearsonCorrelationSimilarity(model);
12        UserNeighborhood neighborhood = new ThresholdUserNeighborhood(0.1, similarity, model);
13        UserBasedRecommender recommender =
14            new GenericUserBasedRecommender(model, neighborhood, similarity);
15        List<RecommendedItem> recommendations = recommender.recommend(2, 3);
16        for (RecommendedItem recommendation : recommendations) {
17            System.out.println(recommendation);
18        }
19    }
20 }
```

## 2) LogLikelihood & ThresholdUserNeignborhood

The screenshot shows the Eclipse IDE interface with the following details:

- Project Explorer:** Shows the project structure for "recommendationApp". It includes a "data" folder containing "movie\_rating.csv" and "music\_rating.csv".
- Code Editor:** Displays the "App.java" file. The code uses LoglikelihoodSimilarity and ThresholdUserNeighborhood.
- Console:** Shows the output of the application execution. It includes SLF4J logs about failed to load class and defaulting to no-operation (NOP) logger implementation. The application then prints recommended items: RecommendedItem[item:1859, value:5.0], RecommendedItem[item:309, value:5.0], and RecommendedItem[item:759, value:5.0].

```
1 package com.bigdata.recommendation.recommendApp;
2
3 import java.io.File;
4
5 public class App
6 {
7     public static void main( String[] args ) throws Exception
8     {
9         //load the data from the file
10        DataModel model = new FileDataModel(new File("data/movie_rating.csv"));
11        UserSimilarity similarity = new LoglikelihoodSimilarity(model);
12        UserNeighborhood neighborhood = new ThresholdUserNeighborhood(0.1, similarity, model);
13        UserBasedRecommender recommender =
14            new GenericUserBasedRecommender(model, neighborhood, similarity);
15        List<RecommendedItem> recommendations = recommender.recommend(2, 3);
16        for (RecommendedItem recommendation : recommendations) {
17            System.out.println(recommendation);
18        }
19    }
20 }
```

```
<terminated> App [Java Application] /Library/Java/JavaVirtualMachines/jdk1.7.0_80.jdk/Contents/Home/bin/java (Oct 21, 2016, 7:48:30 PM)
SLF4J: Failed to load class "org.slf4j.impl.StaticLoggerBinder".
SLF4J: Defaulting to no-operation (NOP) logger implementation
SLF4J: See http://www.slf4j.org/codes.html#StaticLoggerBinder for further details.
RecommendedItem[item:1859, value:5.0]
RecommendedItem[item:309, value:5.0]
RecommendedItem[item:759, value:5.0]
```

## 3) PearsonCorrelationSimilarity&NearestNUserNeighborhood

The screenshot shows the Eclipse IDE interface with the following details:

- Project Explorer:** Shows the project structure for "recommendationApp". It includes a "data" folder containing "dataset.csv", "music\_rating.csv", and "music\_rating2.csv".
- Code Editor:** Displays the "App.java" file. The code uses PearsonCorrelationSimilarity and NearestNUserNeighborhood.
- Console:** Shows the output of the application execution. It includes SLF4J logs about defaulting to no-operation (NOP) logger implementation. The application then prints recommended items: RecommendedItem[item:877, value:5.0], RecommendedItem[item:153, value:4.5874786], and RecommendedItem[item:86, value:4.544808].

```
1 package com.bigdata.recommendation.recommendApp;
2
3 import java.io.File;
4
5 public class App
6 {
7     public static void main( String[] args ) throws Exception
8     {
9         //load the data from the file
10        DataModel model = new FileDataModel(new File("data/music_rating.csv"));
11        UserSimilarity similarity = new PearsonCorrelationSimilarity(model);
12        UserNeighborhood neighborhood = new NearestNUserNeighborhood(10, similarity, model);
13        UserBasedRecommender recommender =
14            new GenericUserBasedRecommender(model, neighborhood, similarity);
15        List<RecommendedItem> recommendations = recommender.recommend(2, 3);
16        for (RecommendedItem recommendation : recommendations) {
17            System.out.println(recommendation);
18        }
19    }
20 }
```

```
<terminated> App [Java Application] /Library/Java/JavaVirtualMachines/jdk1.7.0_80.jdk/Contents/Home/bin/java (Oct 20, 2016, 4:04:39 PM)
SLF4J: Defaulting to no-operation (NOP) logger implementation
SLF4J: See http://www.slf4j.org/codes.html#StaticLoggerBinder for further details.
RecommendedItem[item:877, value:5.0]
RecommendedItem[item:153, value:4.5874786]
RecommendedItem[item:86, value:4.544808]
```

#### 4) Loglikelihood & NearestNUserNeighborhood

The screenshot shows an IDE interface with two tabs: 'App.java' and 'Console'. The 'App.java' tab contains Java code for a recommendation application. The 'Console' tab shows the application's output.

```
1 package com.bigdata.recommendation.recommendationApp;
2
3 import java.io.File;
4
5 public class App
6 {
7     public static void main( String[] args ) throws Exception
8     {
9         //load the data from the file
10        DataModel model = new FileDataModel(new File("data/movie_rating.csv"));
11        UserSimilarity similarity = new LogLikelihoodSimilarity(model);
12        UserNeighborhood neighborhood = new NearestNUserNeighborhood(10, similarity, model);
13        UserBasedRecommender recommender =
14            new GenericUserBasedRecommender(model, neighborhood, similarity);
15        List<RecommendedItem> recommendations = recommender.recommend(2, 3);
16        for (RecommendedItem recommendation : recommendations) {
17            System.out.println(recommendation);
18        }
19    }
20 }
21
22
23
24
25
26
27
28
29
30
31
32
33
34
```

```
<terminated> App [Java Application] /Library/Java/JavaVirtualMachines/jdk1.7.0_80.jdk/Contents/Home/bin/java (Oct 21, 2016, 7:42:23 PM)
SLF4J: Failed to load class "org.slf4j.impl.StaticLoggerBinder".
SLF4J: Defaulting to no-operation (NOP) logger implementation
SLF4J: See http://www.slf4j.org/codes.html#StaticLoggerBinder for further details.
RecommendedItem[item:778, value:5.0]
RecommendedItem[item:475, value:5.0]
RecommendedItem[item:318, value:4.6249423]
```

#### 5) EuclideanDistanceSimilarity & ThresholdUserNeighborhood

The screenshot shows an IDE interface with two tabs: 'Package Explorer' and 'App.java'. The 'App.java' tab contains Java code for a recommendation application. The 'Console' tab shows the application's output.

```
1 package com.bigdata.recommendation.recommendationApp;
2
3 import java.io.File;
4
5 public class App
6 {
7     public static void main( String[] args ) throws Exception
8     {
9         //load the data from the file
10        DataModel model = new FileDataModel(new File("data/movie_rating.csv"));
11        UserSimilarity similarity = new EuclideanDistanceSimilarity(model);
12        UserNeighborhood neighborhood = new ThresholdUserNeighborhood(0.1, similarity, model);
13        UserBasedRecommender recommender =
14            new GenericUserBasedRecommender(model, neighborhood, similarity);
15        List<RecommendedItem> recommendations = recommender.recommend(2, 3);
16        for (RecommendedItem recommendation : recommendations) {
17            System.out.println(recommendation);
18        }
19    }
20 }
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
```

```
<terminated> App [Java Application] /Library/Java/JavaVirtualMachines/jdk1.7.0_80.jdk/Contents/Home/bin/java (Oct 21, 2016, 7:59:20 PM)
SLF4J: Failed to load class "org.slf4j.impl.StaticLoggerBinder".
SLF4J: Defaulting to no-operation (NOP) logger implementation
SLF4J: See http://www.slf4j.org/codes.html#StaticLoggerBinder for further details.
RecommendedItem[item:1859, value:5.0]
RecommendedItem[item:309, value:5.0]
RecommendedItem[item:759, value:5.0]
```

## 6) TanimotoCoefficientSimilarity & NearestNUserNeighborhood

The screenshot shows the Eclipse IDE interface with the following components:

- Package Explorer:** Shows the project structure for "recommendationApp". It includes a "src" folder containing "main/java" and "test/java", a "JRE System Library [J2SE-1.6]", "Maven Dependencies", ".settings", and ".project". Inside "src/main/java", there is a "data" folder containing "movie\_rating.csv" and "music\_rating.csv".
- App.java:** The main Java file. The code imports necessary classes from Apache Mahout's CF library. It defines a public class "App" with a main method. The main method loads a file named "movie\_rating.csv" into a "FileDataModel", creates a "TanimotoCoefficientSimilarity" object, a "NearestNUserNeighborhood" object with a neighborhood size of 15, and a "UserBasedRecommender" object. It then recommends 3 items for user ID 2 and prints them to the console.
- Console:** Displays the output of the application. It shows SLF4J logging messages indicating it failed to load a static logger binder and defaulted to a no-operation (NOP) logger implementation. It also shows the recommended items for user ID 2: RecommendedItem[item:318, value:4.8150177], RecommendedItem[item:307, value:4.4811134], and RecommendedItem[item:608, value:4.382062].

As we can see, we use 2 datasets in the recommender. One is the music\_rating.csv, which includes 425 users and their ratings to different songs. The dataset includes about 10102 user-supplied ratings. Another is movie\_rating.csv, which includes 100,000 ratings from 700 users. And the algorithm generates top 3 recommended item for the user with ID 2.

## Part 2: Clustering

In this part, I use the data from the example tutorial and from the reuters news

### (1) Reuters News:

#### a. K-means

In this part, I use the “cluster-reuters.sh” to do the clustering job.

At first, when I run the command, one error occurred. Then I create the directory using hdfs dfs, then the problem disappeared.

```
16/10/20 23:08:09 INFO JobSubmitter: Cleaning up the staging area file:/tmp/hadoop-wangxucan/mapred/staging/wangxucan1501208986/_staging/job_local1501208986_0001
Exception in thread "main" org.apache.hadoop.mapreduce.lib.input.InvalidInputException: Input path does not exist: hdfs://localhost:9000/tmp/mahout-work-wangxucan/reuters-out-seqdir
    at org.apache.hadoop.mapreduce.lib.input.FileInputFormat.singleThreadedListStatus(FileInputFormat.java:323)
    at org.apache.hadoop.mapreduce.lib.input.FileInputFormat.listStatus(FileInputFormat.java:265)
    at org.apache.hadoop.mapreduce.lib.input.SequenceFileInputFormat.listStatus(SequenceFileInputFormat.java:59
)
```

```
[wangxucan (master) mahout-trunk $ ./examples/bin/cluster-reuters.sh
Discovered Hadoop v2.
Setting dfs command to /usr/local/Cellar/hadoop/2.7.3/bin/dfs dfs, dfs rm to /usr/local/Cellar/hadoop/2.7.3/bin/dfs
fs dfs -rm -r -skipTrash.
Please select a number to choose the corresponding clustering algorithm
1. kmeans clustering (runs from this example script in cluster mode only)
2. fuzzykmeans clustering (may require increased heap space on yarn)
3. lda clustering
4. streamingkmeans clustering
5. clean -- cleans up the work area in /tmp/mahout-work-wangxucan
Enter your choice : 1

bytes written=288000
16/10/21 11:53:57 INFO ClusterEvaluator: Scaled Inter-Cluster Density = 0.45931222023923834
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[384] = 0.6196735100006415
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[316] = NaN
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[963] = NaN
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[788] = NaN
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[1237] = NaN
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[794] = NaN
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[734] = 0.56666666666666664
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[1455] = NaN
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[1283] = 0.56666666666666668
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[932] = NaN
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[389] = NaN
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[452] = 0.6905822623826653
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[1142] = NaN
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[108] = NaN
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[950] = 0.56666666666666668
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[1229] = NaN
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[1449] = NaN
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[1000] = NaN
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[1499] = 0.5906008266391152
16/10/21 11:53:57 INFO ClusterEvaluator: Intra-Cluster Density[19] = NaN
16/10/21 11:53:57 INFO ClusterEvaluator: Average Intra-Cluster Density = 0.6001442665170703
16/10/21 11:53:57 INFO ClusterDumper: Wrote 20 clusters
16/10/21 11:53:57 INFO MahoutDriver: Program took 10854 ms (Minutes: 0.1809)
>{"identifier":"CL-304","r":[{"0.2":0.513}, {"0.5":0.692}, {"0.6":0.768}, {"0.7":0.543}, {"0.75":0.765}], "Top Terms:
he => 3.9134248570128753
said => 3.6165204931806016
would => 3.5078525243105587
have => 3.324054274525676
last => 2.8937940247408993
from => 2.5793220063475668
some => 2.5299734569096066
billion => 2.5207431616483036
which => 2.497113834727894
banks => 2.4121890968376105
pct => 2.3752352791232663
more => 2.371036330493845
its => 2.3143290189596324
market => 2.308005894814338
has => 2.281259171612613
u.s => 2.25774195644408535
year => 2.235641827950111
about => 2.210096547653625
we => 2.197323724106475
been => 2.1645465997549205
Weight : [props - optional]: Point:
>{"identifier":"VL-316","r":[],"c":[{"0":4.031}, {"2":1.951}, {"29.53":7.363}, {"55":4.724}, {"action":4
Top Terms:
ilc => 13.455424308776855
montoya => 10.986308097839355
disbursements => 10.986308097839355
```

```

{"identifier":"VL-1499","r":[{"0":0.438},{"0.01":0.263},{"0.077":0.278},{"0.1":0.456},{"0.2":0.492},
  Top Terms:
    said                      => 1.4483192144182537
    mln                       => 1.278848085702802
    dlrsls                     => 1.1386757567477181
    vs                         => 1.073048673174001
    pct                         => 0.9123013656464698
    cts                         => 0.9051320905642659
    its                         => 0.8622784409898737
    from                        => 0.847106046053479
    year                        => 0.8291803765648226
    net                         => 0.7893499478860058
    inc                         => 0.7836996392513376
    2                           => 0.7777500565611792
    corp                        => 0.7185749607464963
    company                     => 0.7084198830449986
    4                           => 0.6160361142131016
    has                         => 0.6152814222710932
    shr                         => 0.5927909036082842
    u.s                          => 0.5921665951657341
    1986                        => 0.569558723548679
    loss                        => 0.5630149809234345
  Weight : [props - optional]: Point:

{"identifier":"VL-19","r":[],"c":[{"16":2.834}, {"26":2.816}, {"48":4.724}, {"buy":5.512}, {"common":3.5
  Top Terms:
    sterling                  => 8.452377319335938
    software                   => 7.59517240524292
    discretion                 => 7.363028049468994
    consent                     => 7.075345993041992
    ok                          => 7.075345993041992
    holders                    => 6.6149139404296875
    depending                  => 5.822583198547363
    buy                         => 5.511603355407715
    majority                    => 5.242764472961426
    required                    => 5.060442924499512
    debentures                 => 4.965132713317871
    subordinated                => 4.935279846191406
    conditions                  => 4.798078536987305
    48                          => 4.723970890045166
    convertible                 => 4.590439319610596
    received                    => 4.334506034851074
    eight                       => 4.13090705871582
    purchase                    => 4.11783504486084
    its                         => 3.9509592056274414
    common                      => 3.571291208267212
  Weight : [props - optional]: Point:

Inter-Cluster Density: 0.45931222023923834
Intra-Cluster Density: 0.6001442665170703
CDbw Inter-Cluster Density: 0.0
CDbw Intra-Cluster Density: 18.84288640124127
CDbw Separation: 28011.919431283415
wangxucan (master) mahout-trunk $
```

Page 6 of 8 208 Words English (US)

## b. Lda clustering

It takes more time than the k-means clustering method. The K-means takes about 1 minute and the Lda clustering takes about 3 minutes. So the Lda clustering takes more time

**result:**

## (2) Wikipedia article

I changed some of the content in the cluster-reuters.sh and made the cluster-wikipedia.sh to do clustering on the Wikipedia articles, here is the screen shot of my script. At first I messed up with the path thing, so I just make some change to the cluster-reuters.sh without changing the path name but the file we used is the wiki file. You can see from the script.

```

cluster-wikipedia.sh *
1 |if [ "$1" = "--help" ] || [ "$1" = "--?" ]; then
2 |echo "This script Bayes and CBayes classifiers over the last wikipedia dump."
3 |exit
4 |fi
5 |
6 # ensure that MAHOUT_HOME is set
7 if [ -z "$MAHOUT_HOME" ]; then
8 echo "Please set MAHOUT_HOME."
9 |exit
10 |fi
11 |
12 SCRIPT_PATH=$(dirname $0)
13 if [ "$0" != "$SCRIPT_PATH" ] && [ "$SCRIPT_PATH" != "" ]; then
14 cd $SCRIPT_PATH
15 fi
16 START_PATH=`pwd`
17 |
18 # Set commands for dfs
19 source ${START_PATH}/set-dfs-commands.sh
20 |
21 MAHOUT="../../bin/mahout"
22 |
23 if [ ! -e $MAHOUT ]; then
24 echo "Can't find mahout driver in $MAHOUT, cwd `pwd`, exiting..."
25 |exit 1
26 |fi
27 |
28 if [ -z "$MAHOUT_WORK_DIR" ]; then
29 WORK_DIR=/tmp/mahout-work-$USER
30 else
31 WORK_DIR=$MAHOUT_WORK_DIR
32 fi
33 |
34 algorithm=( kmeans fuzzykmeans clean )
35 if [ -n "$1" ]; then
36 choice=$1
37 else
38 echo "Please select a number to choose the corresponding task to run"
39 |exit 1
40 fi
41 |
42 case $choice in
43 |${algorithm[@]}) ;;
44 |*) echo "Unknown algorithm $choice"; exit 1;;
45 esac
46 |
47 echo "Starting $choice algorithm"
48 |
49 if [ -z "$HADOOP_HOME" ] && [ -z "$MAHOUT_LOCAL" ]; then
50 echo "Copying wikipedia data to HDFS"
51 set +
52 DFSRM ${WORK_DIR}/wikixml
53 DFSRM ${WORK_DIR}/reuters-sgm
54 DFSRM ${WORK_DIR}/reuters-out
55 DFS -mkdir -p ${WORK_DIR}/
56 DFS -mkdir ${WORK_DIR}/reuters-sgm
57 DFS -mkdir ${WORK_DIR}/reuters-out
58 DFS -put ${WORK_DIR}/reuters-sgm ${WORK_DIR}/reuters-sgm
59 DFS -put ${WORK_DIR}/reuters-out ${WORK_DIR}/reuters-out
60 set -
61 DFS -put ${START_PATH}/wikixml ${WORK_DIR}/wikixml
62 fi
63 |
64 echo "Creating sequence files from wikiXML"
65 $MAHOUT_HOME/bin/mahout seqwiki -c ${WORK_DIR}/categories.txt \
66 -i ${WORK_DIR}/wikixml/enwiki-latest-pages-articles.xml \
67 -o ${WORK_DIR}/wikipediainput
68 |
69 if [ "x$clustertype" == "xkmeans" ]; then
70 $MAHOUT seq2sparse \
71 -i ${WORK_DIR}/wikipediainput/ \
72 -o ${WORK_DIR}/reuters-out-seqdir-sparse-kmeans --maxDFPercent 85 --namedVector \
73 && \
74 fi
75 echo $START_PATH
76 |
77 set -e
78 |
79 set -x
80 echo "Preparing wikipedia data"
81 rm -rf ${WORK_DIR}/wiki
82 mkdir ${WORK_DIR}/wiki
83 |
84 cp $MAHOUT_HOME/examples/bin/resources/categories.txt ${WORK_DIR}/categories.txt
85 chmod 666 ${WORK_DIR}/categories.txt
86 |
87 if [ -z "$HADOOP_HOME" ] && [ -z "$MAHOUT_LOCAL" ]; then
88 echo "Copying wikipedia data to HDFS"
89 set +
90 DFSRM ${WORK_DIR}/wikixml
91 DFSRM ${WORK_DIR}/reuters-sgm
92 DFSRM ${WORK_DIR}/reuters-out
93 DFS -mkdir -p ${WORK_DIR}/
94 DFS -mkdir ${WORK_DIR}/reuters-sgm
95 DFS -mkdir ${WORK_DIR}/reuters-out
96 DFS -put ${WORK_DIR}/reuters-sgm ${WORK_DIR}/reuters-sgm
97 DFS -put ${WORK_DIR}/reuters-out ${WORK_DIR}/reuters-out
98 set -
99 DFS -put ${START_PATH}/wikixml ${WORK_DIR}/wikixml
100 fi
101 |
102 echo "Creating sequence files from wikiXML"
103 $MAHOUT_HOME/bin/mahout seqwiki -c ${WORK_DIR}/categories.txt \
104 -i ${WORK_DIR}/wikixml/enwiki-latest-pages-articles.xml \
105 -o ${WORK_DIR}/wikipediainput
106 |
107 if [ "x$clustertype" == "xkmeans" ]; then
108 $MAHOUT seq2sparse \
109 -i ${WORK_DIR}/wikipediainput/ \
110 -o ${WORK_DIR}/reuters-out-seqdir-sparse-kmeans --maxDFPercent 85 --namedVector \
111 && \
112 fi

```

```

cluster-wikipedia.sh ●

112 && \
113 $MAHOUT kmeans \
114   -i ${WORK_DIR}/reuters-out-seqdir-sparse-kmeans/tfidf-vectors/ \
115   -c ${WORK_DIR}/reuters-kmeans-clusters \
116   -o ${WORK_DIR}/reuters-kmeans \
117   -dm org.apache.mahout.common.distance.EuclideanDistanceMeasure \
118   -x 10 -k 20 -ow --clustering \
119 && \
120 $MAHOUT clusterdump \
121   -i '$DFS -ls -d ${WORK_DIR}/reuters-kmeans/clusters-*--final | awk '{print $8}'`' \
122   -o ${WORK_DIR}/reuters-kmeans/clusterdump \
123   -d ${WORK_DIR}/reuters-out-seqdir-sparse-kmeans/dictionary.file=0 \
124   -dt sequencefile -b 100 -n 20 --evaluate -dm org.apache.mahout.common.distance.EuclideanDistanceMeasure \
125   --pointsDir ${WORK_DIR}/reuters-kmeans/clusteredPoints \
126 && \
127 cat ${WORK_DIR}/reuters-kmeans/clusterdump
128
129 elif [ "$clustertype" == "xfuzzykmeans" ]; then
130   $MAHOUT seq2sparse \
131     -i ${WORK_DIR}/wikipediainput/ \
132     -o ${WORK_DIR}/reuters-out-seqdir-sparse-fkmeans --maxDFPercent 85 --namedVector \
133 && \
134   $MAHOUT fkmeans \
135     -i ${WORK_DIR}/reuters-out-seqdir-sparse-fkmeans/tfidf-vectors/ \
136     -c ${WORK_DIR}/reuters-fkmeans-clusters \
137     -o ${WORK_DIR}/reuters-fkmeans \
138     -dm org.apache.mahout.common.distance.EuclideanDistanceMeasure \
139     -x 10 -k 20 -ow -m 1.1 \
140 && \
141   $MAHOUT clusterdump \
142     -i '$DFS -ls -d ${WORK_DIR}/reuters-fkmeans/clusters-*--final | awk '{print $8}'`' \
143     -o ${WORK_DIR}/reuters-fkmeans/clusterdump \
144     -d ${WORK_DIR}/reuters-out-seqdir-sparse-fkmeans/dictionary.file=0 \
145     -dt sequencefile -b 100 -n 20 -sp 0 \
146 && \
147   cat ${WORK_DIR}/reuters-fkmeans/clusterdump
148 fi
149

```

## K-means clustering

```

[wangxucan (master *) mahout-trunk $ ./examples/bin/cluster-wikipedia.sh
Discovered Hadoop v2.
Setting dfs command to /usr/local/Cellar/hadoop/2.7.3/bin/hdfs dfs, dfs rm to /usr/local/Cellar/hadoop/2.7.3/bin/hdfs dfs
-rm -r -skipTrash.
Please select a number to choose the corresponding task to run
1. Kmeans (may require increased heap space on yarn)
2. fuzzykmeans
3. clean -- cleans up the work area in /tmp/mahout-work-wangxucan
Enter your choice : 1

```

```

{"identifier": "VL-261", "r": [{"0": 2.786}, {"0.01": 1.229}, {"0.05": 1.255}, {"0.06": 1.5}, {"0.09": 1.464}, {"0.1": 1.464}], "Weight : [props - optional]: Point:
Top Terms:
ipa => 19.673717751222497
h:title => 15.814844468060661
alphabet => 15.066266452564912
languages => 14.69857546862434
consonant => 14.522702974431654
style => 14.51487386226654
language => 13.67295113731833
dotted => 12.629559853497673
arabic => 12.624196529388428
vowels => 12.294184656704173
vowel => 12.192736373228186
transl => 11.714216821333942
script => 11.44367150699391
ar => 11.166593158946318
dialects => 10.96870883773355
consonants => 10.273593537947711
debian => 10.176236320944394
span => 9.897901969797472
unicode => 9.775709572960348
text => 9.628707247621874
Weight : [props - optional]: Point:
{"identifier": "VL-521", "r": [], "c": [{"0": 3.522}, {"0.0117": 14.948}, {"00381": 6.781}, {"03": 3.354}, {"05": 0.0}], "Weight : [props - optional]: Point:
Top Terms:
cyprinids => 34.249671936035156
cyprinidae => 31.70903205871582
carp => 30.4379825592041
ndash => 27.325862884521484
chubs => 25.890316009521484
daces => 24.7880916595459
leuciscinae => 23.634517669677734
minnows => 23.634517669677734
value:rgb => 23.489185333251953
subfamily => 23.09800148010254
fish => 21.207988739013672
minnow => 21.205276489257812
cypriniformes => 21.205276489257812
fishes => 21.14458465576172
barb => 21.144407272338867
barb.23.03 => 21.139354705810547
2.588 => 21.139354705810547
labeoninae => 19.77405548895703
cyprinus => 18.701295852661133
teleostei => 18.547693252563477
Weight : [props - optional]: Point:
Inter-Cluster Density: 0.516279679351974
Intra-Cluster Density: 0.5938834956053567
CDbw Inter-Cluster Density: 0.0
CDbw Intra-Cluster Density: 7.2931502624357964
CDbw Separation: 122048.75037995353

```

### fuzzykmeans:

The fuzzykmeans took more time than the k-means and at last it caused a running out of memory error like the following. It took more heap space and resources. So in this part I only included the error it caused.

```

[wangxucan (master *) mahout-trunk $ ./examples/bin/cluster-wikipedia.sh
Discovered Hadoop v2.
Setting dfs command to /usr/local/Cellar/hadoop/2.7.3/bin/hdfs dfs, dfs rm to /usr/local/Cellar/hadoop/2.7.3/bin/hdfs dfs
-rm -r -skipTrash.
Please select a number to choose the corresponding task to run
1. kmeans (may require increased heap space on yarn)
2. fuzzykmeans
3. clean -- cleans up the work area in /tmp/mahout-work-wangxucan
Enter your choice : 2

```



```

bin — bash — 122x38
university => 2.9158049378909533
:{"identifier":"SC-3308","r": [{"0":2.532}, {"0,0":0.386}, {"0,0,0,0.2":0.273}, {"0,0,0,0.75":0.249}, {"0,
Top Terms:
his => 4.909861443134916
web => 4.672786094358718
accessdate => 4.6444707443407145
he => 4.603584201289596
news => 3.60802931419652074
book => 3.582953779862118
p => 3.57151741608864045
journal => 3.541884710880157
had => 3.529570198296847
0 => 3.471325876773377
isbn => 3.4514516496645946
work => 3.385435091999827
https => 3.1996617975938366
world => 3.172006711929576
nbsp => 3.1148656820878946
3 => 3.1052091515010005
who => 3.060180431723578
2010 => 3.0528499262785544
2011 => 3.046125370201141
university => 2.995477666580426
:{"identifier":"SC-1737","r": [{"0":2.478}, {"0,0":0.371}, {"0,0,0,0.2":0.242}, {"0,0,0,0.75":0.217}, {"0,
Top Terms:
his => 4.369862970140771
he => 4.132471001354284
web => 3.9835628012725675
accessdate => 3.9418508235808916
book => 3.0928718904267023
journal => 3.056672563538523
p => 3.049803999571807
had => 3.0405236406344476
news => 3.003432296933629
0 => 2.9931302252974565
isbn => 2.9797833154461313
work => 2.90850779913192
https => 2.730923577826035

```

### (3)synthetic control.data

#### K-means clustering

```

mahout-trunk — java -Xmx1000m -Djava.net.preferIPv4Stack=true -Djava.security.krb5.realm=-Djava.security.krb5.kdc=-Dhadoop.log.dir=/usr/local/Cellar/hadoop/2.7.3/libexec/logs -Dhadoop.log.file=
[wangxucan (master) mahout-trunk $ hdfs dfs -mkdir -p /user/wangxucan/testdata
16/10/21 20:11:41 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
[wangxucan (master) mahout-trunk $ hdfs dfs -put /Users/wangxucan/Desktop/synthetic_control.data testdata/synthetic_control.data
16/10/21 20:12:08 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
[wangxucan (master) mahout-trunk $ hadoop jar libexec/mahout-examples-0.11.0-job.jar org.apache.mahout.clustering.syntheticcontrol.kmeans.Job
Not a valid JAR: /usr/local/Cellar/mahout-trunk/libexec/mahout-examples-0.11.0-job.jar
[wangxucan (master) mahout-trunk $ hadoop jar libexec/mahout-examples-0.11.0-job.jar org.apache.mahout.clustering.syntheticcontrol.kmeans.Job
16/10/21 20:14:28 INFO Kmeans: Job: Running with default arguments
16/10/21 20:14:28 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable

```

**result:**

| Permission | Owner     | Group      | Size  | Last Modified                      | Replication | Block Size | Name             |
|------------|-----------|------------|-------|------------------------------------|-------------|------------|------------------|
| -rw-r--r-- | wangxucan | supergroup | 194 B | October 21, 2016 at 8:14:49 PM EDT | 1           | 128 MB     | _policy          |
| drwxr-xr-x | wangxucan | supergroup | 0 B   | October 21, 2016 at 8:14:50 PM EDT | 0           | 0 B        | clusteredPoints  |
| drwxr-xr-x | wangxucan | supergroup | 0 B   | October 21, 2016 at 8:14:35 PM EDT | 0           | 0 B        | clusters-0       |
| drwxr-xr-x | wangxucan | supergroup | 0 B   | October 21, 2016 at 8:14:37 PM EDT | 0           | 0 B        | clusters-1       |
| drwxr-xr-x | wangxucan | supergroup | 0 B   | October 21, 2016 at 8:14:38 PM EDT | 0           | 0 B        | clusters-2       |
| drwxr-xr-x | wangxucan | supergroup | 0 B   | October 21, 2016 at 8:14:40 PM EDT | 0           | 0 B        | clusters-3       |
| drwxr-xr-x | wangxucan | supergroup | 0 B   | October 21, 2016 at 8:14:41 PM EDT | 0           | 0 B        | clusters-4       |
| drwxr-xr-x | wangxucan | supergroup | 0 B   | October 21, 2016 at 8:14:43 PM EDT | 0           | 0 B        | clusters-5       |
| drwxr-xr-x | wangxucan | supergroup | 0 B   | October 21, 2016 at 8:14:44 PM EDT | 0           | 0 B        | clusters-6       |
| drwxr-xr-x | wangxucan | supergroup | 0 B   | October 21, 2016 at 8:14:46 PM EDT | 0           | 0 B        | clusters-7       |
| drwxr-xr-x | wangxucan | supergroup | 0 B   | October 21, 2016 at 8:14:48 PM EDT | 0           | 0 B        | clusters-8       |
| drwxr-xr-x | wangxucan | supergroup | 0 B   | October 21, 2016 at 8:14:49 PM EDT | 0           | 0 B        | clusters-9-final |
| drwxr-xr-x | wangxucan | supergroup | 0 B   | October 21, 2016 at 8:14:33 PM EDT | 0           | 0 B        | data             |
| drwxr-xr-x | wangxucan | supergroup | 0 B   | October 21, 2016 at 8:14:34 PM EDT | 0           | 0 B        | random-seeds     |

### Canopy clustering:

```
wangxuan [master] manout-trunk> $ hadoop jar $EXAMPLES/manout-examples-v.1.1-w-100.jar org.apache.manout.clustering.syntheticcontrol.canopy.job  
16/10/21 20:17:18 INFO InputFormat: Job: Beginning with default arguments  
16/10/21 20:17:19 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable  
16/10/21 20:17:20 INFO Configuration: Deleting output  
16/10/21 20:17:20 INFO Configuration.deprecation: session.id is deprecated. Instead, use dfs.metrics.session-id  
16/10/21 20:17:20 INFO jvm.JvmMetrics: Initializing JVM Metrics with processName=JobTracker, sessionId=  
16/10/21 20:17:21 WARN JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.  
16/10/21 20:17:22 INFO InputFormat: Total input paths to process : 1
```

result:

## Browse Directory

/user/wangxucan/output Go!

| Permission | Owner     | Group      | Size | Last Modified                      | Replication | Block Size | Name           |
|------------|-----------|------------|------|------------------------------------|-------------|------------|----------------|
| drwxr-xr-x | wangxucan | supergroup | 0 B  | October 21, 2016 at 8:17:27 PM EDT | 0           | 0 B        | clusteredPoint |
| drwxr-xr-x | wangxucan | supergroup | 0 B  | October 21, 2016 at 8:17:26 PM EDT | 0           | 0 B        | clusters-0-fin |
| drwxr-xr-x | wangxucan | supergroup | 0 B  | October 21, 2016 at 8:17:23 PM EDT | 0           | 0 B        | data           |

Hadoop, 2016.

## Part 3: Classification

Use datasets from online news or Wikipedia articles and do clustering to find related documents.

## 1) 20 news by date

In this part, I use the news dataset to train a Cnaive Bayesian model to do classification

a. Upload file and make directory in HDFS

```
wangxucan (master) mahout-trunk $ hdfs dfs -mkdir /work  
16/10/21 12:06:27 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable  
wangxucan (master) mahout-trunk $ hdfs dfs -put /Users/wangxucan/Downloads/20news-bydate /work  
16/10/21 12:07:50 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable  
wangxucan (master) mahout-trunk $ hdfs dfs -mkdir /work/20news-all  
16/10/21 12:28:10 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable  
wangxucan (master) mahout-trunk $ hdfs dfs -cp -p /work/20news-bydate/*/* /work/20news-all  
16/10/21 12:28:49 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable  
cp: '/work/20news-bydate/*/*': No such file or directory  
wangxucan (master) mahout-trunk $ hdfs dfs -cp -p /work/20news-bydate/*/* /work/20news-all  
16/10/21 12:29:08 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable  
wangxucan (master) mahout-trunk $
```

b. Convert the dataset into <Text,Text> sequence file

```
[wangxucan (master) mahout-trunk $ ./bin/mahout seqdirectory -i /work/20news-all -o /work/20news-seq -ow
Running on hadoop, using /usr/local/Cellar/hadoop/2.7.3/bin/hadoop and HADOOP_CONF_DIR=
MAHOUT-JOB: /usr/local/Cellar/mahout-trunk/examples/target/mahout-examples-0.13.0-SNAPSHOT-job.jar
16/10/21 12:46:37 INFO AbstractJob: Command line arguments: {--charset=[UTF-8], --chunkSize=[64], --endPhase=[2147483647], --fileFilterClass=[org.apache.mahout.text.PrefixAdditionFilter], --input=[/work/20news-all], --keyPrefix=[], --method=[mapreduce], --output=[/work/20news-seq], --overwrite=null, --startPhase=[0], --tempdir=[temp]}
16/10/21 12:46:38 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
16/10/21 12:46:38 INFO deprecation: mapred.input.dir is deprecated. Instead, use mapreduce.input.fileinputformat.inputdir
16/10/21 12:46:38 INFO deprecation: mapred.compress.map.output is deprecated. Instead, use mapreduce.map.output.compress
16/10/21 12:46:38 INFO deprecation: mapred.output.dir is deprecated. Instead, use mapreduce.output.fileoutputformat.outputdir
16/10/21 12:46:39 INFO deprecation: session.id is deprecated. Instead, use dfs.metrics.session-id
```

c. Convert the dataset into a <Text, VectorWrtble> SequenceFile Containing term frequencies for each document



d. Split the preprocessed data into training set and testing set

```
16/10/21 14:02:53 INFO MahoutDriver: Program took 52953 ms (Minutes: 0.88255)
[wangxucan (master) mahout-trunk $ ./bin/mahout split -i /work/20news-vectors/tfidf-vectors --trainingOutput /work/20news-train-vectors --testOut put /work/20news-test-vectors --randomSelectionPct 40 --overwrite --sequenceFiles -xm sequential
Running on hadoop, using /usr/local/Cellar/hadoop/2.7.3/bin/hadoop and HADOOP_CONF_DIR=
MAHOUT-JOB: /usr/local/Cellar/mahout-trunk/examples/target/mahout-examples-0.13_0-SNAPSHOT-job.jar
```

e. Train the classifier

```
[root@centos master]# mahout-trunk$ ./bin/mahout trainbm -i /work/20news-train-vectors -o /work/model -li /work/labelindex -ow -o /work/20news-testing -c
Running on hadoop, using /usr/local/cellar/mahout/hadoop/2.7.3/bin/hadoop and HADOOP_CONF_DIR=
MAHOUT-JOB-0: /usr/local/cellar/mahout-trunk/examples/target/mahout-examples-0.13.0-SNAPSHOT-job.jar
16/10/21 14:19:33 WARN MahoutDriver: No trainbm.props found on classpath, will use command-line arguments only
16/10/21 14:19:33 INFO AbstractJob: Command line arguments: {--alpha=[1.0], --endPhase=[2147483647], --input=[/work/20news-train-vectors], --labelIndex=[/work/labelindex], --output=[/work/20news-testing], --temp=[/tmp], --trainCompress=[true], --useNativeCode=[true]}
16/10/21 14:19:33 INFO NativeCodeLoader: Unable to load native mahout-hadoop library for your platform... using builtin-java classes where applicable
16/10/21 14:19:41 INFO CodecPool: Got brand-new compressor [lzfdeflate]
16/10/21 14:19:41 INFO CodecPool: Got brand-new decompressor [lzfdeflate]
16/10/21 14:19:44 INFO Deprecation: mapred.input.dir is deprecated. Instead, use mapreduce.input.fileinputformat.inputdir
16/10/21 14:19:44 INFO Deprecation: mapred.compress.map.output is deprecated. Instead, use mapreduce.map.output.compress
```

#### f. Test the classifier

```
16/10/21 14:10:57 INFO MahoutDriver: Program took 18088 ms (Minutes: 0.3014666666666666)
[wangxucan [master] mahout-trunk]$ ./bin/mahout testnb -l /work/20news-test/vectors -m /work/model -l /work/labelindex -ow -o /work/20news-testing -c
Running on hadoop, using /usr/local/Cellar/hadoop/2.7.3/bin/hadoop and HADOOP_CONF_DIR=
MAHOUT-JOB: /usr/local/Cellar/mahout-trunk/examples/target/mahout-examples-0.13.0-SNAPSHOT-job.jar
```

result:

```
=====
Statistics
-----
Kappa                      0.8568
Accuracy                   89.1194%
Reliability                84.5424%
Reliability (standard deviation) 0.2181
Weighted precision          0.8928
Weighted recall              0.8912
Weighted F1 score            0.8891

16/10/21 14:18:03 INFO MahoutDriver: Program took 9805 ms (Minutes: 0.1634166666666665)
wangxucan (master) mahout-trunk $
```

## Then I used the existed model in classify-20newsgroups.sh and try different algorithms

### C Naivebayes

```
wangxucan (master) mahout-trunk $ ./examples/bin/classify-20newsgroups.sh
Discovered Hadoop v2.
Setting dfs command to /usr/local/Cellar/hadoop/2.7.3/bin/hdfs dfs, dfs rm to /usr/local/Cellar/hadoop/2.7.3/bin/hdfs dfs -rm -r -skipTrash.
Please select a number to choose the corresponding task to run
1. cnaiivebayes-MapReduce
2. naivebayes-MapReduce
3. cnaiivebayes-Spark
4. naivebayes-Spark
5. sgd
6. clean-- cleans up the work area in /tmp/mahout-work-wangxucan
Enter your choice : 1
```

The result is:

| Summary                          |                 |
|----------------------------------|-----------------|
| Correctly Classified Instances   | : 6740 89.0592% |
| Incorrectly Classified Instances | : 828 10.9408%  |
| Total Classified Instances       | : 7568          |

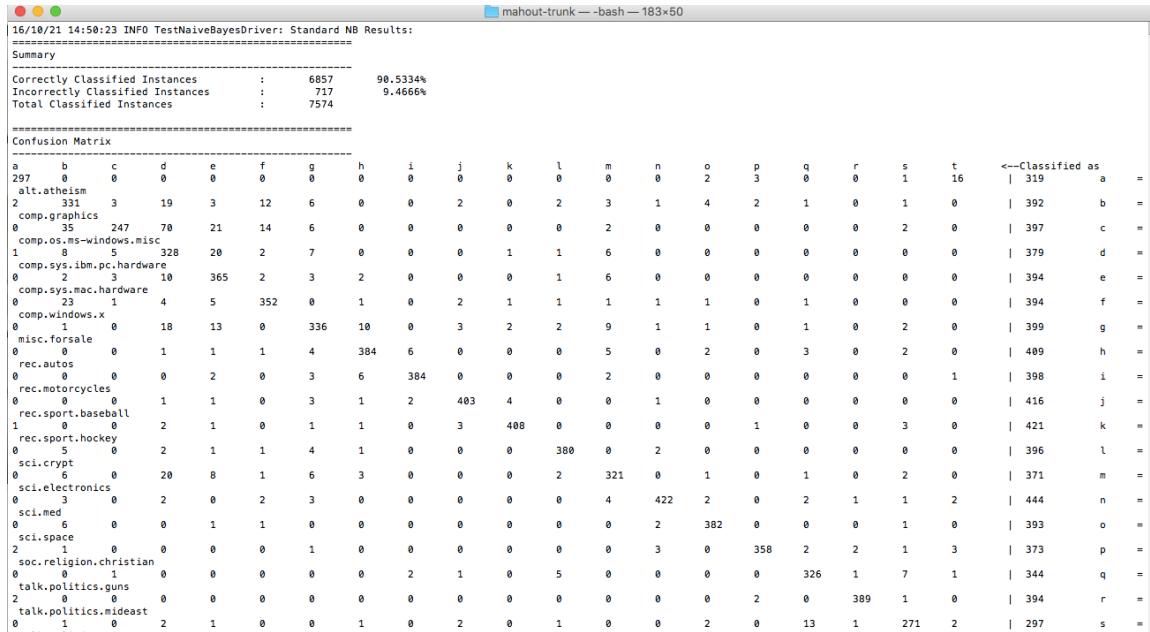
| Confusion Matrix |     |     |     |     |     |     |     |     |     |     |     |     |     |     |    |     |   |     |       |                  |       |       |
|------------------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|----|-----|---|-----|-------|------------------|-------|-------|
| a                | b   | c   | d   | e   | f   | g   | h   | i   | j   | k   | l   | m   | n   | o   | p  | q   | r | s   | t     | <--Classified as |       |       |
| 288              | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 2   | 0   | 1   | 0   | 4  | 2   | 1 | 3   | 3     | 2                | 3     | 359 a |
| 3                | 279 | 3   | 11  | 3   | 16  | 9   | 2   | 1   | 2   | 1   | 8   | 4   | 2   | 1   | 3  | 3   | 3 | 3   | 2     | 3                | 359 b |       |
| 3                | 18  | 278 | 58  | 12  | 9   | 11  | 6   | 3   | 2   | 3   | 7   | 8   | 3   | 7   | 2  | 1   | 0 | 3   | 1     | 435 c            |       |       |
| 1                | 9   | 4   | 341 | 12  | 5   | 9   | 3   | 2   | 0   | 2   | 5   | 7   | 4   | 3   | 2  | 2   | 0 | 3   | 2     | 416 d            |       |       |
| 0                | 2   | 7   | 8   | 349 | 4   | 3   | 3   | 0   | 2   | 3   | 0   | 0   | 1   | 1   | 2  | 0   | 1 | 1   | 0     | 387 e            |       |       |
| 0                | 12  | 0   | 1   | 1   | 334 | 2   | 1   | 3   | 1   | 0   | 4   | 0   | 2   | 3   | 1  | 0   | 0 | 2   | 1     | 368 f            |       |       |
| 1                | 10  | 1   | 27  | 17  | 4   | 280 | 12  | 8   | 3   | 5   | 3   | 11  | 3   | 2   | 3  | 2   | 2 | 3   | 1     | 398 g            |       |       |
| 1                | 0   | 0   | 1   | 2   | 0   | 3   | 366 | 4   | 1   | 2   | 1   | 2   | 2   | 1   | 0  | 4   | 0 | 2   | 0     | 392 h            |       |       |
| 1                | 0   | 0   | 0   | 0   | 0   | 1   | 3   | 408 | 0   | 0   | 0   | 1   | 1   | 0   | 1  | 1   | 0 | 0   | 0     | 417 i            |       |       |
| 0                | 0   | 0   | 1   | 0   | 0   | 1   | 0   | 1   | 371 | 9   | 0   | 2   | 4   | 0   | 0  | 0   | 0 | 0   | 0     | 389 j            |       |       |
| 0                | 0   | 1   | 0   | 1   | 0   | 2   | 0   | 1   | 3   | 412 | 0   | 0   | 0   | 1   | 1  | 0   | 3 | 0   | 0     | 425 k            |       |       |
| 0                | 0   | 0   | 1   | 0   | 0   | 0   | 0   | 0   | 1   | 0   | 408 | 1   | 1   | 0   | 1  | 2   | 0 | 1   | 0     | 416 l            |       |       |
| 0                | 7   | 0   | 9   | 12  | 1   | 6   | 4   | 2   | 2   | 2   | 2   | 344 | 2   | 5   | 4  | 2   | 1 | 2   | 1     | 408 m            |       |       |
| 3                | 1   | 1   | 1   | 2   | 0   | 1   | 2   | 0   | 2   | 1   | 0   | 5   | 359 | 1   | 1  | 0   | 0 | 1   | 2     | 383 n            |       |       |
| 1                | 0   | 0   | 0   | 1   | 0   | 1   | 1   | 0   | 0   | 1   | 0   | 2   | 2   | 378 | 1  | 2   | 0 | 0   | 2     | 392 o            |       |       |
| 4                | 0   | 0   | 1   | 0   | 1   | 0   | 0   | 1   | 0   | 0   | 0   | 3   | 1   | 369 | 1  | 3   | 1 | 4   | 1     | 389 p            |       |       |
| 0                | 0   | 1   | 0   | 0   | 1   | 1   | 0   | 0   | 0   | 0   | 4   | 0   | 0   | 2   | 0  | 360 | 0 | 3   | 0     | 372 q            |       |       |
| 0                | 1   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 2   | 0   | 0   | 0   | 1   | 2   | 0  | 361 | 0 | 1   | 368 r |                  |       |       |
| 2                | 0   | 0   | 0   | 0   | 0   | 2   | 0   | 0   | 2   | 2   | 0   | 0   | 0   | 0   | 0  | 17  | 3 | 276 | 0     | 304 s            |       |       |
| 21               | 0   | 0   | 0   | 0   | 0   | 1   | 0   | 0   | 1   | 0   | 2   | 0   | 1   | 1   | 18 | 8   | 3 | 8   | 179   | 243 t            |       |       |

```
=====
Statistics
-----
Kappa                                0.856
Accuracy                             89.0592%
Reliability                          84.7123%
Reliability (standard deviation)    0.2181
Weighted precision                   0.8913
Weighted recall                      0.8906
Weighted F1 score                    0.8883

16/10/21 14:34:33 INFO MahoutDriver: Program took 10292 ms (Minutes: 0.1715333333333334)
```

## Naivebayes

```
[wangxucan (master) mahout-trunk $ ./examples/bin/classify-20newsgroups.sh
Discovered Hadoop v2.
Setting dfs command to /usr/local/Cellar/hadoop/2.7.3/bin/hdfs dfs, dfs rm to /usr/local/Cellar/hadoop/2.7.3/bin/dfs dfs -rm -r -skipTrash.
Please select a number to choose the corresponding task to run
1. cnaivebayes-MapReduce
2. naivebayes-MapReduce
3. cnaivebayes-Spark
4. naivebayes-Spark
5. sgd
6. clean-- cleans up the work area in /tmp/mahout-work-wangxucan
Enter your choice : 2]
```



```
16/10/21 14:50:23 INFO TestNaiveBayesDriver: Standard NB Results:
=====
Summary
-----
Correctly Classified Instances      :       6657      98.5334%
Incorrectly Classified Instances   :        717      9.4666%
Total Classified Instances         :      7574

=====

Confusion Matrix
-----
a   b   c   d   e   f   g   h   i   j   k   l   m   n   o   p   q   r   s   t   <--Classified as
297  0   0   0   0   0   0   0   0   0   0   0   0   0   0   2   0   0   1   0   16
| 319  a   =  

2   331  3   19   3   12   6   0   0   0   2   0   2   3   1   4   2   1   0   1   0   | 392  b   =  

comp.graphics  0   35   247  70   21   14   6   0   0   0   0   0   2   0   0   0   0   0   2   0   | 397  c   =  

comp.os.ms-windows.misc  1   8   5   328  20   2   7   0   0   0   1   1   6   0   0   0   0   0   0   0   | 379  d   =  

comp.sys.ibm.pc.hardware  0   2   3   10   365  2   3   2   0   0   0   1   6   0   0   0   0   0   0   0   | 394  e   =  

comp.sys.mac.hardware  0   23  1   4   5   352  0   1   0   2   1   1   1   1   1   0   1   0   0   0   | 394  f   =  

comp.windows.x  0   1   0   18   13  0   336  10   0   3   2   2   9   1   1   0   1   0   2   0   | 399  g   =  

misc.forsale  0   0   0   1   1   1   4   384  6   0   0   0   5   0   2   0   3   0   2   0   | 489  h   =  

rec.autos  0   0   0   0   0   2   0   3   6   384  0   0   0   2   0   0   0   0   0   1   | 398  i   =  

rec.motorcycles  0   0   0   1   1   0   3   1   2   403  4   0   0   1   0   0   0   0   0   0   | 416  j   =  

rec.sport.baseball  1   0   0   2   1   0   1   1   0   3   408  0   0   0   0   1   0   0   3   0   | 421  k   =  

rec.sport.hockey  0   5   0   2   1   1   4   1   0   0   0   380  0   2   0   0   0   0   0   0   | 396  l   =  

sci.crypt  0   6   0   20   8   1   6   3   0   0   0   0   2   321  0   1   0   1   0   2   0   | 371  m   =  

sci.electronics  0   3   0   2   0   2   3   0   0   0   0   0   4   422  2   0   2   1   1   2   | 444  n   =  

sci.med  0   6   0   0   1   1   0   0   0   0   0   0   0   2   382  0   0   0   1   0   0   | 393  o   =  

sci.space  2   1   0   0   0   0   1   0   0   0   0   0   0   3   0   0   358  2   2   1   3   | 373  p   =  

soc.religion.christian  0   0   1   0   0   0   0   0   2   1   0   5   0   0   0   0   326  1   7   1   | 344  q   =  

talk.politics.guns  2   0   0   0   0   0   0   0   0   0   0   0   0   0   0   0   0   0   0   0   | 394  r   =  

talk.politics.mideast  0   1   0   2   1   0   0   1   0   2   0   1   0   0   0   2   0   389  1   0   | 297  s   =  

talk.politics.misc  0   0   0   0   0   0   0   0   0   0   0   0   0   0   0   0   0   0   0   0   |
```

```
=====
Statistics
-----
Kappa                                0.8775
Accuracy                             90.5334%
Reliability                          85.8573%
Reliability (standard deviation)    0.2166
Weighted precision                   0.91
Weighted recall                      0.9053
Weighted F1 score                    0.9047
```

```
16/10/21 14:50:23 INFO MahoutDriver: Program took 10648 ms (Minutes: 0.1774666666666666)
wangxucan (master) mahout-trunk $
```

Comparison: The accuracy of naïvebayes is 90.5334%, reliability is 85.8573%.

The accuracy of Cnaivebayes is 89.0592%, the reliability is 84.7123%

And the accuracy of my own trained model is 89.1194% and the reliability is 84.5424% which is a little lower than the existed one.

The running time of are both 5 minutes. So the naivebayes is better.

## 2) Wikipedia

I use the existed classify-wikipedia.sh in the examples

### CBayes:

```
wangxucan (master) mahout-trunk $ ./examples/bin/classify-wikipedia.sh
Discovered Hadoop v2.
Setting dfs command to /usr/local/Cellar/hadoop/2.7.3/bin/dfs, dfs rm to /usr/local/Cellar/hadoop/2.7.3/bin/dfs dfs -rm -r -skipTrash.
Please select a number to choose the corresponding task to run
1. CBayes (may require increased heap space on yarn)
2. BinaryCBayes
3. clean -- cleans up the work area in /tmp/mahout-work-wiki
Enter your choice : 1
```

result:

```
Testing on holdout set: CBayes
+ /usr/local/Cellar/mahout/bin/mahout testnb -i /tmp/mahout-work-wiki/testing -m /tmp/mahout-work-wiki/model -l /tmp/mahout-work-wiki/labelindex -ow -o /tmp/mahout-work-wiki/output -c
Running on hadoop, using /usr/local/Cellar/hadoop/2.7.3/bin/hadoop and HADOOP_CONF_DIR=
MAHOUT-JOB: /usr/local/Cellar/mahout/examples/target/mahout-examples-0.11.0-job.jar
16/10/21 17:59:19 INFO AbstractJob: Command line arguments: {--endPhase=[2147483647], --input=[/tmp/mahout-work-wiki/testing], --labelIndex=[/tmp/mahout-work-wiki/labelindex], --model=[/tmp/mahout-work-wiki/model], --output=[/tmp/mahout-work-wiki/output], --overwrite=null, --runSequential=null, --startPhase=[0], --tempDir=[temp], --testComplementary=null}
16/10/21 17:59:19 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
16/10/21 17:59:25 INFO HadoopUtil: Deleting /tmp/mahout-work-wiki/output
16/10/21 17:59:25 INFO CodecPool: Got brand-new compressor [.deflate]
16/10/21 17:59:25 INFO CodecPool: Got brand-new decompressor [.deflate]
16/10/21 17:59:34 INFO TestNaiveBayesDriver: Complementary Results:
=====
Summary
=====
Correctly Classified Instances : 2013 85.6231%
Incorrectly Classified Instances : 338 14.3769%
Total Classified Instances : 2351
=====
Confusion Matrix
-----
a   b   c   d   e   f   g   h   i   j   <--Classified as
759 8   44  6   6   4   8   9   9   1   | 859   a   = australia
0   107 2   2   1   1   6   0   2   0   | 121   b   = austria
0   0   9   1   2   1   0   0   0   0   | 13    c   = bahamas
4   7   20  441 8   7   2   11  7   2   | 509   d   = canada
0   0   0   0   34  0   0   0   0   0   | 34    e   = colombia
0   0   1   2   1   44  0   0   0   2   | 56    f   = cuba
0   1   1   0   5   0   88  0   0   0   | 67    g   = pakistan
0   1   0   3   3   2   0   6   0   0   | 15    h   = panama
6   17  51  4   6   5   11  14  482  6   | 602   i   = united kingdom
2   2   2   0   2   1   0   1   0   51   | 61    j   = vietnam
=====
Statistics
-----
Kappa                           0.7892
Accuracy                        85.6231%
Reliability                     74.2079%
Reliability (standard deviation) 0.2924
Weighted precision               0.9232
Weighted recall                 0.8562
Weighted F1 score                0.8823
16/10/21 17:59:34 INFO MahoutDriver: Program took 15185 ms (Minutes: 0.2530833333333333)
```

The accuracy is 85.6231%, the reliability is 74.2079%

### Binary-CBayes:

```
wangxucan (master) mahout-trunk $ ./examples/bin/classify-wikipedia.sh
Discovered Hadoop v2.
Setting dfs command to /usr/local/Cellar/hadoop/2.7.3/bin/dfs, dfs rm to /usr/local/Cellar/hadoop/2.7.3/bin/dfs dfs -rm -r -skipTrash.
Please select a number to choose the corresponding task to run
1. CBayes (may require increased heap space on yarn)
2. BinaryCBayes
3. clean -- cleans up the work area in /tmp/mahout-work-wiki
Enter your choice : 2
```

Result:

```

=====
Summary
-----
Correctly Classified Instances      :    2397      87.4498%
Incorrectly Classified Instances   :     344      12.5502%
Total Classified Instances        :    2741

=====
Confusion Matrix
-----
a      b      <--Classified as
609    11      | 620      a      = united kingdom
333    1788    | 2121     b      = united states

=====
Statistics
-----
Kappa                           0.6953
Accuracy                        87.4498%
Reliability                     60.8419%
Reliability (standard deviation) 0.5315
Weighted precision               0.9153
Weighted recall                 0.8745
Weighted F1 score                0.8823

16/10/21 18:13:59 INFO MahoutDriver: Program took 4846 ms (Minutes: 0.08076666666666667)
wangxucan (master) mahout-trunk $ █
1 9 4 343 12 5 9 3 2 8 2 5 7 4 3 2 3

```

The accuracy is 87.4498%, reliability is 60.8419%

Compared to the cBayes model, the binaryCBayes model has higher accuracy but is less reliable.

## Part 4: Install Spark and run word count algorithm

### a. download the pre-built spark

---

#### Download Apache Spark™

Our latest stable version is Apache Spark 2.0.1, released on Oct 3, 2016 ([release notes](#)) ([git tag](#))

1. Choose a Spark release:
2. Choose a package type:
3. Choose a download type:
4. Download Spark: [spark-2.0.1-bin-hadoop2.7.tgz](#)
5. Verify this release using the [2.0.1 signatures and checksums](#) and [project release KEYS](#).

```
# Apache Spark

Spark is a fast and general cluster computing system for Big Data. It provides
high-level APIs in Scala, Java, Python, and R, and an optimized engine that
supports general computation graphs for data analysis. It also supports a
rich set of higher-level tools including Spark SQL for SQL and DataFrames,
MLlib for machine learning, GraphX for graph processing,
and Spark Streaming for stream processing.

<http://spark.apache.org/>

## Online Documentation

You can find the latest Spark documentation, including a programming
guide, on the [project web page](http://spark.apache.org/documentation.html)
and [project wiki](https://cwiki.apache.org/confluence/display/SPARK).
This README file only contains basic setup instructions.

## Building Spark

Spark is built using [Apache Maven](http://maven.apache.org/).
To build Spark and its example programs, run:

    build/mvn -DskipTests clean package

(You do not need to do this if you downloaded a pre-built package.)

You can build Spark using more than one thread by using the -T option with Maven, see ["Parallel builds in M
in+Maven+3].
More detailed documentation is available from the project site, at
["Building Spark"](http://spark.apache.org/docs/latest/building-spark.html).
For developing Spark using an IDE, see [Eclipse](https://cwiki.apache.org/confluence/display/SPARK/Useful+De
and [IntelliJ](https://cwiki.apache.org/confluence/display/SPARK/Useful+Developer+Tools).

## Interactive Scala Shell

The easiest way to start using Spark is through the Scala shell:

    ./bin/spark-shell

Try the following command, which should return 1000:

    scala> sc.parallelize(1 to 1000).count()

## Interactive Python Shell

Alternatively, if you prefer Python, you can use the Python shell:

"README.md" 99L, 3828C
```

b. set the environment variable for spark

```
wangxucan spark-2.0.1-bin-hadoop2.7 $ export SPARK_HOME=/usr/local/Cellar/hadoop/2.7.3/spark-2.0.1-bin-hadoop2.7  
wangxucan spark-2.0.1-bin-hadoop2.7 $
```

c. use the ./bin/pyspark to start it

d1: read the source file License.txt into a resilient distributed dataset

d2: count the number of items in the RDD, in my case, 1562 lines.

d3: First item in RDD which is empty in my case

d4: Count how many lines contains “License”

d5: Count how many lines contains “Apache”

The file looks like this

