

EECS E6893 Homework 2

Zhuxi Cai zc2270

1. Recommendation:

(1) Choose any two datasets from Yahoo Labs Ratings and Classification Data

We choose the following two datasets from Yahoo Labs Ratings and Classification Data:

- ydata-ymusic-rating-study-v1_0-train.txt
- ydata-delicious-popular-urls-and-tags-v1_0.txt

(2) Try various recommendation algorithms provided by Mahout

First we try Mahout working with Eclipse and Maven using dataset ydata-ymusic-rating-study-v1_0-train.txt. We apply three recommendation algorithms: user-based recommendation, item-based recommendation and SVD recommendation. Then we evaluate the three models we built separately.

User-based Recommender Java Code:

```
public class User_based_recommender
{
    public static void main( String[] args ) throws IOException, TasteException
    {
        DataModel model = new FileDataModel(new File("data/ydata-ymusic-rating-study-v1_0-train.txt"));
        UserSimilarity similarity = new PearsonCorrelationSimilarity(model);
        UserNeighborhood neighborhood = new ThresholdUserNeighborhood(0.1, similarity, model);
        UserBasedRecommender recommender = new GenericUserBasedRecommender(model, neighborhood, similarity);
        int userID=2;
        int itemNum=3;
        List<RecommendedItem> recommendations = recommender.recommend(userID, itemNum);
        System.out.println("User-based Recommender");
        System.out.println("For user "+userID+", recommend "+itemNum+" items:");
        for (RecommendedItem recommendation : recommendations) {
            System.out.println(recommendation);
        }
    }
}
```

User-based Recommender Result:

```
User-based Recommender
For user 2, recommend 3 items:
RecommendedItem[item:514, value:4.619235]
RecommendedItem[item:266, value:4.6109624]
RecommendedItem[item:69, value:4.4430842]
```

User-based Recommender Evaluation:

```

public class User_based_recommender_eva {
    public static void main(String[] args) throws IOException, TasteException
    {
        DataModel model = new FileDataModel(new File("data/ydata-ymusic-rating-study-v1_0-train.txt"));
        RecommenderEvaluator evaluator = new AverageAbsoluteDifferenceRecommenderEvaluator();
        RecommenderBuilder builder = new UserRecommenderBuilder();
        double result = evaluator.evaluate(builder, null, model, 0.9, 1.0);
        System.out.println("User-based Recommender Evaluation");
        System.out.println(result);
    }
}

class UserRecommenderBuilder implements RecommenderBuilder{
    public Recommender buildRecommender(DataModel dataModel) throws TasteException {
        UserSimilarity similarity = new PearsonCorrelationSimilarity(dataModel);
        UserNeighborhood neighborhood = new ThresholdUserNeighborhood(0.1, similarity, dataModel);
        return new GenericUserBasedRecommender(dataModel, neighborhood, similarity);
    }
}

```

User-based Recommender Evaluation
1.267003051378381

Item-based Recommender Java Code:

```

public class Item_based_recommender
{
    public static void main( String[] args ) throws IOException, TasteException
    {
        DataModel model = new FileDataModel(new File("data/ydata-ymusic-rating-study-v1_0-train.txt"));
        ItemSimilarity similarity = new PearsonCorrelationSimilarity(model);
        ItemBasedRecommender recommender = new GenericItemBasedRecommender(model, similarity);
        int userID=2;
        int itemNum=3;
        List<RecommendedItem> recommendations = recommender.recommend(userID, itemNum);
        System.out.println("Item-based Recommender");
        System.out.println("For user "+userID+", recommend "+itemNum+" items:");
        for (RecommendedItem recommendation : recommendations) {
            System.out.println(recommendation);
        }
    }
}

```

Item-based Recommender Result:

```

Item-based Recommender
For user 2, recommend 3 items:
RecommendedItem[item:26, value:5.0]
RecommendedItem[item:13, value:5.0]
RecommendedItem[item:5, value:5.0]

```

Item-based Recommender Evaluation:

```

public class Item_based_recommender_eva {
    public static void main(String[] args) throws IOException, TasteException {
        DataModel model = new FileDataModel(new File("data/ydata-ymusic-rating-study-v1_0-train.txt"));
        RecommenderEvaluator evaluator = new AverageAbsoluteDifferenceRecommenderEvaluator();
        RecommenderBuilder builder = new ItemRecommenderBuilder();
        double result = evaluator.evaluate(builder, null, model, 0.9, 1.0);
        System.out.println("Item-based Recommender Evaluation");
        System.out.println(result);
    }
}

class ItemRecommenderBuilder implements RecommenderBuilder{
    public Recommender buildRecommender(DataModel dataModel) throws TasteException {
        ItemSimilarity similarity = new PearsonCorrelationSimilarity(dataModel);
        return new GenericItemBasedRecommender(dataModel, similarity);
    }
}

```

Item-based Recommender Evaluation

0.9814732361074119

SVD Recommender Java Code:

```

public class SVD_recommender {
    public static void main( String[] args ) throws IOException, TasteException {
        DataModel model = new FileDataModel(new File("data/ydata-ymusic-rating-study-v1_0-train.txt"));
        SVDRecommender recommender = new SVDRecommender(model, new ALSWRFactorizer(model, 10, 0.05, 10));
        int userID=2;
        int itemNum=3;
        List<RecommendedItem> recommendations = recommender.recommend(userID, itemNum);
        System.out.println("SVD Recommender");
        System.out.println("For user "+userID+", recommend "+itemNum+" items:");
        for (RecommendedItem recommendation : recommendations) {
            System.out.println(recommendation);
        }
    }
}

```

SVD Recommender Result:

SVD Recommender

For user 2, recommend 3 items:

RecommendedItem[item:773, value:6.2266994]

RecommendedItem[item:763, value:5.7881827]

RecommendedItem[item:629, value:5.784256]

SVD Recommender Evaluation:

```

public class SVD_recommender_eva {

    public static void main(String[] args) throws IOException, TasteException
    {
        DataModel model = new FileDataModel(new File("data/ydata-ymusic-rating-study-v1_0-train.txt"));
        RecommenderEvaluator evaluator = new AverageAbsoluteDifferenceRecommenderEvaluator();
        RecommenderBuilder builder = new SVDRecommenderBuilder();
        double result = evaluator.evaluate(builder, null, model, 0.9, 1.0);
        System.out.println("SVD Recommender Evaluation");
        System.out.println(result);
    }
}

class SVDRecommenderBuilder implements RecommenderBuilder{

    public Recommender buildRecommender(DataModel dataModel) throws TasteException {
        return new SVDRecommender(dataModel, new ALSWRFactorizer(dataModel, 10, 0.05, 10));
    }
}

SVD Recommender Evaluation
0.9719408046658731

```

To compare these three models, we separate the whole dataset into training and testing dataset, each of which is 90% and 10%. Then we calculate the average absolute difference for three recommenders and find that SVD recommender performs best here.

Second we install Mahout in bash and try Mahout locally using dataset ydata-delicious-popular-urls-and-tags-v1_0.txt. We choose item-base recommendation algorithm and get the following results:

Item-based Recommender Code:

```
Zhuxis-MacBook-Pro:mahout-trunk Jovial$ bin/mahout recommenditembased -s SIMILARITY_LOGLIKELIHOOD -i /Users/Jovial/Desktop/Web
scope_R5/ydata-delicious-popular-urls-and-tags-v1_0.txt -o /Users/Jovial/Desktop/HW2/output --numRecommendations 10
```

Item-based Result:

```

15/10/22 21:41:05 INFO Job: Job job_local795347300_0007 completed successfully
15/10/22 21:41:05 INFO Job: Counters: 30
  Tamp File System Counters
variant...ion.pptx FILE: Number of bytes read=1274843590
FILE: Number of bytes written=1318318326
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
  Map-Reduce Framework
    Map input records=11575
    Map output records=104124
    Map output bytes=94154628
    Map output materialized bytes=45301574
    Input split bytes=130
    Combine input records=0
    Combine output records=0
    Reduce input groups=7642
    Reduce shuffle bytes=45301574
    Reduce input records=104124
    Reduce output records=7641
    Spilled Records=312372
    Shuffled Maps =1
    Failed Shuffles=0
    Merged Map outputs=1
    GC time elapsed (ms)=9
    Total committed heap usage (bytes)=2509766656
  Shuffle Errors
    BAD_ID=0
    CONNECTION=0
    IO_ERROR=0
    WRONG_LENGTH=0
    WRONG_MAP=0
    WRONG_REDUCE=0
  File Input Format Counters
    Bytes Read=9843516
  File Output Format Counters
    Bytes Written=1421372
15/10/22 21:41:05 INFO MahoutDriver: Program took 73299 ms (Minutes: 1.22165)

```

```

1 [1800159233:13.0,1800353817:13.0,1800202853:13.0,1800024141:13.0,1800122300:13.0,1807858489:13.0,1804090611:13.0,1800019431:13.0,1800353749:13.0,1800026398:13.0]
2 [1800112450:11.001758,1800071079:10.370626,1804738128:10.332838,1800212379:10.018739,1808405650:10.009023,180024003:9.877988,1803453994:9.799826,1800128578:9.77842,1802771203:9.54708,1800258211:9.54532]
3 [1802953401:13.0,1800185256:13.0,1800090318:13.0,1800444945:13.0,1800021580:13.0,1804090611:13.0,1800082735:13.0,1804857429:13.0,1804383575:13.0,1800019316:13.0]
4 [1800019304:10.5210285,1800379216:10.506128,1800022996:10.503612,1800174173:10.501894,1800019665:10.499486,1800185256:10.498645,1800060404:10.497791,1807733433:10.496107,1800026332:10.490099,1800128578:10.0433445]
5 [1807993019:13.0,1800168311:13.0,18000058561:13.0,1804090611:13.0,1808404659:13.0,1808416901:12.758992,1808465610:12.746555,1800223031:12.67576,1804861086:12.6658535,1808403030:12.661338]
6 [1808412549:13.0,1808481189:13.0,1800112611:13.0,1800361193:12.65157,1808405708:12.530771,1808437167:12.523912,1808443070:12.514433,1808405886:12.511748,1807993019:12.505257,1800018659:12.503544]
7 [1800444945:13.0,1802953401:13.0,1800022145:13.0,1800026205:13.0,1808486011:13.0,1800021539:13.0,180002016:13.0,1800155161:13.0,1800183197:13.0,1800020102:13.0]
8 [1804476879:13.0,1800019191:13.0,1808458391:13.0,1808431618:13.0,1808403836:13.0,1800112450:13.0,1808435651:13.0,1808457961:13.0,1800062915:13.0,1808406133:13.0]
9 [1808467432:13.0,1808404170:12.510724,1808403030:12.503729,1808429384:12.503609,1800018810:12.50152,1807634180:12.501279,1808404010:12.500442,1808421673:12.499154,1808420230:12.496009,1808410707:12.490825]
10 [1808505124:13.0,1807592183:13.0,1807537463:12.67474,1808488292:12.667972,1808471149:12.665307,1808405256:12.505206,1800265585:12.490035,1808442969:12.395251,1808415480:12.321217,1807859434:12.250883]

```

2. Clustering:

(1) Using datasets from Online news and Wikipedia articles

We choose the following two datasets to do clustering:

- reuters21578
- enwiki-20150901-pages-articles1.xml-p000000010p000010000

(2) Do clustering to find related documents

First we apply K-means and Fuzzy K-means on Reuters news dataset:

K-means clustering:

```
15/10/25 12:29:39 INFO Job: Job job_local1473567275_0004 completed successfully
15/10/25 12:29:39 INFO Job: Counters: 15
  File System Counters
    FILE: Number of bytes read=384817170
    FILE: Number of bytes written=307065530
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
  Map-Reduce Framework
    Map input records=21578
    Map output records=21578
    Input split bytes=164
    Spilled Records=0
    Failed Shuffles=0
    Merged Map outputs=0
    GC time elapsed (ms)=12
    Total committed heap usage (bytes)=492830720
  File Input Format Counters
    Bytes Read=17038887
  File Output Format Counters
    Bytes Written=18065095
15/10/25 12:29:39 INFO MahoutDriver: Program took 18344 ms (Minutes: 0.30573333333333336)
```

```
Key: 5646: Value: wt: 1.0 distance: 0.8102479274555142 vec: [{"1697":5.609}, {"2689":3.135}, {"3730":3.051}, {"5997":1.1}, {"6236":5.6}
Key: 1531: Value: wt: 1.0 distance: 0.7224820861423382 vec: [{"2689":3.135}, {"3730":3.051}, {"5997":1.119}, {"6236":5.6}
Key: 371: Value: wt: 1.0 distance: 0.9216730299859334 vec: [{"2969":2.807}, {"3730":3.051}, {"5636":9.593}, {"5703":3.9}
Key: 1531: Value: wt: 1.0 distance: 0.7031398855630863 vec: [{"1697":5.609}, {"2689":3.135}, {"2962":2.807}, {"3730":3.0
Key: 18541: Value: wt: 1.0 distance: 0.5817853698775622 vec: [{"797":8.9}, {"2689":3.135}, {"2962":2.807}, {"3730":3.051}
Key: 1531: Value: wt: 1.0 distance: 0.7343683798371794 vec: [{"2689":3.135}, {"3170":2.829}, {"3310":9.881}, {"3730":3.0
Key: 5646: Value: wt: 1.0 distance: 0.8016823132333883 vec: [{"2540":9.188}, {"2689":3.135}, {"3730":3.051}, {"4397":2.9
Key: 18734: Value: wt: 1.0 distance: 0.8459251196403172 vec: [{"2689":3.135}, {"3730":3.051}, {"5997":1.119}, {"7424":4.5
Key: 15960: Value: wt: 1.0 distance: 0.8350294307372016 vec: [{"933":9.188}, {"2689":3.135}, {"3730":3.051}, {"5997":1.11
Key: 659: Value: wt: 1.0 distance: 0.6897259934760255 vec: [{"2689":3.135}, {"3000":7.514}, {"3730":3.051}, {"5433":7.7
Key: 1531: Value: wt: 1.0 distance: 0.7279302559069704 vec: [{"2689":3.135}, {"3730":3.051}, {"4397":5.024}, {"5997":1.1
Key: 10601: Value: wt: 1.0 distance: 0.7403892433923107 vec: [{"689":9.593}, {"2689":3.135}, {"3170":2.829}, {"3730":3.05
Key: 6533: Value: wt: 1.0 distance: 0.8124119927554354 vec: [{"2689":3.135}, {"2983":9.188}, {"3730":3.051}, {"5848":3.8
Key: 10601: Value: wt: 1.0 distance: 0.6675468748327895 vec: [{"2962":2.807}, {"3730":3.051}, {"3932":7.514}, {"4918":3.6
Key: 659: Value: wt: 1.0 distance: 0.5518133521258943 vec: [{"649":9.37}, {"1509":7.545}, {"1711":7.453}, {"2090":7.088
Key: 10601: Value: wt: 1.0 distance: 0.8691348832521122 vec: [{"2689":3.135}, {"3730":3.051}, {"7027":4.332}, {"30839":3.
Key: 1531: Value: wt: 1.0 distance: 0.7326471852442071 vec: [{"2689":3.135}, {"2962":2.807}, {"3730":3.051}, {"4471":9.1
Key: 659: Value: wt: 1.0 distance: 0.720175488989391 vec: [{"1696":2.643}, {"2689":3.135}, {"2962":2.807}, {"3358":4.31
Key: 659: Value: wt: 1.0 distance: 0.7337654372961773 vec: [{"1982":3.839}, {"2689":3.135}, {"3170":2.829}, {"3730":3.0
Key: 659: Value: wt: 1.0 distance: 0.7574858324075376 vec: [{"1043":4.128}, {"2689":3.135}, {"2962":2.807}, {"3730":3.0
Key: 659: Value: wt: 1.0 distance: 0.7363096348627212 vec: [{"1982":2.714}, {"2273":9.37}, {"2689":3.135}, {"2962":2.80
Key: 659: Value: wt: 1.0 distance: 0.7287894296950976 vec: [{"955":9.37}, {"2689":3.135}, {"3170":2.829}, {"3730":3.051
Key: 1531: Value: wt: 1.0 distance: 0.6777502276330576 vec: [{"616":9.034}, {"2689":3.135}, {"3730":3.051}, {"5997":1.11
Key: 5646: Value: wt: 1.0 distance: 0.7965285548503003 vec: [{"2689":3.135}, {"3730":3.051}, {"4500":4.899}, {"5997":1.1
Key: 3362: Value: wt: 1.0 distance: 0.86023320857273 vec: [{"1696":2.643}, {"2962":2.807}, {"3730":3.051}, {"5081":3.633
Key: 8462: Value: wt: 1.0 distance: 0.7869521142778453 vec: [{"2689":3.135}, {"3730":3.051}, {"5359":9.188}, {"6316":3.1
Key: 659: Value: wt: 1.0 distance: 0.7751149363273682 vec: [{"2689":3.135}, {"3730":3.051}, {"5081":3.633}, {"5241":3.4
Key: 1531: Value: wt: 1.0 distance: 0.8307084824069053 vec: [{"2689":3.135}, {"3730":3.051}, {"5081":3.633}, {"5997":1.1
Key: 659: Value: wt: 1.0 distance: 0.6494662008824592 vec: [{"2460":2.742}, {"2689":3.135}, {"3730":3.051}, {"4918":3.6
Key: 659: Value: wt: 1.0 distance: 0.6428341618646571 vec: [{"1512":7.545}, {"1528":9.602}, {"1890":9.37}, {"2215":9.18
Key: 5798: Value: wt: 1.0 distance: 0.7720094460152325 vec: [{"2689":3.135}, {"3730":3.051}, {"3948":3.954}, {"4747":4.0
Key: 1531: Value: wt: 1.0 distance: 0.850513526595389 vec: [{"2689":3.135}, {"3730":3.051}, {"4481":7.647}, {"4747":4.02
Key: 371: Value: wt: 1.0 distance: 0.8835230467926887 vec: [{"2490":7.129}, {"2689":3.135}, {"3730":3.051}, {"3861":8.4
Key: 864: Value: wt: 1.0 distance: 0.860525819961259 vec: [{"2563":8.677}, {"2689":3.135}, {"3555":3.035}, {"3730":3.05
Key: 659: Value: wt: 1.0 distance: 0.5775528611191603 vec: [{"834":8.677}, {"1555":7.453}, {"2689":3.135}, {"3730":3.05
Count: 21578
15/10/25 12:31:00 INFO MahoutDriver: Program took 2508 ms (Minutes: 0.0418)
```

Fuzzy K-means clustering:

```

15/10/25 12:44:48 INFO Job: Job job_local656430242_0003 completed successfully
15/10/25 12:44:48 INFO Job: Counters: 15
    File System Counters
        FILE: Number of bytes read=421570140
        FILE: Number of bytes written=336995791
        FILE: Number of read operations=0
        FILE: Number of large read operations=0
        FILE: Number of write operations=0
    Map-Reduce Framework
        Map input records=21578
        Map output records=21578
        Input split bytes=164
        Spilled Records=0
        Failed Shuffles=0
        Merged Map outputs=0
        GC time elapsed (ms)=62
        Total committed heap usage (bytes)=374341632
    File Input Format Counters
        Bytes Read=17038887
    File Output Format Counters
        Bytes Written=18064863
15/10/25 12:44:48 INFO MahoutDriver: Program took 37710 ms (Minutes: 0.6285)

```

```

Key: 19713: Value: wt: 1.0 distance: 0.8944488474491454 vec: [{"1697":5.609}, {"2689":3.135}, {"3730":3.051}, {"5997":1.1
Key: 10770: Value: wt: 1.0 distance: 0.8215352424766806 vec: [{"2689":3.135}, {"3730":3.051}, {"5997":1.119}, {"6236":5.6
Key: 3632: Value: wt: 1.0 distance: 0.9433907515767265 vec: [{"2962":2.807}, {"3730":3.051}, {"5636":9.593}, {"5703":3.9
Key: 10770: Value: wt: 1.0 distance: 0.823656911829862 vec: [{"1697":5.609}, {"2689":3.135}, {"2962":2.807}, {"3730":3.05
Key: 10770: Value: wt: 1.0 distance: 0.821178598849267 vec: [{"797":8.9}, {"2689":3.135}, {"2962":2.807}, {"3730":3.051},
Key: 10770: Value: wt: 1.0 distance: 0.839137424563706 vec: [{"2689":3.135}, {"3170":2.829}, {"3310":9.881}, {"3730":3.0
Key: 10770: Value: wt: 1.0 distance: 0.8774963852944634 vec: [{"2540":9.188}, {"2689":3.135}, {"3730":3.051}, {"4397":2.9
Key: 10770: Value: wt: 1.0 distance: 0.9033489384910172 vec: [{"2689":3.135}, {"3730":3.051}, {"5997":1.119}, {"7424":4.5
Key: 18309: Value: wt: 1.0 distance: 0.8904144999396243 vec: [{"933":9.188}, {"2689":3.135}, {"3730":3.051}, {"5997":1.11
Key: 3632: Value: wt: 1.0 distance: 0.8440708899234151 vec: [{"2689":3.135}, {"3000":7.514}, {"3730":3.051}, {"5433":7.7
Key: 10770: Value: wt: 1.0 distance: 0.8252576101160689 vec: [{"2689":3.135}, {"3730":3.051}, {"4397":5.024}, {"5997":1.1
Key: 18309: Value: wt: 1.0 distance: 0.8309546087624643 vec: [{"689":9.593}, {"2689":3.135}, {"3170":2.829}, {"3730":3.05
Key: 18309: Value: wt: 1.0 distance: 0.8434792643576414 vec: [{"2689":3.135}, {"2983":9.188}, {"3730":3.051}, {"5848":3.8
Key: 19713: Value: wt: 1.0 distance: 0.8726591351570927 vec: [{"2962":2.807}, {"3730":3.051}, {"3932":7.514}, {"4918":3.6
Key: 3632: Value: wt: 1.0 distance: 0.8189114506286984 vec: [{"649":9.37}, {"1509":7.545}, {"1711":7.453}, {"2090":7.088
Key: 3632: Value: wt: 1.0 distance: 0.9259366882680182 vec: [{"2689":3.135}, {"3730":3.051}, {"7027":4.332}, {"30839":3.
Key: 10770: Value: wt: 1.0 distance: 0.8270938788222709 vec: [{"2689":3.135}, {"2962":2.807}, {"3730":3.051}, {"4471":9.1
Key: 3632: Value: wt: 1.0 distance: 0.8565607439727175 vec: [{"1696":2.643}, {"2689":3.135}, {"2962":2.807}, {"3358":4.3
Key: 3632: Value: wt: 1.0 distance: 0.8667990170391024 vec: [{"1982":3.839}, {"2689":3.135}, {"3170":2.829}, {"3730":3.0
Key: 3632: Value: wt: 1.0 distance: 0.869694044406639 vec: [{"1043":4.128}, {"2689":3.135}, {"2962":2.807}, {"3730":3.05
Key: 3632: Value: wt: 1.0 distance: 0.8685169304792749 vec: [{"1982":2.7143}, {"2273":9.37}, {"2689":3.135}, {"2962":2.80
Key: 3632: Value: wt: 1.0 distance: 0.8653340016551617 vec: [{"955":9.37}, {"2689":3.135}, {"3170":2.829}, {"3730":3.051
Key: 10770: Value: wt: 1.0 distance: 0.7832506280366477 vec: [{"616":9.034}, {"2689":3.135}, {"3730":3.051}, {"5997":1.11
Key: 3632: Value: wt: 1.0 distance: 0.8973751845574479 vec: [{"2689":3.135}, {"3730":3.051}, {"4500":4.899}, {"5997":1.1
Key: 10770: Value: wt: 1.0 distance: 0.9170032145028807 vec: [{"1696":2.643}, {"2962":2.807}, {"3730":3.051}, {"5081":3.6
Key: 3632: Value: wt: 1.0 distance: 0.9420723425297225 vec: [{"2689":3.135}, {"3730":3.051}, {"5359":9.188}, {"6316":3.1
Key: 3632: Value: wt: 1.0 distance: 0.8613135254723545 vec: [{"2689":3.135}, {"3730":3.051}, {"5081":3.633}, {"5241":3.4
Key: 10770: Value: wt: 1.0 distance: 0.8641642374205118 vec: [{"2689":3.135}, {"3730":3.051}, {"5081":3.633}, {"5997":1.1
Key: 3632: Value: wt: 1.0 distance: 0.828936147581975 vec: [{"2460":2.742}, {"2689":3.135}, {"3730":3.051}, {"4918":3.67
Key: 3632: Value: wt: 1.0 distance: 0.8456475378733378 vec: [{"1512":7.545}, {"1528":9.602}, {"1890":9.37}, {"2215":9.18
Key: 18309: Value: wt: 1.0 distance: 0.8543240999782055 vec: [{"2689":3.135}, {"3730":3.051}, {"3948":3.954}, {"4747":4.0
Key: 18309: Value: wt: 1.0 distance: 0.8763700768290311 vec: [{"2689":3.135}, {"3730":3.051}, {"4481":7.647}, {"4747":4.0
Key: 3025: Value: wt: 1.0 distance: 0.9258367793715158 vec: [{"2490":7.129}, {"2689":3.135}, {"3730":3.051}, {"3861":8.4
Key: 18309: Value: wt: 1.0 distance: 0.8654578762723307 vec: [{"2563":8.677}, {"2689":3.135}, {"3555":3.035}, {"3730":3.0
Key: 3632: Value: wt: 1.0 distance: 0.8398590781420272 vec: [{"834":8.677}, {"1555":7.453}, {"2689":3.135}, {"3730":3.05
Count: 21578
15/10/25 12:47:31 INFO MahoutDriver: Program took 3038 ms (Minutes: 0.050633333333333336)

```

By comparing the two clustering results, we find that the result we get using K-means is different from what we get using Fuzzy K-means. Besides, Fuzzy K-means will take more time to process the data.

Second, we apply Canopy clustering on Wikipedia articles dataset:

```

15/10/25 16:10:24 INFO Job: Job job_local1571132803_0001 completed successfully
15/10/25 16:10:24 INFO Job: Counters: 30
  File System Counters
    FILE: Number of bytes read=349478053
    FILE: Number of bytes written=200849960
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
  Map-Reduce Framework
    Map input records=6273
    Map output records=2
    Map output bytes=4982695
    Map output materialized bytes=4982717
    Input split bytes=316
    Combine input records=0
    Combine output records=0
    Reduce input groups=1
    Reduce shuffle bytes=4982717
    Reduce input records=2
    Reduce output records=1
    Spilled Records=4
    Shuffled Maps =2
    Failed Shuffles=0
    Merged Map outputs=2
    GC time elapsed (ms)=148
    Total committed heap usage (bytes)=1066926080
  Shuffle Errors
    BAD_ID=0
    CONNECTION=0
    IO_ERROR=0
    WRONG_LENGTH=0
    WRONG_MAP=0
    WRONG_REDUCE=0
  File Input Format Counters
    Bytes Read=65602380
  File Output Format Counters
    Bytes Written=6675436
15/10/25 16:10:24 INFO MahoutDriver: Program took 6697 ms (Minutes: 0.11161666666666667)

```

(capiq_ya) 3. Generate vectors from SeqD
arist...
sunsip...
\$bin/mahout seq2sparse -i reuter...
sparse-kmeans
4. Cluster with KMeans...
\$bin/mahout kmeans -i reuter...
reuters/kmeans-clusters -o reuter...
org.apache.mahout.common.distanc...
5. Dump the cluster result to...
\$bin/mahout clusterdump -i reuter...
out-seqdir-sparse-kmeans/dict...
kmeans-dump -n 5 -b 100
6. Show the result
For cluster points:
\$bin/mahout seqdumper -i reuter...
reuters/kmeans-clusters

By observing the running time, we find that Canopy clustering takes more time even though we use a smaller dataset. According to the algorithm, we know that Canopy clustering needs to find the number of clusters that have approximately that size first in order to prevent all points close to an already existing canopy from being the center of a new canopy. This can explain the longer running time of Canopy clustering.

3. Classification:

(1) Using the 20 newspaper data, try various classification algorithms provided by Mahout, and discuss their performance

First we apply the classification algorithms provided by Mahout:

```

[dyn-160-39-173-100:mahout-trunk Jovial$ ./examples/bin/classify-20newsgroups.sh
Discovered Hadoop v2.
Setting dfs command to /usr/local/Cellar/hadoop/2.7.0/bin/hdfs dfs, dfs rm to /usr/local/Cellar/hadoop/2.7.0/bin/hdfs d
fs -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-Jovial
Enter your choice : ]
```

We apply Complement Naïve Bayes, Naïve Bayes and Stochastic Gradient Descent on our 20 newsgroups data:

Complement Naïve Bayes:

Naïve Bayes:

Stochastic Gradient Descent

408
955 407
956 407
957 407
958 406
959 406
960 406
961 405
962 405
963 404
964 403
965 402
966 401
967 401
968 401
969 401
970 400
971 399
972 399
973 399
974 398
975 398
976 397
977 397
978 420 396
979 395
980 395
981 394
982 394
983 394
984 393
985 393
986 Peru 393
987 392
988 391
989 391
990 390
991 390
992 389
993 389
994 Haduo 389
995 389
996 389
997 388
998 388
999 388
1000 388

15/10/25 17:06:27 INFO MahoutDriver: Program took 388303 ms (Minutes: 6.471716666666667)

HW2

capiq_yahoo_comp_arison.csv

OPT

DOCX

RDF

Fall 2014

Windows 7

TXT

sunspots.txt

HW4.pdf

EECS E6893 Homework...huxi Ca!

Strategy Code Local

Spring 2015

useful

Screen Shot 2015-10...8.27 PM

SAS BASE

Visa签证

Screen Shot 2015-10...2.02 PM

Screen Shot 2015-10...59.41 PM

By comparing three classification results, we find that SGD takes much more running time than Complement Naïve Bayes and Naïve Bayes. The reason is that SGD does the same simple operation for each training example and the data size here is quite large. For Complement Naïve Bayes and Naïve Bayes, they have the similar running time while Naïve Bayes has higher accuracy. So Naïve Bayes is the best algorithm here to apply on our 20 newsgroups data.

Second we train our own model and apply our model on the 20 newsgroups data:

By observing the running time and accuracy, we can conclude that our model

performs well.

(2) Do similar experiments on the Wikipedia data that you downloaded

Here we apply CBayes and Binary CBayes classification algorithms on the Wikipedia data:

```
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
```

CBayes:

```
=====
Summary
-----
Correctly Classified Instances      :     1399      83.4228%
Incorrectly Classified Instances   :      278      16.5772%
Total Classified Instances        :    1677

=====
Confusion Matrix
-----
a   b   c   d   e   f   g   h   i   j   <-Classified as
391  9  13  0   3   4   2   21  9   3   |  455  a  = australia
3   114  1   4   0   1   0   2   1   2   |  128  b  = austria
1   0   4   0   0   0   0   0   0   0   |  5   c  = bahamas
1   4   7   342  9   2   1   6   7   2   |  381  d  = canada
0   1   1   0   28  0   0   0   0   2   |  32  e  = colombia
0   0   1   0   1   28  1   1   0   1   |  33  f  = cuba
0   0   1   0   2   1   58  2   0   0   |  64  g  = pakistan
0   0   0   1   2   0   0   3   0   1   |  7   h  = panama
2   19  41   6   15  8   11  23  381  7   |  513  i  = united kingdom
0   0   0   1   1   1   5   1   0   59  |  59  j  = vietnam

=====
Statistics
-----
Kappa                           0.7598
Accuracy                        83.4228%
Reliability                     73.6005%
Reliability (standard deviation) 0.2791
Weighted precision              0.9137
Weighted recall                 0.8342
Weighted F1 score               0.8645

15/10/25 19:57:58 INFO MahoutDriver: Program took 10316 ms (Minutes: 0.1719333333333333)
dyn-160-39-173-100:mahout-trunk Jovial$
```

Binary CBayes:

```
=====
Summary
-----
Correctly Classified Instances      :     2083      86.3599%
Incorrectly Classified Instances   :      329      13.6401%
Total Classified Instances        :    2412

=====
Confusion Matrix
-----
a   b   <-Classified as
497  7   |  504  a  = united kingdom
322  1586 |  1908  b  = united states

=====
Statistics
-----
Kappa                           0.6626
Accuracy                        86.3599%
Reliability                     60.5783%
Reliability (standard deviation) 0.5303
Weighted precision              0.9144
Weighted recall                 0.8636
Weighted F1 score               0.8737

15/10/25 20:09:33 INFO MahoutDriver: Program took 3957 ms (Minutes: 0.06595)
dyn-160-39-173-100:mahout-trunk Jovial$
```

By comparing the results of two classification algorithms, we find that Binary CBayes here takes less time and has higher accuracy. So we conclude that Binary CBayes performs better here.