

Recommendation Engine *(Using graphs)*

DST project



Efforts :

Aashish Singh (2K19/IT/001)

Aman Jain (2K19/IT/013)



Problem Statement

- ❖ **Recommend users related products based on general user data using graphs**
- ❖ A recommendation engine is a system that suggests products, services, information to users based on analysis of data. Notwithstanding, the recommendation can derive from a variety of factors such as the history of the user and the behaviour of similar users



Dataset used

- We have used Amazon's actual dataset for our project
- Dataset consist of product metadata
- Sample dataset-

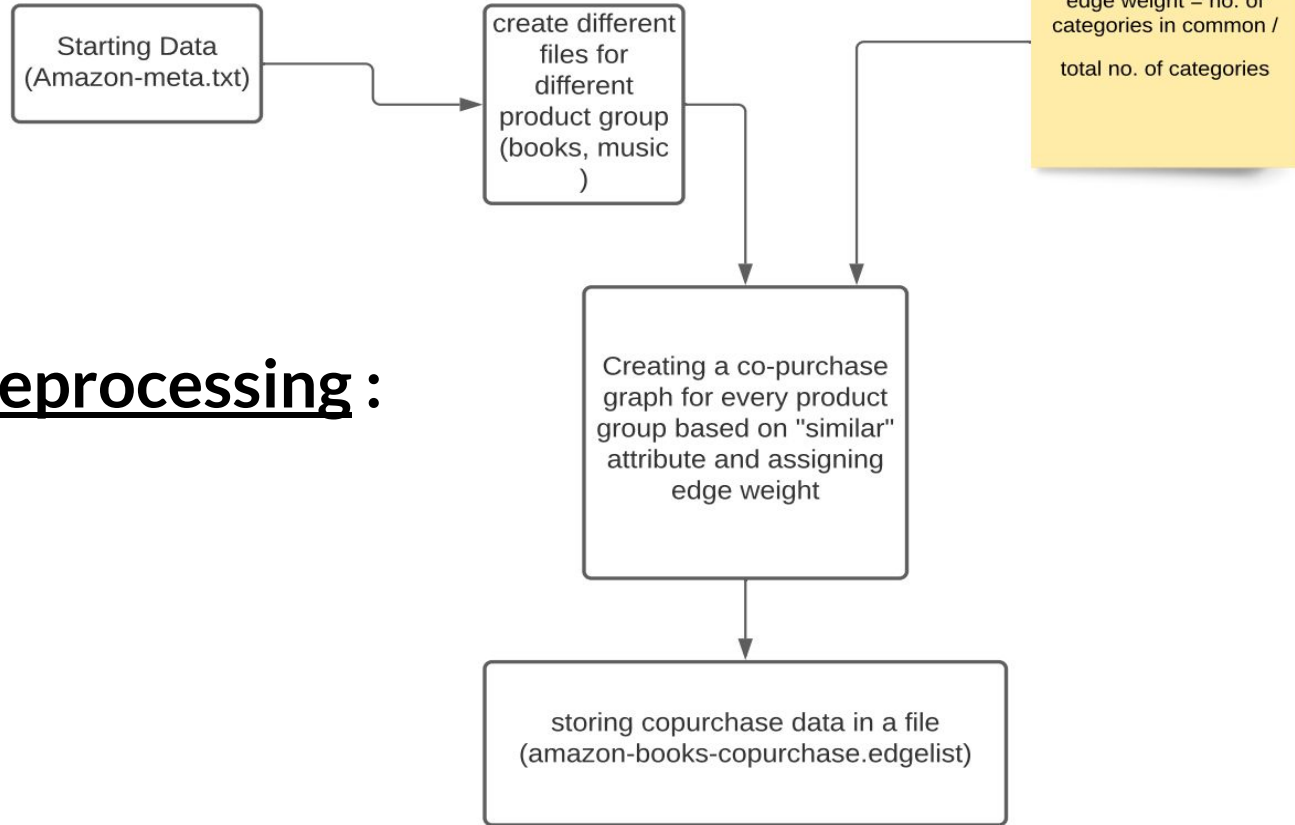
```
Id: 1
ASIN: 0827229534
title: Patterns of Preaching: A Sermon Sampler
group: Book
salesrank: 396585
similar: 5 0804215715 156101074X 0687023955 0687074231 082721619X
categories: 2
|Books[283155]|Subjects[1000]|Religion & Spirituality[22]|Christianity[12290]|Clergy[12360]|Preaching[12368]
|Books[283155]|Subjects[1000]|Religion & Spirituality[22]|Christianity[12290]|Clergy[12360]|Sermons[12370]
reviews: total: 2 downloaded: 2 avg rating: 5
2000-7-28 cutomer: A2JW67OY8U6HHK rating: 5 votes: 10 helpful: 9
2003-12-14 cutomer: A2VE83MZF98ITY rating: 5 votes: 6 helpful: 5
```

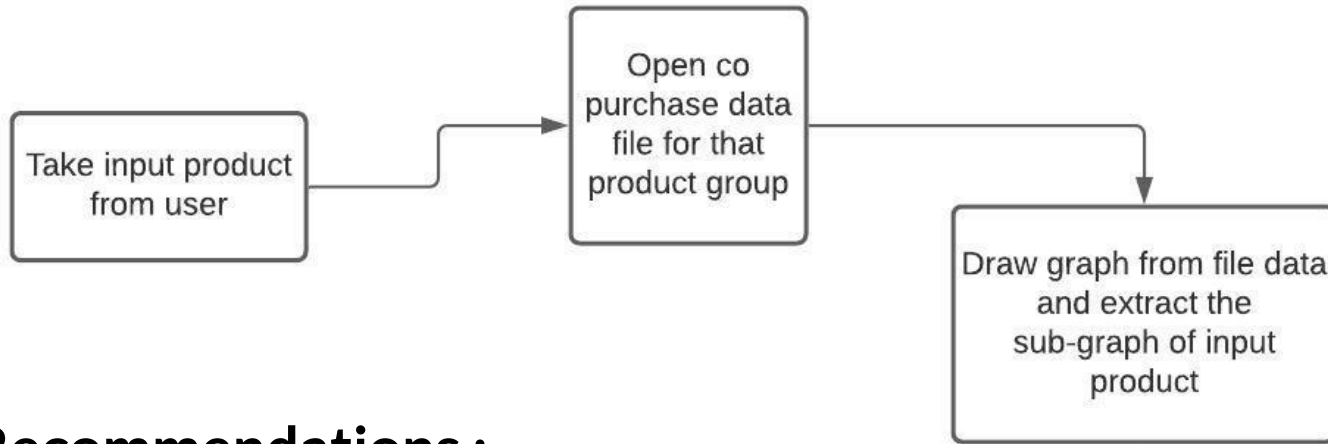


Flow charts

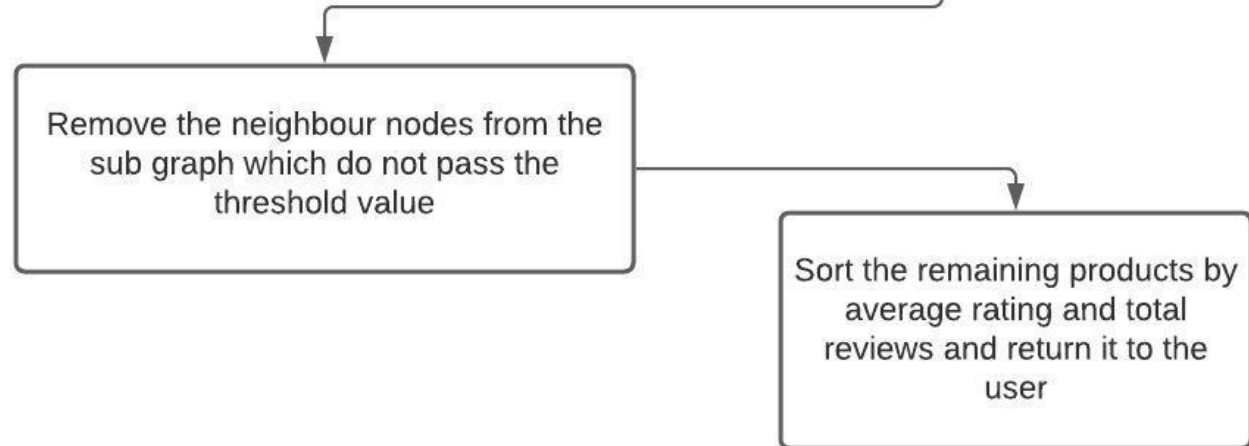
1. Preprocessing
2. Recommendation System

Preprocessing :





Recommendations :





Code snippets

Data formatting

```
8
9 # PREPROCESSING
10 #-----
11 fhr = open('amazon-meta.txt', 'r', encoding='utf-8', errors='ignore')# Reading input file
12 amazonProducts = {} #initialising an empty object
13
14 # Read data from amazon meta files and populate amazon products nested dictionary
15 (Id, ASIN, Title, Categories, Group, Copurchased, SalesRank, TotalReviews, AvgRating) = \
16 | | ("", "", "", "", "", "", 0, 0, 0.0) # initialising values
17 for line in fhr:
18     line = line.strip() #removing whitespaces
19     # a product block started
20     if(line.startswith("Id")):
21         Id = line[3:].strip()
22     elif(line.startswith("ASIN")):
23         ASIN = line[5:].strip()
24     elif(line.startswith("title")):
25         Title = line[6:].strip()
26         Title = ' '.join(Title.split())
27
28     elif(line.startswith("group")):
29         Group = line[6:].strip()
30     elif(line.startswith("salesrank")):
31         SalesRank = line[10:].strip()
32     elif(line.startswith("similar")):
33         ls = line.split()
34         Copurchased = ' '.join([c for c in ls[2:]])
35     elif(line.startswith("categories")):
36         ls = line.split()
37         Categories = ' '.join((fhr.readline()).lower() for i in range(int(ls[1].strip()))))
38         Categories = re.sub('!|!%', '%', Categories) # no escape(string digits + string punctuation)) sub('!', ' ', Categories)
```



```

37     Categories = ' '.join((fhr.readline()).lower() for i in range(int(ls[1].strip())))
38     Categories = re.compile('[%s]' % re.escape(string.digits + string.punctuation)).sub(' ', Categories)
39     Categories = ' '.join(set(Categories.split()) - set(stopwords.words("english")))
40     Categories = ' '.join(stem(word) for word in Categories.split())
41 elif(line.startswith("reviews")):
42     ls = line.split()
43     TotalReviews = ls[2].strip()
44     AvgRating = ls[7].strip()
45 # product block end
46 # write out fields to amazonProducts dictionary
47 elif(line==""):
48     try:
49         MetaData={}
50         if(ASIN != ""):
51             amazonProducts[ASIN] = MetaData
52         MetaData['Id'] = Id
53         MetaData['Title'] = Title
54         MetaData['Categories'] = ' '.join(set(Categories.split()))
55         MetaData['Group'] = Group
56         MetaData['Copurchased'] = Copurchased
57         MetaData['SalesRank'] = int(SalesRank)
58         MetaData['TotalReviews'] = int(TotalReviews)
59         MetaData['AvgRating'] = float(AvgRating)
60         MetaData['DegreeCentrality'] = DegreeCentrality
61         MetaData['ClusteringCoeff'] = ClusteringCoeff
62     except NameError:
63         continue
64     (Id, ASIN, Title, Categories, Group, Copurchased, SalesRank, TotalReviews, AvgRating, DegreeCentrality, ClusteringCoeff) = \
65     ("", "", "", "", "", "", 0, 0, 0.0, 0, 0.0)
66 fhr.close()

```

Data Classification

```
#create specific dictionaries
amazonBooks = {}
amazonMusic = {}

for asin,metadata in amazonProducts.items():
    if (metadata['Group']=='Book'):
        amazonBooks[asin] = amazonProducts[asin]
    elif(metadata['Group']=='Music'):
        amazonMusic[asin] = amazonProducts[asin]
    elif(metadata[asin]=='') \

# Write amazonBooks data to file
fhw = open('amazon-books.txt', 'w', encoding='utf-8', errors='ignore')
fhw.write("Id\t" + "ASIN\t" + "Title\t" + "Categories\t" + "Group\t" + "Copurchased\t" + "SalesRank\t" + "TotalReviews\t" + "AvgRating\n")
for asin, metadata in amazonBooks.items():
    fhw.write(metadata['Id'] + "\t" + \
        asin + "\t" + \
        metadata['Title'] + "\t" + \
        metadata['Categories'] + "\t" + \
        metadata['Group'] + "\t" + \
        metadata['Copurchased'] + "\t" + \
        str(metadata['SalesRank']) + "\t" + \
        str(metadata['TotalReviews']) + "\t" + \
        str(metadata['AvgRating']) + "\n")
fhw.close()
```

Creating co-purchased graph

```
120 # Create a product copurchase graph for analysis where the graph nodes for product ASINs
121 # and graph edge exists if two products were copurchased,
122 # with edge weight being a measure of category similarity between ASINs
123
124 copurchaseGraphBooks = networkx.Graph()
125 for asin,metadata in amazonBooks.items():
126     copurchaseGraphBooks.add_node(asin)
127     for a in metadata['Copurchased'].split():
128         copurchaseGraphBooks.add_node(a.strip())
129         similarity = 0
130         n1 = set((amazonBooks[asin]['Categories']).split())
131         n2 = set((amazonBooks[a]['Categories']).split())
132         n1In2 = n1 & n2      # intersection (Number of words that are common between Categories of connected Nodes)
133         n1Un2 = n1 | n2     # union (Total number of words in both Categories of connected Nodes)
134         if(len(n1Un2)) > 0:
135             similarity = round(len(n1In2)/len(n1Un2), 2)
136         copurchaseGraphBooks.add_edge(asin, a.strip(), weight=similarity)
137
138 fhw = open('amazon-books-copurchase.edgelist', 'wb')
139 networkx.write_weighted_edgelist(copurchaseGraphBooks, fhw)
140 fhw.close()
```

Taking user input

```
16 print("Available groups\n")
17 print("1. Books\n")
18 print("2. Music\n")
19
20 val = input("Choose the group\n")
21 numberToGroupMapping = {
22     "1": "amazon-books",
23     "2": "amazon-music"
24 }
25
26 # key = ASIN; value = MetaData associated with ASIN
27 fhr = open(numberToGroupMapping[val]+".txt", 'r', encoding='utf-8', errors='ignore')
28 amazonProducts = {}
29 display = []
30 fhr.readline() #read first line and skip to second line
31 for line in fhr:
32     cell = line.split('\t')
33     MetaData = {}
34     MetaData['Id'] = cell[0].strip()
35     ASIN = cell[1].strip()
36     MetaData['Title'] = cell[2].strip()
37     MetaData['Categories'] = cell[3].strip()
38     MetaData['Group'] = cell[4].strip()
39     MetaData['Copurchased'] = cell[5].strip()
40     MetaData['SalesRank'] = int(cell[6].strip())
41     MetaData['TotalReviews'] = int(cell[7].strip())
42     MetaData['AvgRating'] = float(cell[8].strip())
43     display.append(cell[1].strip())
44     amazonProducts[ASIN] = MetaData
45 fhr.close()
```

Getting sub-graph of product selected by user

```
60
61 purchasedAsin = input('Enter your product id\n')
62
63 print("Looking for Recommendations for this Product:")
64 print("\n-----")
65
66 # Get the depth-1 ego network of purchasedAsin from copurchaseGraph
67 try:
68     n = purchasedAsin
69     ego = networkx.ego_graph(copurchaseGraph, n, radius=1)
70     purchasedAsinEgoGraph = networkx.Graph(ego)
71 except:
72     print("No similar product available")
73     exit()
74
75
76 # only retain edges with Threshold value
77 threshold = 0.1
78 purchasedAsinEgoTrimGraph = networkx.Graph()
79 purchasedAsinEgoTrimGraph.add_node(purchasedAsin)
80
81 for f,t,e in purchasedAsinEgoGraph.edges(data=True):
82     if e['weight'] >= threshold:
83         purchasedAsinEgoTrimGraph.add_edge(f,t, weight=e['weight'])
84
85 networkx.draw_networkx(purchasedAsinEgoTrimGraph)
86 plt.show()
87
```


Returning recommendations to the user

```
90
91 # Get the list of nodes connected to the purchasedAsin
92 purchasedAsinNeighbours = purchasedAsinEgoTrimGraph.neighbors(purchasedAsin)
93
94 # Get Top Five book recommendations from among the purchasedAsinNeighbours based on one or more of the following data of the
95 # neighbouring nodes: SalesRank, AvgRating, TotalReviews
96
97 # Accessing metadata with ASIN in purchasedAsinNeighbours
98 AsMeta = []
99 for asin in purchasedAsinNeighbours:
100     ASIN = asin
101     Title = amazonProducts[asin]['Title']
102     SalesRank = amazonProducts[asin]['SalesRank']
103     TotalReviews = amazonProducts[asin]['TotalReviews']
104     AvgRating = amazonProducts[asin]['AvgRating']
105     Similarity = str(amazonProducts[asin]['Similarity']*100)+"%"
106     AsMeta.append([ASIN, Title, AvgRating, TotalReviews, Similarity])
107
108 # Sorting the top five nodes in purchasedAsinNeighbour by Average Rating then by TotalReviews
109 T5_byAvgRating_then_byTotalReviews = sorted(AsMeta, key=lambda x: (x[3], x[2]), reverse=True)
110
111 # Print Top Recommendations
112 print('\nTop Recommendations by AvgRating then by TotalReviews for this Product:')
113 print('\n-----')
114 t = PrettyTable(['Product ID', 'Title', 'Average Rating', 'Total Reviews', 'Similarity'])
115 for asin in T5_byAvgRating_then_byTotalReviews:
116     t.add_row(asin)
117
118 print(t)
```

Result

Available groups

1. Books

2. Music

Choose the group

1

Product ID	Title	Average Rating	Total Reviews
0028638506	The Complete Idiot's Guide to American History, Second Edition	3.0	14
1586211943	Swimming Across : A Memoir	4.5	26
1570194424	War of the Worlds	4.5	5
9626341017	Treasure Island (Classic Literature With Classical Music. Junior Classics)	4.0	197
0913321079	Muhammad the Prophet	0.0	0
0252067282	Pistol Packin' Mama: Aunt Molly Jackson and the Politics of Folksong (Music in American Life)	5.0	2
0875885225	Crissy Family Encyclopedia	4.5	7
0595193501	The Claims of Christ: What Jesus Had to Say About Himself	5.0	6
0670877697	The Roald Dahl Treasury	4.5	16
1893342026	Guide to Maryland Trout Fishing: The Catch and Release Streams	5.0	1
0486264815	Chinese Brushwork in Calligraphy and Painting : Its History, Aesthetics, and Techniques	4.0	1
0140388346	The Thief	4.5	87
0805039392	The Evidence of Things Not Seen : Reissued Edition	4.0	3
0764120611	Japanese Grammar (Barron's Grammar Series)	4.0	8
1590790154	Mastering the Rockefeller Habits: What You Must Do to Increase the Value of Your Fast-Growth Firm	5.0	14

Enter your product id

Product ID	Title	Average Rating	Total Reviews
0028638506	The Complete Idiot's Guide to American History, Second Edition	3.0	14
1586211943	Swimming Across : A Memoir	4.5	26
1570194424	War of the Worlds	4.5	5
9626341017	Treasure Island (Classic Literature With Classical Music. Junior Classics)	4.0	197
0913321079	Muhammad the Prophet	0.0	0
0252067282	Pistol Packin' Mama: Aunt Molly Jackson and the Politics of Folksong (Music in American Life)	5.0	2
0875885225	Crissy Family Encyclopedia	4.5	7
0595193501	The Claims of Christ: What Jesus Had to Say About Himself	5.0	6
0670877697	The Roald Dahl Treasury	4.5	16
1893342026	Guide to Maryland Trout Fishing: The Catch and Release Streams	5.0	1
0486264815	Chinese Brushwork in Calligraphy and Painting : Its History, Aesthetics, and Techniques	4.0	1
0140388346	The Thief	4.5	87
0805039392	The Evidence of Things Not Seen : Reissued Edition	4.0	3
0764120611	Japanese Grammar (Barron's Grammar Series)	4.0	8
1590790154	Mastering the Rockefeller Habits: What You Must Do to Increase the Value of Your Fast-Growth Firm	5.0	14

Enter your product id

0140388346

Looking for Recommendations for this Product:

Top Recommendations by AvgRating then by TotalReviews for this Product:

Product ID	Title	Average Rating	Total Reviews	Similarity
0140376410	The Ear, the Eye, and the Arm	4.5	265	65.0%
0140386351	A Girl Named Disaster	4.0	83	47.0%
0531095398	A Girl Named Disaster	4.0	83	41.0%
068982033X	The Moorchild (Aladdin Fantasy)	4.5	54	43.0%
068817423X	The Queen of Attolia	4.5	31	31.0%
0380733048	The Queen of Attolia (rpkg)	4.5	31	32.0%
0374410828	Celine (Sunburst Book)	4.5	12	24.0%
0446522481	A Kind of Grace : The Autobiography of the World's Greatest Female Athlete	5.0	9	10.0%

Example of a sub-graph

