

```

1 import urllib
2 import random
3 import gzip
4 import math
5 import numpy
6 import string
7 import random
8 import pandas as pd
9 import seaborn as sns
10 import matplotlib.pyplot as plt
11 from collections import defaultdict
12 import scipy.optimize
13 from sklearn import svm
14 from sklearn import linear_model
15 from sklearn.metrics import mean_absolute_percentage_error
16 from sklearn.metrics import confusion_matrix
17 from sklearn.metrics import mean_squared_error

```

```

1 from google.colab import drive
2 drive.mount('/content/drive')

```

## ▼ PreProcessing Methods

```

1 #PreProcessing methods
2 userDict = {}
3 userIndex=0
4 bookDict = {}
5 bookIndex=0
6 def getUserId(key):
7     if key not in userDict.keys():
8         global userIndex
9         userIndex=userIndex+1
10        userDict[key] = userIndex
11        return userIndex
12    return userDict[key]
13 def getBookId(key):
14     if key not in bookDict.keys():
15         global bookIndex
16         bookIndex=bookIndex+1
17         bookDict[key] = bookIndex
18         return bookIndex
19    return bookDict[key]
20 def preProcess(val):
21     val["user_id"] = getUserId(val["user_id"])
22     val["book_id"] = getBookId(val["book_id"])
23     del val["review_id"]
24     del val["date_added"]
25     del val["date_updated"]
26     del val["read_at"]
27     del val["started_at"]
28     del val["n_votes"]
29     del val["n_comments"]
30     return val

```

```

1 f = gzip.open("/content/drive/MyDrive/Final Project/young_adult_10000.json.gz")
2 dataset = []
3 for l in f:
4     dataset.append(preProcess(eval(l)))
5 random.shuffle(dataset)
6 dataTrain = dataset[:9000]
7 dataTest = dataset[9000:]
8 df = pd.DataFrame (dataset, columns = ['user_id', 'book_id', 'rating'])

```

```

1 df

```

	user_id	book_id	rating
0	494	3222	3
1	480	2174	5
2	477	924	0
3	243	2354	5
4	66	949	5
...	...	...	...

```
1 df.describe()
```

	user_id	book_id	rating
count	10000.000000	10000.000000	10000.000000
mean	349.337300	1728.146700	3.755500
std	204.133854	1443.718744	1.245662
min	1.000000	1.000000	0.000000
25%	162.000000	389.000000	3.000000
50%	367.000000	1396.500000	4.000000
75%	506.000000	2824.000000	5.000000
max	715.000000	4980.000000	5.000000

```
1 df.info()
```

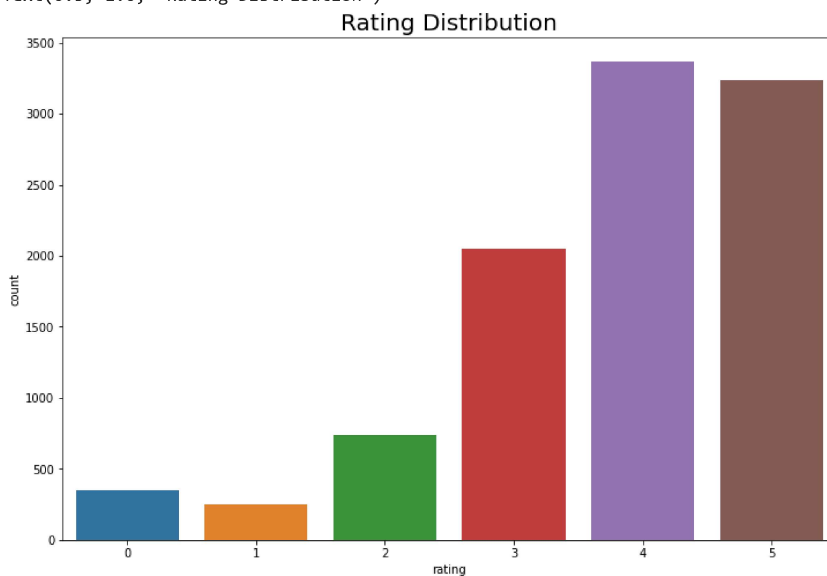
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 3 columns):
#   Column  Non-Null Count  Dtype
---  -
0   user_id  10000 non-null     int64
1   book_id  10000 non-null     int64
2   rating   10000 non-null     int64
dtypes: int64(3)
memory usage: 234.5 KB
```

```
1 df.isnull().sum()
```

```
user_id    0
book_id    0
rating     0
dtype: int64
```

```
1 plt.figure(figsize=(12,8))
2 sns.countplot(x='rating',data=df)
3 plt.title('Rating Distribution',size=20)
```

```
Text(0.5, 1.0, 'Rating Distribution')
```



## ▼ Bag of Words Approach

```

1 wordCount = defaultdict(int)
2 punctuation = set(string.punctuation)
3
4 for d in dataset:
5     r = ''.join([c for c in d['review_text'].lower() if not c in punctuation])
6     for w in r.split():
7         wordCount[w] += 1
8
9 sorted_counts = sorted(wordCount.items(), key=lambda x:x[1])[:-1]
10
11 counts = [(i[1], i[0]) for i in sorted_counts]
12
13 words = [x[1] for x in counts[:1000]]
14 wordId = dict(zip(words, range(len(words))))
15 wordSet = set(words)

```

```

1 def feature(datum):
2     feat = [0]*len(words)
3     # removing punctuation from review
4     r = ''.join([c for c in datum['review_text'].lower() if not c in punctuation])
5     # adding word counts to feature
6     for w in r.split():
7         if w in words:
8             feat[wordId[w]] += 1
9     feat.append(1)
10    return feat
11 def accuracy(y, predictions):
12     correct = predictions == y
13     acc = sum(correct) / len(correct)
14     return acc

```

```

1 X = [feature(d) for d in dataset]
2 y = [d['rating'] for d in dataset]
3
4 Xtrain = X[:9*len(X)//10]
5 ytrain = y[:9*len(y)//10]
6 Xtest = X[9*len(X)//10:]
7 ytest = y[9*len(y)//10:]

```

```

1 # Training a model for Logistic Regression
2 mod = linear_model.LogisticRegression(C=1)
3 mod.fit(Xtrain,ytrain)
4 predictions = mod.predict(Xtest)
5 acc = accuracy(ytest, predictions)

```

```
1 print('The accuracy for Logistic Regression for Bag of Words approach - \n', acc*100, '%')
```

The accuracy for Logistic Regression for Bag of Words approach -  
45.5 %

```

1 # Training a model for Linear Regression
2 mod1 = linear_model.LinearRegression()
3 mod1.fit(Xtrain,ytrain)
4 predictions1 = mod1.predict(Xtest)
5 mse1 = mean_squared_error(predictions1, ytest)
6 mape_1 = mean_absolute_percentage_error(predictions1, ytest)*100

```

```

1 print('The Mean Squared Error of LInear Regression for Bag of Words approach - \n', mse1, '\n')
2 print('The Mean absolute percentage error of Linear Regression for Bag of Words approach - \n', mape_1)

```

The Mean Squared Error of LInear Regression for Bag of Words approach -  
1.5754416572189098

The Mean absolute percentage error of Linear Regression for Bag of Words approach -  
36.60850676472431

## ▼ Final Approach - Jaccard Similarity

```

1 usersPerItem = defaultdict(set) # Maps an item to the users who rated it
2 itemsPerUser = defaultdict(set) # Maps a user to the items that they rated

```

```

3 reviewsPerUser = defaultdict(list)
4 reviewsPerItem = defaultdict(list)
5 ratingDict = {} # To retrieve a rating for a specific user/item pair
6
7 for d in dataTrain:
8     user,item = d['user_id'], d['book_id']
9     usersPerItem[item].add(user)
10    itemsPerUser[user].add(item)
11    ratingDict[(user,item)] = d['rating']

```

```

1 userAverages = {}
2 itemAverages = {}
3
4 for u in itemsPerUser:
5     rs = [ratingDict[(u,i)] for i in itemsPerUser[u]]
6     userAverages[u] = sum(rs) / len(rs)
7
8 for i in usersPerItem:
9     rs = [ratingDict[(u,i)] for u in usersPerItem[i]]
10    itemAverages[i] = sum(rs) / len(rs)
11
12 for d in dataTrain:
13     user,item = d['user_id'], d['book_id']
14     reviewsPerUser[user].append(d)
15     reviewsPerItem[item].append(d)
16

```

```

1 # From the pseudo code for jaccard similarity
2 def Jaccard(s1, s2):
3     number = len(s1.intersection(s2))
4     denom = len(s1.union(s2))
5     if denom == 0:
6         return 0
7     return number / denom
8
9 def predictRating(user,item):
10    ratings = []
11    similarities = []
12    for d in reviewsPerUser[user]:
13        i2 = d['book_id']
14        if i2 == item: continue
15        ratings.append(d['rating'] - itemAverages[i2])
16        similarities.append(Jaccard(usersPerItem[item],usersPerItem[i2]))
17    if (sum(similarities) > 0):
18        weightedRatings = [(x*y) for x,y in zip(ratings,similarities)]
19        return itemAverages[item] + sum(weightedRatings) / sum(similarities)
20    else:
21        # User hasn't rated any similar items
22        ratingMean = sum([d['rating'] for d in dataTrain]) / len(dataTrain)
23        return ratingMean
24

```

```

1 predictions2 = [predictRating(d['user_id'], d['book_id']) for d in dataTest]
2 labels = [d['rating'] for d in dataTest]
3 mse2 = mean_squared_error(predictions2, labels)
4 mape_2 = mean_absolute_percentage_error(predictions2, labels)*100

```

```

1 print('The Mean Squared Error for Jaccard Similarities - \n', mse2, '\n')
2 print('The Mean absolute percentage error for Jaccard Similarities - \n', mape_2)

```

The Mean Squared Error for Jaccard Similarities -  
1.8375712166843186

The Mean absolute percentage error for Jaccard Similarities -  
31.011538168802392

✓ 0s completed at 8:52 PM

● ✕