

Open in app



Neha Patel

13 Followers

About



Recommendation system in python using ALS algorithm and Apache Spark



Neha Patel Apr 7, 2019 · 4 min read

Recommendation system has become integral part of many online platforms given it's use in increasing sales and customer retention. Amazon, Netflix, Facebook, Linkedin and many more make use of recommendation system on their platform to recommend products, movies, people, posts respectively to users. Ever wondered how these companies provides recommendations.

RECOMMENDATION SYSTEM

Given it's popularity there are three main techniques used for providing recommendations online-Collaborative filtering, Content-based and Hybrid



We are going to build a recommendation system in python using Apache spark and Jupyter Notebook. You can make this simple recommendation model as a mini project in college or as a personal project.

Requirements:

Your machine should have latest version of Python, Apache Spark and Jupyter Notebook installed. You also need to connect pyspark with Jupyter Notebook and many tutorials are available out there to do the same.

Dataset used is downloaded from the link given below:

Amazon review data

This dataset contains product reviews and metadata from Amazon, including 142.8 million reviews spanning May 1996 ...

jmcauley.ucsd.edu

The data used here is of Musical Instruments for training the model.

Terminologies:

There are certain terminologies which needs to be understood before moving forward.

- 1. **Apache Spark:** Apache Spark is an open-source distributed general-purpose cluster-computing framework.It can be used with Hadoop too.
- 2. **Collaborative filtering:** Collaborative filtering is a method of making automatic predictions (filtering) about the interests of a user by collecting preferences or taste information from many users. Consider example if a person A likes item 1, 2, 3 and B like 2,3,4 then they have similar interests and A should like item 4 and B should like item 1.
- 3. **Alternating least square(ALS) matrix factorization:** The idea is basically to take a large (or potentially huge) matrix and factor it into some smaller representation of the original matrix through alternating least squares. We end up with two or more



4. **PySpark**: **PySpark** is the collaboration of Apache Spark and Python. PySpark is the Python API for Spark.

So let's start making our recommendation model in jupyter notebook.

1.Initialize spark session:

```
from pyspark.sql import SparkSession
spark = SparkSession.builder.appName('Recommendation_system')
.getOrCreate()
```

2. Load Dataset in Apache Spark

```
df = spark.read.json("Musical_Instruments_5.json")
df.show(100,truncate=True)
```

Your dataframe will look like following:

```
asin| helpful|overall|
                              reviewText| reviewTime| reviewerID|
                                                                reviewerName
iewTime
+------
|1384719342| [0, 0]| 5.0|Not much to write...|02 28, 2014|A2IBPI20UZIR0U|cassandra tu "Yea...|
                                                                                                139
|1384719342|[13, 14]| 5.0|The product does ... |03 16, 2013|A14VAT5EAX3D9S|
                                                                        Jake
                                                                                                136
3392666
                  5.0|The primary job o... | 08 28, 2013 | A195EZSQDW3E21 | Rick Bennette "Ri... | It Does The Job Well |
|1384719342| [1, 1]|
                                                                                                137
7648999
|1384719342| [0, 0]|
                  5.0|Nice windscreen p...|02 14, 2014|A2C00NNG1ZQQG2|RustyBill "Sunday...|GOOD WINDSCREEN F...|
2336860
                                                                SEAN MASLANKA No more pops when...
                  5.0|This pop filter i...|02 21, 2014| A94QU4C90B1AX|
|1384719342| [0, 0]|
2946860
|B00004Y2UT| [0, 0]|
                  5.0|So good that I bo...|12 21, 2012|A2A039TZMZHH9Y| Bill Lewey "blewey"|
6948999
|B00004Y2UT| [0, 0]|
                   5.0|I have used monst...|01 19, 2014|A1UPZM995ZAH90|
                                                                       Brian | Monster Standard ...
0089600
|B00004Y2UT| [0, 0]|
                   3.0|I now use this ca...|11 16, 2012| AJNFQI3YR6XJ5| Fender Guy "Rick"|Didn't fit my 199...|
                                                                                                135
3024000
```

Original dataframe

3. Select appropriate columns



We do not need all the columns present in the dataframe .Only asin which is ProductID, reviewerID and overall(rating given by users to each product) is required.

reviewerID	overall	asin
A2IBPI20UZIR0U	5.0	1384719342
A14VAT5EAX3D95		1384719342
A195EZSQDW3E21	200 000	1384719342
A2C00NNG1ZQQG2	5.0	1384719342
A94QU4C90B1AX	5.0	1384719342
A2A039TZMZHH9Y	5.0	B00004Y2UT
A1UPZM995ZAH90	5.0	B00004Y2UT
AJNFQI3YR6XJ5	3.0	B00004Y2UT
A3M1PLEYNDEY08	5.0	B00004Y2UT
AMNTZU1YQN1TH	5.0	B00004Y2UT
A2NYK9KWFMJV4Y	5.0	B00004Y2UT
A35QFQI0M46LW0	4.0	B00005ML71
A2NIT6BKW11XJQ	3.0	B00005ML71
A1C0009L0LVI39	5.0	B00005ML71
A17SLR18TUMULM	5.0	B00005ML71
A2PD27UKAD3Q00	2.0	B00005ML71
AKSFZ4G1AXYFC	4.0	B000068NSX
A670JZLHBBUQ9	5.0	B000068NSX
A2EZWZ8MBEDOLN	5.0	B000068NSX
A1CL807EOUPVP1	5.0	B000068NSX

4. Importing important modules

from pyspark.ml.evaluation import RegressionEvaluator
from pyspark.ml.recommendation import ALS

5. Converting String to index

Before making an ALS model it needs to be clear that ALS only accepts integer value as parameters. Hence we need to convert asin and reviewerID column in index form.



```
indexer = [StringIndexer(inputCol=column, outputCol=column+"_index")
for column in list(set(nd.columns)-set(['overall'])) ]

pipeline = Pipeline(stages=indexer)
transformed = pipeline.fit(nd).transform(nd)
transformed.show()
```

asin	overall	reviewerID	reviewerID_index	asin_index
1384719342	5.0	A2IBPI20UZIR0U	72.0	781.0
1384719342	5.0	A14VAT5EAX3D9S	359.0	781.0
1384719342	5.0	A195EZSQDW3E21	436.0	781.0
1384719342	5.0	A2C00NNG1ZQQG2	1216.0	781.0
1384719342	5.0	A94QU4C90B1AX	1137.0	781.0
B00004Y2UT	5.0	A2A039TZMZHH9Y	54.0	629.6
B00004Y2UT	5.0	A1UPZM995ZAH90	348.0	629.6
B00004Y2UT	3.0	AJNFQI3YR6XJ5	324.0	629.6
B00004Y2UT	5.0	A3M1PLEYNDEY08	12.0	629.6
B00004Y2UT	5.0	AMNTZU1YQN1TH	185.0	629.6
B00004Y2UT	5.0	A2NYK9KWFMJV4Y	4.0	629.6
B00005ML71	4.0	A35QFQI0M46LWO	425.0	870.6
B00005ML71	3.0	A2NIT6BKW11XJQ	652.0	870.6
B00005ML71	5.0	A1C0009L0LVI39	55.0	870.6
B00005ML71	5.0	A17SLR18TUMULM	651.0	870.6
B00005ML71	2.0	A2PD27UKAD3Q00	290.0	870.6
B000068NSX	4.0	AKSFZ4G1AXYFC	93.0	538.6
8000068NSX	5.0	A670JZLHBBUQ9	258.0	538.6
B000068NSX	5.0	A2EZWZ8MBEDOLN	3.0	538.6
B000068NSX	5.0	A1CL807EOUPVP1	31.0	538.6

String to index conversion

6. Creating training and test data

```
(training,test)=transformed.randomSplit([0.8, 0.2])
```

7. Creating ALS model and fitting data



```
gative=True;
model=als.fit(training)
```

8. Generate predictions and evaluate rmse

```
evaluator=RegressionEvaluator(metricName="rmse",labelCol="overall",pr
edictionCol="prediction")

predictions=model.transform(test)
rmse=evaluator.evaluate(predictions)

print("RMSE="+str(rmse))
predictions.show()
```

RMSE=1.168435412679992

asin o	verall	reviewerID	reviewerID_index	asin_index	prediction
 B000CCJP4I	5.0		858.0	148.0	3.0042644
B000CCJP4I		A3J8U952XAL34Z	502.0	148.0	
B000CCJP4I	41,537,170	A1YR3RVSBZK8CW	30.0	148.0	
B000KIPTE4	1000000	A3F49ZMUC1GSRP	1,000,000	463.0	
B000KIPTE4	70000	A2EZWZ8MBEDOLN	3.0	463.0	
B000KIPTE4	5.0	A2AH7HRHDTQENH		463.0	
B0010LZYUU	3.0		530.0	471.0	
B0010LZYUU	5.0	A2IBPI20UZIR0U	72.0	471.0	5.033518
B0002GZ052	5.0	A26HM2R5529NYY	822.0	496.0	4.2211905
B0002GZ052	5.0	A1T4U9CAQ25IBR	750.0	496.0	4.59568
B000BKY8CU	4.0	A3PGQWCSJPCYDH	64.0	243.0	3.317696
B000RPUMII	5.0	A223S6N0DBQBHP	236.0	392.0	3.5661318
B000RPUMII	1.0	ANAKK5KNUAP17	663.0	392.0	4.1016426
B000RPUMII	5.0	A30JØRGAECAGH8	381.0	392.0	2.6469836
B000KUCQXY	4.0	A2Q6KC2KU2TØOL	1315.0	540.0	3.1362884
B000KUCQXY	5.0	A208BAXJPDSVØM	365.0	540.0	4.13888
B001GGYF4E	4.0	A3S737ZGWE1GKY	1346.0	623.0	3.714615
B001GGYF4E	3.0	A8DCZN408QYKC	953.0	623.0	3.6110063
B009EOKTCM	5.0	A1SD1C8XK3Z3V1	6.0	858.0	4.6142297
B0002GLCRC	5.0	A24VCDADYAIHAM	481.0	31.0	4.6885085

only showing top 20 rows



user recs=model.recommendForAllUsers(20).show(10)

```
| reviewerID_index| recommendations|
| 471|[[746, 5.981454],...|
| 1342|[[412, 6.348332],...|
| 463|[[898, 6.0316505]...|
| 833|[[426, 6.379773],...|
| 496|[[426, 5.978744],...|
| 148|[[898, 6.291563],...|
| 1088|[[426, 4.7776713]...|
| 1238|[[579, 5.840673],...|
| 540|[[898, 4.944048],...|
| 392|[[710, 5.531966],...|
```

Recommendations

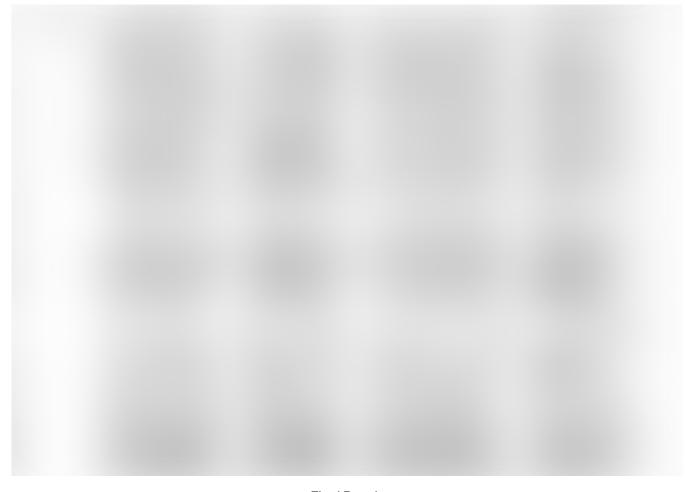
10. Converting back to string form

As seen in above image the results are in integer form we need to convert it back to its original name. The code is little bit longer given so many conversions.

```
import pandas as pd
recs=model.recommendForAllUsers(10).toPandas()
nrecs=recs.recommendations.apply(pd.Series) \
            .merge(recs, right index = True, left index = True) \
            .drop(["recommendations"], axis = 1) \
            .melt(id vars = ['reviewerID index'], value name =
"recommendation") \
            .drop("variable", axis = 1) \
            .dropna()
nrecs=nrecs.sort values('reviewerID index')
nrecs=pd.concat([nrecs['recommendation'].apply(pd.Series),
nrecs['reviewerID index']], axis = 1)
nrecs.columns = [
        'ProductID index',
        'Rating',
        'UserID index'
```



```
dict1 =dict(zip(md['reviewerID_index'],md['reviewerID']))
dict2=dict(zip(md['asin_index'],md['asin']))
nrecs['reviewerID']=nrecs['UserID_index'].map(dict1)
nrecs['asin']=nrecs['ProductID_index'].map(dict2)
nrecs=nrecs.sort_values('reviewerID')
nrecs.reset_index(drop=True, inplace=True)
new=nrecs[['reviewerID','asin','Rating']]
new['recommendations'] = list(zip(new.asin, new.Rating))
res=new[['reviewerID','recommendations']]
res_new=res['recommendations'].groupby([res.reviewerID]).apply(list).
reset_index()
```



Final Result

Voila!!

You just made an recommendation application.

Get started

Open in app



Data Science

Recommendation System

Apache Spark

Python

About Write Help Legal

Get the Medium app



