```
pip install pyspark
→ Collecting pyspark
      Downloading pyspark-3.5.2.tar.gz (317.3 MB)
                                               - 317.3/317.3 MB 2.8 MB/s eta 0:00:00
      Preparing metadata (setup.py) ... done
    Requirement already satisfied: py4j==0.10.9.7 in /usr/local/lib/python3.10/dist-packages (from pyspark) (0.10.9.7)
    Building wheels for collected packages: pyspark
      Building wheel for pyspark (setup.py) ... done
      Stored in directory: /root/.cache/pip/wheels/34/34/bd/03944534c44b677cd5859f248090daa9fb27b3c8f8e5f49574
    Successfully built pyspark
     Installing collected packages: pyspark
    Successfully installed pyspark-3.5.2
import pandas as pd
from pyspark.sql import SparkSession
spark = SparkSession.builder.master("local").appName(
    "Movie_recommendation").config(
        "spark.some.config.option","some-value").getOrCreate()
rating=spark.read.csv(
    '/content/ratings_small.csv',header=True,inferSchema=True)
rating=rating.drop('timestamp')
movies=spark.read.csv(
    '<u>/content/movies_metadata.csv</u>',header=True,inferSchema=True)
#Rename the 'id' column to 'movieId' in the movies DataFrame.
movies = movies.withColumnRenamed('id','movieId')
movie_data_vishnav=rating.join(movies,on='movieId')
Columns=len(movie_data_vishnav.columns)
Rows=movie_data_vishnav.count()
print('Number of Columns: {}\nNumber of Rows: {}'.format(Columns,Rows))
movie_data_vishnav.columns
    Number of Columns: 26
    Number of Rows: 44925
    ['movieId'
      'userId',
      'rating',
     'adult',
      'belongs_to_collection',
     'budget',
      genres',
     'homepage',
     'imdb_id',
      'original_language',
      'original_title',
      'overview',
      'popularity'
      'poster_path'
      'production_companies',
      'production_countries',
      'release_date',
      'revenue',
     'runtime'
      'spoken_languages',
      'status',
      'tagline',
      'title',
     'video',
      'vote average',
     'vote_count']
```

Data Cleaning: Replace all the zeros in the above mentioned fields with NaN

```
import numpy as np
from pyspark.sql.functions import when
movie_data_vishnav=movie_data_vishnav.withColumn(
    "userId",when(movie_data_vishnav.userId==0,np.nan).otherwise(movie_data_vishnav.userId))
movie_data_vishnav=movie_data_vishnav.withColumn(
    "movieId",when(movie_data_vishnav.movieId==0,np.nan).otherwise(movie_data_vishnav.movieId))
movie_data_vishnav=movie_data_vishnav.withColumn(
    "rating",when(movie_data_vishnav.rating==0,np.nan).otherwise(movie_data_vishnav.rating))
movie_data_vishnav=movie_data_vishnav.withColumn(
    "title",when(movie_data_vishnav.title==0,np.nan).otherwise(movie_data_vishnav.title))
from pyspark.sql.types import IntegerType
movie_data_vishnav=movie_data_vishnav.withColumn(
    "budget", movie_data_vishnav["budget"].cast(IntegerType()))
movie_data_vishnav.show()
```

```
genres|
|movieId|userId|rating|adult|belongs_to_collection| budget|
                                                                                                      homepage| imdb_id|original_language|
     ---+----+
   949.0| 647.0|
                    4.0|False|
                                                   NULL | 60000000 | [{'id': 28, 'name...|
                                                                                                            NULL|tt0113277|
                                                                                                                                              en|
   949.0 564.0
                   3.0|False|
                                                   NULL | 60000000 | [{'id': 28, 'name...|
                                                                                                            NULL|tt0113277|
                                                                                                                                              en|
                                                   NULL | 60000000 | [{'id': 28, 'name...|
NULL | 60000000 | [{'id': 28, 'name...|
   949.01 558.01
                    4.0|False|
                                                                                                            NULL | tt0113277 |
                                                                                                                                              enl
   949.0 547.0
                    4.0|False|
                                                                                                            NULL | tt0113277 |
                                                                                                                                              enl
                                                   NULL|60000000|[{'id': 28, 'name...|
NULL|60000000|[{'id': 28, 'name...|
   949.0 537.0
                    3.0|Falsel
                                                                                                            NULL | tt0113277 |
                                                                                                                                              en|
   949.0 509.0
                     2.0|False|
                                                                                                            NULL|tt0113277|
                                                                                                                                              en|
   949.0 505.0
                    3.5|False|
                                                   NULL|60000000|[{'id': 28, 'name...|
                                                                                                            NULL|tt0113277|
                                                                                                                                              en
                                                   NULL|60000000|[{'id': 28, 'name...|
NULL|60000000|[{'id': 28, 'name...|
   949.0 | 452.0 |
                                                                                                            NULL|tt0113277|
                     4.5|False|
                                                                                                                                              en|
   949.0 387.0
                     5.0 False
                                                                                                            NULL | tt0113277 |
                                                                                                                                              enl
   949.0 363.0
                                                   NULL|60000000|[{'id': 28,
                     4.0|False|
                                                                                 'name...
                                                                                                            NULL|tt0113277|
                                                                                                                                              en l
                                                  NULL | 60000000 | [{'id': 28, 'name...|
   949.0 | 311.0
                    3.0|False|
                                                                                                            NULL|tt0113277|
                                                                                                                                              en|
   949.0 263.0
                                                   NULL|60000000|[{'id': 28, 'name...|
NULL|60000000|[{'id': 28, 'name...|
                                                                                                            NULL|tt0113277|
                     3.0|False|
                                                                                                                                              enl
   949.0 242.0
                     5.0|False|
                                                                                                            NULL|tt0113277|
                                                                                                                                              enl
                                                   NULL | 60000000 | [{'id': 28, 'name...|
NULL | 60000000 | [{'id': 28, 'name...|
   949.01 232.01
                     2.0|Falsel
                                                                                                            NULL|tt0113277|
                                                                                                                                              enl
   949.0 102.0
                     4.0|False|
                                                                                                            NULL | ++0113277 |
                                                                                                                                              enl
                                                   NULL|60000000|[{'id': 28, 'name...|
   949.0 23.0
                     3.5|False|
                                                                                                            NULL | tt0113277 |
                                                                                                                                              en|
                     2.0|False| {'id': 645, 'name...|58000000|[{'id': 12, 'name...|http://www.mgm.co...|tt0113189|
   710.0 | 390.0
                                                                                                                                              en|
   710.0 | 358.0 |
                     1.0|False| {'id': 645, 'name...|58000000|[{'id': 12, 'name...|http://www.mgm.co...|tt0113189|
                                                                                                                                               en|
                                                   NULL|98000000|[{'id': 28,
  1408.0 | 665.0
                     3.0|False|
                                                                                 'name...
                                                                                                            NULL|tt0112760|
                                                                                                                                              enlCu
                                                   NULL|98000000|[{'id': 28, 'name...|
  1408.0 | 658.0
                    5.0|False|
                                                                                                            NULL|tt0112760|
only showing top 20 rows
```

## Building the Recommendatin model using ALS on the training data

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

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

# we set cold strategy to 'drop' to ensure we dont get NaN evaluation
als_vishnav=ALS(maxIter=5,regParam=0.09,rank=25,userCol='userId',itemCol='movieId',coldStartStrategy="drop",nonnegative=True)
model_vishnav=als_vishnav.fit(training)

evaluator=RegressionEvaluator(metricName="rmse",labelCol="rating",predictionCol="prediction")
predictions=model_vishnav.transform(test)
rmse=evaluator.evaluate(predictions)
print("RMSE="+str(rmse))
predictions.show()
```

## RMSE=0.9195242773432019

Ė	homepage	genres	budget	ongs_to_collection	adult b	gl	rating	userId	movieId u
/hPObJrJhUw92[	marking the star	the gravitas of	is charged	ne perfect scap	t when Ahmed	-+ 0  But	2.0	15	1
/hPObJrJhUw920	marking the star	the gravitas of $\ldots$	is charged	ne perfect scap	t when Ahmed	0  But	5.0	89	1
/hPObJrJhUw920	marking the star	the gravitas of $\ldots$	is charged	ne perfect scap	t when Ahmed	0  But	4.0	90	1
/hPObJrJhUw920	marking the star	the gravitas of	is charged	ne perfect scap	t when Ahmed	0  But	5.0	91	1
/hPObJrJhUw920	marking the star	the gravitas of $\ldots$	is charged	ne perfect scap	t when Ahmed	0  But	5.0	112	1
/hPObJrJhUw920	marking the star	the gravitas of	is charged	ne perfect scap	t when Ahmed	0  But	5.0	128	1
/hPObJrJhUw920	marking the star	the gravitas of	is charged	ne perfect scap	t when Ahmed	0  But	3.0	130	1
/hPObJrJhUw92[	marking the star	the gravitas of	is charged	ne perfect scap	t when Ahmed	0  But	2.0	138	1
/hPObJrJhUw920	marking the star	the gravitas of	is charged	ne perfect scap	t when Ahmed	0  But	5.0	168	1
/hPObJrJhUw920	marking the star	the gravitas of	is charged	ne perfect scap	t when Ahmed	0  But	4.0	175	1
/hPObJrJhUw920	marking the star	the gravitas of	is charged	ne perfect scap	t when Ahmed	0 But	2.0	176	1
/hPObJrJhUw920	marking the star	the gravitas of	is charged	ne perfect scap	t when Ahmed	0 But	5.0	179	1
/hPObJrJhUw920	marking the star	the gravitas of	is charged	ne perfect scap	t when Ahmed	0 But	5.0	185	1
/hPObJrJhUw920	marking the star	the gravitas of	is charged	ne perfect scap	t when Ahmed	5  But	4.5	205	1
/hPObJrJhUw920	marking the star	the gravitas of	is charged	ne perfect scap	t when Ahmed	0 But	3.0	212	1
/hPObJrJhUw920	marking the star	the gravitas of	is charged	ne perfect scap	t when Ahmed	0 But	3.0	213	1
/hPObJrJhUw920	marking the star	the gravitas of	is charged	ne perfect scap	t when Ahmed	0  But	5.0	219	1
/hPObJrJhUw920	marking the star	the gravitas of		ne perfect scap	t when Ahmed	0  But	5.0	242	1
/hPObJrJhUw920	marking the star	the gravitas of	is charged	ne perfect scap	t when Ahmed	5  But	1.5	261	1
	marking the star				t when Ahmed			283	1

only showing top 20 rows

#Now we will use model.transform() function in order to generate recommended movies along with their predicted features.

```
recomendations = model_vishnav.transform(single_user)
recomendations.orderBy('prediction',ascending=False).show(truncate = False )
```

<del>_</del>	++					genres	prediction			
	2459  58559	29  29	An Elephant Can  Confession of a	Be Extremely Child of the	Deceptive Century	[{'id': 35,  [{'id': 18,	'name': 'name':	'Comedy'}] 'Drama'}]	3.742302  3.169405	   