In [1]: | pip install pyspark Collecting pyspark Downloading pyspark-3.1.2.tar.gz (212.4 MB) | 212.4 MB 63 kB/s Collecting py4j==0.10.9Downloading py4j-0.10.9-py2.py3-none-any.whl (198 kB) | 198 kB 49.5 MB/s Building wheels for collected packages: pyspark o/o o-o o\o o|o o/o o-o o\o Created wheel for pyspark: filename=pyspark-3.1.2-py2.py3-none-any.whl size=212880768 sha256=1cd545 a48bdd2824ae2e2c57fcf5048b6e39f1df540ccbe8a12b9f79e7d8dde4 Stored in directory: /root/.cache/pip/wheels/a5/0a/c1/9561f6fecb759579a7d863dcd846daaa95f598744e71b 02c77 Successfully built pyspark Installing collected packages: py4j, pyspark Successfully installed py4j-0.10.9 pyspark-3.1.2 WARNING: Running pip as root will break packages and permissions. You should install packages reliabl y by using venv: https://pip.pypa.io/warnings/venv Note: you may need to restart the kernel to use updated packages. In [2]: from pyspark.sql import functions as f from pyspark.sql import DataFrameNaFunctions as DFna from pyspark.sql.functions import udf, col, when import matplotlib.pyplot as plt import pyspark as ps import os, sys, requests, json from pyspark.sql.functions import col,size,regexp replace,lit from pyspark.ml.evaluation import RegressionEvaluator from pyspark.ml.recommendation import ALS from pyspark.ml.evaluation import RegressionEvaluator from pyspark.ml.recommendation import ALS from pyspark.ml.tuning import CrossValidator, ParamGridBuilder from pyspark.ml import Pipeline from pyspark.sql import Row import numpy as np import math from pyspark.sql.functions import regexp replace from pyspark.sql import spark = ( SparkSession.builder.appName("ModelTraining") .config("spark.executor.memory","16g") .getOrCr eate()) sc = spark.sparkContext In [3]: import numpy as np import pandas as pd from operator import add from pyspark.ml.feature import RegexTokenizer, CountVectorizer, Tokenizer from pyspark.ml.feature import StopWordsRemover, VectorAssembler from pyspark.ml.feature import Word2Vec, Word2VecModel from pyspark.ml.feature import IDF from pyspark.ml import Pipeline, PipelineModel from pyspark.sql.functions import \* from pyspark.sql.types import \* import folium import html In [4]: pro=spark.read.csv('../input/data-science-for-good-careervillage/professionals.csv', header=True, quote= '"', sep=", ", multiLine=True) email=spark.read.csv('../input/data-science-for-good-careervillage/emails.csv', header=True,quote='"',s ep=",",multiLine=True) ques=spark.read.csv('../input/data-science-for-good-careervillage/questions.csv', header=True, quote='"' , sep=",", multiLine=True) match=spark.read.csv('../input/data-science-for-good-careervillage/matches.csv', header=True, quote='"', sep=",",multiLine=True) ans=spark.read.csv('../input/data-science-for-good-careervillage/answers.csv', header=True, quote='"', se p=",",multiLine=**True**) ans score=spark.read.csv('../input/data-science-for-good-careervillage/answer scores.csv', header=True, quote='"', sep=",", multiLine=True) In [5]: | #match ques=match.join(ques, match.matches question id==ques.questions id) #match ques=match ques.join(ans,match ques.questions id==ans.answers question id) #match ques=match ques.join(email, match ques.matches email id==email.emails id) #match ques=match ques.join(pro,match ques.emails recipient id==pro.professionals id) #ques ans=match ques.select('questions id','questions body','answers id','answers body','professionals In [6]: In [7]: #ques ans=ques ans.join(ans score, ques ans.answers id == ans score.id).select('questions id', 'questions b ody','professionals id','score') In [8]: #ques ans=ques ans.select('questions id','questions body','professionals id','score') In [9]: #Data cleaning question def datacleaning que(test): user regex=r"( $@\{1,15\}$ )" clean test=test.withColumn( 'user mentioned', f.array remove( f.array( f.regexp extract(f.col("Full questions"), user regex, 1), f.regexp extract( f.col("Full\_questions"), "".join([f"{user\_regex}.\*?" for i in range(0,2)]),2), f.regexp extract( f.col("Full questions"), "".join([f"{user regex}.\*?" for i in range(0,3)]),3), f.regexp extract( f.col("Full\_questions"), "".join([f"{user\_regex}.\*?" for i in range(0,4)]),4), f.regexp extract( f.col("Full questions"), "".join([f"{user regex}.\*?" for i in range(0,5)]),5), f.regexp extract( f.col("Full questions"), "".join([f"{user regex}.\*?" for i in range(0,6)]),6),), "",).a lias('user mentioned')) clean test=clean test.withColumn( 'original text',f.col("Full questions")).withColumn('Full question ns',f.regexp replace(f.col("Full questions"), user regex,"").alias("Full questions")) url regex=r"((https?|ftp|file):\/ $\{2,3\}$ )+([-\w+&@#/\%-~|\$?!:,.]\*)|(www.)+([-\w+&@#/\%-~|\$?!:,.]\*)" clean test=clean test.withColumn( 'Full questions', f.regexp replace(f.col("Full questions"), url re #clean html clean test = clean test.withColumn('Full questions', f.regexp replace(f.col("Full questions"), r'<.</pre> \*?>', '')) email regex=r" $[\w.-]+@[\w.-]+\.[a-zA-Z]{1,}$ " clean test=clean test.withColumn('Full questions',f.regexp replace(f.col("Full questions"), email r egex,"")) test\_cleaned=(clean\_test.withColumn( 'Full\_questions', f.regexp\_replace(f.col("Full\_questions"),"[^a -zA-Z]"," ")).withColumn( 'Full questions',f.regexp replace(f.col("Full questions")," +"," ")).withColu mn( 'Full questions',f.trim(f.col("Full questions"))).filter("Full questions !=''")) test data=test cleaned.select('questions id','Full questions','questions author id').coalesce(3).ca che() return test data In [10]: #a=f.concat(f.col("questions title"),f.lit(' '),f.col("questions body")) #ques=ques.withColumn("Full\_questions",a).alias("Full") In [11]: | ques=ques.withColumnRenamed("questions body", "Full questions") ques.show() questions id| questions author id|questions date added| questions title| Full quest ions |332a511f156944448...|8f6f374ffd834d258...|2016-04-26 11:14:...|Teacher career ...|What is a mat |eb80205482e4424ca...|acccbda28edd4362a...|2016-05-20 16:48:...|I want to become ...|I am Priyanka fr |4ec31632938a40b98...|f2c179a563024ccc9...|2017-02-08 19:13:...|Will going abroad...|I'm planning on |2f6a9a99d9b24e5ba...|2c30ffba444e40eab...|2017-09-01 14:05:...|To become a speci...|i hear business |5af8880460c141dbb...|aa9eb1a2ab184ebbb...|2017-09-01 02:36:...|Are there any sch...|I'm trying to fi |7c336403258f4da3a...|d1e4587c0e784c62b...|2017-03-01 04:27:...|How many years of...|To be an enginee |be3c5edfdb07423e9...|71b4554d4a824253a...|2017-09-01 04:59:...|I want to become ...|I am a musician |0f1d6a4f276c4a058...|585ac233015447cc9...|2016-05-19 22:16:...|what kind of col...|I like soccer be |d4999cdc470049a1a...|654e1d6fd5b947249...|2017-08-31 19:20:...|What are the coll...|I'm asking becau |e214acfbe6644d65b...|16908136951a48ed9...|2012-09-09 05:33:...|what does it take...|I am a sophomore |6351c23f2e144b359...|4627ecfaf1814d2e8...|2018-08-14 20:10:...|What major is bes...|#majors #occupat |a978f83b826c4fba8...|47be51b6231e4bbf8...|2016-10-24 17:25:...|How does the gend...|As a woman of co |63171d8c2d3f4a8c9...|240a5e4e4acd469b9...|2016-10-25 21:09:...|As a civil engine...|I would love to |fecd4c7f68144042a...|3acd97cf60da4b239...|2018-01-15 22:03:...|what is the harde...|Throughout the p |baa937b4cd184a22a...|43411aa04f6041b0a...|2018-03-12 19:38:...|How to have a boa...|I've been accept |eb0027b3dcd04d88b...|aa425c019c7b49af9...|2017-03-10 01:05:...|To be a dental hy...|I'm a sophomore |be0bac33e0594b21a...|8801c375ce6a439a9...|2018-07-09 00:01:...|Is it better to s...|I have a lot of |5cf8f5929e1c4357a...|7b7ad31c59f74b3fb...|2016-05-13 01:49:...|Who/what do you t...|Most people clai m... |f156171f9c4944de9...|454f348c03a14cb3a...|2018-01-17 17:26:...|How can I get my ...|I have developed |fb615cf6770d4b7b9...|70a65013686d490eb...|2016-03-01 18:52:...|Which topics in c...|Hi, I am a high only showing top 20 rows In [12]: | ques = datacleaning que(ques) ques.show(20) questions id| Full questions | questions author id | +----+ |332a511f156944448...|What is a maths t...|8f6f374ffd834d258...| |eb80205482e4424ca...|I am Priyanka fro...|acccbda28edd4362a...| |4ec31632938a40b98...|I m planning on g...|f2c179a563024ccc9...| |2f6a9a99d9b24e5ba...|i hear business m...|2c30ffba444e40eab...| |5af8880460c141dbb...|I m trying to fin...|aa9eb1a2ab184ebbb...| |7c336403258f4da3a...|To be an engineer...|d1e4587c0e784c62b...| |be3c5edfdb07423e9...|I am a musician a...|71b4554d4a824253a...| |Of1d6a4f276c4a058...|I like soccer bec...|585ac233015447cc9...| |d4999cdc470049a1a...|I m asking becaus...|654e1d6fd5b947249...| |e214acfbe6644d65b...|I am a sophomore ...|16908136951a48ed9...| |6351c23f2e144b359...|majors occupation...|4627ecfaf1814d2e8...| |a978f83b826c4fba8...|As a woman of col...|47be51b6231e4bbf8...| |63171d8c2d3f4a8c9...|I would love to b...|240a5e4e4acd469b9...| |fecd4c7f68144042a...|Throughout the pr...|3acd97cf60da4b239...| |baa937b4cd184a22a...|I ve been accepte...|43411aa04f6041b0a...| |eb0027b3dcd04d88b...|I m a sophomore i...|aa425c019c7b49af9...| |be0bac33e0594b21a...|I have a lot of i...|8801c375ce6a439a9...| |5cf8f5929e1c4357a...|Most people claim...|7b7ad31c59f74b3fb...| |f156171f9c4944de9...|I have developed ...|454f348c03a14cb3a...| |fb615cf6770d4b7b9...|Hi I am a high sc...|70a65013686d490eb...| +----+ only showing top 20 rows In [13]: | ques ans final = ques.dropDuplicates().select('questions id', 'Full questions') ques ans final=ques ans final.na.drop(subset=["Full questions"]) In [14]: from pyspark.sql.functions import lit,row number,col from pyspark.sql.window import Window w = Window().partitionBy(lit('a')).orderBy(lit('a')) ques ans final = ques ans final.withColumn("row num", row number().over(w)) In [15]: # Build the pipeline tokenizer = RegexTokenizer(gaps = False, pattern = '\w+',inputCol = 'Full questions', outputCol = 'toke stopWordsRemover = StopWordsRemover(inputCol = 'token', outputCol = 'nostopwrd') countVectorizer = CountVectorizer(inputCol="nostopwrd", outputCol="rawFeature") iDF = IDF(inputCol="rawFeature", outputCol="idf vec") word2Vec = Word2Vec(vectorSize = 100, minCount = 5, inputCol = 'nostopwrd', outputCol = 'word vec', see vectorAssembler = VectorAssembler(inputCols=['idf vec', 'word vec'], outputCol='comb vec') pipeline = Pipeline(stages=[tokenizer, stopWordsRemover, countVectorizer, iDF, word2Vec, vectorAssemble r]) # fit the model pipeline mdl = pipeline.fit(ques ans final) In [16]: ques\_ans\_final.show(5) +----+ questions id| Full questions|row num| +----+ |16bae94c380446b89...|Mainly asking thi...| |cf0fcf3eea9849a59...|I ts really rare ...| 3 | |e16b6cde5e3f4e5d8...|How do you become...| |a30fd1fe06f74fd4b...|I want to be able...| 4 | |4f487e3644c8441c8...|I want to find a ...| +----+ only showing top 5 rows In [17]: # transform the question data ques pipeline df = pipeline mdl.transform(ques ans final) In [18]: def CosineSim(vec1, vec2): return np.dot(vec1, vec2) / np.sqrt(np.dot(vec1, vec1)) / np.sqrt(np.dot(vec2, vec2)) In [19]: all ques vecs = ques pipeline df.select('questions id', 'word vec').rdd.map(lambda x: (x[0], x[1])).col lect() In [20]: def getQuestionDetails(in ques): a = in\_ques.alias("a") b = ques.alias("b") return a.join(b, col("a.question id") == col("b.questions id"), 'inner') \ .select([col('a.'+xx) for xx in a.columns] + [col('b.Full questions')]) In [21]: def getKeyWordsRecoms(key words, sim bus limit): #print('\nQuestion Relevent "' + key words + '"') input words df = sc.parallelize([(0, key words)]).toDF(['question id', 'Full questions']) # transform the the key words to vectors input words df = pipeline mdl.transform(input words df) # choose word2vec vectors input key words vec = input words df.select('word vec').collect()[0][0] # get similarity sim\_bus\_byword\_rdd = sc.parallelize((i[0], float(CosineSim(input\_key\_words\_vec, i[1]))) for i in al l\_ques\_vecs) sim bus byword df = spark.createDataFrame(sim bus byword rdd) \ .withColumnRenamed('\_1', 'question\_id') \ .withColumnRenamed(' 2', 'score') \ .orderBy("score", ascending = False) sim\_bus\_byword\_df=sim\_bus\_byword\_df.na.drop(subset=["score"]) # return top 10 similar a = sim bus byword df.limit(sim bus limit) return getQuestionDetails(a).sort('score', ascending=False) In [22]: ans score=ans score.withColumnRenamed("score", "vote") In [23]: **def** acc 1ques(id ques): key words = str(np.array(ques.filter(ques.questions id==id ques).select('Full questions').collect ())[0][0]) keywords recom df = getKeyWordsRecoms(key words, 20) list truth=np.array(ans.filter(ans.answers question id==id ques).select('answers author id').collec t()) rs1 = keywords recom df.join(ans, keywords recom df.question id==ans.answers question id).select("qu estion id", "answers id", "answers author id", "score") # top vote 20 rs2 = rs1.join(ans score,rs1.answers id==ans score.id).select("answers author id","question id","vo rs2 = rs2.withColumn("ts", rs2.vote\*rs2.score) rs2 = rs2.sort('ts', ascending=False).limit(20) # new top 20 : 0.9807 #rs2 = rs1.sort('score', ascending=False).limit(20) # new top 20 vote >0: 0.4219 #rs2 = rs1.join(ans score,rs1.answers id==ans score.id).select("answers author id","question id","s core", "vote") #rs2=rs2.filter(rs2.vote>0) #rs2 = rs2.sort('score', ascending=False).limit(20) # Remove 1.0 score #rs1 = rs1.filter(rs1.question id!=id ques) #rs2 = rs1.sort('score', ascending=False).limit(20) rs3 = rs2.join(pro,rs2.answers\_author\_id==pro.professionals\_id).select("professionals\_id","professi onals industry", "professionals headline") list\_recom=np.array(rs3.select('professionals\_id').dropDuplicates().collect()) sum count=ans.filter(ans.answers question id==id ques).count() for i in list recom: for j in list truth: if i==j: count=count+1 count=float(count/sum count) print(count) return count In [24]: **def** score(test): sum=0for t in test: sum=sum+acc 1ques(t) score = float(sum/len(test)) return score In [25]: test=['01352c4d67fe435ca59e745ff2520d2a', '03eee1ca07174470b160717027ab46d6', '04a979f4e7fd49b9a07b6fae7a5727ee', '062f49f153de4b8793e4e669ec5b5331', '083965c88d894a9f9e4e71e521641338', '09e3bdc69a6149aa8656bbc18162ac37', '0d7fab391dc145a384da4af0a078b77f', '0db6ed5d24df42f18d19958ccb32cd6e', 'la039cb9f3064f76b386f84f303edc43', 'la444e5e5824446eaf37f31effd72ce0'] In [26]: a = score(test)/opt/conda/lib/python3.7/site-packages/ipykernel launcher.py:2: RuntimeWarning: invalid value encount ered in double scalars 1.0 1.0 0.11111111111111111 0.181818181818182 0.16666666666666666 0.23529411764705882 0.13333333333333333 0.7 0.7692307692307693 0.15384615384615385 Out[26]: 0.44513003336532747