In [1]: pip install pyspark Collecting pyspark Downloading pyspark-3.1.2.tar.gz (212.4 MB) | 212.4 MB 64 kB/s s eta 0:00:01 | 36.8 MB 13.2 MB/s eta 0:00:14 41.5 MB 13.2 MB/s eta 0: | 46.0 MB 41.7 MB/s eta 0:00:04 00:13 | 61.2 MB 41.7 MB/s e | 47.8 MB 41.7 MB/s eta 0:00:04 | 68.8 MB 41.7 MB/s eta 0:00:04 ta 0:00:04 Collecting py4j == 0.10.9Downloading py4j-0.10.9-py2.py3-none-any.whl (198 kB) | 198 kB 45.0 MB/s eta 0:00:01 Building wheels for collected packages: pyspark Building wheel for pyspark (setup.py) ... done Created wheel for pyspark: filename=pyspark-3.1.2-py2.py3-none-any.whl size=212880768 sha256=b3aa16 0fa72c372b2a090e940fc974b5a7a135712c2d3f317cd4bfa371f1e4ca 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 import pandas as pd 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,reqexp replace,lit from pyspark.ml.evaluation import RegressionEvaluator from pyspark.ml.recommendation import ALS from pyspark.ml.evaluation import RegressionEvaluator, MulticlassClassificationEvaluator 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 reqexp replace from pyspark.sql import SparkSession spark = SparkSession.builder.master("local[\*]").config("spark.executor.memory", "70g").config("spark.dr iver.memory", "50g").config("spark.memory.offHeap.enabled", True).config("spark.memory.offHeap.size", "20 g").appName("sampleCodeForReference").getOrCreate() sc = spark.sparkContext In [3]: import numpy as np import pandas as pd from pyspark.sql import functions as f 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]: from pyspark.sql import SQLContext sqlContext = SQLContext(sc) In [5]: #Data cleaning question def datacleaning ques(test): col name='questions body' user regex=r"( $@\{1,15\}$ )" clean test=test.withColumn( 'user mentioned', f.array remove( f.regexp\_extract(f.col(col\_name), user\_regex, 1), f.regexp extract( f.col(col\_name), "".join([f"{user\_regex}.\*?" for i in range(0,2)]),2), f.regexp extract( f.col(col name), "".join([f"{user regex}.\*?" for i in range(0,3)]),3), f.regexp extract( f.col(col name), "".join([f"{user regex}.\*?" for i in range(0,4)]),4), f.regexp extract( f.col(col\_name), "".join([f"{user\_regex}.\*?" for i in range(0,5)]),5), f.regexp extract( f.col(col\_name), "".join([f"{user\_regex}.\*?" for i in range(0,6)]),6),), "",).alias('user\_mentione) **d'**)) clean test=clean test.withColumn( 'original text', f.col(col name)).withColumn(col name, f.regexp rep lace(f.col(col name), user regex,"").alias(col name)) # remove hastag hashtag\_user\_regex="#(\w{1,})" clean\_test=clean\_test.withColumn( 'hashtags',f.array\_remove( f.regexp extract(f.col(col name), hashtag user regex, 1), f.regexp extract( f.col(col name), "".join([f"{hashtag user regex}.\*?" for i in range(0,2)]),2), f.regexp extract( f.col(col name), "".join([f"{hashtag user regex}.\*?" for i in range(0,3)]),3), f.regexp\_extract( f.col(col\_name), "".join([f"{hashtag\_user\_regex}.\*?" for i in range(0,4)]),4), f.regexp extract( f.col(col\_name), "".join([f"{hashtag\_user\_regex}.\*?" for i in range(0,5)]),5), f.regexp\_extract( f.col(col name), "".join([f"{hashtag user regex}.\*?" for i in range(0,6)]),6),), "",).alias('hashta clean\_test=clean\_test.withColumn(col\_name, f.regexp\_replace(f.col(col\_name), hashtag\_user\_regex, "\$1" ).alias(col name)) url regex=r"((https?|ftp|file):\/ $\{2,3\}$ )+([-\w+&@#/%-~|\$?!:,.]\*)|(www.)+([-\w+&@#/%-~|\$?!:,.]\*)" clean test=clean test.withColumn( col name, f.regexp replace(f.col(col name), url regex,"")) email\_regex=r"[\w.-]+\.[a-zA-Z]{1,}" clean test=clean test.withColumn(col name, f.regexp replace(f.col(col name), email regex,"")) #clean test=clean test.withColumn('text', html unescape("text")) test cleaned=(clean test.withColumn( col name, f.regexp replace(f.col(col name), "[^a-zA-Z]", " ")).wi thColumn(col name, f.regexp replace(f.col(col name), " +", " ")).withColumn(col name, f.trim(f.col(col name) e))).filter(col(col name) !='')) test data=test cleaned.select('professionals id',col name).coalesce(2).cache() return test data In [6]: pro=spark.read.csv('../input/data-science-for-good-careervillage/professionals.csv', header=True, quote= ""', sep=", ", multiLine=True) ques=spark.read.csv('../input/data-science-for-good-careervillage/questions.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, inferSchema=True) **Collaborative** In [7]: from pyspark.sql.functions import lit,row\_number,col from pyspark.sql.window import Window w = Window().partitionBy(lit('a')).orderBy(lit('a')) ques = ques.withColumn("ques\_id", row\_number().over(w)) pro = pro.withColumn("pro id", row number().over(w)) In [8]: ans\_score\_new=ans\_score.join(ans, ans\_score.id==ans.answers\_id,'left').select('answers\_author\_id','answ ers\_question\_id','score') ans\_score\_new=ans\_score\_new.join(pro,ans\_score\_new.answers\_author\_id==pro.professionals\_id,'left').sele ct('pro\_id', 'answers\_question\_id', 'score') ans\_score\_new=ans\_score\_new.join(ques,ans\_score\_new.answers\_question\_id==ques.questions\_id,'left').sele ct('pro\_id','ques\_id','score') pro ques collar=ans score new pro\_ques\_collar=pro\_ques\_collar.na.drop("any") In [9]: (training, test) = pro ques collar.randomSplit([0.8, 0.2]) In [11]: als = ALS(rank=40, maxIter=15, regParam=0.01, coldStartStrategy="drop", implicitPrefs=False, userCol="pro i d", itemCol="ques\_id", ratingCol="score") In [12]: model = als.fit(pro\_ques\_collar) In [13]: | predictions = model.transform(pro\_ques\_collar) In [14]: def recommendations\_for\_pro(pro\_id): print('\nInfo of Professional: ') pro.filter(pro.pro\_id==pro\_id).drop('pro\_id').show() df=userRecs.filter(userRecs.pro id==pro id) x=df.select('recommendations').collect() df = sc.parallelize(x[0][0]).toDF(['id ques','value']) df=df.join(ques,ques\_id==df.id\_ques).select('questions\_id','questions\_title','questions\_body') print('\nRecommendations: ') return df In [15]: **def** add recommendations ques(ques id): proRecs=model.recommendForAllItems(10) id=np.array(ques.filter(ques.questions\_id==ques\_id).select('ques\_id').collect())[0][0] df=proRecs.filter(proRecs.ques\_id==int(id)) x=df.select('recommendations').collect() df = sc.parallelize(x[0][0]).toDF(['id\_pro','value']) df=df.join(pro,pro.pro\_id==df.id\_pro).select('professionals\_id','professionals\_location','profes nals\_industry','professionals\_headline','professionals\_date\_joined') return df In [29]: def recommendations\_for\_ques(ques\_id,top): id=np.array(ques.filter(ques.questions id==ques id).select('ques\_id').collect())[0][0] df=predictions.filter(predictions.ques\_id==int(id)).orderBy("prediction", ascending = False).select ('pro\_id','prediction') df=df.filter(df.prediction>=0) df=df.join(pro,pro.pro\_id==df.pro\_id).select('professionals\_id','professionals\_location','profes nals industry','professionals headline','professionals date joined').distinct() if len(np.array(df.select('professionals\_id').collect()))<15:</pre> df\_add=add\_recommendations\_ques(ques\_id) df=df.union(df add) df=df.limit(top) return df def acc 1ques collar(id ques): keywords\_recom\_df = recommendations\_for\_ques(id\_ques,20) list\_truth=np.array(ans.filter(ans.answers\_question\_id==id\_ques).select('answers\_author\_id').collec t()) list\_recom=np.array(keywords\_recom\_df.select('professionals\_id').collect()) count=0 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 break count=float(count/sum count) return count Content In [18]: pro ans=pro.join(ans,pro.professionals id==ans.answers author id).select('professionals id', 'answers qu estion id', 'answers id') pro\_ans=pro\_ans.join(ans\_score,pro\_ans.answers\_id==ans\_score.id).select('professionals\_id','answers\_que stion\_id','score') pro ans=pro ans.join(ques,pro ans.answers question id==ques.questions id).select('professionals id','qu estions body','questions title','score') pro\_ans=pro\_ans.withColumn('combined\_ques',concat(pro\_ans.questions\_title,lit(' '),pro\_ans.questions\_bo pro\_ques\_content=pro\_ans.groupBy('professionals\_id').agg(f.collect\_list('combined\_ques').alias('questio ns\_body')) In [19]: from pyspark.sql.functions import udf, col join udf = udf(lambda x: ",".join(x)) pro ques content=pro ques content.withColumn("questions body", join udf(col("questions body"))) In [20]: pro ques content=datacleaning ques(pro ques content) pro ques content=pro ques content.na.drop(subset=["questions body"]) In [21]: # Build the pipeline tokenizer = RegexTokenizer(gaps = False, pattern = '\w+',inputCol = 'questions\_body', outputCol = 'toke stopWordsRemover = StopWordsRemover(inputCol = 'token', outputCol = 'nostopwrd') word2Vec = Word2Vec(vectorSize = 100, minCount = 5, inputCol = 'nostopwrd', outputCol = 'word\_vec', see pipeline = Pipeline(stages=[tokenizer, stopWordsRemover, word2Vec]) # fit the model pipeline mdl = pipeline.fit(pro ques content) In [22]: # transform the question data ques pipeline df = pipeline mdl.transform(pro ques content) In [23]: def CosineSim(vec1, vec2): return np.dot(vec1, vec2) / np.sqrt(np.dot(vec1, vec1)) / np.sqrt(np.dot(vec2, vec2)) In [24]: all ques vecs = ques pipeline df.select('professionals id', 'word vec').rdd.map(lambda x: (x[0], x[1])) .collect() In [25]: def getQuestionDetails(in ques): a = in\_ques.alias("a") b = pro.alias("b") return a.join(b, col("a.professionals id") == col("b.professionals id"), 'inner') \ .select([col('a.score')]+[col('b.'+xx) for xx in b.columns] ) In [26]: def getKeyWordsRecoms(key words, sim bus limit): input\_words df = sc.parallelize([(0, key\_words)]).toDF(['professionals\_id', 'questions\_body']) # 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 1 ques vecs) sim bus byword df = spark.createDataFrame(sim bus byword rdd) \ .withColumnRenamed(' 1', 'professionals 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).select('professionals id', 'professionals location', 'professionals indu stry','professionals headline','professionals date joined') In [301: def acc 1ques final(id\_ques): key words = str(np.array(ques.filter(ques.questions id==id ques).select('questions body').collect ())[0][0]) recom\_df1 = getKeyWordsRecoms(key words, 5) recom df2 = recommendations for ques(id ques, 15) keywords recom df=recom df1.union(recom df2).distinct() list truth=np.array(ans.filter(ans.answers question id==id ques).select('answers author id').collec t()) list recom=np.array(keywords recom df.select('professionals id').collect()) count=0 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 break count=float(count/sum count) return count test=['01352c4d67fe435ca59e745ff2520d2a', In [31]: '03eee1ca07174470b160717027ab46d6', '04a979f4e7fd49b9a07b6fae7a5727ee', '062f49f153de4b8793e4e669ec5b5331', '083965c88d894a9f9e4e71e521641338', '09e3bdc69a6149aa8656bbc18162ac37', '0d7fab391dc145a384da4af0a078b77f', '0db6ed5d24df42f18d19958ccb32cd6e', 'la039cb9f3064f76b386f84f303edc43', '1a444e5e5824446eaf37f31effd72ce0'] sum=0 for t in test: x=acc\_1ques final(t) sum=sum+x print(x) score=float(sum/len(test)) score 1.0 1.0 0.9090909090909091 0.6666666666666666 0.5882352941176471 0.75 1.0 0.5769230769230769 Out[31]: 0.8179804835687188

In [ ]: