Komal Suryan: NUID 002747707

Big Data System Engineering with Scala Spring 2023 Assignment No. 2 (Spark)



TASKS:

Url: https://www.kaggle.com/competitions/titanic/dataLinks to an external site.

For this assignment you will use the training and testing datasets.(train.csv & test.csv)

You are to load the dataset using Spark and perform the operations below:

Exploratory Data Analysis- Follow up on the previous spark assignment 1 and explain a few statistics. (20 pts)

Feature Engineering - Create new attributes that may be derived from the existing attributes. This may include removing certain columns in the dataset. (30 pts)

Prediction - Use the train.csv to train a Machine Learning model of your choice & test it on the test.csv. You are required to predict if the records in test.csv survived or not. Note(1 = Survived, 0 = Dead) (50 pts)

Please note: Do not include the test.csv while training the model. Also do not use the 'Survived' column during training. Doing these would defeat the purpose of the entire model.

The classifier must have an accuracy of 70 % for 100 pts.

The submission must be on a new GitHub repository or data bricks public notebook along with the Assignment pdf. Please provide Screenshots of the results in the Assignment.pdf for submission.

Screenshots

```
↑ ↓ © 目 $ ♬ î :
import org.apache.spark.sql.SparkSession
       val spark = SparkSession
         .builder()
          .appName("spark2")
          .getOrCreate()
      Intitializing Scala interpreter ...

Spark Web UI available at <a href="http://lo.110.48.55:4040">http://lo.110.48.55:4040</a>

SparkContext available as 'sc' (version = 3.3.2, master = local[*], app id = local-1680803218993)

SparkSession available as 'spark'
23/04/06 13:47:01 WARN SparkSession: Using an existing Spark session; only runtime SQL configurations will take effect.
      import org.apache.spark.sql.SparkSession
spark: org.apache.spark.sql.SparkSession = org.apache.spark.sql.SparkSession@5f30e7a1
val sc = spark.sparkContext
 sc: org.apache.spark.SparkContext = org.apache.spark.SparkContext@46e24902
[ ] val test = spark.read.option("header", "true").csv("test.csv")
val train = spark.read.option("header", "true").csv("train.csv")
      test: org.apache.spark.sql.DataFrame = [PassengerId: string, Pclass: string ... 9 more fields]
train: org.apache.spark.sql.DataFrame = [PassengerId: string, Survived: string ... 10 more fields]
EDA
[ ] train.show(5)
       |PassengerId|Survived|Pclass|
                                                                                       Sex|Age|SibSp|Parch|
                                                                                                                                                    Fare|Cabin|Embarked|
                      11
                                    01
                                             3|Braund. Mr. Owen ...| male| 22| 1| 0|
                                                                                                                                A/5 21171| 7.25| null|
```

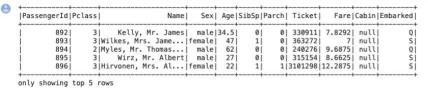
▼ EDA

[] train.show(5)

Passe	ngerId Sur	vived Pc	lass	Name	Sex	Age	SibSp	Parch		Ticket	Fare	Cabin E	Embarked
	1	0	3 Braund, Mr.	0wen	male	22	1	0	A/5	21171	7.25	null	S
İ	2	1	1 Cumings, Mrs	. Joh	female	38	1	0	PC	17599	71.2833	C85	C
į.	3	1	3 Heikkinen, M	liss	female	26	0	0	STON/02. 3	101282	7.925	null	S
ĺ	4 j	1	1 Futrelle, Mr	s. Ja	female	35	1	0		113803	53.1	C123	S
į.	5	0	3 Allen, Mr. W	/illia	male	35	0	0		373450	8.05	null	S

only showing top 5 rows

test.show(5)



[] import org.apache.spark.sql.functions.count

```
[ ] import org.apache.spark.sql.functions.count
     train.groupBy("Sex", "Survived")
          .agg(count("Survived").alias("count"))
          .show()
         Sex|Survived|count|
                         81|
      female|
                    1|
                       233 |
109 |
468 |
      female
        malel
     import org.apache.spark.sql.functions.count
     train.groupBy("Pclass", "Survived")
          .agg(count("Survived").alias("count"))
.show()
     |Pclass|Survived|count|
                        136|
                       97|
                    0|
1|
0|
1|
0|
           2|
3|
3|
2|
                                                                                                                                                    T madil-44
                       372
                                                                                                                                                    Testattutatt
                         801
[ ]
     train.groupBy("Pclass", "Sex")
           .agg(count("Sex").alias("count"))
           .show()
     |Pclass|
                  Sex | count |
            2|female|
                          76
            2i malei
                         108
               male
                         347
            1| male|
                         122
            1|female|
                          94
            3|female|
                         144
Feature Engineering
```

```
[ ] val test_new = test
    .withColumn("Age", col("Age").cast("Double"))
    .withColumn("Fare", col("Fare").cast("Double"))
    .withColumn("Pclass", col("Pclass").cast("Integer"))
    .withColumn("PassengerId", col("PassengerId").cast("Integer"))

test_new: org.apache.spark.sql.DataFrame = [PassengerId: int, Pclass: int ... 9 more fields]
```

```
[ ] val train_new = train
    .withColumn("Age", col("Age").cast("Double"))
    .withColumn("Fare", col("Fare").cast("Double"))
    .withColumn("Pclass", col("Pclass").cast("Integer"))
    .withColumn("Survived", col("Survived").cast("Integer"))
    .withColumn("PassengerId", col("PassengerId").cast("Integer"))
```

Prediction

```
[ ] val women = train.filter(col("Sex") === "female").select("Survived")
     val rate_women = women.agg(avg("Survived")).as[Double].head()
println("% of women survived: " + rate_women)
     % of women survived: 0.7420382165605095
    women: org.apache.spark.sql.DataFrame = [Survived: string] rate_women: Double = 0.7420382165605095
val women = train.filter(col("Sex") === "male").select("Survived")
     val rate_men = women.agg(avg("Survived")).as[Double].head()
println("% of men survived: " + rate_men)
% of men survived: 0.18890814558058924
women: org.apache.spark.sql.DataFrame = [Survived: string]
rate_men: Double = 0.18890814558058924
[ ] import org.apache.spark.ml.Pipeline
     import org.apache.spark.ml.feature.{StringIndexer, OneHotEncoder, VectorAssembler}
     import org apache spark ml classification LogisticRegression
     val sexIndexer = new StringIndexer()
       .setInputCol("Sex")
        .setOutputCol("SexIndex")
     val embarkedIndexer = new StringIndexer()
       .setInputCol("Embarked")
       setOutputCol("EmbarkedIndex")
     val sexEncoder = new OneHotEncoder()
setInputCol("SexIndex")
```

```
import org.apache.spark.ml.Pipeline
 import org.apache.spark.ml.feature.{StringIndexer, OneHotEncoder, VectorAssembler}
 import org.apache.spark.ml.classification.LogisticRegression
val sexIndexer = new StringIndexer()
   .setInputCol("Sex")
   .setOutputCol("SexIndex")
val embarkedIndexer = new StringIndexer()
   .setInputCol("Embarked")
   .setOutputCol("EmbarkedIndex")
val sexEncoder = new OneHotEncoder()
   .setInputCol("SexIndex")
   .setOutputCol("SexVec")
val embarkedEncoder = new OneHotEncoder()
   .setInputCol("EmbarkedIndex")
   .setOutputCol("EmbarkedVec")
val assembler = new VectorAssembler()
   .setInputCols(Array("Pclass", "Age", "Fare"))
   .setOutputCol("features")
val lr = new LogisticRegression()
  .setLabelCol("Survived")
  .setFeaturesCol("features")
   .setMaxIter(10)
val pipeline = new Pipeline()
   .setStages(Array(sexIndexer, embarkedIndexer, sexEncoder, embarkedEncoder, assembler, lr))
val pipelineModel = pipeline.fit(train_final)
val predictions = pipelineModel.transform(test_final)
```

[] 23/04/06 14:48:22 WARN InstanceBuilder\$JavaBLAS: Failed to load implementation from:dev.ludovic.netlib.blas.VectorBLAS
import org.apache.spark.ml.Pipeline
import org.apache.spark.ml.feature.{StringIndexer, OneHotEncoder, VectorAssembler}
import org.apache.spark.ml.feature.StringIndexer = strIdx_1le1809eeda4
embarkedIndexer: org.apache.spark.ml.feature.StringIndexer = strIdx_17316ab1fba1
sexEncoder: org.apache.spark.ml.feature.OneHotEncoder = oneHotEncoder_34b6c30464cd
embarkedEncoder: org.apache.spark.ml.feature.OneHotEncoder = oneHotEncoder_38baf25c479b
assembler: org.apache.spark.ml.feature.VectorAssembler = VectorAssembler: uid=vecAssembler_8e210e4ce38e, handleInvalid=error, numInputCols=3
lr: org.apache.spark.ml.feature.VectorAssembler = logreg_a35e622db1da
pipeline: org.apache.spark.ml.Pipeline = pipeline_7c06efe6a676
pipel...

predictions.show()

0

SexVe	EmbarkedIndex	SexIndex	Embarked	Cabin	Fare	Ticket	Parch	SibSp	Age	Sex	Name	assengerId Pclass
[(1,[0],[1.0]	2.0	0.0	Q	null	7.8292	330911	0	0	34.5	male	Kelly, Mr. James	892 3
[(1,[],[]	0.0	1.0	S	null	7.0	363272	0	1	47.0	female	Wilkes, Mrs. Jame	893 3
[(1,[0],[1.0]	2.0	0.0	Q	null	9.6875	240276			62.0	male	Myles, Mr. Thomas	894 2
[(1,[0],[1.0]	0.0	0.0	S	null	8.6625	315154	0	0	27.0	male	Wirz, Mr. Albert	895 3
[(1,[],[]	0.0	1.0	S	null	12.2875	3101298	1		22.0	female	Hirvonen, Mrs. Al	896 3
[(1,[0],[1.0]	0.0	0.0	S	null	9.225	7538	0	0	14.0	male	Svensson, Mr. Joh	897 3
[(1,[],[]	2.0	1.0	Q	null	7.6292	330972	0	0	30.0	female	Connolly, Miss. Kate	898 3
[(1,[0],[1.0]	0.0	0.0	S	null	29.0	248738	1	1	26.0	male	Caldwell, Mr. Alb	899 2
[(1,[],[]	1.0	1.0	C	null	7.2292	2657	0	0	18.0	female	Abrahim, Mrs. Jos	900 3
[(1,[0],[1.0]	0.0	0.0	S	null	24.15	A/4 48871	0	2	21.0	male	Davies, Mr. John	901 3
[(1,[0],[1.0]	0.0	0.0	S	null	7.8958	349220	0	0	30.272590361445783	male	Ilieff, Mr. Ylio	902 3
[(1,[0],[1.0]	0.0	0.0	S	null	26.0	694	0	0	46.0	male	Jones, Mr. Charle	903 1
(1.[].[]	0.0	1.0	S	B45	82.2667	21228	0	1	23.0	female	Snyder, Mrs. John	904 1
	0 ==	0.0	S	null	26.0	24065	0	1	63.0	male	Howard, Mr. Benjamin	905 2
HILL Commerce	0	1.0	S	E31	61.175	W.E.P. 5734	0	1	47.0	female	Chaffee, Mrs. Her	906 1
	1	1.0	C	null	27.7208	SC/PARIS 2167	0	1	24.0	female	del Carlo, Mrs. S	907 2
	2	0.0	Q	null	12.35	233734	0	0	35.0	male	Keane, Mr. Daniel	908 2
	1	0.0	C	null	7.225	2692	0	0	21.0	male	Assaf, Mr. Gerios	909 3
(1,1,1)	0.0	1.0	S	null	7.925	STON/02. 3101270	0	1	27.0	female	Ilmakangas, Miss	910 3
[(1,[],[]	1.0	1.0	C	null	7.225	2696	0	0	45.0	female	"Assaf Khalil, Mr	911 3