```python
from pyspark.sql import SparkSession
from pyspark.ml.feature import StringIndexer, VectorAssembler
from pyspark.ml.classification import LogisticRegression
from pyspark.ml.evaluation import BinaryClassificationEvaluator

spark = SparkSession.builder.getOrCreate()

marksDF = spark.read.csv("teach_scores_1.csv",header = True, inferSchema = True)

marksDF.printSchema()

"""
oot
 |-- subject_1_gp: double (nullable = true)
 |-- subject_2_gp: double (nullable = true)
 |-- subject_3_gp: double (nullable = true)
 |-- subject_4_gp: double (nullable = true)
 |-- subject_5_gp: double (nullable = true)
 |-- grade: string (nullable = true)
"""
```

```
root
 |-- subject_1_gp: double (nullable = true)
 |-- subject_2_gp: double (nullable = true)
 |-- subject_3_gp: double (nullable = true)
 |-- subject_4_gp: double (nullable = true)
 |-- subject_5_gp: double (nullable = true)
 |-- grade: string (nullable = true)
```

```
'\noot\n |-- subject_1_gp: double (nullable = true)\n |-- subject_2_gp: double (nullable = true)\n |-- subject_3_gp: double (nullabl
```

```python
marksDF.show(10,False)

"""
+------------+------------+------------+------------+------------+-----+
|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_gp|subject_5_gp|grade|
+------------+------------+------------+------------+------------+-----+
|2.0         |2.1         |3.5         |2.4         |3.0         |F    |
|2.0         |2.0         |2.0         |3.0         |3.0         |F    |
|2.1         |2.0         |2.4         |3.5         |3.0         |F    |
|9.0         |9.0         |9.0         |9.0         |9.0         |A+   |
|8.0         |8.0         |8.0         |8.0         |8.0         |A    |
|7.0         |7.0         |7.0         |7.0         |7.0         |B    |
|6.0         |6.0         |6.0         |6.0         |6.0         |C    |
|5.0         |5.0         |5.0         |5.0         |5.0         |D    |
|4.0         |4.0         |4.0         |4.0         |4.0         |E    |
|3.0         |3.0         |3.0         |3.0         |3.0         |F    |
+------------+------------+------------+------------+------------+-----+
"""

marksDF.describe("subject_1_gp").show()

"""
+-------+------------------+
|summary|      subject_1_gp|
+-------+------------------+
|  count|               172|
|   mean| 6.6639534883720930|
| stddev|2.1988786504628766|
|    min|               0.0|
|    max|               9.9|
+-------+------------------+
"""
```

```
+------------+------------+------------+------------+------------+-----+
|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_gp|subject_5_gp|grade|
+------------+------------+------------+------------+------------+-----+
```

```
|2.0        |2.1        |3.5        |2.4        |3.0        |F    |
|2.0        |2.0        |2.0        |3.0        |3.0        |F    |
|2.1        |2.0        |2.4        |3.5        |3.0        |F    |
|9.0        |9.0        |9.0        |9.0        |9.0        |A+   |
|8.0        |8.0        |8.0        |8.0        |8.0        |A    |
|7.0        |7.0        |7.0        |7.0        |7.0        |B    |
|6.0        |6.0        |6.0        |6.0        |6.0        |C    |
|5.0        |5.0        |5.0        |5.0        |5.0        |D    |
|4.0        |4.0        |4.0        |4.0        |4.0        |E    |
|3.0        |3.0        |3.0        |3.0        |3.0        |F    |
+-----------+-----------+-----------+-----------+-----------+-----+
only showing top 10 rows

+-------+------------------+
|summary|       subject_1_gp|
+-------+------------------+
|  count|               172|
|   mean| 6.663953488372093|
| stddev|2.1988786504628766|
|    min|               0.0|
|    max|               9.9|
+-------+------------------+


'\n+-------+------------------+\n|summary|       subject_1_gp|\n+-------+------------------+\n|  count|               172|\n|   mean|
```

```python
marksDF.describe("grade").show()
"""
+-------+-----+
|summary|grade|
+-------+-----+
|  count|  172|
|   mean| null|
| stddev| null|
|    min|    A|
|    max|    F|
+-------+-----+
"""
```

```
+-------+-----+
|summary|grade|
+-------+-----+
|  count|  172|
|   mean| null|
| stddev| null|
|    min|    A|
|    max|    F|
+-------+-----+


'\n+-------+-----+\n|summary|grade|\n+-------+-----+\n|  count|  172|\n|   mean| null|\n| stddev| null|\n|    min|    A|\n|    max|
```

```python
inputCols = ["subject_1_gp","subject_2_gp","subject_3_gp","subject_4_gp","subject_5_gp"]

outputCol = "features"

marksDF_assembler = VectorAssembler(inputCols = inputCols,outputCol = outputCol)

featuresDf = marksDF_assembler.transform(marksDF)

print("featuresDF printSchema")

featuresDf.printSchema()

"""
root
 |-- subject_1_gp: double (nullable = true)
 |-- subject_2_gp: double (nullable = true)
 |-- subject_3_gp: double (nullable = true)
 |-- subject_4_gp: double (nullable = true)
 |-- subject_5_gp: double (nullable = true)
 |-- grade: string (nullable = true)
 |-- features: vector (nullable = true)

"""
```

```
featuresDF printSchema
root
 |-- subject_1_gp: double (nullable = true)
 |-- subject_2_gp: double (nullable = true)
 |-- subject_3_gp: double (nullable = true)
 |-- subject_4_gp: double (nullable = true)
 |-- subject_5_gp: double (nullable = true)
 |-- grade: string (nullable = true)
 |-- features: vector (nullable = true)


'\nroot\n |-- subject_1_gp: double (nullable = true)\n |-- subject_2_gp: double (nullable = true)\n |-- subject_3_gp: double (nullab
```

```python
featuresDf.show(10,False)

print("featureDf show")
```

```
+------------+------------+------------+------------+------------+-----+--------------------+
|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_gp|subject_5_gp|grade|features            |
+------------+------------+------------+------------+------------+-----+--------------------+
|2.0         |2.1         |3.5         |2.4         |3.0         |F    |[2.0,2.1,3.5,2.4,3.0]|
|2.0         |2.0         |2.0         |3.0         |3.0         |F    |[2.0,2.0,2.0,3.0,3.0]|
|2.1         |2.0         |2.4         |3.5         |3.0         |F    |[2.1,2.0,2.4,3.5,3.0]|
|9.0         |9.0         |9.0         |9.0         |9.0         |A+   |[9.0,9.0,9.0,9.0,9.0]|
|8.0         |8.0         |8.0         |8.0         |8.0         |A    |[8.0,8.0,8.0,8.0,8.0]|
|7.0         |7.0         |7.0         |7.0         |7.0         |B    |[7.0,7.0,7.0,7.0,7.0]|
|6.0         |6.0         |6.0         |6.0         |6.0         |C    |[6.0,6.0,6.0,6.0,6.0]|
|5.0         |5.0         |5.0         |5.0         |5.0         |D    |[5.0,5.0,5.0,5.0,5.0]|
|4.0         |4.0         |4.0         |4.0         |4.0         |E    |[4.0,4.0,4.0,4.0,4.0]|
|3.0         |3.0         |3.0         |3.0         |3.0         |F    |[3.0,3.0,3.0,3.0,3.0]|
+------------+------------+------------+------------+------------+-----+--------------------+
only showing top 10 rows

featureDf show
```

```python
grade_indexer = StringIndexer(inputCol = "grade", outputCol = "label")


label_df = grade_indexer.fit(featuresDf).transform(featuresDf)

print("after adding label")

label_df.printSchema()

label_df.createOrReplaceTempView()

"""
root
 |-- subject_1_gp: double (nullable = true)
 |-- subject_2_gp: double (nullable = true)
 |-- subject_3_gp: double (nullable = true)
 |-- subject_4_gp: double (nullable = true)
 |-- subject_5_gp: double (nullable = true)
 |-- grade: string (nullable = true)
 |-- features: vector (nullable = true)
 |-- label: double (nullable = false)

"""
```

```
after adding label
root
 |-- subject_1_gp: double (nullable = true)
 |-- subject_2_gp: double (nullable = true)
 |-- subject_3_gp: double (nullable = true)
 |-- subject_4_gp: double (nullable = true)
 |-- subject_5_gp: double (nullable = true)
 |-- grade: string (nullable = true)
 |-- features: vector (nullable = true)
 |-- label: double (nullable = false)


'\nroot\n |-- subject_1_gp: double (nullable = true)\n |-- subject_2_gp: double (nullable = true)\n |-- subject_3_gp: double (nullab
```

```python
print("label included df")


label_df.show(10,False)

"""
+------------+------------+------------+------------+------------+-----+---------------------+-----+
|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_gp|subject_5_gp|grade|features             |label|
+------------+------------+------------+------------+------------+-----+---------------------+-----+
|2.0         |2.1         |3.5         |2.4         |3.0         |F    |[2.0,2.1,3.5,2.4,3.0]|0.0  |
|2.0         |2.0         |2.0         |3.0         |3.0         |F    |[2.0,2.0,2.0,3.0,3.0]|0.0  |
|2.1         |2.0         |2.4         |3.5         |3.0         |F    |[2.1,2.0,2.4,3.5,3.0]|0.0  |
|9.0         |9.0         |9.0         |9.0         |9.0         |A+   |[9.0,9.0,9.0,9.0,9.0]|6.0  |
|8.0         |8.0         |8.0         |8.0         |8.0         |A    |[8.0,8.0,8.0,8.0,8.0]|1.0  |
|7.0         |7.0         |7.0         |7.0         |7.0         |B    |[7.0,7.0,7.0,7.0,7.0]|2.0  |
|6.0         |6.0         |6.0         |6.0         |6.0         |C    |[6.0,6.0,6.0,6.0,6.0]|3.0  |
|5.0         |5.0         |5.0         |5.0         |5.0         |D    |[5.0,5.0,5.0,5.0,5.0]|4.0  |
|4.0         |4.0         |4.0         |4.0         |4.0         |E    |[4.0,4.0,4.0,4.0,4.0]|5.0  |
|3.0         |3.0         |3.0         |3.0         |3.0         |F    |[3.0,3.0,3.0,3.0,3.0]|0.0  |
+------------+------------+------------+------------+------------+-----+---------------------+-----+
"""
```

```
label included df
+------------+------------+------------+------------+------------+-----+---------------------+-----+
|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_gp|subject_5_gp|grade|features             |label|
+------------+------------+------------+------------+------------+-----+---------------------+-----+
|2.0         |2.1         |3.5         |2.4         |3.0         |F    |[2.0,2.1,3.5,2.4,3.0]|0.0  |
|2.0         |2.0         |2.0         |3.0         |3.0         |F    |[2.0,2.0,2.0,3.0,3.0]|0.0  |
|2.1         |2.0         |2.4         |3.5         |3.0         |F    |[2.1,2.0,2.4,3.5,3.0]|0.0  |
|9.0         |9.0         |9.0         |9.0         |9.0         |A+   |[9.0,9.0,9.0,9.0,9.0]|6.0  |
|8.0         |8.0         |8.0         |8.0         |8.0         |A    |[8.0,8.0,8.0,8.0,8.0]|1.0  |
|7.0         |7.0         |7.0         |7.0         |7.0         |B    |[7.0,7.0,7.0,7.0,7.0]|2.0  |
|6.0         |6.0         |6.0         |6.0         |6.0         |C    |[6.0,6.0,6.0,6.0,6.0]|3.0  |
|5.0         |5.0         |5.0         |5.0         |5.0         |D    |[5.0,5.0,5.0,5.0,5.0]|4.0  |
|4.0         |4.0         |4.0         |4.0         |4.0         |E    |[4.0,4.0,4.0,4.0,4.0]|5.0  |
```

```
|3.0         |3.0         |3.0         |3.0         |3.0         |F    |[3.0,3.0,3.0,3.0,3.0]|0.0  |
+------------+------------+------------+------------+------------+-----+--------------------+-----+
only showing top 10 rows


'\n+------------+------------+------------+------------+------------+-----+--------------------+-----+\n|subject_1_gp|subject_2_gp|
```

```python
trainingData,testdata = label_df.randomSplit([0.7,0.3],seed = 42)

print("display training data")

trainingData.show(10,False)

"""
+------------+------------+------------+------------+------------+-----+--------------------+-----+
|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_gp|subject_5_gp|grade|features            |label|
+------------+------------+------------+------------+------------+-----+--------------------+-----+
|0.0         |0.0         |0.0         |0.0         |0.0         |F    |(5,[],[])           |0.0  |
|0.0         |0.0         |0.0         |0.0         |0.0         |F    |(5,[],[])           |0.0  |
|1.0         |1.0         |1.0         |1.0         |1.0         |F    |[1.0,1.0,1.0,1.0,1.0]|0.0  |
|2.0         |2.0         |2.0         |2.0         |2.0         |F    |[2.0,2.0,2.0,2.0,2.0]|0.0  |
|2.0         |2.0         |2.0         |2.0         |2.0         |F    |[2.0,2.0,2.0,2.0,2.0]|0.0  |
|2.0         |2.0         |2.0         |3.0         |3.0         |F    |[2.0,2.0,2.0,3.0,3.0]|0.0  |
|2.1         |2.0         |2.4         |3.5         |3.0         |F    |[2.1,2.0,2.4,3.5,3.0]|0.0  |
|2.1         |2.0         |2.4         |3.5         |3.0         |F    |[2.1,2.0,2.4,3.5,3.0]|0.0  |
|3.0         |3.0         |3.0         |3.0         |3.0         |F    |[3.0,3.0,3.0,3.0,3.0]|0.0  |
|4.0         |4.0         |4.0         |4.0         |4.0         |E    |[4.0,4.0,4.0,4.0,4.0]|5.0  |
+------------+------------+------------+------------+------------+-----+--------------------+-----+
"""
```

```
display training data
+------------+------------+------------+------------+------------+-----+--------------------+-----+
|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_gp|subject_5_gp|grade|features            |label|
+------------+------------+------------+------------+------------+-----+--------------------+-----+
|0.0         |0.0         |0.0         |0.0         |0.0         |F    |(5,[],[])           |0.0  |
|0.0         |0.0         |0.0         |0.0         |0.0         |F    |(5,[],[])           |0.0  |
|1.0         |1.0         |1.0         |1.0         |1.0         |F    |[1.0,1.0,1.0,1.0,1.0]|0.0  |
|2.0         |2.0         |2.0         |2.0         |2.0         |F    |[2.0,2.0,2.0,2.0,2.0]|0.0  |
|2.0         |2.0         |2.0         |2.0         |2.0         |F    |[2.0,2.0,2.0,2.0,2.0]|0.0  |
|2.0         |2.0         |2.0         |3.0         |3.0         |F    |[2.0,2.0,2.0,3.0,3.0]|0.0  |
|2.1         |2.0         |2.4         |3.5         |3.0         |F    |[2.1,2.0,2.4,3.5,3.0]|0.0  |
|2.1         |2.0         |2.4         |3.5         |3.0         |F    |[2.1,2.0,2.4,3.5,3.0]|0.0  |
|3.0         |3.0         |3.0         |3.0         |3.0         |F    |[3.0,3.0,3.0,3.0,3.0]|0.0  |
|4.0         |4.0         |4.0         |4.0         |4.0         |E    |[4.0,4.0,4.0,4.0,4.0]|5.0  |
+------------+------------+------------+------------+------------+-----+--------------------+-----+
only showing top 10 rows


'\n+------------+------------+------------+------------+------------+-----+--------------------+-----+\n|subject_1_gp|subject_2_gp|
```

```python
logisticRegression = LogisticRegression().setMaxIter(100).setRegParam(0.02).setElasticNetParam(0.8)

logisticRegressionModel = logisticRegression.fit(trainingData)

predictionDf = logisticRegressionModel.transform(testdata)


print("logisticregession model prediction")

predictionDf.show(10,False)

"""
+-----------+-----------+-----------+-----------+-----------+-----+-------------------+-----+----------------------------
|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_gp|subject_5_gp|grade|features           |label|rawPrediction
+-----------+-----------+-----------+-----------+-----------+-----+-------------------+-----+----------------------------
|1.0        |1.0        |1.0        |1.0        |1.0        |F    |[1.0,1.0,1.0,1.0,1.0]|0.0  |[4681.43057975775,-4092.0005704376
|2.0        |2.0        |2.0        |3.0        |3.0        |F    |[2.0,2.0,2.0,3.0,3.0]|0.0  |[2730.5696233099634,-3009.84817828
|2.0        |2.1        |3.5        |2.4        |3.0        |F    |[2.0,2.1,3.5,2.4,3.0]|0.0  |[2679.6994520823405,-2898.39258868
|2.0        |2.1        |3.5        |2.4        |3.0        |F    |[2.0,2.1,3.5,2.4,3.0]|0.0  |[2679.6994520823405,-2898.39258868
|3.0        |3.0        |3.0        |3.0        |3.0        |F    |[3.0,3.0,3.0,3.0,3.0]|0.0  |[2668.323731176173,-2643.263258509
|4.0        |4.0        |4.0        |4.0        |4.0        |E    |[4.0,4.0,4.0,4.0,4.0]|5.0  |[1661.7703068853848,-1918.89460254
|4.0        |4.0        |4.0        |4.0        |4.0        |E    |[4.0,4.0,4.0,4.0,4.0]|5.0  |[1661.7703068853848,-1918.89460254
|4.1        |4.1        |4.1        |4.1        |4.1        |E    |[4.1,4.1,4.1,4.1,4.1]|5.0  |[1561.1149644563066,-1846.45773694
|4.2        |4.2        |4.2        |4.2        |4.2        |E    |[4.2,4.2,4.2,4.2,4.2]|5.0  |[1460.459622027227,-1774.020871353
|4.3        |4.3        |4.3        |4.3        |4.3        |E    |[4.3,4.3,4.3,4.3,4.3]|5.0  |[1359.8042795981482,-1701.58400575
+-----------+-----------+-----------+-----------+-----------+-----+-------------------+-----+----------------------------
"""
```

```
logisticregession model prediction
+-----------+-----------+-----------+-----------+-----------+-----+-------------------+-----+----------------------------
|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_gp|subject_5_gp|grade|features           |label|rawPrediction
+-----------+-----------+-----------+-----------+-----------+-----+-------------------+-----+----------------------------
|1.0        |1.0        |1.0        |1.0        |1.0        |F    |[1.0,1.0,1.0,1.0,1.0]|0.0  |[8.677062451954509,-5.6366376319
|2.0        |2.0        |2.0        |3.0        |3.0        |F    |[2.0,2.0,2.0,3.0,3.0]|0.0  |[5.405825633464345,-3.8271204090
|2.0        |2.1        |3.5        |2.4        |3.0        |F    |[2.0,2.1,3.5,2.4,3.0]|0.0  |[5.458506640504615,-3.5359403936
|2.0        |2.1        |3.5        |2.4        |3.0        |F    |[2.0,2.1,3.5,2.4,3.0]|0.0  |[5.458506640504615,-3.5359403936
|3.0        |3.0        |3.0        |3.0        |3.0        |F    |[3.0,3.0,3.0,3.0,3.0]|0.0  |[5.405825633464345,-3.1098862217
|4.0        |4.0        |4.0        |4.0        |4.0        |E    |[4.0,4.0,4.0,4.0,4.0]|5.0  |[3.7702072242192637,-1.846511116
|4.0        |4.0        |4.0        |4.0        |4.0        |E    |[4.0,4.0,4.0,4.0,4.0]|5.0  |[3.7702072242192637,-1.846511116
|4.1        |4.1        |4.1        |4.1        |4.1        |E    |[4.1,4.1,4.1,4.1,4.1]|5.0  |[3.6066453832947563,-1.720173566
|4.2        |4.2        |4.2        |4.2        |4.2        |E    |[4.2,4.2,4.2,4.2,4.2]|5.0  |[3.443083542370247,-1.5938360156
|4.3        |4.3        |4.3        |4.3        |4.3        |E    |[4.3,4.3,4.3,4.3,4.3]|5.0  |[3.2795217014457396,-1.467498465
+-----------+-----------+-----------+-----------+-----------+-----+-------------------+-----+----------------------------
only showing top 10 rows
```

```
'\n+-----------+-----------+-----------+-----------+-----------+-----+-------------------+-----+----------------------------
```

```python
evaluator = BinaryClassificationEvaluator() .setLabelCol("label").setRawPredictionCol("prediction").setMetricName("areaUnderROC")

accuracy = evaluator.evaluate(predictionDf)

print("accuracy of the model")

print(accuracy * 100)
```

```
accuracy of the model
95.23809523809523
```

```
df1 = spark.createDataFrame(
    [
        (9.1,9.2,9.3,9.4,9.5),
        (9.0,9.0,9.0,9.0,9.0),
        (2.1,2.0,2.4,3.5,3.0),
        (8.0,8.1,8.2,8.3,8.4),
        (7.0,7.1,7.2,7.3,7.35),
        (6.0,6.1,6.2,6.3,6.4),
        (5.0,5.1,5.2,5.3,5.4)
    ],
    ["subject_1_gp","subject_2_gp","subject_3_gp","subject_4_gp","subject_5_gp"]
)

print("new values for prediction")

df1.printSchema()

df1.show(10,False)

"""
root
 |-- subject_1_gp: double (nullable = true)
 |-- subject_2_gp: double (nullable = true)
 |-- subject_3_gp: double (nullable = true)
 |-- subject_4_gp: double (nullable = true)
 |-- subject_5_gp: double (nullable = true)
"""
"""
+------------+------------+------------+------------+------------+
|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_gp|subject_5_gp|
+------------+------------+------------+------------+------------+
|9.1         |9.2         |9.3         |9.4         |9.5         |
|9.0         |9.0         |9.0         |9.0         |9.0         |
|2.1         |2.0         |2.4         |3.5         |3.0         |
|8.0         |8.1         |8.2         |8.3         |8.4         |
|7.0         |7.1         |7.2         |7.3         |7.35        |
|6.0         |6.1         |6.2         |6.3         |6.4         |
|5.0         |5.1         |5.2         |5.3         |5.4         |
+------------+------------+------------+------------+------------+
"""
```

```
new values for prediction
root
 |-- subject_1_gp: double (nullable = true)
 |-- subject_2_gp: double (nullable = true)
 |-- subject_3_gp: double (nullable = true)
 |-- subject_4_gp: double (nullable = true)
 |-- subject_5_gp: double (nullable = true)

+------------+------------+------------+------------+------------+
|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_gp|subject_5_gp|
+------------+------------+------------+------------+------------+
|9.1         |9.2         |9.3         |9.4         |9.5         |
|9.0         |9.0         |9.0         |9.0         |9.0         |
|2.1         |2.0         |2.4         |3.5         |3.0         |
|8.0         |8.1         |8.2         |8.3         |8.4         |
|7.0         |7.1         |7.2         |7.3         |7.35        |
|6.0         |6.1         |6.2         |6.3         |6.4         |
|5.0         |5.1         |5.2         |5.3         |5.4         |
+------------+------------+------------+------------+------------+


'\n+------------+------------+------------+------------+------------+\n|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_gp|subject_!
```

```python
df2 = marksDF_assembler.transform(df1)

df3 = logisticRegressionModel.transform(df2)

df3.createOrReplaceTempView("input_marks_view")

print("prediction of given data")

df3.show()

"""
+-----------+-----------+-----------+-----------+-----------+--------------------+--------------------+--------------------+------
|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_gp|subject_5_gp|            features|       rawPrediction|         probability|predi
+-----------+-----------+-----------+-----------+-----------+--------------------+--------------------+--------------------+------
|        9.1|        9.2|        9.3|        9.4|        9.5|[9.1,9.2,9.3,9.4,...|[-3854.1141599748...|[0.0,1.0426740666...|
|        9.0|        9.0|        9.0|        9.0|        9.0|[9.0,9.0,9.0,9.0,...|[-3370.9968145685...|[0.0,6.8300184358...|
|        2.1|        2.0|        2.4|        3.5|        3.0|[2.1,2.0,2.4,3.5,...|[2554.92106006113...|[0.99999999296771...|
|        8.0|        8.1|        8.2|        8.3|        8.4|[8.0,8.1,8.2,8.3,...|[-2746.9053932549...|[0.0,1.0,1.656005...|
|        7.0|        7.1|        7.2|        7.3|       7.35|[7.0,7.1,7.2,7.3,...|[-1706.3408951380...|[0.0,1.3146646688...|
|        6.0|        6.1|        6.2|        6.3|        6.4|[6.0,6.1,6.2,6.3,...|[-733.79854467341...|[0.0,0.0,1.174650...|
|        5.0|        5.1|        5.2|        5.3|        5.4|[5.0,5.1,5.2,5.3,...|[272.754879617376...|[0.0,0.0,0.0,6.66...|
+-----------+-----------+-----------+-----------+-----------+--------------------+--------------------+--------------------+------
"""
```

prediction of given data
```
+-----------+-----------+-----------+-----------+-----------+--------------------+--------------------+--------------------+----
|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_gp|subject_5_gp|            features|       rawPrediction|         probability|pre
+-----------+-----------+-----------+-----------+-----------+--------------------+--------------------+--------------------+----
|        9.1|        9.2|        9.3|        9.4|        9.5|[9.1,9.2,9.3,9.4,...|[-5.2169138587886...|[1.31392309736027...|
|        9.0|        9.0|        9.0|        9.0|        9.0|[9.0,9.0,9.0,9.0,...|[-4.4078848220061...|[4.96397730443946...|
|        2.1|        2.0|        2.4|        3.5|        3.0|[2.1,2.0,2.4,3.5,...|[5.36192479426412...|[0.49945872236119...|
|        8.0|        8.1|        8.2|        8.3|        8.4|[8.0,8.1,8.2,8.3,...|[-3.4177336086190...|[3.74339932298724...|
|        7.0|        7.1|        7.2|        7.3|       7.35|[7.0,7.1,7.2,7.3,...|[-1.7047243628317...|[0.00399493576025...|
|        6.0|        6.1|        6.2|        6.3|        6.4|[6.0,6.1,6.2,6.3,...|[-0.1464967901288...|[0.02170122208900...|
|        5.0|        5.1|        5.2|        5.3|        5.4|[5.0,5.1,5.2,5.3,...|[1.48912161911619...|[0.08683248402132...|
+-----------+-----------+-----------+-----------+-----------+--------------------+--------------------+--------------------+----
```

'\n+-----------+-----------+-----------+-----------+-----------+--------------------+--------------------+--------------------+-

```python
spark.sql("select subject_1_gp,subject_2_gp,subject_3_gp,subject_4_gp,subject_5_gp,prediction from input_marks_view").show()

"""
+-----------+-----------+-----------+-----------+-----------+----------+
|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_gp|subject_5_gp|prediction|
+-----------+-----------+-----------+-----------+-----------+----------+
|        9.1|        9.2|        9.3|        9.4|        9.5|       6.0|
|        9.0|        9.0|        9.0|        9.0|        9.0|       6.0|
|        2.1|        2.0|        2.4|        3.5|        3.0|       0.0|
|        8.0|        8.1|        8.2|        8.3|        8.4|       1.0|
|        7.0|        7.1|        7.2|        7.3|       7.35|       2.0|
|        6.0|        6.1|        6.2|        6.3|        6.4|       3.0|
|        5.0|        5.1|        5.2|        5.3|        5.4|       4.0|
+-----------+-----------+-----------+-----------+-----------+----------+
"""
```

```
+-----------+-----------+-----------+-----------+-----------+----------+
|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_gp|subject_5_gp|prediction|
+-----------+-----------+-----------+-----------+-----------+----------+
|        9.1|        9.2|        9.3|        9.4|        9.5|       6.0|
|        9.0|        9.0|        9.0|        9.0|        9.0|       6.0|
|        2.1|        2.0|        2.4|        3.5|        3.0|       0.0|
|        8.0|        8.1|        8.2|        8.3|        8.4|       1.0|
|        7.0|        7.1|        7.2|        7.3|       7.35|       2.0|
|        6.0|        6.1|        6.2|        6.3|        6.4|       2.0|
|        5.0|        5.1|        5.2|        5.3|        5.4|       4.0|
+-----------+-----------+-----------+-----------+-----------+----------+
```

'\n+-----------+-----------+-----------+-----------+-----------+----------+\n|subject_1_gp|subject_2_gp|subject_3_gp|subject_4_g

```
final_out =spark.sql ("SELECT main_df.subject_1_gp,main_df.subject_2_gp,main_df.subject_3_gp," +
    "main_df.subject_4_gp,main_df.subject_5_gp,main_df.grade,main_df.label,input_marks_view.prediction FROM main_df  " +
    "JOIN input_marks_view  ON main_df.subject_1_gp = input_marks_view.subject_1_gp AND main_df.subject_2_gp = input_marks_view.subj
    "AND main_df.subject_3_gp = input_marks_view.subject_3_gp AND main_df.subject_4_gp = input_marks_view.subject_4_gp AND " +
    "main_df.subject_5_gp = input_marks_view.subject_5_gp  GROUP BY main_df.subject_1_gp,main_df.subject_2_gp," +
    "main_df.subject_3_gp,main_df.subject_4_gp,main_df.subject_5_gp,main_df.grade,input_marks_view.prediction,main_df.label")
```

```
AnalysisException: Table or view not found: main_df; line 1 pos 173;
'Aggregate ['main_df.subject_1_gp, 'main_df.subject_2_gp, 'main_df.subject_3_gp, 'main_df.subject_4_gp, 'main_df.subject_5_gp, 'main
+- 'Join Inner, (((('main_df.subject_1_gp = 'input_marks_view.subject_1_gp) AND ('main_df.subject_2_gp = 'input_marks_view.subject_2
   :- 'UnresolvedRelation [main_df], [], false
   +- 'UnresolvedRelation [input_marks_view], [], false
```

```
final_out.describe()
```

```
NameError: name 'final_out' is not defined
```