

```

from pyspark.sql import SparkSession
from pyspark.ml.feature import StringIndexer, VectorAssembler
from pyspark.ml.classification import RandomForestClassifier
from pyspark.ml.evaluation import BinaryClassificationEvaluator

spark = SparkSession.builder.getOrCreate()

suvDF = spark.read.csv("suv_data.csv", header = True, inferSchema = True)

suvDF.printSchema()

"""
root
 |-- User_ID: integer (nullable = true)
 |-- Gender: string (nullable = true)
 |-- Age: integer (nullable = true)
 |-- EstimatedSalary: integer (nullable = true)
 |-- Purchased: integer (nullable = true)
"""

root
 |-- User_ID: integer (nullable = true)
 |-- Gender: string (nullable = true)
 |-- Age: integer (nullable = true)
 |-- EstimatedSalary: integer (nullable = true)
 |-- Purchased: integer (nullable = true)

'\nroot\n |-- User_ID: integer (nullable = true)\n |-- Gender: string (nullable = true)\n |-- Age: integer (nullable = true)\n |-- E:

```

```

suvDF.show(10, False)

"""
+-----+-----+-----+-----+-----+
|User_ID|Gender|Age|EstimatedSalary|Purchased|
+-----+-----+-----+-----+-----+
|15624510|Male  |19 |19000          |0        |
|15810944|Male  |35 |20000          |0        |
|15668575|Female|26 |43000          |0        |
|15603246|Female|27 |57000          |0        |
|15804002|Male  |19 |76000          |0        |
|15728773|Male  |27 |58000          |0        |
|15598044|Female|27 |84000          |0        |
|15694829|Female|32 |150000         |1        |
|15600575|Male  |25 |33000          |0        |
|15727311|Female|35 |65000          |0        |
+-----+-----+-----+-----+-----+
"""

+-----+-----+-----+-----+-----+
|User_ID|Gender|Age|EstimatedSalary|Purchased|
+-----+-----+-----+-----+-----+
|15624510|Male  |19 |19000          |0        |
|15810944|Male  |35 |20000          |0        |
|15668575|Female|26 |43000          |0        |
|15603246|Female|27 |57000          |0        |
|15804002|Male  |19 |76000          |0        |
|15728773|Male  |27 |58000          |0        |
|15598044|Female|27 |84000          |0        |
|15694829|Female|32 |150000         |1        |
|15600575|Male  |25 |33000          |0        |
|15727311|Female|35 |65000          |0        |
+-----+-----+-----+-----+-----+

only showing top 10 rows

'\n+-----+-----+-----+-----+-----+\n|User_ID|Gender|Age|EstimatedSalary|Purchased|\n+-----+-----+-----+-----+-----+

```

```
suvDF.describe().show()

suvDF.createOrReplaceTempView("first_view")
```

```
"""
+-----+-----+-----+-----+-----+
|summary|      User_ID|Gender|      Age| EstimatedSalary|      Purchased|
+-----+-----+-----+-----+-----+
|  count|         400|   400|         400|           400|           400|
|   mean| 1.56915397575E7| null|    37.655|    69742.5|    0.3575|
| stddev|71658.32158119006| null|10.482876597307927|34096.9602824248|0.4798639635968691|
|   min|    15566689|Female|    18|    15000|           0|
|   max|    15815236| Male|    60|   150000|           1|
+-----+-----+-----+-----+-----+
"""
```

```
+-----+-----+-----+-----+-----+
|summary|      User_ID|Gender|      Age| EstimatedSalary|      Purchased|
+-----+-----+-----+-----+-----+
|  count|         400|   400|         400|           400|           400|
|   mean| 1.56915397575E7| null|    37.655|    69742.5|    0.3575|
| stddev|71658.32158119006| null|10.482876597307927|34096.9602824248|0.4798639635968691|
|   min|    15566689|Female|    18|    15000|           0|
|   max|    15815236| Male|    60|   150000|           1|
+-----+-----+-----+-----+-----+
```

```
'\n+-----+-----+-----+-----+-----+-----+\n|summary|      User_ID|Gender|
```

```
data_set_suv= spark.sql("select Gender,Age,EstimatedSalary,Purchased from first_view ")
data_set_suv.show(5,False)
```

```
"""
+-----+-----+-----+-----+
|Gender|Age|EstimatedSalary|Purchased|
+-----+-----+-----+-----+
|Male  |19 |19000           |0        |
|Male  |35 |20000           |0        |
|Female|26 |43000           |0        |
|Female|27 |57000           |0        |
|Male  |19 |76000           |0        |
+-----+-----+-----+-----+
"""
```

```
+-----+-----+-----+-----+
|Gender|Age|EstimatedSalary|Purchased|
+-----+-----+-----+-----+
|Male  |19 |19000           |0        |
|Male  |35 |20000           |0        |
|Female|26 |43000           |0        |
|Female|27 |57000           |0        |
|Male  |19 |76000           |0        |
+-----+-----+-----+-----+
```

only showing top 5 rows

```
'\n+-----+-----+-----+-----+-----+-----+\n|Gender|Age|EstimatedSalary|Purchased|\n+-----+-----+-----+-----+-----+-----+\n|Male  |19 |
```

```

suv_indexer = StringIndexer(inputCol = "Gender", outputCol = "gen_label")

gen_label_df = suv_indexer.fit(data_set_suv).transform(data_set_suv)

print("after adding label")

gen_label_df.printSchema()

gen_label_df.show(5,False)

```

```

"""
|-- Gender: string (nullable = true)
|-- Age: integer (nullable = true)
|-- EstimatedSalary: integer (nullable = true)
|-- Purchased: integer (nullable = true)
|-- gen_label: double (nullable = false)
"""

```

```

"""
+-----+-----+-----+-----+
|Gender|Age|EstimatedSalary|Purchased|gen_label|
+-----+-----+-----+-----+
|Male  |19 |19000          |0        |1.0      |
|Male  |35 |20000          |0        |1.0      |
|Female|26 |43000          |0        |0.0      |
|Female|27 |57000          |0        |0.0      |
|Male  |19 |76000          |0        |1.0      |
+-----+-----+-----+-----+

```

```

"""

```

after adding label

```

root
 |-- Gender: string (nullable = true)
 |-- Age: integer (nullable = true)
 |-- EstimatedSalary: integer (nullable = true)
 |-- Purchased: integer (nullable = true)
 |-- gen_label: double (nullable = false)

```

```

+-----+-----+-----+-----+
|Gender|Age|EstimatedSalary|Purchased|gen_label|
+-----+-----+-----+-----+
|Male  |19 |19000          |0        |1.0      |
|Male  |35 |20000          |0        |1.0      |
|Female|26 |43000          |0        |0.0      |
|Female|27 |57000          |0        |0.0      |
|Male  |19 |76000          |0        |1.0      |
+-----+-----+-----+-----+

```

only showing top 5 rows

```

'\n+-----+-----+-----+-----+-----+\n|Gender|Age|EstimatedSalary|Purchased|gen_label|\n+-----+-----+-----+-----+

```

```

inputCols = ["Age","EstimatedSalary","gen_label"]

```

```

outputCol = "features"

```

```

suvDF_assembler = VectorAssembler(inputCols = inputCols,outputCol = outputCol)
featuresDf = suvDF_assembler.transform(gen_label_df)

```

```

print("featuresDF printSchema")

```

```

featuresDf.printSchema()

```

```

"""
root
 |-- Gender: string (nullable = true)
 |-- Age: integer (nullable = true)
 |-- EstimatedSalary: integer (nullable = true)
 |-- Purchased: integer (nullable = true)
 |-- gen_label: double (nullable = false)
 |-- features: vector (nullable = true)
"""

```

```
featuresDF printSchema
root
 |-- Gender: string (nullable = true)
 |-- Age: integer (nullable = true)
 |-- EstimatedSalary: integer (nullable = true)
 |-- Purchased: integer (nullable = true)
 |-- gen_label: double (nullable = false)
 |-- features: vector (nullable = true)

'\nroot\n |-- Gender: string (nullable = true)\n |-- Age: integer (nullable = true)\n |-- EstimatedSalary: integer (nullable = true)'
```

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

```
print("featureDf show")
```

```
"""
+-----+-----+-----+-----+-----+-----+
|Gender|Age|EstimatedSalary|Purchased|gen_label|features|
+-----+-----+-----+-----+-----+-----+
|Male  |19 |19000          |0        |1.0      |[19.0,19000.0,1.0]|
|Male  |35 |20000          |0        |1.0      |[35.0,20000.0,1.0]|
|Female|26 |43000          |0        |0.0      |[26.0,43000.0,0.0]|
|Female|27 |57000          |0        |0.0      |[27.0,57000.0,0.0]|
|Male  |19 |76000          |0        |1.0      |[19.0,76000.0,1.0]|
|Male  |27 |58000          |0        |1.0      |[27.0,58000.0,1.0]|
|Female|27 |84000          |0        |0.0      |[27.0,84000.0,0.0]|
|Female|32 |150000         |1        |0.0      |[32.0,150000.0,0.0]|
|Male  |25 |33000          |0        |1.0      |[25.0,33000.0,1.0]|
|Female|35 |65000          |0        |0.0      |[35.0,65000.0,0.0]|
+-----+-----+-----+-----+-----+-----+
"""
```

```
+-----+-----+-----+-----+-----+-----+
|Gender|Age|EstimatedSalary|Purchased|gen_label|features|
+-----+-----+-----+-----+-----+-----+
|Male  |19 |19000          |0        |1.0      |[19.0,19000.0,1.0]|
|Male  |35 |20000          |0        |1.0      |[35.0,20000.0,1.0]|
|Female|26 |43000          |0        |0.0      |[26.0,43000.0,0.0]|
|Female|27 |57000          |0        |0.0      |[27.0,57000.0,0.0]|
|Male  |19 |76000          |0        |1.0      |[19.0,76000.0,1.0]|
|Male  |27 |58000          |0        |1.0      |[27.0,58000.0,1.0]|
|Female|27 |84000          |0        |0.0      |[27.0,84000.0,0.0]|
|Female|32 |150000         |1        |0.0      |[32.0,150000.0,0.0]|
|Male  |25 |33000          |0        |1.0      |[25.0,33000.0,1.0]|
|Female|35 |65000          |0        |0.0      |[35.0,65000.0,0.0]|
+-----+-----+-----+-----+-----+-----+
```

only showing top 10 rows

```
featureDf show
```

```
'\n+-----+-----+-----+-----+-----+-----+\n|Gender|Age|EstimatedSalary|Purchased|gen_label|features'
```

```
suv_indexer = StringIndexer(inputCol = "Purchased", outputCol = "label")

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

print("after adding pur_label")

label_df.createOrReplaceTempView("main_df")

label_df.printSchema()

"""
root
 |-- Gender: string (nullable = true)
 |-- Age: integer (nullable = true)
 |-- EstimatedSalary: integer (nullable = true)
 |-- Purchased: integer (nullable = true)
 |-- gen_label: double (nullable = false)
 |-- features: vector (nullable = true)
 |-- label: double (nullable = false)
"""
```

```
after adding pur_label
root
 |-- Gender: string (nullable = true)
 |-- Age: integer (nullable = true)
 |-- EstimatedSalary: integer (nullable = true)
 |-- Purchased: integer (nullable = true)
 |-- gen_label: double (nullable = false)
 |-- features: vector (nullable = true)
 |-- label: double (nullable = false)
```

```
label_df.show(10,False)
```

```
"""
+-----+-----+-----+-----+-----+-----+-----+
|Gender|Age|EstimatedSalary|Purchased|gen_label|features          |label|
+-----+-----+-----+-----+-----+-----+-----+
|Male  |19 |19000          |0        |1.0      |[19.0,19000.0,1.0]|0.0  |
|Male  |35 |20000          |0        |1.0      |[35.0,20000.0,1.0]|0.0  |
|Female|26 |43000          |0        |0.0      |[26.0,43000.0,0.0]|0.0  |
|Female|27 |57000          |0        |0.0      |[27.0,57000.0,0.0]|0.0  |
|Male  |19 |76000          |0        |1.0      |[19.0,76000.0,1.0]|0.0  |
|Male  |27 |58000          |0        |1.0      |[27.0,58000.0,1.0]|0.0  |
|Female|27 |84000          |0        |0.0      |[27.0,84000.0,0.0]|0.0  |
|Female|32 |150000         |1        |0.0      |[32.0,150000.0,0.0]|1.0  |
|Male  |25 |33000          |0        |1.0      |[25.0,33000.0,1.0]|0.0  |
|Female|35 |65000          |0        |0.0      |[35.0,65000.0,0.0]|0.0  |
+-----+-----+-----+-----+-----+-----+-----+
"""
```

```
+-----+-----+-----+-----+-----+-----+-----+
|Gender|Age|EstimatedSalary|Purchased|gen_label|features          |label|
+-----+-----+-----+-----+-----+-----+-----+
|Male  |19 |19000          |0        |1.0      |[19.0,19000.0,1.0]|0.0  |
|Male  |35 |20000          |0        |1.0      |[35.0,20000.0,1.0]|0.0  |
|Female|26 |43000          |0        |0.0      |[26.0,43000.0,0.0]|0.0  |
|Female|27 |57000          |0        |0.0      |[27.0,57000.0,0.0]|0.0  |
|Male  |19 |76000          |0        |1.0      |[19.0,76000.0,1.0]|0.0  |
|Male  |27 |58000          |0        |1.0      |[27.0,58000.0,1.0]|0.0  |
|Female|27 |84000          |0        |0.0      |[27.0,84000.0,0.0]|0.0  |
|Female|32 |150000         |1        |0.0      |[32.0,150000.0,0.0]|1.0  |
|Male  |25 |33000          |0        |1.0      |[25.0,33000.0,1.0]|0.0  |
|Female|35 |65000          |0        |0.0      |[35.0,65000.0,0.0]|0.0  |
+-----+-----+-----+-----+-----+-----+-----+
```

only showing top 10 rows

```
'\n+-----+-----+-----+-----+-----+-----+-----+\n|Gender|Age|EstimatedSalary|Purchased|gen_label|featu
```

Gender	Age	EstimatedSalary	Purchased	gen_label	features	label	rawPrediction	probability
--------	-----	-----------------	-----------	-----------	----------	-------	---------------	-------------

```
only showing top 10 rows
```

11/11/2016

```
'\naccuracy of the model\n83.09386973180077\n'
```

```
'\nroot\n |-- User_ID: integer (nullable = true)\n |-- Gender: string (nullable = true)\n |-- Age: integer (nullable = true)\n |-- E:
```

```
df1= spark.sql("select Gender,Age,EstimatedSalary from input_view ")
df1.show(5,False)
```

```
"""
+-----+-----+
|Gender|Age|EstimatedSalary|
+-----+-----+
|Male  |19 |19000          |
|Male  |35 |20000          |
|Female|26 |43000          |
|Female|27 |57000          |
|Male  |19 |76000          |
+-----+-----+
"""
```

```
+-----+-----+
|Gender|Age|EstimatedSalary|
+-----+-----+
|Male  |19 |19000          |
|Male  |35 |20000          |
|Female|26 |43000          |
|Female|27 |57000          |
|Male  |19 |76000          |
+-----+-----+
only showing top 5 rows
```

```
'\n+-----+-----+\n|Gender|Age|EstimatedSalary|\n+-----+-----+\n|Male  |19 |19000          |\n|Male  |35 |20000          |'
```

```
input_indexer = StringIndexer(inputCol = "Gender", outputCol = "gen_label")
```

```
gen_label_input_df = input_indexer.fit(data_set_suv).transform(df1)
```

```
print("after adding label")
```

```
gen_label_input_df.printSchema()
```

```
gen_label_input_df.show(5,False)
```

```
"""
root
 |-- Gender: string (nullable = true)
 |-- Age: integer (nullable = true)
 |-- EstimatedSalary: integer (nullable = true)
 |-- gen_label: double (nullable = false)
```

```
+-----+-----+-----+
|Gender|Age|EstimatedSalary|gen_label|
+-----+-----+-----+
|Male  |19 |19000          |1.0      |
|Male  |35 |20000          |1.0      |
|Female|26 |43000          |0.0      |
|Female|27 |57000          |0.0      |
|Male  |19 |76000          |1.0      |
+-----+-----+-----+
"""
```

```
after adding label
root
 |-- Gender: string (nullable = true)
 |-- Age: integer (nullable = true)
 |-- EstimatedSalary: integer (nullable = true)
 |-- gen_label: double (nullable = false)
```

```
+-----+-----+-----+
|Gender|Age|EstimatedSalary|gen_label|
+-----+-----+-----+
|Male  |19 |19000          |1.0      |
|Male  |35 |20000          |1.0      |
|Female|26 |43000          |0.0      |
|Female|27 |57000          |0.0      |
|Male  |19 |76000          |1.0      |
+-----+-----+-----+
only showing top 5 rows
```



```
inputCols = ["Age","EstimatedSalary","gen_label"]

outputCol = "features"

input_assembler = VectorAssembler(inputCols = inputCols,outputCol = outputCol)
featuresDf = input_assembler.transform(gen_label_input_df)

print("featuresDF printSchema")

featuresDf.printSchema()

"""
root
 |-- Gender: string (nullable = true)
 |-- Age: integer (nullable = true)
 |-- EstimatedSalary: integer (nullable = true)
 |-- gen_label: double (nullable = false)
 |-- features: vector (nullable = true)
"""

featuresDf.show(5,False)
```

```
"""
+-----+-----+-----+-----+-----+
|Gender|Age|EstimatedSalary|gen_label|features      |
+-----+-----+-----+-----+-----+
|Male  |19 |19000          |1.0      |[19.0,19000.0,1.0]|
|Male  |35 |20000          |1.0      |[35.0,20000.0,1.0]|
|Female|26 |43000          |0.0      |[26.0,43000.0,0.0]|
|Female|27 |57000          |0.0      |[27.0,57000.0,0.0]|
|Male  |19 |76000          |1.0      |[19.0,76000.0,1.0]|
+-----+-----+-----+-----+-----+
"""
```

```
featuresDF printSchema
root
 |-- Gender: string (nullable = true)
 |-- Age: integer (nullable = true)
 |-- EstimatedSalary: integer (nullable = true)
 |-- gen_label: double (nullable = false)
 |-- features: vector (nullable = true)

+-----+-----+-----+-----+-----+
|Gender|Age|EstimatedSalary|gen_label|features      |
+-----+-----+-----+-----+-----+
|Male  |19 |19000          |1.0      |[19.0,19000.0,1.0]|
|Male  |35 |20000          |1.0      |[35.0,20000.0,1.0]|
|Female|26 |43000          |0.0      |[26.0,43000.0,0.0]|
|Female|27 |57000          |0.0      |[27.0,57000.0,0.0]|
|Male  |19 |76000          |1.0      |[19.0,76000.0,1.0]|
+-----+-----+-----+-----+-----+

only showing top 5 rows
```

```
print("prediction_input")
input_pre = ran_for_Model.transform(featuresDf)
input_pre.show(5,False)

"""
+-----+-----+-----+-----+-----+-----+-----+
|Gender|Age|EstimatedSalary|gen_label|features      |rawPrediction      |probability      |
+-----+-----+-----+-----+-----+-----+-----+
|Male  |19 |19000          |1.0      |[19.0,19000.0,1.0]| [20.0,0.0]        |[1.0,0.0]        |
|Male  |35 |20000          |1.0      |[35.0,20000.0,1.0]| [20.0,0.0]        |[1.0,0.0]        |
|Female|26 |43000          |0.0      |[26.0,43000.0,0.0]| [20.0,0.0]        |[1.0,0.0]        |
|Female|27 |57000          |0.0      |[27.0,57000.0,0.0]| [20.0,0.0]        |[1.0,0.0]        |
|Male  |19 |76000          |1.0      |[19.0,76000.0,1.0]| [19.977272727272727,0.0227272727272728]| [0.9988636363636363,0.0011363636363636]|
+-----+-----+-----+-----+-----+-----+-----+
"""
```

prediction_input

Gender	Age	EstimatedSalary	gen_label	features	rawPrediction	probability
Male	19	19000	1.0	[19.0,19000.0,1.0]	[20.0,0.0]	[1.0,0.0]
Male	35	20000	1.0	[35.0,20000.0,1.0]	[20.0,0.0]	[1.0,0.0]
Female	26	43000	0.0	[26.0,43000.0,0.0]	[20.0,0.0]	[1.0,0.0]
Female	27	57000	0.0	[27.0,57000.0,0.0]	[20.0,0.0]	[1.0,0.0]
Male	19	76000	1.0	[19.0,76000.0,1.0]	[19.9772727272727,0.0227272727272728]	[0.998863636363636,0.0011363636363636]

only showing top 5 rows

