

```
[1]:
import findspark
findspark.init()
# pyspark is expecting a zero indexed integer for the label column
import pyspark
findspark.find()

Out[1]:
'C:\Spark\spark-3.2.1-bin-hadoop3.2'

In [2]:
from pyspark.sql import SparkSession

spark = SparkSession.builder.config("spark.driver.memory", "150g").appName("Titanic").getOrCreate()

spark

Out[2]:
SparkSession - in-memory

SparkContext

Spark UI

Version
V3.2.1

Master
local[*]

AppName
Titanic

In [3]:
import pandas as pd
df = (spark.read.format("csv").options(header="True")
     .load("C:\Users\Rajesh Kumar\Documents\Train.csv"))

In [4]:
df.limit(20).toPandas()

Out[4]:
PassengerId  Survived  Pclass      Name  Sex  Age  SibSp  Parch      Ticket  Fare  Cabin  Embarked
0            1         0      3   Braund, Mr. Owen Harris  male  22  1  0      A/5 21171  7.25  None  S
1            2         1      1  Cumings, Mrs. John Bradley (Florence Briggs Th...  female  38  1  0      PC 17599  71.2833  C85  C
2            3         1      3   Heikkinen, Miss. Laina  female  26  0  0      STON/O2 3101282  7.925  None  S
3            4         1      1  Futrelle, Mrs. Jacques Heath (Lily May Peel)  female  35  1  0      113803  53.1  C123  S
4            5         0      3   Allen, Mr. William Henry  male  35  0  0      373450  8.05  None  S
5            6         0      3   Moran, Mr. James  male  None  0  0      330877  8.4683  None  Q
6            7         0      1  McCarthy, Mr. Timothy J  male  54  0  0      17463  51.8625  E46  S
7            8         0      3  Palsson, Master. Gosta Leonard  male  2  3  1      349909  21.075  None  S
8            9         1      3  Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)  female  27  0  2      34742  11.133  None  S
9           10         1      2  Nasser, Mrs. Nicholas (Adele Achem)  female  14  1  0      23776  30.0708  None  C
10           11         1      3  Sandstrom, Miss. Marguerite Rut  female  4  1  1      PP 9549  16.7  G6  S
11           12         1      1  Bonnell, Miss. Elizabeth  female  58  0  0      113783  26.56  C103  S
12           13         0      3  Saudercock, Mr. William Henry  male  20  0  0      A/5 2151  8.05  None  S
13           14         0      3  Anderson, Mr. Andrew Johan  male  39  1  5      347082  31.275  None  S
14           15         0      3  Vestrom, Miss. Hulda Amanda Adolfina  female  14  0  0      350406  7.8542  None  S
15           16         1      2  Hewlett, Mrs. (Mary D Kingcome)  female  55  0  0      248706  16  None  S
16           17         0      3  Rice, Master. Eugene  male  2  4  1      382652  29.125  None  Q
17           18         1      2  Williams, Mr. Charles Eugene  male  None  0  0      244373  13  None  S
18           19         0      3  Vander Planke, Mrs. Julius (Emelia Maria Vander...  female  31  1  0      345763  18  None  S
19           20         1      3  Masselmani, Mrs. Fatma  female  None  0  0      2649  7.225  None  C

In [5]:
df.count()

Out[5]:
891

In [6]:
df.printSchema()

root
 |-- PassengerId: string (nullable = true)
 |-- Survived: string (nullable = true)
 |-- Pclass: string (nullable = true)
 |-- Name: string (nullable = true)
 |-- Sex: string (nullable = true)
 |-- Age: string (nullable = true)
 |-- SibSp: string (nullable = true)
 |-- Parch: string (nullable = true)
 |-- Ticket: string (nullable = true)
 |-- Fare: string (nullable = true)
 |-- Cabin: string (nullable = true)
 |-- Embarked: string (nullable = true)

In [7]:
df.groupBy("Sex").count().show()

df.groupBy("Survived").count().show()

df.groupBy("Sex", "Survived").count().orderBy("Sex").show()

df.groupBy("Pclass", "Survived").count().orderBy("Pclass").show()

-----+-----+
| Sex|count|
+----+-----+
|female| 314|
|male| 577|
+----+-----+

-----+-----+
| Survived|count|
+-----+-----+
| 0| 81|
| 1| 342|
+-----+-----+

-----+-----+
| Sex|Survived|count|
+----+-----+-----+
|male| 0| 81|
|female| 1| 233|
|male| 1| 109|
|male| 0| 109|
+-----+-----+

-----+-----+
| Column Name|Null Values|Count| Null Value_Percent|
+-----+-----+-----+-----+
| Age| 277| 19.865319865319865|
| Cabin| 687| 77.104371043711|
| Embarked| 218| 244680.1392783|
+-----+-----+

In [8]:
df.select(mean("age")).show()

-----+-----+
| avg(age)|
+-----+
| 29.69911764785882|
+-----+

In [10]:
df.select("name").show(10, False)

-----+-----+
| name|
+-----+
| Braund, Mr. Owen Harris|
| Cumings, Mrs. John Bradley (Florence Briggs Thayer)|
| Heikkinen, Miss. Laina|
| Futrelle, Mrs. Jacques Heath (Lily May Peel)|
| Allen, Mr. William Henry|
| Moran, Mr. James|
| McCarthy, Mr. Timothy J|
| Palsson, Master. Gosta Leonard|
| Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)|
| Nasser, Mrs. Nicholas (Adele Achen)|
+-----+
only showing top 10 rows

In [11]:
df = df.withColumn("Initial", regexp_extract(col("name"), "([a-zA-Z-]+)", 1))

df = df.replace(("Miss", "Mie", "Ms", "Dr", "Major", "Lady", "Countess", "Jonkheer", "Col", "Rev", "Capt", "Sir", "Don"),
               ["Miss", "Miss", "Mr", "Mr", "Mrs", "Mrs", "Other", "Other", "Other", "Mr", "Mr", "Mr", "Don"])

In [13]:
df = df.withColumn("Age", when(df["Initial"] == "Miss") & (df["Age"].isNull()), 22) otherwise(df["Age"])
df = df.withColumn("Age", when(df["Initial"] == "Other") & (df["Age"].isNull()), 45) otherwise(df["Age"])
df = df.withColumn("Age", when(df["Initial"] == "Master") & (df["Age"].isNull()), 6) otherwise(df["Age"])
df = df.withColumn("Age", when(df["Initial"] == "Mr") & (df["Age"].isNull()), 32) otherwise(df["Age"])
df = df.withColumn("Age", when(df["Initial"] == "Mrs") & (df["Age"].isNull()), 33) otherwise(df["Age"])

df = df.na.fill({"Embarked": 'S'})
df = df.drop("cabin")

df.select("initial", "name", "age").filter("initial == 'Miss'").show(100, False)

-----+-----+-----+-----+
| initial|name|age|
+-----+-----+-----+-----+
| Miss|Heikkinen, Miss. Laina|16|
| Miss|Sandstrom, Miss. Marguerite Rut|14|
| Miss|Bonnell, Miss. Elizabeth|158|
| Miss|Vestrom, Miss. Hulda Amanda Adolfina|14|
| Miss|"McGowan, Miss. Anna""Annie""|15|
| Miss|Palsson, Miss. Torborg Danira|18|
| Miss|"O'Dowd, Miss. Ellen""Mellie""|122|
| Miss|Glynn, Miss. Mary Agatha|122|
| Miss|Vander Planke, Mrs. Augusta Maria|118|
| Miss|Nicola-Yarred, Miss. Jamila|114|
| Miss|Laroche, Miss. Simone Marie Anne Andree|13|
| Miss|Devenay, Miss. Margarete Delia|119|
| Miss|O'Driscoll, Miss. Bridget|122|
| Miss|Rugg, Miss. Emily|121|
| Miss|West, Miss. Constance Mirium|122|
| Miss|Icard, Miss. Amelie|138|
| Miss|Andersson, Miss. Elna Alexandra|15|
| Miss|Gossewin, Miss. Lillian Amy|116|
| Miss|Dowdell, Miss. Elizabeth|130|
| Miss|McDermott, Miss. Bridget Delia|117|
| Miss|Ilett, Miss. Bertha|117|
| Miss|Fortune, Miss. Helen Helen|124|
| Miss|Petranec, Miss. Matilda|128|
| Miss|Salakelavik, Miss. Anna Kristine|121|
| Miss|Moran, Miss. Bertha|122|
| Miss|Zabour, Miss. Heleni|114,5|
| Miss|Jussila, Miss. Mattina|130|
| Miss|Atkalah, Miss. Melae|117|
| Miss|Andersson, Miss. Ellis Anna Maria|122|
| Miss|Webber, Miss. Sarah|125,5|
| Miss|Peter, Miss. Anna|122|
| Miss|Newson, Miss. Helen Morypney|119|
| Miss|Pyster, Miss. Anna Sofia|135|
| Miss|"Ford, Miss. Robina Maggie""Ruby""|19|
| Miss|"Lingham, Miss. Katherine""Katie""|12|
| Miss|Johnson, Miss. Eleanor Ileen|11|
| Miss|Isham, Miss. Anna Elizabeth|160|
| Miss|Sage, Miss. Constance Gladys|122|
| Miss|Kink-Hellmann, Miss. Luisa Gretchen|14|
| Miss|Andersen Jensen, Miss. Carla Christine Nielsine|158|
| Miss|Lurette, Miss. Elise|158|
| Miss|"Madigan, Miss. Margaret""Maggie""|122|
| Miss|"Vroily, Miss. Henriette""Mrs. Harbeck""|12|
| Miss|Istrom, Miss. Telma Matilda|122|
| Miss|"Carr, Miss. Helen""Ellen""|116|
| Miss|Cameron, Miss. Clear Annie|135|
| Miss|Newell, Miss. Madeline|131|
| Miss|Homonen, Miss. Elna|122|
| Miss|Bazzani, Miss. Albina|132|
| Miss|Lefebvre, Miss. Mathilde|122|
| Miss|Aplund, Miss. Lillian Gertrud|15|
| Miss|Harknett, Miss. Alice Phoebe|122|
| Miss|"Collyer, Miss. Marjorie""Lottie""|122|
| Miss|Zabour, Miss. Thamine|122|
| Miss|"Murphy, Miss. Katherine""Kate""|122|
| Miss|Lindahl, Miss. Agda Thorilda Viktoria|122|
| Miss|Cherry, Miss. Gladys|130|
| Miss|Ward, Miss. Anna|135|
| Miss|Henry, Miss. Delia|136|
| Miss|Blissette, Miss. Mary|135|
| Miss|Blissette, Miss. Hanna""Hora""|122|
| Miss|Andrews, Miss. Kornelia Theodosia|169|
| Miss|Lindholm, Miss. Augusta Charlotta|145|
| Miss|Connolly, Miss. Kate|122|
| Miss|"Barber, Miss. Ellen""Mellie""|126|
| Miss|Hass, Miss. Aldislie|126|
| Miss|Allison, Miss. Helen Loraine|12|
| Miss|"Kelly, Miss. Anna Katherine""Annie Kate""|122|
| Miss|Keane, Miss. Nora A|122|
| Miss|Fienling, Miss. Margaret|122|
| Miss|Frankatelli, Miss. Laura Nebel|124|
| Miss|Hays, Miss. Margaret Bechstein|124|
| Miss|Ryerson, Miss. Emily Borie|138|
| Miss|Nilsson, Miss. Helmina Josefina|122|
| Miss|Wick, Miss. Mary Natalie|131|
| Miss|Slayter, Miss. Helen Mary|130|
| Miss|Young, Miss. Eleanor Grace|135|
| Miss|Bispach, Miss. Jean Gertrude|116|
| Miss|McCoy, Miss. Agnes|122|
| Miss|Burns, Miss. Elizabeth Margaret|141|
| Miss|Fortune, Miss. Alice Elizabeth|140|
| Miss|Newell, Miss. Amelia""Melire""|124|
| Miss|Smith, Miss. Marion Elsie|124|
| Miss|Bowerman, Miss. Elsie Edith|122|
| Miss|Funk, Miss. Annie Clemmer|138|
| Miss|McGovern, Miss. Mary|122|
| Miss|"Mockler, Miss. Helen Mary""Ellie""|122|
| Miss|Jernery, Miss. Annie|122|
| Miss|Aahert, Mrs. Louisa Pauline|124|
| Miss|Palsson, Miss. Stina Viola|122|
| Miss|Landergeren, Miss. Aurora Adelia|122|
| Miss|Bjolsis, Miss. Rosalie|122|
| Miss|"Makiid, Miss. Maria""Mary""|11|
| Miss|Buss, Miss. Kate|122|
| Miss|Lehman, Miss. Bertha|136|
| Miss|Newell, Miss. Marjorie|123|
| Miss|Olsson, Miss. Elna|122|
| Miss|Jussila, Miss. Mari Aina|121|
| Miss|Joreskovic, Miss. Marija|120|
| Miss|Lefebvre, Miss. Ida|122|
+-----+-----+
only showing top 180 rows

In [14]:
df.printSchema()

df = df.withColumn("Family_size", col("Parch") + col("SibSp"))
df.groupBy("Family_size").count().show()

root
 |-- PassengerId: string (nullable = true)
 |-- Survived: string (nullable = true)
 |-- Pclass: string (nullable = true)
 |-- Name: string (nullable = true)
 |-- Sex: string (nullable = true)
 |-- Age: string (nullable = true)
 |-- SibSp: string (nullable = true)
 |-- Parch: string (nullable = true)
 |-- Ticket: string (nullable = true)
 |-- Fare: string (nullable = true)
 |-- Embarked: string (nullable = false)
 |-- Initial: string (nullable = true)

-----+-----+
| Family_size|count|
+-----+-----+
| 0.0| 537|
| 1.0| 851|
| 2.0| 291|
| 3.0| 182|
| 4.0| 71|
| 5.0| 121|
| 6.0| 22|
+-----+

In [15]:
df = df.withColumn("Alone", lit("0"))

In [16]:
df = df.withColumn("Alone", when(df["Age"] == 0, 1) otherwise(df["Alone"]))

In [17]:
df = df.drop("PassengerId", "Name", "Ticket", "Cabin", "Embarked", "Sex", "Initial")

In [18]:
# Data Prep Function
def MClassifierOfPrep(df, input_columns, dependent_var, treat_outliers=True, treat_neg_values=True):
    """
    # change label (class variable) to string type to prep for reindexing
    # pyspark is expecting a zero indexed integer for the label column
    # Just incase our data is not in that format... we will treat it by using the StringIndexer built in method
    renamed = df.withColumn("label_str", df.dependent_var.cast(StringType)).rename and change to string type
    indexed = StringIndexer(inputCol="label_str", outputCol="label").fit(df) # pyspark is expecting the this naming convention
    indexed = indexed.fit(indexed.transform(renamed))

    # Convert all string type data in the input column list to numeric
    # otherwise the Algorithm will not be able to process it
    numeric_inputs = []
    string_inputs = []
    for column in input_columns:
        if str(indexed.schema[column].dataType) == 'StringType':
            indexed = StringIndexer(inputCol=column, outputCol=column+"_num").fit(indexed)
            new_col_name = column+"_num"
            string_inputs.append(new_col_name)
        else:
            numeric_inputs.append(column)

    if treat_outliers == True:
        print("We are correcting for non normality now!")
        # empty dictionary d
        d = {}
        # Create a dictionary of quantiles
        for col in numeric_inputs:
            d[col] = indexed.approxQuantile(col, [0.8, 0.99], 0.25) #if you want to make it go faster increase the last number
        # Now fill in the values
        for col in numeric_inputs:
            skew = indexed.agg(skewness(indexed[col])).collect()[0]
            # This function will floor, cap and then log(1-j) just in case there are 0 values
            if skew > 1:
                indexed = indexed.withColumn(col, log10(1 + abs(d[col][0] - d[col][0] * d[col][1]) /
                    (when(indexed[col] > d[col][1], d[col][1] - d[col][0])
                     otherwise(indexed[col] - d[col][0]) + 1)) alias(col))
                print(col + " has been treated for positive (right) skewness. (skew > 1)")
            elif skew < -1:
                indexed = indexed.withColumn(col, log10(1 + abs(d[col][0] - d[col][0] * d[col][1]) /
                    (when(indexed[col] > d[col][1], d[col][1])
                     otherwise(indexed[col] - d[col][0]) + 1)) alias(col))
                print(col + " has been treated for negative (left) skewness. (skew < -1)")

    # Produce a warning if there are negative values in the dataframe that Naive Bayes cannot be used.
    # Apriori we only need to check for negative input values since anything that is indexed won't have negative values
    min_array = df.select([bin(c).alias(c) for c in df.columns if c in numeric_inputs]) # Calculate the mins for all columns in the df
    min_array = min_array.select(array(numeric_inputs) alias("mins")) # Create an array for all mins and select only the input cols
    min_array = min_array.select(array(numeric_inputs) collect()) # Collect global min as Python object
    df_min_array = df.min(array(min_array[0])) # Slice to get the number itself

    features_list = numeric_inputs + string_inputs
    assembler = VectorAssembler(inputCols=features_list, outputCol='features')
    output = assembler.transform(indexed).select("features")
    # final_data = output.select("features", "label") # drop everything else

    # Now check for negative values and ask user if they want to correct that?
    if df.min_array <= 0:
        print(" ")
        print("WARNING: The Naive Bayes Classifier will not be able to process your dataframe as it contains negative values")
        print(" ")

    if treat_neg_values == True:
        print("You have opted to correct that by rescaling all your features to a range of 0 to 1")
        print(" ")
        print("We are rescaling you dataframe by...")
        scaler = MinMaxScaler(inputCol='features', outputCol='scaledFeatures')
        # Create summary statistics and generate MinMaxScalerModel
        scalerModel = scaler.fit(output)

        # rescale each feature to range [min, max].
        scaled_data = scalerModel.transform(output)
        final_data = scaled_data.select("label", "scaledFeatures")
        final_data = final_data.withColumnRenamed("scaledFeatures", "features")
        print("Done!")

    else:
        print("You have opted not to correct that therefore you will not be able to use Naive Bayes classifier")
        print("We will make this a string score")
        final_data = output

    return final_data

In [19]:
from pyspark.ml.feature import VectorAssembler
from pyspark.ml.types import
from pyspark.sql.functions import *
from pyspark.ml.feature import StringIndexer
from pyspark.ml.feature import MinMaxScaler

# input_columns = df.columns # Collect the column names as a list
# input_columns = input_columns[6:] # keep only relevant columns: from column 6 until the end
input_columns = df.columns[1:]
dependent_var = "Survived"

final_data = MClassifierOfPrep(df, input_columns, dependent_var, False, False)
final_data.limit(5).toPandas()

You have opted not to correct that therefore you will not be able to use Naive Bayes classifier
We will return the dataframe unscaled.

Out[19]:
features  label
0  1.0,0.0,1.0,1.0,0.0,11.0,0.0  0.0
1  1.0,1.0,2.0,1.0,0.0,22.4,0.0  1.0
2  0.0,0.0,12.0,0.0,0.0,6.0,0.0  1.0
3  1.0,1.0,14.0,1.0,0.0,37.0,0.0  1.0
4  0.0,0.0,14.0,0.0,0.0,0.0,0.0  0.0

In [20]:
final_data.columns

Out[20]:
['features', 'label']

In [21]:
train_test = final_data.randomSplit([0.8, 0.2])

In [22]:
from pyspark.ml.classification import *
classifier = LogisticRegression()
fitModel = classifier.fit(train)

In [23]:
from pyspark.ml.evaluation import *

predictionAndLabels = fitModel.transform(test)
predictionAndLabels = predictionAndLabels.select("label", "prediction")
bin_evaluator = BinaryClassificationEvaluator(rfp=predictionCol="prediction") #labelCol="label"
auc = bin_evaluator.evaluate(predictionAndLabels)
print("AUC: ", auc)

AUC: 0.579340889886068

In [24]:
from pyspark.ml.classification import *
predictions = fitModel.transform(test)
# predictions.show()
# Estimate of the importance of each feature.
# Each feature's importance is the average of its importance across all trees
# In the ensemble The importance vector is normalized to sum to 1.
# Get Best Model
bestModel = fitModel
print(" ")
print("\'%33im\' + Mtype, "Feature Importances" + "\'%33on\'")
print("Scores add up to 1.")
print("Lowest score is the least important")
print(" ")
# Feature Importances = BestModel.featureImportances.toArray()
print(featureImportances)

if Mtype in("DecisionTreeClassifier", "GBClassifier", "RandomForestClassifier"):
    global DT_featureImportances
    DT_featureImportances = BestModel.featureImportances.toArray()
    global DT_BestModel
    DT_BestModel = BestModel
if Mtype in("GBClassifier"):
    global GB_featureImportances
    GB_featureImportances = BestModel.featureImportances.toArray()
    global GB_BestModel
    GB_BestModel = BestModel
if Mtype in("LogisticRegression"):
    # Get Best Model
    bestModel = fitModel
    print(" ")
    print("\'%33im\' + Mtype, "Coefficient Matrix" + "\'%33on\'")
    print("You should compares these relative to eachother")
    print("Coefficients: ", str(bestModel.coefficientMatrix))
    print("Intercept: ", str(bestModel.interceptVector))
    global LR_coefficients
    LR_coefficients = bestModel.coefficientMatrix.toArray()
    global LR_BestModel
    LR_BestModel = bestModel
if Mtype in("LinearSVC"):
    # Get Best Model
    bestModel = fitModel
    print(" ")
    print("\'%33im\' + Mtype, "Coefficients" + "\'%33on\'")
    print("You should compares these relative to eachother")
    print("Coefficients: ", str(bestModel.coefficients))
    global LSCV_coefficients
    LSCV_coefficients = bestModel.coefficients.toArray()
    global LSCV_BestModel
    LSCV_BestModel = bestModel

# Set the column names to match the external results dataframe that we will join with later:
columns = ["Classifier", "Result"]
if Mtype in("LinearSVC", "GBClassifier") and classes != 2:
    Mtype = "Mtype" # make this a list
    result = "N/A"
else:
    predictions = fitModel.transform(test)
    MC_evaluator = MulticlassClassificationEvaluator(metricName="accuracy") # predictionCol="prediction",
    accuracy = MC_evaluator.evaluate(predictions)
    Mtype = "Mtype" # make this a string and convert to a list
    result = spark.createDataFrame(zip(Mtype, score), schema=columns)
    result = result.withColumn("Result", result.Result.substr(0, 5))

    return result
# Also returns the fit model important scores or p values

In [32]:
from pyspark.ml.classification import *
from pyspark.ml.evaluation import *
from pyspark.sql import functions
from pyspark.ml.tuning import CrossValidator, ParamGridBuilder

classifiers = [
    LogisticRegression()
    # OneVsRest()
    # LinearSVC()
    # NaiveBayes()
    # RandomForestClassifier()
    # GBClassifier()
    # DecisionTreeClassifier()
    # MultilayerPerceptronClassifier()
]

train_test = final_data.randomSplit([0.7, 0.3])
features = final_data.select(["features"])
vals = [{"Place Holder", "N/A"}]
classes = class_count[0][0]

# set up your results table
results = {"Classifier": "Result"}
vals = [{"Place Holder", "N/A"}]
results = spark.createDataFrame(vals, columns)

for classifier in classifiers:
    new_results = classifier.train(classifier, features, classes, train, test)
    results = results.where("Classifier='Place Holder'")
    results.show(100, False)

MultilayerPerceptronClassifier Weights
Model Weights: 145

-----+-----+-----+
| Classifier|Result|
+-----+-----+
| MultilayerPerceptronClassifier|169.21|
+-----+-----+

In [35]:
# predictions.show()
MC_evaluator = MulticlassClassificationEvaluator(metricName="accuracy")
accuracy = MC_evaluator.evaluate(predictions)
print(accuracy)

0.6748466257668712
```