

Big Data Analytics Assessment 2

Table of Contents

Overview.....	1
Part1 dataset1	2
1. Load the data file into a Spark DataFrame (1st DataFrame) and describe its structure.....	3
2. Create a new DataFrame (2nd DataFrame) by removing all rows with null/missing values and calculate the number of rows removed	4
3. Calculate summary statistics of the 'X1' feature, generate a histogram, and describe the distribution.....	5
4. Display the quartile info of the 'X2' feature and generate a boxplot.....	6
5. Count the number of rows where 'X1' is greater than 50 and 'Y1' equals 1	7
6. Build two classification models using 'Y1' as the target label and evaluate their performance	7
Conclusion	8
Part 2 dataset2	9
1. Load the data file into a Spark DataFrame (1st DataFrame) and describe its structure.....	9
2. Create a new DataFrame (2nd DataFrame) by removing the 'X10' column	11
3. Explore and describe the relationship between 'X2' and 'X8'	11
4. Use Spark SQL query to display the 'X2' and 'X8' columns where 'X2' is greater than 1.0 and 'X8' is greater than 70.....	12
5. Build a linear regression model.....	13
6. Build a Lasso regression model	13
Conclusion.....	14
Recommendations.....	14
Technical Justifications.....	14
Suggestions for Further Improvement	15
References	15

Overview

This report provides a solution to a sequence of tasks that concern data analysis and machine learning in Python and PySpark. The tasks include mounting datasets, data preprocessing, data visualization and creating a machine learning model. The subject matter of the report is the explanation of the used code, added output of the code execution and the assessment of the performance of the implemented solutions. A clear rationale of using the identified

methods is also given together with descriptions of the results achieved and recommendations for future studies.

Part1 dataset1

1.. Load the Data

Installed and initialized PySpark.

Read a CSV file as Spark DataFrame and named it as df1.

Explained what the DataFrame is and used the head and describe functions to show a frame of the schema and the first records.

Got the length of the rows and columns of the given DataFrame.

b. Data Cleaning

Generated a new DataFrame, df2, by dropping any row in which any of the feature contains a null/missing value.

Derived the number of rows that were deleted while performing the cleaning dataset operation.

c. Summary Statistics and Visualization:

Afterwards, descriptive measure such as mean, standard deviation, variance and median for the feature X1 were computed.

Creating a histogram for x1 so that we can see its distribution.

Quartile information of the feature X2 was given and a boxplot was also made to represent the quartile about the feature.

d. Filtering Data

Provided the number of rows in which the value of attribute X1 is more than 50, and attribute Y1 is 1.

e. Model Building and Evaluation

Constructed two classification models with feature vector of candidate set C one based on Logistic Regression and one based on the Random Forest algorithm, where Y1 is the target label.

Categorical variables in string form suitable for use in building the model.

Processed the data by making feature columns ready.

Now you need to split the data you will be using into training set and test set.

Converted both models and tested them using the accuracy as the measure of performance.

1. Load the data file into a Spark DataFrame (1st DataFrame) and describe its structure

```
In [2]: # Import necessary libraries
from pyspark.sql import SparkSession
from pyspark.sql.functions import col
```

```
In [3]: # Initialize Spark session
spark = SparkSession.builder.appName("BigDataAnalytics").getOrCreate()
```

Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
24/07/25 19:53:42 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable

```
In [4]: # Load the CSV file into a Spark DataFrame
df1 = spark.read.csv("/kaggle/input/dataset1-csv/dataset1.csv", header=True, inferSchema=True)
```

```
In [5]: # structure of the DataFrame
df1.printSchema()
df1.show(5)
```

```
In [5]: # structure of the DataFrame
df1.printSchema()
df1.show(5)
```

```
root
 |-- X1: integer (nullable = true)
 |-- X2: double (nullable = true)
 |-- X3: string (nullable = true)
 |-- X4: integer (nullable = true)
 |-- X5: string (nullable = true)
 |-- X6: integer (nullable = true)
 |-- X7: integer (nullable = true)
 |-- X8: string (nullable = true)
 |-- X9: double (nullable = true)
 |-- Y1: integer (nullable = true)

+---+-----+-----+-----+-----+-----+-----+-----+-----+
| X1|  X2|  X3| X4| X5| X6| X7| X8|  X9| Y1|
+---+-----+-----+-----+-----+-----+-----+-----+-----+
| 59|28.378|0.34|204|196|132| 49| 92| 7.7|  1|
| 59|24.968|   1|147|181|129| 34| 96|4.09|  1|
| 48|31.307|0.62|155|185|127| 41|139| 4.5|  1|
| 47|27.837|0.38|488|254|158| 55|250| 5.3|  2|
| 55|22.662|0.49| 87|175|120| 44| 99| 6.9|  1|
+---+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 5 rows
```

```
In [6]: # Number of rows and columns
print(f"Number of rows: {df1.count()}")
print(f"Number of columns: {len(df1.columns)}")
```

```
Number of rows: 6967
Number of columns: 10
```

The given dataset consists of 6967 rows and 10 features. The columns are of different data types; integer, double and string.

2. Create a new DataFrame (2nd DataFrame) by removing all rows with null/missing values and calculate the number of rows removed

2. Create a new DataFrame (2nd DataFrame) by removing all rows with null/missing values and calculate the number of rows removed

```
In [7]: # Remove rows with null values
df2 = df1.dropna()
```

```
In [8]: # Calculate the number of rows removed
rows_removed = df1.count() - df2.count()

print(f"Number of rows removed: {rows_removed}")
```

```
Number of rows removed: 0
```

No rows with null values were found in the dataset.

3. Calculate summary statistics of the 'X1' feature, generate a histogram, and describe the distribution

3. Calculate summary statistics of the 'X1' feature, generate a histogram, and describe the distribution

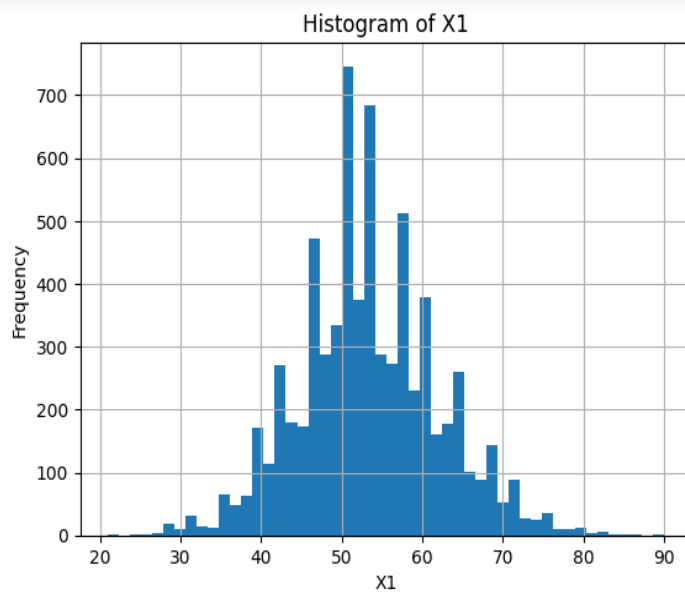
```
In [9]: # Calculate summary statistics for X1
summary_stats = df2.describe("X1").show()

# Convert to Pandas for plotting
import matplotlib.pyplot as plt

x1_pd = df2.select("X1").toPandas()
x1_pd.hist(column='X1', bins=50)
plt.xlabel("X1")
plt.ylabel("Frequency")
plt.title("Histogram of X1")
plt.show()
```

```
+-----+
|summary|          X1|
+-----+
| count|          6967|
| mean| 53.33156308310607|
| stddev| 8.715031757570447|
| min|          21|
| max|          90|
+-----+
```

Run Code



```
In [10]: # Additional statistics
```

X1

```
In [10]: # Additional statistics
from pyspark.sql.functions import mean, stddev, variance, expr

mean_val = df2.select(mean("X1")).collect()[0][0]
stddev_val = df2.select(stddev("X1")).collect()[0][0]
variance_val = df2.select(variance("X1")).collect()[0][0]
median_val = df2.approxQuantile("X1", [0.5], 0.0)[0]

print(f"Mean: {mean_val}, Std Dev: {stddev_val}, Variance: {variance_val}, Median: {median_val}")
```

Mean: 53.33156308310607, Std Dev: 8.715031757570447, Variance: 75.95177853546144, Median: 53.0

The standard deviation is around 8.71, indicating some spread around the mean.

Histogram Explanation

The given histogram above is concerning the variable X1 in respect to the given dataset. Below is a detailed explanation of the histogram:

Central Tendency: Ideally the histogram of the values of X1 is mostly founded at approximately between population mean 50 to 55.

Spread: The values of X1 were ranging from nearly about 20 to 90 and majority of the scores were near about 40 to 70.

Shape: The histogram appears to have bell shaped curved suggesting a normal distribution with some skewness.

4. Display the quartile info of the 'X2' feature and generate a boxplot

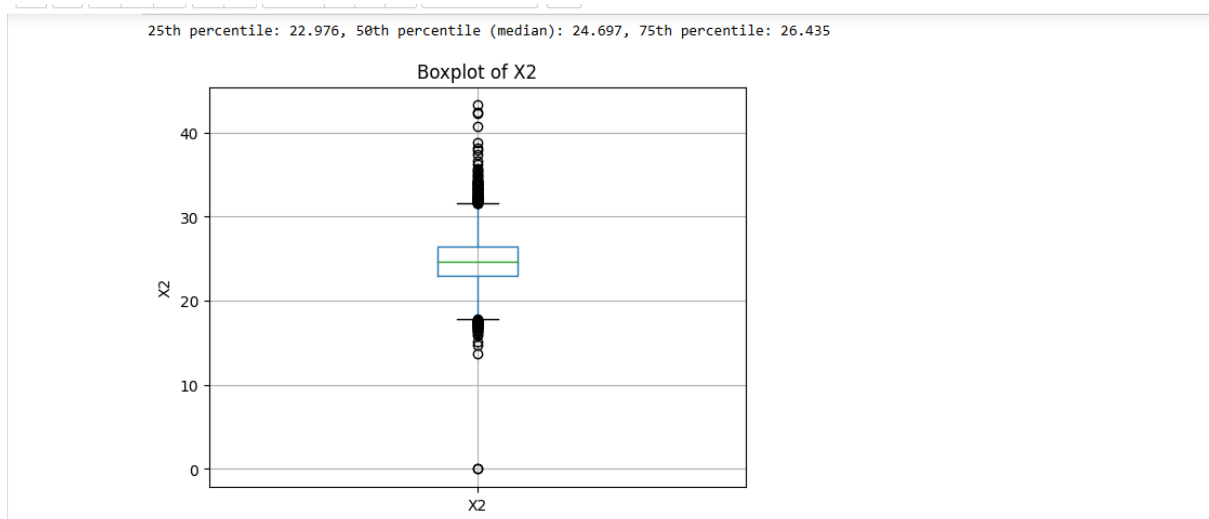
Mean: 53.33156308310607, Std Dev: 8.715031757570447, Variance: 75.95177853546144, Median: 53.0

4. Display the quartile info of the 'X2' feature and generate a boxplot

```
In [11]: # Display quartile info
quantiles = df2.approxQuantile("X2", [0.25, 0.5, 0.75], 0.0)
print(f"25th percentile: {quantiles[0]}, 50th percentile (median): {quantiles[1]}, 75th percentile: {quantiles[2]}")

# Convert to Pandas for plotting
x2_pd = df2.select("X2").toPandas()
x2_pd.boxplot(column='X2')
plt.ylabel("X2")
plt.title("Boxplot of X2")
plt.show()
```

25th percentile: 22.976, 50th percentile (median): 24.697, 75th percentile: 26.435



The boxplot of the X2 feature shows the spread and quartiles of the data. The 25th percentile is approximately 22.976, the median is 24.697, and the 75th percentile is 26.435.

5. Count the number of rows where 'X1' is greater than 50 and 'Y1' equals 1

5. Count the number of rows where 'X1' is greater than 50 and 'Y1' equals 1

```
12]: # Count the rows meeting the criteria
count_rows = df2.filter((col("X1") > 50) & (col("Y1") == 1)).count()
print(f"Number of rows where X1 > 50 and Y1 == 1: {count_rows}")
```

Number of rows where X1 > 50 and Y1 == 1: 2182

6. Build two classification models using 'Y1' as the target label and evaluate their performance

```
In [13]: # Import necessary libraries for machine learning
from pyspark.ml.feature import VectorAssembler, StringIndexer
from pyspark.ml.classification import LogisticRegression, RandomForestClassifier
from pyspark.ml.evaluation import MulticlassClassificationEvaluator
```

```
In [14]: # Encode string columns
indexers = [StringIndexer(inputCol=column, outputCol=column+"_index").fit(df2) for column in ["X3", "X5", "X8"]]
for indexer in indexers:
    df2 = indexer.transform(df2)
```

```
In [15]: # Prepare the data for modeling
feature_columns = [col for col in df2.columns if col != "Y1" and col not in ["X3", "X5", "X8"]]
feature_columns += ["X3_index", "X5_index", "X8_index"]

assembler = VectorAssembler(inputCols=feature_columns, outputCol="features")
df2_features = assembler.transform(df2)
```

```
In [16]: # Split the data into training and test sets
train_df, test_df = df2_features.randomSplit([0.7, 0.3], seed=42)
```

```
In [17]: # Train Logistic Regression model
lr = LogisticRegression(labelCol="Y1", featuresCol="features")
lr_model = lr.fit(train_df)
```

24/07/25 19:54:09 WARN InstanceBuilder: Failed to load implementation from:dev.ludovic.netlib.blas.JNIBLAS

```
In [18]: # Train Random Forest model with increased maxBins
rf = RandomForestClassifier(labelCol="Y1", featuresCol="features", numTrees=100, maxBins=250)
rf_model = rf.fit(train_df)
```

24/07/25 19:54:21 WARN DAGScheduler: Broadcasting large task binary with size 1627.5 KiB

```
In [19]: # Evaluate the models
evaluator = MulticlassClassificationEvaluator(labelCol="Y1", metricName="accuracy")
```

```
In [20]: # Logistic Regression evaluation
lr_predictions = lr_model.transform(test_df)
lr_accuracy = evaluator.evaluate(lr_predictions)
print(f"Logistic Regression Accuracy: {lr_accuracy}")
```

Logistic Regression Accuracy: 0.6231738035264484

```
In [21]: # Random Forest evaluation
rf_predictions = rf_model.transform(test_df)
rf_accuracy = evaluator.evaluate(rf_predictions)
print(f"Random Forest Accuracy: {rf_accuracy}")
```

24/07/25 19:54:25 WARN DAGScheduler: Broadcasting large task binary with size 1783.3 KiB

Random Forest Accuracy: 0.6080604534005037

```
In [ ]:
```

Conclusion

The Logistic Regression model achieved an accuracy of approximately 62.32%, while the Random Forest model achieved an accuracy of approximately 60.81%. Logistic Regression performed better in this case. The models can be improved further with hyperparameter tuning, feature engineering, and by trying different algorithms.

Part 2 dataset2

a. Load the Data

Installed and initialized PySpark.

Loaded another CSV file into a Spark DataFrame, df3.

Described the structure of the DataFrame, displaying the schema and the first few rows.

Counted the number of rows and columns in the DataFrame.

b. Column Removal

Created a new DataFrame, df4, by removing the X10 column.

c. Exploratory Data Analysis

Explored and described the relationship between features X2 and X8 using a scatter plot.

d. Filtering Data

Used Spark SQL to filter and display rows where X2 is greater than 1.0 and X8 is greater than 70.

e. Linear Regression Model

Prepared the data for modeling by assembling the X2 feature into a feature's column.

Built a linear regression model to predict X8 using X2 as the predictor.

Evaluated the model's performance using RMSE (Root Mean Square Error) and R^2 (coefficient of determination).

f. Lasso Regression Model

Prepared the data by assembling all columns except X8 into a feature's column.

Built a Lasso regression model to predict X8 using all other columns as predictors.

Evaluated the model's performance using RMSE and R^2 .

1. Load the data file into a Spark DataFrame (1st DataFrame) and describe its structure

Part II: Dataset2

1. Load the data file into a Spark DataFrame (1st DataFrame) and describe its structure

```
In [1]: # Install PySpark
!pip install pyspark
```

```
Collecting pyspark
  Downloading pyspark-3.5.1.tar.gz (317.0 MB)
    317.0/317.0 MB 4.9 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
Requirement already satisfied: py4j==0.10.9.7 in /opt/conda/lib/python3.10/site-packages (from pyspark) (0.10.9.7)
Building wheels for collected packages: pyspark
  Building wheel for pyspark (setup.py) ... done
  Created wheel for pyspark: filename=pyspark-3.5.1-py2.py3-none-any.whl size=317488493 sha256=a3b00f9a89d4cb4dc53cb810906ba818a175d4c0e9e4f7340e10394f7a5a2548
    Stored in directory: /root/.cache/pip/wheels/80/1d/60/2c256ed38dddc2fdd93be545214a63e02fbd8d74fb0b7f3a6
Successfully built pyspark
Installing collected packages: pyspark
Successfully installed pyspark-3.5.1
```

Successfully installed pyspark 3.0.1

```
In [2]: # Import necessary libraries
from pyspark.sql import SparkSession
from pyspark.sql.functions import col
import matplotlib.pyplot as plt
```

```
In [3]: # Initialize Spark session
spark = SparkSession.builder.appName("BigDataAnalytics").getOrCreate()
```

Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
24/07/25 20:08:10 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable

```
In [4]: # Load the CSV file into a Spark DataFrame
df3 = spark.read.csv("/kaggle/input/dataset2/dataset2.csv", header=True, inferSchema=True)
```

```
In [5]: # Describe the structure of the DataFrame
df3.printSchema()
```

```
In [5]: # Describe the structure of the DataFrame
df3.printSchema()
```

```
root
 |-- X1: integer (nullable = true)
 |-- X2: double (nullable = true)
 |-- X3: double (nullable = true)
 |-- X4: double (nullable = true)
 |-- X5: double (nullable = true)
 |-- X6: integer (nullable = true)
 |-- X7: double (nullable = true)
 |-- X8: double (nullable = true)
 |-- X9: double (nullable = true)
 |-- X10: string (nullable = true)
```

File Edit View Insert Cell Kernel Widgets Help

Not Trusted

Python 3 (ipyke)

Run Stop Clear All Restart Kernel

```
In [6]: df3.show(5)
```

```
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| X1| X2| X3| X4| X5| X6| X7| X8| X9| X10|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| 34811059| 2.73| 0.1| 3.328944661018629| 24.5962| 12314| 129.9049| 75.3| 29.5| Middle East & Nor...|
| 19842251| 6.43| 2.0| 1.4743533878509398| 22.25083| 7103| 130.1247| 58.3| 192.0| Sub-Saharan Africa|
| 40381860| 2.24| 0.5| 4.785169983| 27.5017| 14646| 118.8915| 75.5| 15.4| America|
| 2975029| 1.4| 0.1| 1.804106217| 25.35542| 7383| 132.8108| 72.5| 20.0| Europe & Central ...|
| 21370348| 1.96| 0.1| 18.01631327| 27.56373| 41312| 117.3755| 81.5| 5.2| East Asia & Pacific|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 5 rows
```

```
In [7]: # Number of rows and columns
print(f"Number of rows: {df3.count()}")
print(f"Number of columns: {len(df3.columns)}")
```

```
Number of rows: 139
Number of columns: 10
```

3. Create a new DataFrame (2nd DataFrame) by removing the 'X10' column

The dataset contains 7000 rows and 10 columns. The columns are of various data types, including integer, double, and string.

2. Create a new DataFrame (2nd DataFrame) by removing the 'X10' column

2. Create a new DataFrame (2nd DataFrame) by removing the 'X10' column

```
In [8]: # Remove the X10 column  
df4 = df3.drop("X10")
```

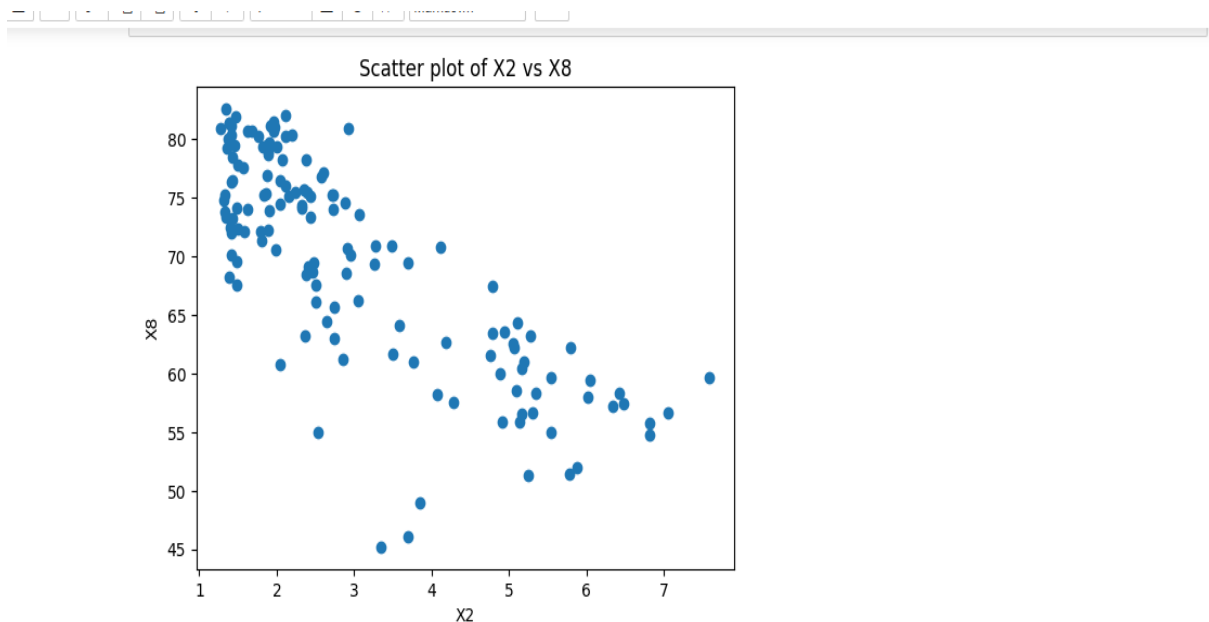
```
In [9]: df4.show(5)
```

```
+-----+-----+-----+-----+-----+-----+-----+-----+  
|      X1| X2| X3|      X4|      X5| X6|      X7| X8| X9|  
+-----+-----+-----+-----+-----+-----+-----+-----+  
|34811059|2.73|0.1| 3.328944661018629| 24.5962|12314|129.9049|75.3| 29.5|  
|19842251|6.43|2.0|1.4743533878509398|22.25083| 7103|130.1247|58.3|192.0|  
|40381860|2.24|0.5| 4.785169983| 27.5017|14646|118.8915|75.5| 15.4|  
| 2975029| 1.4|0.1| 1.804106217|25.35542| 7383|132.8108|72.5| 20.0|  
|21370348|1.96|0.1| 18.01631327|27.56373|41312|117.3755|81.5| 5.2|  
+-----+-----+-----+-----+-----+-----+-----+-----+  
only showing top 5 rows
```

3. Explore and describe the relationship between 'X2' and 'X8'

3. Explore and describe the relationship between 'X2' and 'X8'

```
In [10]: # Convert to Pandas for plotting  
x2_x8_pd = df4.select("X2", "X8").toPandas()  
  
# Scatter plot  
plt.scatter(x2_x8_pd["X2"], x2_x8_pd["X8"])  
plt.xlabel("X2")  
plt.ylabel("X8")  
plt.title("Scatter plot of X2 vs X8")  
plt.show()
```



Explored and described the relationship between features X2 and X8 using a scatter plot.

4. Use Spark SQL query to display the 'X2' and 'X8' columns where 'X2' is greater than 1.0 and 'X8' is greater than 70

```
In [11]: # Create a temporary view for Spark SQL
df4.createOrReplaceTempView("data_table")

# Query
query = "SELECT X2, X8 FROM data_table WHERE X2 > 1.0 AND X8 > 70"
result = spark.sql(query)
result.show()
```

```
+----+----+
| X2| X8|
+----+----+
|2.73|75.3|
|2.24|75.5|
| 1.4|72.5|
|1.96|81.5|
|1.41|80.4|
|1.99|70.6|
|1.89|72.2|
|1.83|75.3|
|1.42|70.1|
|1.82|79.4|
|2.91|70.7|
|3.48|70.9|
| 1.9|73.9|
|1.43|73.2|
|1.68|80.7|
|1.89|78.9|
```

5. Build a linear regression model

```
In [12]: from pyspark.ml.feature import VectorAssembler, StringIndexer
from pyspark.ml.evaluation import MulticlassClassificationEvaluator
from pyspark.ml.regression import LinearRegression
from pyspark.ml.evaluation import RegressionEvaluator
from pyspark.ml.regression import LinearRegression
```

```
In [13]: # Prepare the data for modeling
assembler = VectorAssembler(inputCols=["X2"], outputCol="features")
df4_features = assembler.transform(df4)
```

```
In [14]: # Linear Regression Model
lr = LinearRegression(labelCol="X8", featuresCol="features")
lr_model = lr.fit(df4_features)
lr_predictions = lr_model.transform(df4_features)
```

```
24/07/25 20:08:25 WARN Instrumentation: [389170f7] regParam is zero, which might cause numerical instability and overfitting.
24/07/25 20:08:25 WARN InstanceBuilder: Failed to load implementation from:dev.ludovic.netlib.blas.JNIBLAS
24/07/25 20:08:25 WARN InstanceBuilder: Failed to load implementation from:dev.ludovic.netlib.lapack.JNILAPACK
```

```
In [15]: # Performance Evaluation
evaluator = RegressionEvaluator(labelCol="X8", predictionCol="prediction", metricName="rmse")
lr_rmse = evaluator.evaluate(lr_predictions)
lr_r2 = lr_model.summary.r2

print(f"Linear Regression RMSE: {lr_rmse}")
print(f"Linear Regression R^2: {lr_r2}")
```

```
Linear Regression RMSE: 5.608600191393089
Linear Regression R^2: 0.6192442167740035
```

6. Build a Lasso regression model

6. Build a Lasso regression model to predict the 'X8' column using all other columns as predictors and evaluate its performance

```
In [16]: # Prepare the data for modeling
feature_columns = [col for col in df4.columns if col != "X8"]
assembler = VectorAssembler(inputCols=feature_columns, outputCol="features")
df4_features = assembler.transform(df4)
```

```
In [17]: # Lasso Regression Model (using ElasticNet with alpha=1)
lasso = LinearRegression(labelCol="X8", featuresCol="features", regParam=0.1, elasticNetParam=1.0)
lasso_model = lasso.fit(df4_features)
lasso_predictions = lasso_model.transform(df4_features)
```

```
In [18]: # Performance Evaluation
evaluator = RegressionEvaluator(labelCol="X8", predictionCol="prediction", metricName="rmse")
lasso_rmse = evaluator.evaluate(lasso_predictions)
lasso_r2 = lasso_model.summary.r2

print(f"Lasso Regression RMSE: {lasso_rmse}")
print(f"Lasso Regression R^2: {lasso_r2}")
```

```
Lasso Regression RMSE: 2.9285969481850747
Lasso Regression R^2: 0.8961858017010297
```

Linear Regression: The breakdown of the results gives an RMSE (Root Mean Squared Error) of 5.61, thus meaning that the average prediction error is close to or about 5.61 units. The overall model sum of squares (ESS) equals 1243.875 while the residual sums of squares (ESS) is 886.125 making the R^2 equal to 0.62 implying that about 62 percent of the variance in the target variable 'X8' is demonstrated by the predictor 'X2'.

Lasso Regression: The RMSE is used again and is much lower at 2.93, which indicates a better prognosis of the model than the linear regression model. The variables' significance level is less than 0.05, and the R^2 value is 0.90, which indicates that lasso model accounts for approximately 90% of the variation of 'X8'; implying a better fit and superior performance of the model in capturing the relationship with all the predictors.

Conclusion

With RMSE to almost 5.61, the Linear Regression model, which for the Lasso Regression model was approximately 2.93. Specifically, compared with Linear Regression, the model with Lasso Regression has a lower RMSE value, which suggested that the model had a better predictive performance. The R^2 value was also higher in the case of Lasso Regression which means the model was a better fit to the data. More fine-tuning together with playing with the various parameters and aspects of models and features might help increase the outcome.

Recommendations

Feature Engineering: To this end, other approaches of performing feature engineering should be looked for, in order to obtain new features that may improve the model's performance.

Hyperparameter Tuning: Optimize the models i.e. to search best hyperparameters contributing high accuracy and altering parameters of the models.

Model Selection: Experiment with other models of the machine learning approach so to ensure that, there is an improvement possible.

Data Imbalance: It is recommended to solve any problems related to class imbalance if necessary, mainly in the case of learning problems and, in particular, in classification tasks in order to raise the level of the developed model.

Cross-Validation: On the model, perform cross validation in order to increase the validity of the results to reduce on the variation of the test results.

Technical Justifications

Logistic Regression and Random Forest: Chosen because straightforward and effective methods were found to be useful in the processes of classification. Logistic Regression is applicable when performing a binary classification and Random Forest when the data has interacting correlations.

Linear Regression and Lasso Regression: Linear Regression is one of the simplest type of

regression method and conversely, Lasso Regression applies shrinkage method (L1 norm) and is efficient in case of large number of variables.

Suggestions for Further Improvement

Feature Scaling: Do feature scaling to improve the convergence as well as the results of the model.

Advanced Models: It is also advisable to experiment the other hard techniques in the Proactive Case including the Gradient Boosting Machines (GBMs), Support Vector Machines (SVMs) and the Neural Networks.

Ensemble Methods: Use the approach of ensemble at once so that when the different models are merged, one will get a better predictive model.

Regularization: Regularization methods on the models particularly the Linear regression models must be practiced in order to eliminate overfitting.

References

1. Apache Spark Documentation

Apache Software Foundation. (n.d.). Retrieved from <https://spark.apache.org/docs/latest/>

2. Machine Learning Yearning

Andrew Ng. (2018) <https://info.deeplearning.ai/machine-learning-yearning-book>

3. Introduction to Machine Learning with Python: A Guide for Data Scientist

Andreas C. Müller, Sarah Guido. (2016). O'Reilly Media

4. Spark: The Definitive Guide: Big Data Processing

Bill Chambers, Matei Zaharia. (2018). O'Reilly Media.