

Faculty of Mathematical, Physical and Natural Sciences of Tunis Faculté des Sciences Mathématiques, Physiques et Naturelles de Tunis كلية العلوم للرياضيات والفيزياء والطبيعيات بتونس

Academic Project Report

Machine Learning Pipeline in PySpark

Realized By : Aziz Ayadi

Louai Azzouni

Rayen Ben Fathallah

Supervised By: Manel Zekri

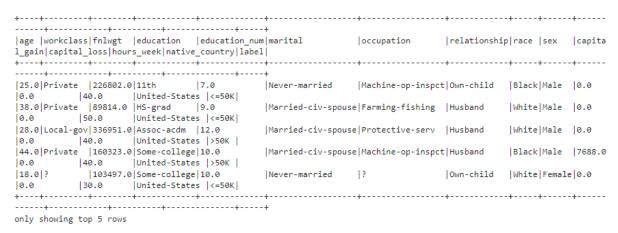
Academic Year : 2024 - 2025

Objective

The primary objective of this project is to build a machine learning model in PySpark to predict whether an individual's income exceeds \$50,000 per year based on census data. This guide uses the "Census Income" dataset, incorporating data preprocessing, model building, evaluation, and hyperparameter tuning.

Dataset Overview

The dataset, referred to as "adult.csv," includes features such as age, work class, education level, occupation, marital status, and capital gain. The target variable, income, is a binary classification label that indicates whether a person's income is above or below \$50,000.



Steps

1. Initialize SparkContext and SQLContext

- SparkContext (sc): Provides the core functionality for Spark jobs. In PySpark, SparkContext is required to establish the initial connection with Spark.
- SQLContext (sqlContext): Enables Spark to interface with SQL-based data sources. It's initialized to read the dataset and perform SQL-like transformations.

```
# initialize spark context
import pyspark
from pyspark import SparkContext
sc =SparkContext()

# initialize sql context
from pyspark.sql import SQLContext
sqlContext = SQLContext(sc)
```

2. Load and Inspect Dataset

- Load Dataset: The CSV file is read into a DataFrame using sqlContext.read.csv(), with header=True and inferSchema=True to automatically infer column types.
- Schema Inspection: The .printSchema() function provides an overview of data types and nullability, showing features such as age, workclass, fnlwgt, education, marital-status, and income.

```
#Todo: using SQLContext to read csv and assign to dataframe
df = sqlContext.read.csv("adult.csv", header=True, inferSchema= True)
#Todo:printSchema
df.printSchema()
root
|-- age: integer (nullable = true)
 |-- workclass: string (nullable = true)
 |-- fnlwgt: integer (nullable = true)
 -- education: string (nullable = true)
 -- educational-num: integer (nullable = true)
 -- marital-status: string (nullable = true)
 -- occupation: string (nullable = true)
 |-- relationship: string (nullable = true)
 -- race: string (nullable = true)
 -- gender: string (nullable = true)
 |-- capital-gain: integer (nullable = true)
 -- capital-loss: integer (nullable = true)
 |-- hours-per-week: integer (nullable = true)
 |-- native-country: string (nullable = true)
 |-- income: string (nullable = true)
```

3. Data Preprocessing

 Column Renaming: Renaming columns for simplicity and consistency in naming conventions.

```
# Run the cell to rename the columns properly:
cols = ['age','workclass','fnlwgt','education','education_num','marital','occupation','relationship','race','sex','capi
#note income -renamed-> as label
df=df.toDF(*cols)
```

 Type Conversion: Conversion of continuous features to FloatType using a custom function. Important continuous features include age, fnlwgt, capital_gain, education_num, capital_loss, and hours_week.

```
# Import all from `sql.types`
from pyspark.sql.types import *

# Write a custom function to convert the data type of DataFrame columns
def convertColumn(df, names, newType):
    for name in names:
        df = df.withColumn(name, df[name].cast(newType))
    return df

# List of continuous features
CONTI_FEATURES = ['age', 'fnlwgt', 'capital_gain', 'education_num', 'capital_loss', 'hours_week']
# Convert the type
df = convertColumn(df, CONTI_FEATURES, FloatType())
# Check the dataset
df.printSchema()
```

```
root
 |-- age: float (nullable = true)
 |-- workclass: string (nullable = true)
 |-- fnlwgt: float (nullable = true)
 |-- education: string (nullable = true)
 -- education num: float (nullable = true)
  -- marital: string (nullable = true)
 -- occupation: string (nullable = true)
 -- relationship: string (nullable = true)
 |-- race: string (nullable = true)
 -- sex: string (nullable = true)
 -- capital_gain: float (nullable = true)
 -- capital_loss: float (nullable = true)
 -- hours_week: float (nullable = true)
 |-- native country: string (nullable = true)
 |-- label: string (nullable = true)
```

 Descriptive Statistics: Use .describe() to generate basic statistical summaries such as mean, standard deviation, min, and max for each column.

 Grouping and Aggregation: For example, computing the average capital_gain by marital status provides insights into capital gains across different marital groups.

```
#example
df.groupby('marital').agg({'capital_gain': 'mean'}).show()

***The state of the state o
```

 Crosstab Analysis: Creates a cross-tabulation between age and income (<=50K or >50K), providing a count of individuals in each age group by income category.

```
#todo crosstab computation
 df.crosstab('age', 'label').sort("age_label").show()
+----+
|age_label|<=50K|>50K|
+----+
    17.0 595 0
    18.0 862 0
    19.0 | 1050 | 3 |
               1
    20.0 1112
    21.0 | 1090 | 6 |
    22.0 1161 17
    23.0 | 1307 | 22 |
    24.0 1162 44
    25.0 1119 76
    26.0 | 1068 | 85 |
    27.0 | 1117 | 115 |
    28.0 | 1101 | 179 |
    29.0 | 1025 | 198 |
```

| 30.0| 1031| 247| | 31.0| 1050| 275| | 32.0| 957| 296| | 33.0| 1045| 290| | 34.0| 949| 354| | 35.0| 997| 340| | 36.0| 948| 400| +-----+----+

> Null Handling and Row Exclusion: Removes rows with missing values using .dropna(). For example, records from native_country with only one entry (like "Holand-Netherlands") are excluded to avoid issues in crossvalidation.

```
#Drop null vals
df = df.dropna()

df_remove = df.filter(df.native_country!='Holand-Netherlands')
```

```
#TODO : follow the above instruction
 df.filter(df.native_country == 'Holand-Netherlands').count()
df.groupby('native_country').agg({'native_country': 'count'}).sort(asc("count(native_country)")).show()
  -----
     native_country|count(native_country)|
 Holand-Netherlands
            Hungary
           Honduras
                                      20
           Scotland
                                      21
Outlying-US(Guam-...
                                      23
        Yugoslavia
                                      23
                                      23
     Trinadad&Tobago
                                      27
           Cambodia
                                      28
               Hong
                                      30
           Thailand
                                      30
            Ireland
                                      37
             France
                                      38
             Ecuador
                                      45
                                      46
             Greece
                                      49
           Nicaragua
                                      49
               Iran
                                      59
              Taiwan
                                      65
                                                                          Activer Windows
           Portugal|
                                                                          Accédez aux paramètres pour activer Windows
only showing top 20 rows
```

4. Feature Engineering

 Age Transformation: Create a new feature, age_square, to capture potential non-linear relationships between age and income.

```
#TODO: # Add age square
  from pyspark.sql.functions import *
  # 1 Select the column
  age square = df.select(col("age")**2)
  # 2 Apply the transformation and add it to the DataFrame
  df = df.withColumn("age_square", col("age")**2)
  df.printSchema()
root
 |-- age: float (nullable = true)
 |-- workclass: string (nullable = true)
 |-- fnlwgt: float (nullable = true)
 -- education: string (nullable = true)
 |-- education_num: float (nullable = true)
 -- marital: string (nullable = true)
 |-- occupation: string (nullable = true)
 |-- relationship: string (nullable = true)
 -- race: string (nullable = true)
 -- sex: string (nullable = true)
 -- capital_gain: float (nullable = true)
 |-- capital_loss: float (nullable = true)
 |-- hours_week: float (nullable = true)
 |-- native_country: string (nullable = true)
 |-- label: string (nullable = true)
 |-- age_square: double (nullable = true)
```

5. Building the Data Processing Pipeline

 Categorical Feature Encoding: Use StringIndexer and OneHotEncoderEstimator to encode categorical features, including workclass, education, marital, occupation, relationship, race, sex, and native country.

• Label Encoding: The target variable income is converted into a binary label (label) using StringIndexer.

```
# 2. Index the label feature
# Convert label into label indices using the StringIndexer
label_stringIdx = StringIndexer(inputCol="label", outputCol="newlabel")
stages += [label_stringIdx]
```

 Feature Vectorization: Use VectorAssembler to combine all transformed categorical and continuous features into a single vector, features.

```
# 3. Add continuous variable
assemblerInputs = [c + "classVec" for c in CATE_FEATURES] + CONTI_FEATURES

# 4. Assemble the steps
assembler = VectorAssembler(inputCols=assemblerInputs, outputCol="features")
stages += [assembler]
```

• **Pipeline Creation**: Instantiate a Pipeline that will execute all the stages sequentially.

```
# Create a Pipeline.
pipeline = Pipeline(stages=stages)
pipelineModel = pipeline.fit(df_remove)
model = pipelineModel.transform(df_remove)
model
```

DataFrame[age: float, workclass: string, fnlwgt: float, education: string, education_num: float, marital: string, occup ation: string, relationship: string, race: string, sex: string, capital_gain: float, capital_loss: float, hours_week: f loat, native_country: string, label: string, age_square: double, workclassIndex: double, workclassClassVec: vector, edu cationIndex: double, educationclassVec: vector, maritalIndex: double, maritalclassVec: vector, occupationIndex: double, occupationclassVec: vector, relationshipIndex: double, relationshipclassVec: vector, raceIndex: double, raceclassVec: vector, newlabel: do uble, features: vector]

6. Model Building - Logistic Regression

 Dataset Preparation: Split the dataset into an 80/20 traintest split for model validation.

```
# Split the data into train and test sets
train_data, test_data = df_train.randomSplit([.8,.2],seed=1234)
```

 Logistic Regression Initialization: Initialize Logistic Regression with labelCol="label" and featuresCol="features".

```
# Print the coefficients and intercept for Logistic regression
print("Coefficients: " + str(linearModel.coefficients))
print("Intercept: " + str(linearModel.intercept))
```

Coefficients: [-0.06629458791643376,-0.15218373640701383,-0.053913615606758065,-0.16967312430449774,-0.12115314684082322,0.13 25974961176333,0.1943887659566995,-0.6386553259560794,-0.20168892525823268,-0.06643613691435478,0.22587144074752571,0.3784635 752251114,-0.0044245321889490345,-0.2958940967195082,-0.011315453541718349,-0.3286032800269526,-0.4220383458346703,0.57488818 55407, -0.3496608365399488, -0.20341828378529217, -0.21097882986337838, -0.16094910982880872, -0.10229984961820654, 0.1935414055955, -0.10414065955, -0.10414065965, -0.1041406696, -0.10414066, -0.1041406, -0.694. - 0.1707732587338623. - 0.1265004208587058. - 0.3046016757809549. - 0.32401360966564163. 0.11297915494496764. 0.12425063820676888.5433645026419642,-0.1908446548558885,-0.06449555808439705,-0.269348178043779,0.1680087916932683,-0.1253754753434063,-0.386858 1630578742, -0.18369433515375264, 0.02912182143548542, -0.08956223920758943, -0.2708959392074213, -0.03487727693207883, -0.32413919, -0.014136443644, -0.01414143644, -0.0141414364, -0.01414144, -0.0141414364, -0.01414144, -0.01414144, -0.01414144, -0.01414144, -0.01414144, -0.01414144, -0.01414144, -0.01414144, -0.0141414, -0.014144, -0.014415336244, -0.04758250791733353, -0.1173475422607489, 0.08015333921255102, -0.36180606787408326, -0.33618210864537906, -0.1663353505126, -0.16633505126, -0.16633505126, -0.16633505126, -0.16633505126, -0.16633505126, -0.16633505126, -0.16633505126, -0.16633505126, -0.16633505126, -0.16633505126, -0.16633505126, -0.16633505126, -0.16633505126, -0.1663505126, -0.1663505126, -0.1663505126, -0.1663505126, -0.1663505126, -0.1663505126, -0.1663505126, -0.1663505126, -0.1663505126, -0.1663505126, -0.1663505126, -0.1663505126, -0.1663505126, -0.1663505126, -0.1663505126, -0.1663505126, -0.166505126, -0.19376112,-0.006091174986221248,-0.45442144073339164,-0.05071249102702611,-0.3521378742759532,-0.2426920862206138,-0.2172135389 7728612,-0.5821780298068088,-0.17111959024994547,-0.00021772139996622131,-0.15792241012962724,-0.13336347352772165,-0.1244484 9686729739, -0.4299324718213068, -0.4802401343425621, -0.3509031990149178, 0.15554349625772512, 0./94484256440286982, -0.2892702852575532,-0.19176517240952143,0.35513884661883627,-0.47121791074628866,0.16608265034986552,-0.6627801214856749,-0.40786931479756 95,-0.22434028566679146,-0.21468806117681136,0.007013088722353848,1.1642146321728235e-07,2.1Å996£359910857e-0550102852355900|| 4261892,0.00022176143684426227,0.008705658401723285] Intercept: -2.0325289129477486

7. Model Evaluation

- Accuracy Calculation: Compute accuracy by comparing predicted labels to true labels in the test set.
- ROC and Area Under ROC: Evaluate the model's performance using the ROC curve, where an area under the ROC of 0.89 indicates strong model performance.

```
#We need to look at the accuracy metric to see how well (or bad) the model performs.

def accuracy_m(model):
    predictions = model.transform(test_data)
    cm = predictions.select("label", "prediction")
    acc = cm.filter(cm.label == cm.prediction).count() / cm.count()
    print("Model accuracy: %.3f%%" % (acc * 100))

accuracy_m(model = linearModel)
```

Model accuracy: 81.887%

```
### Use ROC
from pyspark.ml.evaluation import BinaryClassificationEvaluator

# Evaluate model
evaluator = BinaryClassificationEvaluator(rawPredictionCol="rawPrediction")
print(evaluator.evaluate(predictions))
print(evaluator.getMetricName())
```

8. Hyperparameter Tuning with Cross-Validation

 Parameter Grid: Construct a ParamGridBuilder to test two values for the regParam parameter.

 CrossValidator: Set up a 5-fold cross-validation to identify the optimal regularization parameter and evaluate model accuracy.

Time to train model: 1079.451 seconds

• **Best Model Accuracy**: The cross-validated model achieves an accuracy of 84.82%, demonstrating the effectiveness of tuning.

```
#accuracy of cv selected model
accuracy_m(model = cvModel)
```

Model accuracy: 84.820%

Conclusion

The completed PySpark ML pipeline demonstrates efficient handling of categorical data, preprocessing of features, and model evaluation in Spark. With a final accuracy of around 84.82% on the test set, this model effectively predicts income levels based on census data.

This pipeline can be expanded with additional feature engineering, model evaluation metrics, and alternative algorithms if needed. The entire process is scalable and easily adaptable for similar classification tasks in large datasets using Spark.