# Lab 5 - Assignment

## MM16B023[1]

[a]*Indian Institute of Technology Madras*

*Keyword:* BigQuery, StandardScaler, RandomForest, XGBoost

*Abstract:* This paper presents the solutions to third assignment of the Big Data Laboratory course (CS4830) at *IIT Madras.* All the notations used are as according with the textbook Mining of massive data sets by Anand Rajaraman

## Problem 1

Count using BigQuery the number of Iris Virginica flowers which have sepal width greater than 3 cm and petal length smaller than 2 cm

**Solution:**



$\implies$ There are no flowers with sepal width greater than 3 cm and petal length smaller than 2 cm

## Problem 2

Train a classification model on the dataset and report the accuracy for different preprocessing techniques and models. Provide the details of data exploration and feature engineering steps

**Solution:**

```
### Written by H.Vishal MM16B023 at 1:43 PM on 29th Feb, 2020 ###

from __future__ import print_function
from pyspark.context import SparkContext
```

```python
from pyspark.ml.feature import VectorAssembler
from pyspark.ml.regression import LinearRegression
from pyspark.sql.session import SparkSession
from pyspark.ml import Pipeline
from pyspark.ml.feature import PCA
from pyspark.ml.linalg import Vectors
from pyspark.mllib.regression import LabeledPoint
from pyspark.mllib.util import MLUtils
import numpy as np
from pyspark.ml.feature import StandardScaler
import pyspark.sql.functions as f
import pyspark.sql.types
from pyspark.sql import Row
from pyspark.sql.types import DoubleType
from pyspark.ml import Pipeline
from pyspark.ml.classification import DecisionTreeClassifier
from pyspark.ml.feature import StringIndexer, VectorIndexer
from pyspark.ml.evaluation import MulticlassClassificationEvaluator
from pyspark.ml.classification import RandomForestClassifier
from pyspark.ml.feature import IndexToString, StringIndexer, VectorIndexer
from pyspark.ml.evaluation import MulticlassClassificationEvaluator
from pyspark.ml.feature import PCA

sc = SparkContext()
spark = SparkSession(sc)

## Reading Dataframe

spark_df = spark.read.format("bigquery").option("table", "lab5.table_iris").load().toDF("sl", "sw", "pl
clean_data = spark_df.withColumn("label", spark_df["labclass"])

cols = spark_df.drop('labclass').columns

assembler = VectorAssembler(inputCols=cols, outputCol = 'features')
labelIndexer = StringIndexer(inputCol="labclass", outputCol="indexedLabel").fit(spark_df)

## Standardize the columns

scaler = StandardScaler(inputCol="features", outputCol="scaledFeatures", withStd=False, withMean=True)

## Principal component analysis

pca = PCA(k=3, inputCol='scaledFeatures', outputCol='pcaFeature')

(trainingData, testData) = spark_df.randomSplit([0.8, 0.2])

## Training a RandomForest model

rf = RandomForestClassifier(labelCol="indexedLabel", featuresCol="pcaFeature", numTrees=10)

## Retrieve orginal labels from indexed labels

labelConverter = IndexToString(inputCol="prediction", outputCol="predictedLabel",
                               labels=labelIndexer.labels)

## Modying indexers and forest in a Pipeline
```

```
pipeline = Pipeline(stages=[labelIndexer, assembler,scaler,pca, rf, labelConverter])

## Train the ML model

model = pipeline.fit(trainingData)

## Predictions

predictions = model.transform(testData)
evaluator = MulticlassClassificationEvaluator(
    labelCol="indexedLabel", predictionCol="prediction", metricName="accuracy")
accuracy = evaluator.evaluate(predictions)
print("Test set error fraction = %g" % (1.0 - accuracy))
```

**Output**

| Model | No. of comp | Test set error |
|---|---|---|
| Standardization & | 2 | 16.129 % |
| RandomForest | 3 | 7.407 % |
| Normalization & | 2 | 12.424 % |
| xgboost | 3 | 4.975% |

**Discussion**

- As the number of principal components included increases, better the fit is which indeed results in better performance on the given test data set.

- Ideally, one would keep increasing the number of principal components until the error converges. This would give the optimal number of the components to be included.

- Comparison is made for the above 2 different techniques preprocessing techniques and models. In this case, the $2^{nd}$ choice seems to perform better. In general, one needs try out at least a few different combinations as there in no set rule as to which is going to perform better for a particular problem.



**Fig:** PCA $(k = 3)$ with StandardScaler using RandomForest