Machine Learning With Apache Spark

By: Nikhilesh Vashisht

Table Of Contents

| Machine Learning with Apache Spark: | 4 |
|---|----|
| 1: PREDICT DIABETES | 5 |
| Data Dictionary: | 5 |
| Selecting Features: | 6 |
| Model Selection: | 6 |
| Hyper Parameters | 6 |
| maxDepth: | 6 |
| maxBins: | 6 |
| Accuracy: | 8 |
| 2: PREDICT CANCER | 9 |
| Data Dictionary: | 9 |
| Selecting features: | 10 |
| Model Selection | 10 |
| Hyper Parameters | 10 |
| Tolerance: | 10 |
| MaxIter: | 11 |
| RegParam: | 11 |
| Accuracy: | 12 |
| 3: KAGGLE COMPETITION -1 QUORA QUESTION PAIRS | 13 |
| Data Dictionary: | 13 |
| Feature Engineering: | 13 |
| Model Selection | 14 |
| Hyper Parameters | 14 |
| maxDepth: | 14 |
| maxBins: | 14 |
| Accuracy: | 15 |
| 4: KAGGLE COMPETITION - 2: TWEET SENTIMENT EXTRACTION | 16 |
| Data Dictionary: | 16 |
| Feature Engineering: | 16 |
| Model Selection | 17 |
| Hyper Parameters | 18 |
| maxDepth: | 18 |
| maxBins: | 18 |
| Accuracy: | 18 |
| Appendix: Screenshot of the Code | 19 |
| Q1: PREDICT DIABETES | 19 |
| Q2: PREDICT CANCER | 21 |
| KAGGLE COMPETITION -1 QUORA QUESTION PAIRS | 28 |
| KAGGLE COMPETITION - 2: TWEET SENTIMENT EXTRACTION | 29 |

Machine Learning with Apache Spark:

Machine learning is making the computer learn by studying data and statistics. Machine Learning helps solve problems by analyzing the data and learning to predict the outcome. Instead of addressing the challenges posed by dispersed data, we can concentrate on the data issues and models thanks to Apache Spark machine learning such as infrastructure, configurations, and so on.

1: PREDICT DIABETES

Predict whether a person is at risk of becoming diabetic based on the individual's data. Before analyzing it is very important to understand what we are going to do.

Type 2 Diabetes is caused when there is not enough insulin generated by the pancreas.

Data Dictionary:

| Field Name | Definition | Туре |
|--------------------------|---|---------|
| Pregnancies | Number of times pregnant | Feature |
| Glucose | Plasma glucose concentration, 2 hours in an oral glucose tolerance test | Feature |
| BloodPressure | Diastolic blood pressure (mm Hg) | Feature |
| SkinThickness | Triceps skin fold thickness (mm) | Feature |
| Insulin | 2-Hour serum insulin (mu U/ml) | Feature |
| ВМІ | Body mass index (weight in kg/(height in m)^2) | Feature |
| DiabetesPedigreeFunction | Diabetes pedigree function, a score based on a genetic factor of a person (diabetes has a close relation to family history) | Feature |
| Age | Age (years) | Feature |
| Outcome | Outcome whether the patient has diabetes or not | Label |

Selecting Features:

The major features in our dataset which cause diabetes are:

```
Correlation between Pregnancies and Outcome: 0.22189815303398636
Correlation between Insulin and Outcome: 0.13054795488404794
Correlation between Glucose and Outcome: 0.4665813983068737
Correlation between BloodPressure and Outcome: 0.06506835955033274
Correlation between SkinThickness and Outcome: 0.07475223191831945
Correlation between BMI and Outcome: 0.2926946626444454
Correlation between DiabetesPedigreeFunction and Outcome: 0.17384406565296
Correlation between Age and Outcome: 0.23835598302719757
```

Since there's a weak correlation between SkinThickness, BloodPressure, and Outcome, these features can be neglected. Even though there is a weak correlation between insulin and outcome, we have to consider it as it is one of the major factors of type 2 diabetes.

Model Selection:

I have chosen Random Forest. Random Forest algorithm avoids overfitting by using multiple trees.

Hyper Parameters

Parameters of an algorithm that may be tweaked to improve performance. The optimum hyperparameters are typically hard to predict in advance, and tweaking a model is when machine learning transitions from science to trial-and-error engineering. Because hyperparameter tuning is based on experimental findings rather than theory, the easiest way to identify the ideal settings is to test many different combinations and evaluate the performance of each model.

maxDepth:

Max depth can be defines as the longest path between the root node and the leaf node or max number of levels in each decision tree.

maxBins:

Max bins is a decision tree parameter. Number of bins set while continuous features are made separate or discreet. The default value of maxBin in spark is 32. and the value needs to be set >=2

```
Vector(m, a, x, D, e, p, t, h,
                                , 2 maxBin: 32)
accuracy of the model is :75.36231884057972
Vector(m, a, x, D, e, p, t, h, , 2 maxBin: 33)
accuracy of the model is :77.53623188405797
Vector(m, a, x, D, e, p, t, h, , 2 maxBin: 34)
accuracy of the model is :74.27536231884058
Vector(m, a, x, D, e, p, t, h, , 2 maxBin: 35)
accuracy of the model is :74.63768115942028
Vector(m, a, x, D, e, p, t, h, , 3 maxBin: 32)
accuracy of the model is :77.17391304347827
Vector(m, a, x, D, e, p, t, h, , 3 maxBin: 33)
accuracy of the model is :76.08695652173914
Vector(m, a, x, D, e, p, t, h, , 3 maxBin: 34)
accuracy of the model is :76.81159420289855
Vector(m, a, x, D, e, p, t, h, , 3 maxBin: 35)
accuracy of the model is :76.44927536231883
Vector(m, a, x, D, e, p, t, h, , 4 maxBin: 32)
accuracy of the model is :75.72463768115942
Vector(m, a, x, D, e, p, t, h, , 4 maxBin: 33)
accuracy of the model is :76.08695652173914
Vector(m, a, x, D, e, p, t, h, , 4 maxBin: 34)
accuracy of the model is :75.0
Vector(m, a, x, D, e, p, t, h, , 4 maxBin: 35)
accuracy of the model is :77.17391304347827
Vector(m, a, x, D, e, p, t, h, , 5 maxBin: 32)
accuracy of the model is :75.0
Vector(m, a, x, D, e, p, t, h, , 5 maxBin: 33)
accuracy of the model is :75.36231884057972
Vector(m, a, x, D, e, p, t, h, , 5 maxBin: 34)
accuracy of the model is :75.36231884057972
Vector(m, a, x, D, e, p, t, h, , 5 maxBin: 35)
accuracy of the model is :75.36231884057972
Vector(m, a, x, D, e, p, t, h, , 6 maxBin: 32)
accuracy of the model is :76.44927536231883
Vector(m, a, x, D, e, p, t, h, , 6 maxBin: 33)
accuracy of the model is :74.63768115942028
Vector(m, a, x, D, e, p, t, h, , 6 maxBin: 34)
accuracy of the model is :73.55072463768117
Vector(m, a, x, D, e, p, t, h, , 6 maxBin: 35)
accuracy of the model is :74.63768115942028
Vector(m, a, x, D, e, p, t, h, , 7 maxBin: 32)
accuracy of the model is :73.18840579710145
Vector(m, a, x, D, e, p, t, h, , 7 maxBin: 33)
accuracy of the model is :75.0
Vector(m, a, x, D, e, p, t, h, , 7 maxBin: 34)
accuracy of the model is :73.18840579710145
Vector(m, a, x, D, e, p, t, h, , 7 maxBin: 35)
accuracy of the model is :73.55072463768117
```

As we can see in the results above which parameters are yielding better results. When the maxdepth is 3,4 and 6. And when the maxBin is 32,33,34 & 35.

Accuracy:

On top of this, I wanted to check the accuracies of each submodel in the cross-validator. So, I had passed sedcollectSubModels to true, and ran it through a nested for loop to check the accuracies.

```
for (models <- subModels nikhilesh)
       for (s <- models) {val each model = s.transform(testData)
       val accuracy_nikhilesh_subModel = evaluator_nikhilesh.evaluate(each_model)
println("accuracy of the model is :"+(accuracy_nikhilesh_subModel*100))}
accuracy of the model is :72.46376811594203
accuracy of the model is :72.46376811594203
accuracy of the model is :71.73913043478261
accuracy of the model is :75.0
accuracy of the model is :72.82608695652173
accuracy of the model is :75.36231884057972
accuracy of the model is :74.63768115942028
accuracy of the model is :75.36231884057972
accuracy of the model is :76.81159420289855
accuracy of the model is :76.81159420289855
accuracy of the model is :73.91304347826086
accuracy of the model is :76.44927536231883
accuracy of the model is :73.91304347826086
accuracy of the model is :76.44927536231883
accuracy of the model is :74.27536231884058
accuracy of the model is :75.36231884057972
accuracy of the model is :75.36231884057972
accuracy of the model is :76.08695652173914
accuracy of the model is :76.81159420289855
accuracy of the model is :77.17391304347827
accuracy of the model is :76.08695652173914
accuracy of the model is :75.0
accuracy of the model is :77.89855072463769
accuracy of the model is :75.36231884057972
```

The cross validator choses the model which yields the best accuracy or result, and in this case the highest accuracy is chosen as the best model.

The highest accuracy in which this model can predict whether the person has diabetes or not is: **77.174%**

2: PREDICT CANCER

Cancer is a condition in which some cells in the body develop uncontrollably and spread to other parts of the body. In the millions of cells that make up the human body, cancer can develop practically anywhere. Human cells often divide to create new cells as the body requires them. New cells replace old ones when they die as a result of ageing or damage. Malignant (cancerous) or benign (non-cancerous) tumours can both occur in the body. It is typical for non-cancerous tumours to grow slowly and not spread. Cancerous tumours have the ability to spread throughout the body, grow quickly, infiltrate neighbouring normal tissues, and do great damage. Here we are predicting whether a person's tumour is cancerous in order to decide whether surgery is necessary or not

Data Dictionary:

| Field Name | Value Range | Туре |
|-----------------------------|-------------------------------|---------|
| ID | Sample Code Number | ID |
| Clump Thickness | 1-10 | Feature |
| Uniformity of Cell Size | 1-10 | Feature |
| Uniformity of Cell Shape | 1-10 | Feature |
| Single Epithelial Cell Size | 1-10 | Feature |
| Bare Nuclei | 1-10 | Feature |
| Bland Chromatin | 1-10 | Feature |
| Normal Nucleol | 1-10 | Feature |
| Mitoses | 1-10 | Feature |
| Class | 2 for benign, 4 for malignant | Label |

Selecting features:

Since id is a just a sample code number, we can ignore it and consider all other features. All the other features have no null values. Are within the range fo 1-10 and the class variable which is our predictor variable has binomial values either 2 or 4.

```
// Entering paste mode (ctrl-D to finish)
println("Correlation between UofCSize and Class: "+df_cacner_nikhilesh.stat.corr("UofCSize", "Class", "pear
+"\nCorrelation between UofCShape and Class: "+df cacner nikhilesh.stat.corr("UofCShape", "Class", "pearso
n").toString+
 "\nCorrelation between Marginal Adhesion and Class: "+df_cacner_nikhilesh.stat.corr("Marginal Adhesion", "
Class", "pearson").toString+
 "\nCorrelation between SECSize and Class: "+df_cacner_nikhilesh.stat.corr("SECSize", "Class", "pearson").t
 "\nCorrelation between Bare Nuclei and Class: "+df_cacner_nikhilesh.stat.corr("Bare Nuclei", "Class", "pea
rson").toString+
 "\nCorrelation between Mitoses and Class: "+df_cacner nikhilesh.stat.corr("Mitoses", "Class", "pearson").t
oString+
 "\nCorrelation between Normal Nucleoli and Class: "+df_cacner_nikhilesh.stat.corr("Normal Nucleoli", "Clas
s", "pearson").toString)
// Exiting paste mode, now interpreting.
Correlation between UofCSize and Class: 0.8208014428258776
Correlation between UofCShape and Class: 0.821890947688868
Correlation between Marginal Adhesion and Class: 0.7062941354660839
Correlation between SECSize and Class: 0.6909581590873204
Correlation between Bare Nuclei and Class: 0.8226958729964622
Correlation between Mitoses and Class: 0.4234479212952124
Correlation between Normal Nucleoli and Class: 0.7186771878756356
       П
```

None of the correlations is weak. Hence I have considered all the features except id

Model Selection

I have chosen Logistic Regression. Logistic Regression is usually used in binary classification problems. It is a supervised machine-learning approach that can be used to model the likelihood of a given class or occurrence. According to Kaggle, Logistic Regression is the most used classification algorithm. I have tried Naive-Bayes and Random Forest but Logistic Regression gave the best result.

Hyper Parameters

Parameters of an algorithm that may be tweaked to improve performance. The optimum hyperparameters are typically hard to predict in advance, and tweaking a model is when machine learning transitions from science to trial-and-error engineering. Because hyperparameter tuning is based on experimental findings rather than theory, the easiest way to identify the ideal settings is to test many different combinations and evaluate the performance of each model.

Tolerance:

The dual gap is tested for optimality and the optimization code continues until it is smaller than tolerance only if the updates are smaller than tolerance.

MaxIter:

The most iterations necessary for the solvers to coincide.

RegParam:

RegParam represents the compromise between the two objectives of minimising training error and minimising the model to prevent overfitting.

```
Vector(R, e, g, P, a, r, a, m, :, , 0.1 MaxIter:
                                                 10)
best accuracy of the model is :99.42942942943
Vector(R, e, g, P, a, r, a, m, :, , 0.1 MaxIter:
                                                 20)
best accuracy of the model is :99.96996996998
Vector(R, e, g, P, a, r, a, m, :, , 0.1 MaxIter:
                                                 30)
best accuracy of the model is :99.97997997997
Vector(R, e, g, P, a, r, a, m, :, , 0.1 MaxIter:
                                                 40)
best accuracy of the model is :99.97997997997
Vector(R, e, g, P, a, r, a, m, :, , 0.2 MaxIter:
                                                 10)
best accuracy of the model is :99.64964964964965
Vector(R, e, g, P, a, r, a, m, :, , 0.2 MaxIter:
                                                 20)
best accuracy of the model is :99.91991991992
Vector(R, e, g, P, a, r, a, m, :, , 0.2 MaxIter:
                                                 30)
best accuracy of the model is :99.94994994994995
Vector(R, e, g, P, a, r, a, m, :, , 0.2 MaxIter:
                                                 40)
best accuracy of the model is :99.94994994994995
Vector(R, e, g, P, a, r, a, m, :, , 0.3 MaxIter:
                                                 10)
best accuracy of the model is :99.70970970970971
Vector(R, e, g, P, a, r, a, m, :, , 0.3 MaxIter:
best accuracy of the model is :99.92992992992994
Vector(R, e, g, P, a, r, a, m, :, , 0.3 MaxIter:
best accuracy of the model is :99.92992992992994
Vector(R, e, g, P, a, r, a, m, :, , 0.3 MaxIter:
best accuracy of the model is :99.92992992992994
```

As we can see with all the combinations of RegParam and MaxIter, the model was still overfitting in some way or the other contributing to more than 99% accuracy when we were using BinaryClassificationEvaluator. To avoid overfitting, I changed the evaluator function to MultiClassificationEvaluator and the metric as "accuracy". This resulted in avoiding overfitting.

After multiple trial and error methods, in order to avoid overfitting and minimizing the loss. I have taken RegParam as 0.1 and MaxIter as 10.

Accuracy:

BinaryClassficationEvaluator was overfitting the model. So I have chosen MultiClassEvaluator as my evaluator function. The cross validator chooses the model which yields the best accuracy or result, and in this case, the highest accuracy is chosen as the best model.

```
// Exting paste mode (ctrl-D to finish)

val paramGrid_nikhilesh = new ParamGridBuilder().addGrid(logistic_nikhilesh.regParam, Array(0.1))
.addGrid(logistic_nikhilesh.elasticMctParam, ElasticMctParam1)
.addGrid(logistic_nikhilesh.btol, tolerance)
.addGrid(logistic_nikhilesh.makter, Array(10)).build()

val cross_validator_nikhilesh = new CrossValidator()
.setEstimatorParamMaps(paramGrid_nikhilesh)
.setEstimatorParamMaps(paramGrid_nikhilesh)
.setEvaluator(cvaluator_nikhilesh)
.setEvaluator(cvaluator_nik
```

The highest accuracy with which this model can predict whether the person has diabetes or not is: **94.258**%

3: KAGGLE COMPETITION -1 QUORA QUESTION PAIRS

Link: https://www.kaggle.com/c/guora-guestion-pairs

"Over 100 million people visit Quora every month, so it's no surprise that many people ask similarly worded questions. Multiple questions with the same intent can cause seekers to spend more time finding the best answer to their question and make writers feel they need to answer multiple versions of the same question. Quora values canonical questions because they provide a better experience to active seekers and writers, and offer more value to both of these groups in the long term."

The goal of this competition is to predict which of the provided pairs of questions contain two questions with the same meaning.

Data Dictionary:

| Field Name | Deinition | Туре |
|--------------|--|---------|
| ID | The id of a training set question pair | ID |
| qid1 | Unique id of question1 | ID |
| qid2 | Unique id of question2 | ID |
| question1 | First Question | Feature |
| question2 | Second Question | Feature |
| is_duplicate | Whether the two questions are duplicate or not | Label |

Feature Engineering:

- Since id,qid1 and qid2 are unique identification numbers, these doesnt hold any value over the outcome.
- In the question1 and question2, there were lot of that there was a lot of special characters like double quotes, so in order to remove them from the features, I have used the regular expression "[^a-zA-Z0-9]/g" to remove the special characters.
- The "check" function in the code is used to remove the accent from the words. To do so I have used **stringAccents** in the **StringUtils** class.
- The first step of the NLP process is gathering the data and breaking it into understandable parts. Tokenizer splits paragraphs and sentences into smaller units that can be more easily assigned meaning.
- The StopWordsRemover function removes all stop words from input sequences of strings (such as the result of a tokenizer). The stopWords option specifies the list of

stopwords. Calling the StopWordsRemover will provide access to some languages' default stop words. The possible possibilities for DefaultStopWords(language) are "danish," "dutch," "english," "finnish," "french," "german," "hungarian," "italian," "norwegian," "portuguese," "russian," "spanish," "swedish," and "turkish." a case for boolean parameters If the matches should be case sensitive, sensitive lets you know that.

• Further, had to check how many of the words are the same in both the sentences. For this, the function "getRatio_bw_question" is coded.

Model Selection

To predict whether the 2 questions convey the same meaning, I have chosen Random Forest.

Hyper Parameters

Parameters of an algorithm that may be tweaked to improve performance. The optimum hyperparameters are typically hard to predict in advance, and tweaking a model is when machine learning transitions from science to trial-and-error engineering. Because hyperparameter tuning is based on experimental findings rather than theory, the easiest way to identify the ideal settings is to test many different combinations and evaluate the performance of each model.

maxDepth:

Max depth can be defined as the longest path between the root node and the leaf node or max number of levels in each decision tree.

maxBins:

Max bins is a decision tree parameter. A number of bins are set while continuous features are made separate or discreet. The default value of maxBin in spark is 32. and the value needs to be set >=2

After multiple trials and error methods, in order to avoid overfitting and minimize the loss. I have taken maxdepth is 3,4 and 6. And when the maxBin is 32,33,34 & 35.

Accuracy:

BinaryClassficationEvaluator was overfitting the model. So I have chosen MultiClassEvaluator as my evaluator function. The cross-validator chooses the model which yields the best accuracy or result, and in this case, the highest accuracy is chosen as the best model.

The highest accuracy with which this model can predict whether the pair of questions are duplicates or not is: **72.76%**

4: KAGGLE COMPETITION - 2: TWEET SENTIMENT EXTRACTION

Source: https://www.kaggle.com/competitions/tweet-sentiment-extraction

"The competition is about predicting the word or phrase from the tweet that exemplifies the provided sentiment into 3 categories Positive, Negative & Neutral."

The goal of this competition is to predict the sentiment of the tweet into any of the three categories Positive, Neutral, or Negative.

Data Dictionary:

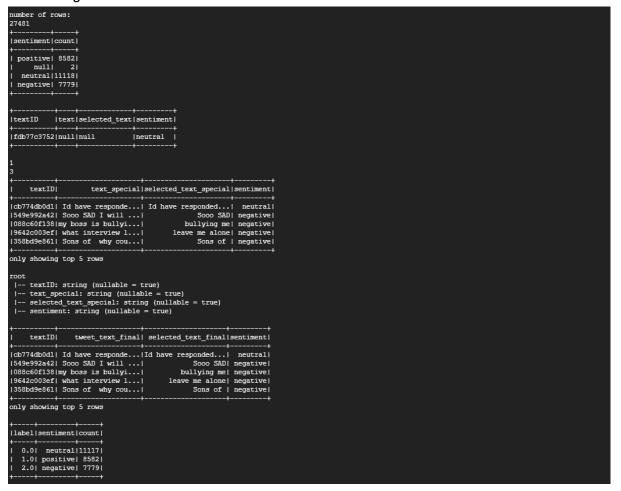
| Field Name | Definition | Туре |
|---------------|---|---------|
| textID | Unique ID for each piece of text. | ID |
| text | The text of the tweet | Feature |
| selected_text | Text that supports the tweet's sentiment | Feature |
| sentiment | Whether the two questions are duplicates or not | Label |

Feature Engineering:

- Since textID is a unique identification number, these doesn't hold any value over the outcome.
- To remove special characters like Asterix or colon. I have used the regular expression "[^a-zA-Z0-9]/g" in the function "check_special" to remove the special characters.
- The "check" function in the code is used to remove the accent from the words. To do so I have used stringAccents in the StringUtils class.
- The first step of the NLP process is gathering the data and breaking it into
 understandable parts. RegexTokenizer splits paragraphs and sentences into smaller
 units that can be more easily assigned meaning. RegexTokenizer also considers only
 the alphabets by using the method .setPattern in which the regular expression needs
 to be passed. The regular expression passed here is "[a-zA-Z']+".
- The StopWordsRemover function removes all stop words from input sequences of strings (such as the result of a tokenizer). The stopWords option specifies the list of stopwords. Calling the StopWordsRemover will provide access to some languages' default stop words. The possible possibilities for DefaultStopWords(language) are "danish," "dutch," "english," "finnish," "french," "german," "hungarian," "italian," "norwegian," "portuguese," "russian," "spanish," "swedish," and "turkish." a case for

- boolean parameters If the matches should be case sensitive, sensitive lets you know that.
- N-grams are continuous word, symbol, or token sequences in a document. They can
 be described technically as the adjacent groups of items in a document. When
 working with text data in NLP, they come into play. Here the N in NGram is 2 to split
 into 2 or to make them as bigrams.
- Count Vectoriser converts a collection of strings to a matrix of token counts. Count Vectoriser has been used on top of N-Gram
- IDF is a method for determining the meaning of word-based sentences and eliminates the shortcomings of the Bag of Words technique, which is effective for classifying texts or assisting a machine to interpret words into numbers. Input for the IDF is the outcome of Count Vectoriser.
- String indexer converts the given set of strings to index. Like in this case there are 3
 different labels Positive, Negative and Neutral. So, the string indexer has converted
 Neutral: 0

Positive: 1 Negative: 2



Model Selection

Since LinearSVC offers only binary classification, I had to consider other classification techniques. After testing Logistic Regression, Decision Trees, and Random Forest. I have used Random Forest because it was giving the best result out of all of them.

Hyper Parameters

Parameters of an algorithm that may be tweaked to improve performance. The optimum hyperparameters are typically hard to predict in advance, and tweaking a model is when machine learning transitions from science to trial-and-error engineering. Because hyperparameter tuning is based on experimental findings rather than theory, the easiest way to identify the ideal settings is to test many different combinations and evaluate the performance of each model.

maxDepth:

Max depth can be defined as the longest path between the root node and the leaf node or max number of levels in each decision tree.

maxBins:

Max bins are a decision tree parameter. A number of bins are set while continuous features are made separate or discreet. The default value of maxBin in spark is 32. and the value needs to be set >=2

After multiple trial and error methods, in order to avoid overfitting and minimize the loss. I have taken maxdepth is 3,4 and 6. And when the maxBin is 32,33,34 & 35.

Accuracy:

So I have chosen MultiClassEvaluator as my evaluator function. The cross-validator chooses the model which yields the best accuracy or result, and in this case, the highest accuracy is chosen as the best model.

```
scale> val predictions_nikhilesh = model_nikhilesh.transform(testData)
predictions_nikhilesh: org.apache.spark.sql.DataFrame = [textID: string, tweet_text_final: string ... 12 more fields]
scale> val accuracy_nikhilesh = evaluator_nikhilesh.evaluate(predictions_nikhilesh)
22/11/23 07:13:50 WARN org.apache.spark.scheduler.DAGScheduler: Broadcasting large task binary with size 1791.0 KiB
accuracy_nikhilesh: Double = 0.7203810369307343
scale>
```

The highest accuracy with which this model can predict the sentiment of the tweet is: **72.039**%

References:

- https://blog.exploratory.io/exploratory-weekly-update-12-3-d4b1d0f620b9
- https://www.kaggle.com/c/quora-question-pairs
- https://www.kaggle.com/competitions/tweet-sentiment-extraction
- https://stackoverflow.com/questions/38874546/spark-crossvalidatormodel-access-oth-er-models-than-the-bestmodel/38874828#38874828
- https://stats.stackexchange.com/questions/530745/what-happens-if-a-random-forest-max-bins-is-set-much-higher-than-the-number-of-c#:~:text=It%20is%20a%20decision%20tree.also%20increases%20computation%20and%20communication
- https://spark.apache.org/docs/2.2.1/api/java/org/apache/spark/ml/classification/Rand-omForestClassifier.html
- https://spark.apache.org/docs/3.1.1/api/scala/org/apache/spark/ml/classification/RandomForestClassifier.html
- https://spark.apache.org/docs/latest/mllib-decision-tree.html
- https://spark.apache.org/docs/1.6.2/api/java/org/apache/spark/ml/feature/RegexTokenizer.html
- https://spark.apache.org/docs/3.0.2/api/scala/org/apache/spark/ml/tuning/CrossValid-ator.html
- https://spark.apache.org/docs/2.2.0/api/java/org/apache/spark/ml/evaluation/BinaryClassificationEvaluator.html

Appendix: Screenshot of the Code

Q1: PREDICT DIABETES

```
ecala> val data nikhilesh = spark.read.format("csv").option("inferSchema", "true").option("header", "true").load("hdfs://10.128.0.19/BigData/diabetes.csv")
iata_nikhilesh: org.apache.spark.sql.DataFrame = [Pregnancies: int, Glucose: int ... 9 more fields]
                                                        data_nikhilesh.na.drop().show(false)
        ssla> var df_nikhilesh = data_nikhilesh.drop("_e9","_c10")
df_nikhilesh: org.apache.spark.sql.DataFrame = [Pregnancies: int, Glucose: int ... 7 more fields]
        | Pregnancies|Glucose|BloodPressure|SkinThickness|Insulin|BMI|DiabetesPedigreeFunction|Age|Outcome
                                                                                                                                                                                                                                                        | 129
| 123
| 10
| 10
| 10
| 10
| 135
| 119
| 10
| 126
| 111
| 10
| 131
| 133
| 10
| 125
| 10
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   | 131 | 10
| 121 | 10
| 129 | 10
| 130 | 10
| 139 | 10
| 157 | 10
| 157 | 10
| 150 | 10
| 152 | 10
| 160 | 10
| 122 | 10
| 128 | 10
| 145 | 10
| 133 | 10
| 133 | 10
| 135 | 10
| 136 | 10
| 137 | 10
|4
|10
|1
|3
|8
|1
|13
|5
|5
|3
|6
|10
|4
|11
                                                                                                                           on between Pregnancies and Outcome: "4df_milhilesh.stat.corr("Fregnancies", "Outcome"), foottripp" information between Insulin and Outcome: "4df_milhilesh.stat.corr("Insulin", "Outcome"), foottripp" information between Bloodfressure and Outcome
                                    val Arroy(crainingData, testData) = df_mikhilesh.randomSplit(Array(0.65, 5.35), 512)
spData; org.spacko.spark.spl. Dataset[org.spacko.spark.spl.Nov] - [Prespanacies: int, Glucose: int ... 7 more fields)
sr org.spacko.spark.spl. Split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(split(
```

```
Vector(m, a, x, D, e, p, t, h,
                                , 2 maxBin: 32)
accuracy of the model is :75.36231884057972
Vector(m, a, x, D, e, p, t, h, , 2 maxBin: 33)
accuracy of the model is :77.53623188405797
Vector(m, a, x, D, e, p, t, h, , 2 maxBin: 34)
accuracy of the model is :74.27536231884058
Vector(m, a, x, D, e, p, t, h, , 2 maxBin: 35)
accuracy of the model is :74.63768115942028
Vector(m, a, x, D, e, p, t, h, , 3 maxBin: 32)
accuracy of the model is :77.17391304347827
Vector(m, a, x, D, e, p, t, h, , 3 maxBin: 33)
accuracy of the model is :76.08695652173914
Vector(m, a, x, D, e, p, t, h, , 3 maxBin: 34)
accuracy of the model is :76.81159420289855
Vector(m, a, x, D, e, p, t, h, , 3 maxBin: 35)
accuracy of the model is :76.44927536231883
Vector(m, a, x, D, e, p, t, h, , 4 maxBin: 32)
accuracy of the model is :75.72463768115942
Vector(m, a, x, D, e, p, t, h, , 4 maxBin: 33)
accuracy of the model is :76.08695652173914
Vector(m, a, x, D, e, p, t, h, , 4 maxBin: 34)
accuracy of the model is :75.0
Vector(m, a, x, D, e, p, t, h, , 4 maxBin: 35)
accuracy of the model is :77.17391304347827
Vector(m, a, x, D, e, p, t, h, , 5 maxBin: 32)
accuracy of the model is :75.0
Vector(m, a, x, D, e, p, t, h, , 5 maxBin: 33)
accuracy of the model is :75.36231884057972
Vector(m, a, x, D, e, p, t, h, , 5 maxBin: 34)
accuracy of the model is :75.36231884057972
Vector(m, a, x, D, e, p, t, h, , 5 maxBin: 35)
accuracy of the model is :75.36231884057972
Vector(m, a, x, D, e, p, t, h, , 6 maxBin: 32)
accuracy of the model is :76.44927536231883
Vector(m, a, x, D, e, p, t, h, , 6 maxBin: 33)
accuracy of the model is :74.63768115942028
Vector(m, a, x, D, e, p, t, h, , 6 maxBin: 34)
accuracy of the model is :73.55072463768117
Vector(m, a, x, D, e, p, t, h, , 6 maxBin: 35)
accuracy of the model is :74.63768115942028
Vector(m, a, x, D, e, p, t, h, , 7 maxBin: 32)
accuracy of the model is :73.18840579710145
Vector(m, a, x, D, e, p, t, h, , 7 maxBin: 33)
accuracy of the model is :75.0
Vector(m, a, x, D, e, p, t, h, , 7 maxBin: 34)
accuracy of the model is :73.18840579710145
Vector(m, a, x, D, e, p, t, h, , 7 maxBin: 35)
accuracy of the model is :73.55072463768117
```

```
scala> val cvModel_nikhilesh = cross_validator_nikhilesh.fit(trainingData)
cvModel_nikhilesh: org.apache.spark.ml.tuning.CrossValidatorModel = CrossValidatorModel: uid=cv_ea1356f3d6a9, bestModel=pipeline_d0a
7ab9cf0b4, numFolds=3

scala> val predictions_nikhilesh = cvModel_nikhilesh.transform(testData)
predictions_nikhilesh: org.apache.spark.sql.DataFrame = [Pregnancies: int, Glucose: double ... 9 more fields]

scala> val accuracy_nikhilesh = evaluator_nikhilesh.evaluate(predictions_nikhilesh)
accuracy_nikhilesh: Double = 0.7717391304347826

scala> println("accuracy of the model is :"+(accuracy_nikhilesh*100))
accuracy of the model is :77.17391304347827

scala>
```

Q2: PREDICT CANCER

П

```
val data_nikhilesh = spark.read.format("csv").option("inferSchema", "true").option("header", "true").load("hdfs://l0.128.0.19/BigData/cancer.csv")
 // Exiting paste mode, now interpreting.
import org.apache.spark.sql.functions_
import org.apache.spark.ml.feature.VectorAssembler
import org.apache.spark.ml.feature.VectorAssembler
import org.apache.spark.ml.repression.LinearRegression
import org.apache.spark.ml.repression.LinearRegression
import org.apache.spark.ml.repression.DinearRegression
import org.apache.spark.ml.repression.PegressionEvaluator
import org.apache.spark.ml.repression.PegressionEvaluator
import org.apache.spark.ml.param.ParamMap
import org.apache.spark.ml.linalg.(Matrix, Vectors)
import org.apache.spark.ml.stat.Correlation
import org.apache.spark.ml.stat.Correlation
import org.apache.spark.sql.functions.
import org.apache.spark.sql.functions.
import org.apache.spark.sql.functions.
import org.apache.spark.sql.functions.
import org.apache.spark.ml.feature.(VectorAssembler, StringIndexer)
import org.apache.spark.ml.feature.(VectorAssembler, StringIndexer)
import org.apache.spark.ml.feature.(VectorAssembler, StringIndexer)
import org.apache.spark.ml.feature.(RandomForestClassi...
 osla> var df_cacner_nikhilesh = data_nikhilesh
df_cacner_nikhilesh: org.apache.spark.sql.DataFrame = [id: int, Clump Thickness: int ... 9 more fields]
       df_cacner_nikhilesh.printSchema
  oot
|-- id: integer (mullable = true)
|-- clump Thickness: integer (mullable = true)
|-- UofcSize: integer (mullable = true)
|-- UofcShape: integer (mullable = true)
|-- Warginal Adhesion: integer (mullable = true)
|-- Marginal Adhesion: integer (mullable = true)
|-- SECSize: integer (mullable = true)
|-- Bland Chomasin: integer (mullable = true)
|-- Normal Nucleoli: integer (mullable = true)
|-- Mitoges: integer (mullable = true)
|-- Class: integer (mullable = true)
|-- Class: integer (mullable = true)
        П
             :paste
// Entering paste mode (ctrl-D to finish)
println("Correlation between UofCSize and Class: "+df_cacner_nikhilesh.stat.corr("UofCSize", "Class", "pear
 son").toString
  +"\nCorrelation between UofCShape and Class: "+df_cacner_nikhilesh.stat.corr("UofCShape", "Class", "pearso
n").toString+
"\nCorrelation between Marginal Adhesion and Class: "+df_cacner_nikhilesh.stat.corr("Marginal Adhesion", "
Class", "pearson").toString+
   "\nCorrelation between SECSize and Class: "+df_cacner_nikhilesh.stat.corr("SECSize", "Class", "pearson").t
oString+
  "\nCorrelation between Bare Nuclei and Class: "+df cacner nikhilesh.stat.corr("Bare Nuclei", "Class", "pea
  "\nCorrelation between Mitoses and Class: "+df_cacner_nikhilesh.stat.corr("Mitoses", "Class", "pearson").t
 oString+
   "\nCorrelation between Normal Nucleoli and Class: "+df_cacner_nikhilesh.stat.corr("Normal Nucleoli", "Clas
 s", "pearson").toString)
 // Exiting paste mode, now interpreting.
Correlation between UofCSize and Class: 0.8208014428258776
Correlation between UofCShape and Class: 0.821890947688868
Correlation between Marginal Adhesion and Class: 0.7062941354660839
Correlation between SECSize and Class: 0.6909581590873204
Correlation between Bare Nuclei and Class: 0.8226958729964622
Correlation between Mitoses and Class: 0.4234479212952124
Correlation between Normal Nucleoli and Class: 0.7186771878756356
```

```
5| 5| 34|
10| 1| 30|
1| 43| 1|
6| 1| 39|
9| 0| 2|
2|355| 21|
7| 2| 9|
3| 28| 43|
8| 2| 19|
4| 7| 41|
    ala> df_cacner_nikhilesh.stat.crosstab("Bare Nuclei","Class").show()
5 | 10 | 20 |

5 | 10 | 20 |

10 | 3 | 129 |

1 | 387 | 15 |

6 | 0 | 4 |

9 | 0 | 9 |

2 | 21 | 9 |

7 | 1 | 7 |

3 | 14 | 14 |

8 | 2 | 19 |

4 | 6 | 13 |
 |Bland Chromatin_Class|
                                    df_cacner_nikhilesh.stat.crosstab("Bland Chromatin", "Class").show()
|Bland Chromatin_Class| 2| 4|
                                    5| 4| 30|

10| 0| 20|

1|148| 2|

6| 1| 8|

9| 0| 11|

2|153| 7|

7| 6| 65|

3|125| 36|

8| 0| 28|

4| 7| 32|
scala> df_cacner_nikhilesh.stat.crosstab("Mitoses","Class").show()
+------
|Mitoses_Class| 2| 4|
                     5| 1| 5|

10| 0| 14|

1|431|132|

6| 0| 3|

2| 8| 27|

7| 1| 8|

3| 2| 31|

8| 1| 7|

4| 0| 12|
    ala> df_cacner_nikhilesh.stat.crosstab("Normal Nucleoli","Class").show()
|Normal Nucleoli_Class| 2| 4|
                                    5| 2|

5| 2|

10| 0|

1|391|

6| 4|

9| 0|

2| 30|

7| 2|

3| 11|

8| 3|

4| 1|
                                                  17|
60|
41|
18|
15|
6|
14|
31|
20|
```

df cacner nikhilesh.stat.crosstab("SECSize", "Class").show()

|SECSize_Class| 2| 4|

```
| Color | Colo
```

```
|Normal Nucleoli|count|
    5| 0|
10| 0|
1|369|
6| 0|
9| 1|
2| 37|
7| 1|
3| 27|
8| 1|
4| 8|
                                                                                                                               30 |
67 |
4 |
25 |
8 |
18 |
25 |
27 |
30 |
                   val assembler_cancer_mikhilesh = new VectorAssembler().setInpurCols(Array("Clump Thickness", "GotCSlape", "Morginal Adhesion", "SSCSIze", "Bare Nuclei", "Bland Chromatin", "Normal Nucleoli", "Mitoses", "Class")).setOutputCol("features" ler_cancer_mikhilesh: org.apache.spark.ml.feature.VectorAssembler = VectorAssembler = Understanding Thickness = "Class").setOutputCol("features" ler_cancer_mikhilesh: org.apache.spark.ml.feature.VectorAssembler = VectorAssembler = Understanding Thickness = "Class").setOutputCol("features" ler_cancer_mikhilesh: org.apache.spark.ml.features VectorAssembler = VectorAssem
                   val indexer_nikhilesh = new StringIndexer().setInputCol("Class").setOutputCol("label")
r_nikhilesh: org.apache.spark.ml.feature.StringIndexer = strIdx_9a30ea86cc0b
                  val Array(training, test) = df_cacmer_nikhilesh.randomSplit(Array(0.7, 0.3))
mg; org.apsche.spark.sgl.Natasef[org.apsche.spark.sgl.Now] = [id: int, Clump Thickness: int ... 9 more fields]
org.apsche.spark.sgl.Dataset[org.apsche.spark.sgl.Now] = [id: int, Clump Thickness: int ... 9 more fields]
                    val logistic_nikhilesh = new LogisticRegression().setMaxIter(10).setLabelCol("label").setFeaturesCol("featic_nikhilesh: org.apache.spark.ml.classification.logisticRegression = logreg_9333a8d59d25
                    val evaluator_mikhilesh = new BinaryClassificationEvaluator[).setlabelCol["label").setEawFredictionCol["rawFrediction"].setHetricHass("aresUnderHOC")
for mikhilesh: org.apache.spark.sl.evaluation.BinaryClassificationEvaluator = BinaryClassificationEvaluator: uid=binEval_c5046ceT4543, metricHass=aresI
 val MaxIter_Array = Array(10,20,30,40)
val pipeline_nikhilesh = new Pipeline()
    .setStages(Array(assembler_cancer_nikhilesh,indexer_nikhilesh, logistic_nikhilesh))
for (RegParami <- RegParam Array)
    for(MaxIter1 <- MaxIter_Array)</pre>
  {
val paramGrid_nikhilesh = new ParamGridBuilder().addGrid(logistic_nikhilesh.regParam, Array(RegParam1))
.addGrid(logistic_nikhilesh.elasticNetParam, ElasticNetParam1)
.addGrid(logistic_nikhilesh.tol, tolerance)
.addGrid(logistic_nikhilesh.maxIter, Array(MaxIter1)).build()
   val cross_validator_nikhilesh = new CrossValidator()
.setEstimatorParamMaps(paramGrid_nikhilesh)
   .setEstimator(pipeline_nikhilesh)
.setEstimator(pipeline nikhilesh)
.setEvaluator(evaluator_nikhilesh)
.setNumFolds(5)
val cvModel_nikhilesh = cross_validator_nikhilesh.fit(training)
val predictions_nikhilesh = cvModel_nikhilesh.transform(test)
val accuracy_nikhilesh = evaluator_nikhilesh.evaluate(predictions_nikhilesh)
println("RegParam: ":+RegParaml.toString+" MaxIter: "+MaxIter1.toString)
println("best accuracy_of the model is :"+(accuracy_nikhilesh*100)))
   // Exiting paste mode, now interpreting.
```

```
Vector(R, e, g, P, a, r, a, m, :, , 0.1 MaxIter: 10)
best accuracy of the model is :99.42942942943
Vector(R, e, g, P, a, r, a, m, :, , 0.1 MaxIter: 20)
best accuracy of the model is :99.96996996996998
Vector(R, e, g, P, a, r, a, m, :, , 0.1 MaxIter: 30)
best accuracy of the model is :99.97997997997997
Vector(R, e, g, P, a, r, a, m, :, , 0.1 MaxIter: 40)
best accuracy of the model is :99.97997997997997
Vector(R, e, g, P, a, r, a, m, :, , 0.2 MaxIter: 10)
best accuracy of the model is :99.64964964965
Vector(R, e, g, P, a, r, a, m, :, , 0.2 MaxIter: 20)
best accuracy of the model is :99.91991991991992
Vector(R, e, g, P, a, r, a, m, :, , 0.2 MaxIter: 30)
best accuracy of the model is :99.94994994994995
Vector(R, e, g, P, a, r, a, m, :, , 0.2 \text{ MaxIter: } 40)
best accuracy of the model is :99.94994994995
Vector(R, e, g, P, a, r, a, m, :, , 0.3 MaxIter: 10)
best accuracy of the model is :99.70970970971
Vector(R, e, g, P, a, r, a, m, :, , 0.3 MaxIter: 20)
best accuracy of the model is :99.92992992994
Vector(R, e, g, P, a, r, a, m, :, , 0.3 MaxIter: 30)
best accuracy of the model is :99.92992992992994
Vector(R, e, g, P, a, r, a, m, :, , 0.3 MaxIter: 40)
best accuracy of the model is :99.92992992992994
```

```
res7: evaluator nikhilesh.type = MulticlassClassificationEvaluator: uid-mcEval_7b87a8a40f71, metricName=f1, metricLabel=0.0, beta=1.0, eps=1.0E-15
   ola> .setFredictionCol("prediction")
es8: res7.type = MulticlassClassificationEvaluator: uid=mcEval_7b87a8a40f71, metricName=f1, metricLabel=0.0, beta=1.0, eps=1.0E-15
  .setMetricName("accuracy")
res9: res8.type = MulticlassClassificationEvaluator: uid=mcEval_7b87a8a40f71, metricName=accuracy, metricLabel=0.0, beta=1.0, eps=1.0E-15
 scala> :paste
// Entering paste mode (ctrl-D to finish)
 for (RegParam1 <- RegParam_Array)
  for(MaxIter1 <- MaxIter_Array)</pre>
     {
    l paramGrid_nikhilesh = new ParamGridBuilder().addGrid(logistic_nikhilesh.regFaram, Array(RegFaram1))
    .addGrid(logistic_nikhilesh.elasticNetFaram, ElasticNetFaram1)
    .addGrid(logistic_nikhilesh.tol, tolerance)
    .addGrid(logistic_nikhilesh.maxHater, Array(MaxHer1)).build()
    l cross_validator_nikhilesh = new CrossValidator()
    retEstimator(pipeline_nikhilesh)
    .setEstimator(pipeline_nikhilesh)
    .setEstimator(pipeline_nikhilesh)
      .setEvaluator(evaluator nikhilesh)
 .setEvaluator(evaluator_mixhiresh)
.setNumFolds(5)
val cvModel_nikhilesh = cross_validator_nikhilesh.fit(training)
val cvModel_nikhilesh = cvModel_nikhilesh.transform(test)
val accuracy_nikhilesh = evaluator_nikhilesh.evaluate(predictions_nikhilesh)
println("RegFaram: ":+RegFaram1+" MaxIter: "+MaxIter1.toString)
println("best accuracy of the model is :"+(accuracy_nikhilesh*100))
 // Exiting paste mode, now interpreting.
Vector(R, e, g, P, a, r, a, m, :, , 0.1 MaxIter: 10)
best accuracy of the model is :94.25837320574163
Vector(R, e, g, P, a, r, a, m, :, , 0.1 MaxIter: 20)
best accuracy of the model is :98.56459330443541
Vector(R, e, g, P, a, r, a, m, :, , 0.1 MaxIter: 30)
best accuracy of the model is :99.04306220095694
 Weccouracy of the model is :99.04306220095694
Wector(R, e, g, P, a, r, a, m, :, , 0.1 MaxIter: 40)
best accuracy of the model is :99.04306220095694
Wector(R, e, g, P, a, r, a, m, :, , 0.2 MaxIter: 10)
best accuracy of the model is :98.08612440191388
   ector(R, e, g, P, a, r, a, m, :, , 0.2 MaxIter: 20)
est accuracy of the model is :98.56459330143541
 vector(R, e, g, P, a, r, a, m, :, , 0.2 Maxiter: 30)
best accuracy of the model is :98.56459330143541
Vector(R, e, g, P, a, r, a, m, :, , 0.2 Maxiter: 40)
best accuracy of the model is :98.56459330143541
  Pector(R, e, g, P, a, r, a, m, :, , 0.3 MaxIter: 10)
Dest accuracy of the model is :97.1291866028708
           pr(R, e, g, P, a, r, a, m, :, , 0.3 MaxIter: 20) accuracy of the model is :98.56459330143541
    ector(R, e, g, P, a, r, a, m, :, , 0.3 MaxIter
est accuracy of the model is :98.56459330143541
   ector(R, e, g, P, a, r, a, m, :, , 0.3 MaxIter: 40) est accuracy of the model is :98.56459330143541
```

```
### Parting parts and (ext-D to finish)

### Parting parts and (ext-D to finis
```

KAGGLE COMPETITION -1 QUORA QUESTION PAIRS

```
|question2|is_duplicate
                                                                                                                                                                                                                                                                                                                         |null
  number of rows after removing duplicates:
                 1|149132|
5| 6|
-36| 1|
0|253343|
        1|149132|
5| 6|
-36| 1|
0|253343|
     pora mikhilesh: org.apache.spark.sql.DataFrame = [id: string, qid1: string ... 4 mor
_mikhilesh: org.apache.spark.sql.DataFrame = [id: int, qid1: int ... 4 more fields]
_uikhilesh: org.apache.spark.sql.DataFrame = [id: int, qid1: int ... 4 more fields]
_heck: org.apache.spark.sql.copressions.UserDefinedPunction = SparkUserDefinedPunction
                                                                                                                                    n($Lambda$5588/496688719$62f4a010,StringType,List(Some(class[value[0]: string])),Some(class[value[0]: string]),None,true,true)
iPunction($Lambda$5589/13607366608fbef88d,StringType,List(Some(class[value[0]: string])),Some(class[value[0]: string]),None,true,tr
               training.groupBy($"label").count().show
               val cwModel_nikhilesh = cross_validator_nikhilesh.fit(training)
1 nikhilesh: org.apache.spark.ml.tuning.CrossValidatorModel = CrossValidatorModel: uid=cv_9842a7f748d0, bestModel=pipeline_e76d678a1983, numFolds=4
   scala> :paste
// Entering paste mode (ctrl-D to finish)
  val predictions_nikhilesh = cvModel_nikhilesh.transform(test)
val accuracy_nikhilesh = evaluator_nikhilesh.evaluate(predictions_nikhilesh)
println("accuracy of the model is :"+(accuracy_nikhilesh*100))
  // Exiting paste mode, now interpreting.
  accuracy of the model is :72.76868009688272
predictions nikhilesh: org.apache.spark.sql.DataFrame = [qid2: int, qid1: int ... 11 more fields]
accuracy_nikhilesh: Double = 0.7276868009688272
```

KAGGLE COMPETITION - 2: TWEET SENTIMENT EXTRACTION

```
scala> :paste
// Entering paste mode (ctrl-D to finish)
 val first_dataframe_twitter_nikhilesh = spark.read.format("csv").option("inferSc
println("number of rows: \n"*first_dataframe_twitter_nikhilesh.count().toString)
first_dataframe_twitter_nikhilesh.groupBy(col("sentiment")).count().show
    orintin(first_dataframe_twitter_nikhilesh.filter(first_dataframe_twitter_nikhilesh("cext").isWull || col("text")=="").count())
printin(first_dataframe_twitter_nikhilesh.filter(first_dataframe_twitter_nikhilesh("selected_text").isWull || col("selected_text")=="").count())
wil deke udf((colValue) String) => {
Stringfills_stripkcounts(colValue)
           check_special = udf((colValue: String) => {
clValue.replaceAll("[^\p{Alpha}\\p{Digit}] | +","");
                nikhienh.nbov(s)

inkhienh.printSchess
inkhienh.pri
 number of rows:
27481
  |sentiment|count|
 | positive| 8582|
| null| 2|
    null| 2|
neutral|11118|
negative| 7779|
|textID |text|selected_text|sentiment|
 |fdb77c3752|null|null |neutral
                                                                            text_special|selected_text_special|sentiment|
only showing top 5 rows
    |-- textID: string (nullable = true)
|-- text_special: string (nullable = true)
|-- selected_text_special: string (nullable = true)
|-- sentiment: string (nullable = true)
             textID| tweet_text_final| selected_text_final|sentiment|
only showing top 5 rows
|label|sentiment|count|
        0.0| neutral|11117|
1.0| positive| 8582|
2.0| negative| 7779|
```

```
val tokenizer_nikhilesh = new RegexTokenizer()
.setFattern("[a-zh-2']+")
.setGaps(false)
.setInputCol("tweet_text_final")
.setOutputCol("words")
val remover_nikhilesh = new StopHordsRemover()
.setInputCol("words")
.setOutputCol("filtered")
val ngram_nikhilesh = new NGram()
.setNiputCol("filtered")
   val ngram nikhilesh = new NGram()
.setN(2)
.setInputCol("filtered")
.setOutputCol("ngram-2")
val cv2_nikhilesh: CountVectorizer
.setInputCol("ngram-2")
.setOutputCol("ngram-2-features")
val cv2idf_nikhilesh = new IDF()
.setInputCol("ngram-2-features")
.setOutputCol("gram-2-features")
.setOutputCol("cv2-idf-features")
                                                                      ctorizer = new CountVectorizer()
   // Exiting paste mode, now interpreting.
   tokenizer_nikhilesh: org.apache.spark.ml.feature.RegexTokenizer = RegexTokenizer: uid=regexTok_cc785446e514, minTokenLength=1, gaps=false, pattern=[a-zh-Z*]+, toLowercase=true remover_nikhilesh: org.apache.spark.ml.feature.StopWordsRemover = StopWordsRemover: uid=stopWords_438627df0911, numStopWords=181, locale=en, caseSensitive=false gyram nikhilesh: org.apache.spark.ml.feature.NGram = NGram: uid=ngram_e73cdd0fda21, n=2 ovz_nikhilesh: org.apache.spark.ml.feature.CountVectorizer = cnuVec_bds=8b98fde cvzldf_nikhilesh: org.apache.spark.ml.feature.UDF = idf_07d48420bfd1
       ala>
scala> :paste
// Entering paste mode (ctrl-D to finish)
val randomForestModel_nikhilesh = new RandomForestClassifier()
.setFeaturesCol("cv2-idf-features")
.setLabelCol("label")
val pipeline_nikhilesh = new Pipeline()
.setStages(Array(cokenizer_nikhilesh, remover_nikhilesh, ngram_nikhilesh, cv2_nikhilesh, cv2idf_nikhilesh, randomForestModel_nikhilesh))
val paramGrid_nikhilesh = new ParamGridBullder()
.addGrid(randomForestModel_nikhilesh.maxBins, Array(3,4,5))
.addGrid(randomForestModel_nikhilesh.maxBins, Array(32,33,34,35))
.addGrid(randomForestModel_nikhilesh.maxBins, Array(32,33,34,35))
.addGrid(randomForestModel_nikhilesh.maxBins, Array(32,33,34,35))
.addGrid(randomForestModel_nikhilesh.impurity, Array("entropy", "gini")).build()
val evaluator_nikhilesh = new MulticlassGlassificationEvaluator()
.setLabelCol("label")
.setTeatinstend("prediction")
.setMetricName("accuracy")
val cross_validator_nikhilesh = new CrossValidator()
.setExtunator (pipeline_nikhilesh)
.setEvaluator (evaluator_nikhilesh)
.setEvaluator (evaluator_nikhilesh)
.setEvaluator (evaluator_nikhilesh)
.setEvaluator (pipeline_nikhilesh)
     setNumFolds(3)
   val Array(trainingData, testData) = df_5_nikhilesh.randomSplit(Array(0.8, 0.2), 754)
  // Exiting paste mode, now interpreting.
  randomForestModel_nikhilesh: org.apache.spark.ml.classification.RandomForestClassifier = rfc_c81dlae75155
pipeline nikhilesh: org.apache.spark.ml.Pipeline = pipeline_d966218e488f
paramGrid_nikhilesh: Array[org.apache.spark.ml.param.ParamMap] =
                      rfc_c8ldlae75155-impurity: entropy,
rfc_c8ldlae75155-maxBins: 32,
rfc_c8ldlae75155-maxDepth: 3
}, {
                      rfc_c81dlae75155-impurity: entropy,
rfc_c81dlae75155-maxBins: 33,
rfc_c81dlae75155-maxDepth: 3
}, {
                      rfc_c81dlae75155-impurity: entropy,
rfc_c81dlae75155-maxBins: 34,
rfc_c81dlae75155-maxDepth: 3
}, {
                      rfc_c81dlae75155-impurity: entropy,
rfc_c81dlae75155-maxBins: 35,
rfc_c81dlae75155-maxDepth: 3
                      rfc_c81dlae75155-impurity: gini, rfc_c81dlae75155-max...
      ala>
    /11/23 05:05:12 When org.apach: spaint.opach plants request. Exit code 137:022-11-23 05:05:12.069]Container exited with a non-zero exit code 137:022-11-23 05:05:12.069]Killed by external signal
                            val predictions_nikhilesh = model_nikhilesh.transform(testData)
   predictions_nikhilesh: org.apache.spark.sql.DataFrame = [textID: string, tweet_text_final: string ... 12 more fields]
   scala> val accuracy_nikhilesh = evaluator_nikhilesh.evaluate(predictions_nikhilesh)
22/11/23 07:13:50 WARN org.apache.spark.scheduler.DAGScheduler: Broadcasting large task binary with size 1791.0 KiB
    accuracy_nikhilesh: Double = 0.7203810369307343
```