In [0]:	Opinion Mining Project on Amazon Product Review Dataset Import all the libraries and function from SparkSQL and SparkMLib from pyspark.sql import SparkSession import pyspark.sql.functions as F from pyspark.sql import SparkSession from pyspark.sql.types import * #import wget from pyspark.ml.feature import Bucketizer,RegexTokenizer,StopWordsRemover,CountVectorizer,IDF, NGram, HashingTF, StringIndexer from pyspark.sql.functions import * from pyspark.sql.functions import LogisticRegression, RandomForestClassificationModel, RandomForestClassifier, GBTClassificationModel, GBTClass from pyspark.ml import Pipeline,PipelineModel from pyspark.ml.evaluation import BinaryClassificationEvaluator
In [0]:	<pre>spark = SparkSession.builder.appName("Opinion mining on Amazon Fashion product reviews").getOrCreate() spark.sparkContext.setLogLevel("WARN") #Create UnionAll function to merge dataframs from functools import reduce # For Python 3.x from pyspark.sql import DataFrame def unionAll(*dfs):</pre>
In [0]: In [0]:	<pre># File location and type for Clothing , Shoes and Jewelry data file_location = "/FileStore/tables/reviews_Clothing_Shoes_and_Jewelry_5_json.gz" file_type = "gz" # The applied options are for CSV files. For other file types, these will be ignored. clsj = spark.read.json(file_location)</pre> # File location and type for beauty data file_location = "/FileStore/tables/reviews_Beauty_5_json.gz" file_type = "gz" # File location = "/FileStore/tables/reviews_Beauty_5_json.gz" file_type = "gz"
In [0]: In [0]:	<pre>file_type = "gz" # The applied options are for CSV files. For other file types, these will be ignored. bt = spark.read.json(file_location) #Execute union to merge bt and clsj dataframes df = unionAll(bt,clsj) #Drop the columns that are not required part 1 df = df.drop("image")\</pre>
In [0]:	<pre>df.columns Out[49]: ['asin', 'helpful', 'overall', 'reviewText', 'reviewTime', 'reviewerID', 'reviewerName', 'summary', 'unixReviewTime']</pre> df.printSchema()
	root asin: string (nullable = true) helpful: array (nullable = true)
In [0]:	<pre># Drop and columns and concatenate summary and reviewtext as Text column df = df.withColumn("text",concat(F.col("summary"), lit(" "),F.col("reviewText")))\ .drop("helpful")\ .drop("reviewerID")\ .drop("reviewerName")\ .drop("reviewTime")</pre> #Print the Count of number of instances df.count() Out[53]: 477179
In [0]:	#Describe the Distribution of the ratings column 'overall' df.describe("overall").show() + summary overall + count 477179 mean 4.2223610007984425 stddev 1.1306299814841152 min 1.0 max 5.0 +
<pre>In [0]:</pre>	<pre>#Data Cleaning : Finding Null Values in each column for col in df.columns: print(col,":",df[df[col].isNull()].count()) asin : 0 overall : 0 reviewText : 0 summary : 0 unixReviewTime : 0 text : 0 #Filter out the rows with neutral overall ratings</pre>
	<pre>df1 = df.filter("overall !=3") #Create a split such that from -inf to 4 will be placed in bucket 0 and 4 to inf will be placed in bucket 1 splits = [-float("inf"), 4.0, float("inf")] #Bucketize data and create labels 0 if overall rating is in (1.0,2.0), #otherwise 1 bucketizer = Bucketizer(splits=splits,\</pre>
In [0]:	overall label count ++ 2.0 0.0 26919 5.0 1.0 277771 1.0 0.0 21718 4.0 1.0 98098 ++ #Diplay the dataframe after droping the columns and creating a new target variable/column 'label' based on overall ratings df2.show() ++ asin overall reviewText summary unixReviewTime text label
In [0]:	1.0 Very oily and cre Don't waste your 1391040000 Don't waste your 0.0 7806397051 4.0 The texture of th great quality 1376425600 great quality The 1.0 7806397051 2.0 I really can't te Do not work on my 1386460800 Don to work on my 0.0 7806397051 5.0 I was very happy Very nice palette! 136594000 Very nice palette 1.0 7806397051 1.0 PLEASE DONT DO IT Smh!!! 1376611200 Smh!! PLEASE DON 0.0 7950931051 2.0 Chalky, Not Pigmen Chalky, Not Pigmen 1387252800 Chalky, Not Pigmen 0.0 975991062 2.0 Did nothing for m Nothing 1405209600 no Lightening, no 0.0 9759991062 1.0 Did nothing for m Nothing 1390435200 This works 1390435200 This works 1390435200 This works 10.0 9759991062 5.0 I bought this pro This works 1390435200 This works 10.0 975991062 5.0 I institute of this companies 1.0 975991062 5.0 I used it for ana Durns 1392631600 Nothing 10.0 975991062 5.0 I order this crea Did work for me 1392631600 Did work for me 1.0 975991062 4.0 Good product. Use excellent 1382634400 Excellent Good pr 1.0 9758972216 5.0 I haven't been a Love the smell of 1382634400 Excellent Good pr 1.0 9788972216 5.0 To his is the first Very good 13790432600 Lurrrrrrrv 1.0 9788972216 5.0 I his product has Great Scent 1394496000 Spring Garden in 1.0 9788972216 5.0 I his product has Great Scent 1395763200 Great Gent Liver 1.0 9788972216 5.0 I his product has Great Scent 1395763200 Great Gent Liver 1.0 9788972216 5.0 I his product has Great Scent 1395763200 Great Gent Liver 1.0 9788972216 5.0 I his product has Great Scent 1394496000 Spring Garden in 1.0 978909061 5.0 I'm very picky who Spring Garden in 139449
In [0]: In [0]:	#Split data as 80-20% Train and Test dataset splitSeed = 5043 trainingData, testData = df3.randomSplit([.80, 0.20], splitSeed) NLP Text Preprocessing, Pipelines and Model Build #Tokenize the sentence based on the regex pattern
	tokenizer = RegexTokenizer(inputCol="text", outputCol="reviewTokensUf", pattern="\\s+ [,.()\\"]") #Remove Stop Words that do not contribute in any way to our analysis stopwords_remover = StopWordsRemover(stopWords=StopWordsRemover.loadDefaultStopWords("english"), inputCol="reviewTokensUf", outputCol="reviewTokens") #converts word documents to vectors of token counts cv = CountVectorizer(inputCol="reviewTokens", outputCol="cv", vocabSize=296337) #IDF model idf = IDF(inputCol="cv", outputCol="features") #Logistic Boosted Classifier lr = LogisticRegression(maxIter=100, regParam=0.02, elasticNetParam=0.3)
In [0]:	<pre>#Create a pipeline by combining all the functions we defined above - tokenizer , stopwords_remover, cv, idf, gbtc steps = [tokenizer, stopwords_remover, cv, idf, lr] lr_pipeline = Pipeline(stages=steps) #fit the training dataset dataset into the pipeline model = lr_pipeline.fit(trainingData) #Obtain the predictions from the model predictions = model.transform(testData) #Call the Binary Classification Evaluator function evaluator = BinaryClassificationEvaluator() areaUnderROC = evaluator.evaluate(lr_predictions) print('Test Area Under ROC for Linear Regression model with Count Vectorizer is ', areaUnderROC)</pre>
In [0]:	Test Area Under ROC for Linear Regression model with Count Vectorizer is 0.9398334608052311 Build a Vocabulary of all the relevant tokens and display the co-efficients #Building the vocabulary to explore the coefficients vocabulary = model.stages[2].vocabulary weights = model.stages[-1].coefficients.toArray() weights = [float(weight) for weight in weights] schema = StructType([StructField('word', StringType()),
In [0]:	#Shows the top 10 positive sentiment tokens cdf.orderBy(desc("weight")).show(10) ++ word weight +
In [0]:	comfortable 0.24460426 excellent 0.22886382 great! 0.22694735 compliments 0.21397798 highly 0.20946273 +
	disappointed -0.4081317 returned -0.31633556 waste -0.29761943 disappointing -0.26205418 poor -0.23398674 returning -0.22233732 unfortunately -0.22159778 return -0.22129418 disappointment -0.20847303 cheap -0.19569387 +
In [0]:	<pre>#model evaluation and metrics for Logistic Regression vanila version lp = predictions.select("label", "prediction") counttotal = predictions.count() correct = lp.filter(F.col("label") == F.col("prediction")).count() wrong = lp.filter(~(F.col("label") == F.col("prediction"))).count() ratioWrong = float(wrong) / float(counttotal) ratioCorrect=correct/counttotal trueneg =(lp.filter(F.col("label") == 0.0).filter(F.col("label") == F.col("prediction")).count()) /counttotal truepos = (lp.filter(F.col("label") == 1.0).filter(F.col("label") == F.col("prediction")).count())/counttotal falseneg = (lp.filter(F.col("label") == 0.0).filter(~(F.col("label") == F.col("prediction"))).count())/counttotal falsepos = (lp.filter(F.col("label") == 1.0).filter(~(F.col("label") == F.col("prediction"))).count())/counttotal</pre>
	<pre>precision = truepos / (truepos + falsepos) recall = truepos / (truepos + falseneg) #fmeasure= 2 precision recall / (precision + recall) accuracy=(truepos + trueneg) / (truepos + trueneg + falsepos + falseneg) print('counttotal :', counttotal) print('correct :', correct) print('wrong :', wrong) print('ratioWrong :', ratioWrong) print('ratioCorrect :', ratioCorrect) print('truen :', trueneg) print('truep :', truepos) print('falsen :', falseneg)</pre>
	<pre>print('falsep</pre>
In [0]:	recall : 0.8731763686263181 accuracy : 0.8694618418000464 # Display the columns of predictions dataframe after transformation predictions.printSchema() display(predictions) root asin: string (nullable = true) overall: double (nullable = true) reviewText: string (nullable = true) summarry: string (nullable = true) unixReviewTime: long (nullable = true) text: string (nullable = true) text: string (nullable = true)
	label: double (nullable = true) reviewTokensUf: array (nullable = true) element: string (containsNull = true) reviewTokens: array (nullable = true) element: string (containsNull = true) reviewTokens: array (nullable = true) rewiewTokens: array (nullable = true) rewiewTokens: array (nullable = true) rewiewTokens: array (nullable = true) reviewTokens: array (nul
In [0]:	<pre>#Model Logistic Regrssion with Bigrams and Hashing TF with IDF vectorization #Create a pipeline by combining all the functions we defined above - tokenizer , stopwords_remover, cv, idf, gbtc #Tokenize the sentence based on the regex pattern tokenizer = RegexTokenizer(inputCol="text",outputCol="reviewTokensUf",pattern="\\s+ [,.()\\"]") bigram = NGram(inputCol = "reviewTokensUf", outputCol = "bigrams", n = 2)</pre>
In [0]:	tfs = HashingTF(inputCol="bigrams", outputCol="h_features") #IDF model idf = IDF(inputCol="h_features", outputCol="features") lr = LogisticRegression(maxIter=20) steps = [tokenizer, stopwords_remover, bigram, tfs, idf, lr] bigrams_pipeline = Pipeline(stages=steps) model = bigrams_pipeline.fit(trainingData) bi_predictions = model.transform(testData) evaluator = BinaryClassificationEvaluator() areaUnderROC = evaluator.evaluate(bi_predictions) print('Test Area Under ROC for Bigrams linear regression', areaUnderROC) Test Area Under ROC for Bigrams linear regression 0.9467960427315577 #model evaluation lp = bi_predictions.select("label", "prediction") counttotal = bi_predictions.count() correct = lp.filter(F.col("label")= F.col("prediction")).count()
	<pre>wrong = lp.filter(-[c.col("label") == F.col("prediction"))).count() ratioWrong = float(wrong) / float(counttotal) ratioCorrect=correct/counttotal trueneg = (lp.filter(F.col("label") == 0.0).filter(F.col("label") == F.col("prediction")).count()) / counttotal truepos = (lp.filter(F.col("label") == 0.0).filter(F.col("label") == F.col("prediction")).count()) / counttotal falseneg = (lp.filter(F.col("label") == 0.0).filter(-(F.col("label") == F.col("prediction"))).count()) / counttotal falseneg = (lp.filter(F.col("label") == 0.0).filter(-(F.col("label") == F.col("prediction"))).count()) / counttotal falseneg = (lp.filter(F.col("label") == 0.0).filter(-(F.col("label") == F.col("prediction"))).count()) / counttotal precision = truepos / (truepos + falseneg) precision = truepos / (truepos + falseneg) #fmeasure= 2 precision recall / (precision + recall) accuracy=(truepos + trueneg) / (truepos + trueneg + falseneg) print('counttotal :', counttotal) print('correct :', correct) print('ratioWrong :', wrong) print('ratioWrong :', wrong) print('ratioWrong :', ratioWrong) print('rtuen :', truepos) print('truen :', truepos) print('falsen :', falseneg) print('falsen :', falseneg) print('falsen :', falseneg) print('recall :', recall) #print('fmeasure :', fmeasure) print('fmeasure :', couracy)</pre>
In [0]:	counttotal : 17244 correct : 15314 wrong : 1930 ratioWorng : 0.11192298770586871 ratioCorrect : 0.8880770122941313 truen : 0.5028995592669914 truep : 0.38517745302713985 falsen : 0.06692182788216192 falsep : 0.045001159823706796 precision : 0.8953895928821785 recall : 0.8519753719856337 accuracy : 0.8880770122941313
	bi_predictions.select(F.col("bigrams")).show(5) +
In [0]:	<pre>#RandomForest Implementation #RandomForest Implementation regexTokenizer = RegexTokenizer(inputCol="text", outputCol="words", pattern="\\w") #Remove Stop Words that do not contribute in any way to our analysis stopwords_remover_rf = \ StopWordsRemover(stopWords=StopWordsRemover.loadDefaultStopWords("english"),inputCol="words",outputCol="filtered") # bag of words count countVectors = CountVectorizer(inputCol="filtered", outputCol="features", vocabSize=10000, minDF=5) label_stringIdx = StringIndexer(inputCol = "label", outputCol = "new_label")</pre>
	<pre>rf = RandomForestClassifier(numTrees=3, maxDepth=2, labelCol="new_label", seed=42, \</pre>
In [0]:	<pre>rf_predictions.printSchema() rf_predictions.show(5) root asin: string (nullable = true) overall: double (nullable = true) reviewText: string (nullable = true) summary: string (nullable = true) unixReviewTime: long (nullable = true) text: string (nullable = true) text: string (nullable = true) reviewTokensUf: array (nullable = true)</pre>
	element: string (containsNull = true) reviewTokens: array (nullable = true) element: string (containsNull = true) bigrams: array (nullable = true) element: string (containsNull = true) h_features: vector (nullable = true) h_features: vector (nullable = true) rawPrediction: vector (nullable = true) probability: vector (nullable = true) prediction: double (nullable = false) +
In [0]:	1.0 Did nothing for m Nothing 1392681600 Nothing Did nothi 0.0 Inothing, did, no Inothing, nothing
In [0]:	<pre>cf = classification_report(y_pred, y_orig) print(cf) precision recall f1-score</pre>
	<pre>lp = rf_predictions.select("label", "prediction") counttotal = rf_predictions.count() correct = lp.filter(F.col("label") == F.col("prediction")).count() wrong = lp.filter(~(F.col("label") == F.col("prediction"))).count() ratioWrong = float(wrong) / float(counttotal) ratioCorrect=correct/counttotal trueneg =(lp.filter(F.col("label") == 0.0).filter(F.col("label") == F.col("prediction")).count()) /counttotal truepos = (lp.filter(F.col("label") == 1.0).filter(F.col("label") == F.col("prediction")).count())/counttotal falseneg = (lp.filter(F.col("label") == 0.0).filter(~(F.col("label") == F.col("prediction"))).count())/counttotal falsepos = (lp.filter(F.col("label") == 1.0).filter(~(F.col("label") == F.col("prediction"))).count())/counttotal precision = truepos / (truepos + falsepos)</pre>
	<pre>precision = truepos / (truepos + falsepos) recall = truepos / (truepos + falseneg) #fmeasure= 2 precision recall / (precision + recall) accuracy=(truepos + trueneg) / (truepos + trueneg + falsepos + falseneg) print('counttotal :', counttotal) print('correct :', correct) print('wrong :', wrong) print('ratioWrong :', ratioWrong) print('ratioCorrect :', ratioCorrect) print('truen :', trueneg) print('truep :', truepos) print('falsen :', falseneg) print('falsep :', falsepos) print('precision :', precision) print('recall :', recall)</pre>
In [0]: In [0]:	recall : 0.9764150943396226 accuracy : 0.5815356065877987 Naive Bayes nb = NaiveBayes(smoothing=1) steps = [regexTokenizer, stopwords_remover_rf, countVectors, label_stringIdx, nb] nb_pipeline = Pipeline(stages=steps) model = nb_pipeline.fit(trainingData) nb_predictions = model.transform(testData) evaluator = BinaryClassificationEvaluator()
In [0]:	<pre>areaUnderROC = evaluator.evaluate(nb_predictions) print('Test Area Under ROC for Naive Bayes Classification is ', areaUnderROC) Test Area Under ROC for Naive Beyes Classification is 0.5399622753791408 #model evaluation lp = nb_predictions.select("label", "prediction") counttotal = nb_predictions.count() correct = lp.filter(F.col("label")== F.col("prediction")).count() wrong = lp.filter(-(F.col("label") == F.col("prediction"))).count() ratioWrong = float(wrong) / float(counttotal) ratioCorrect=correct/counttotal</pre>
	<pre>trueneg =(lp.filter(F.col("label") == 0.0).filter(F.col("label") == F.col("prediction")).count()) /counttotal truepos = (lp.filter(F.col("label") == 1.0).filter(F.col("label") == F.col("prediction")).count())/counttotal falseneg = (lp.filter(F.col("label") == 0.0).filter(~(F.col("label") == F.col("prediction"))).count())/counttotal falsepos = (lp.filter(F.col("label") == 1.0).filter(~(F.col("label") == F.col("prediction"))).count())/counttotal precision = truepos / (truepos + falsepos) recall = truepos / (truepos + falseneg) #fmeasure= 2 precision recall / (precision + recall) accuracy=(truepos + trueneg) / (truepos + trueneg + falsepos + falseneg) print('counttotal :', counttotal) print('correct :', correct) print('wrong :', wrong)</pre>
	<pre>print('ratioWrong :', ratioWrong) print('ratioCorrect :', ratioCorrect) print('truen :', trueneg) print('truen :', truepos) print('falsen :', falseneg) print('falsep :', falsepos) print('precision :', precision) print('recall :', recall) #print('fmeasure :', fmeasure) print('daccuracy :', accuracy)</pre>
In [0]:	ratioWrong : 0.116214335421016 ratioCorrect : 0.883785664578984 truen : 0.5045813036418464 truep : 0.37920436093713755 falsen : 0.0652400835073069 falsep : 0.050974251913709114 precision : 0.881504448638447 recall : 0.8532098121085595 accuracy : 0.883785664578984 #GBTClassifier : Gradiennt Boosted Trees #Tokenize the sentence based on the regex pattern tokenizer = RegexTokenizer(inputCol="text",outputCol="reviewTokensUf",pattern="\\s+ [,.()\"]")
	<pre>#Remove Stop Words that do not contribute in any way to our analysis stopwords_remover = StopWordsRemover(stopWordsRemover.loadDefaultStopWords("english"),inputCol="reviewTokensUf",outputCol="reviewTokens") #converts word documents to vectors of token counts cv = CountVectorizer(inputCol="reviewTokens",outputCol="cv",vocabSize=296337) steps = [tokenizer, stopwords_remover, cv] Pipeline_mo = Pipeline(stages=steps) transformed_model = Pipeline_mo.fit(trainingData) train_tr = transformed_model.transform(trainingData) test_tr = transformed_model.transform(testData)</pre>
In [0]:	<pre>train_tr.printSchema() root asin: string (nullable = true) overall: double (nullable = true) reviewText: string (nullable = true) summary: string (nullable = true) unixReviewTime: long (nullable = true) text: string (nullable = true) label: double (nullable = true) element: string (containsNull = true) element: string (containsNull = true)</pre>
In [0]:	
	·
In [0]:	<pre># training gbtc = GBTClassifier(featuresCol="cv", labelCol="label", maxIter=20) gbtc = gbtc.fit(train_tr) # prediction pred = gbtc.transform(test_tr) pred.show(3) </pre>
In [0]:	1.0 Did nothing for m Nothing 1392681600 Nothing Did nothi 0.0 Inothing, did, no Inothing, nothing 0.0
In [0]:	GBTClassificationModel: uid = GBTClassifier_b7bb4030176d, numTrees=20, numClasses=2, numFeatures=87097 evaluator = BinaryClassificationEvaluator() areaUnderROC = evaluator.evaluate(pred) print('Test Area Under ROC for Gradient Boost Classification is ', areaUnderROC) Test Area Under ROC for Naive Bayes Classification is 0.8692959846983227 pred.show(5) +
	asin overall reviewText summary unixReviewTime text label reviewTokensUf reviewTokens cv rawPrediction probability prediction
In [0]:	from sklearn.metrics import confusion_matrix, classification_report from pyspark.ml.evaluation import MulticlassClassificationEvaluator evaluator=MulticlassClassificationEvaluator(metricName="accuracy") #e acc = evaluator.evaluate(pred) print("Prediction Accuracy: ", acc) y_pred=pred.select("prediction").collect() y_orig=pred.select("label").collect() cm = confusion_matrix(y_orig, y_pred) print("Confusion Matrix:") prediction Accuracy: 0.7811992577128276 Confusion Matrix: [[8332 1494]
In [0]:	[2279 5139]] from sklearn.metrics import confusion_matrix, classification_report cf = classification_report(y_orig, y_pred) print(cf) precision recall f1-score support 0.0 0.79 0.85 0.82 9826 1.0 0.77 0.69 0.73 7418 accuracy 0.78 17244 macro avg 0.78 0.77 0.77 17244
In [0]:	macro avg 0.78 0.77 0.77 17244 weighted avg 0.78 0.78 17244