<pre># File location file_location = file_type = "gz" # The applied op bt = spark.read. df = unionAll(b</pre>	B://projectdatabia678/reviews_Beauty_5.json.gz ons are for CSV files. For other file types, ton(file_location) clsj)		
root asin: string helpful: arr	<pre>"overall', 'reviewText', 'reviewTime', 'reviewTime',</pre>	ewerID', 'reviewerName', 'summary', 'u	nixReviewTime']
df = df.withColu .drop("helpful" .drop("reviewer .drop("reviewer .drop("reviewTi # COMMAND #Print the Count df.count() 477179 #COMMAND #Describe the Di df.describe("ove	<pre>: long (nullable = true) ("text", concat(col("summary"), lit(" "), col("r ')\ ne")\ ') f number of instances ribution of the ratings column 'overall' ll").show()</pre>	reviewText")))\	
<pre>for col in df.co print(col,": # COMMAND #Filter out the df1 = df.filter(splits = [-float #Bucketize data bucketizer = Buc df2= bucketizer.</pre>	<pre>inding Null Values mns: df[df[col].isNull()].count()) ows with neutral overall ratings verall !=3") inf"), 4.0, float("inf")] d create labels 0 if overall rating is in (1.6) tizer(splits=splits,\</pre>	0,2.0), #otherwise 1	
<pre>fractions = {1.0 df3 = df2.stat.s df3.groupBy("lab # COMMAND #Split data as 8 splitSeed = 5043</pre>	eate train and test dataset .1, 0.0 : 1.0} oleBy("label", fractions, 36) ').count().show() 20% Train and Test dataset Oata = df3.randomSplit([0.80, 0.20], splitSeed		
count mean 4.222361 stddev 1.130629 min max +	477179 07984425 14841152 1.0 5.0 + nt + 71 18 98 19		
asin overal 7806397051 1 7806397051 4 7806397051 5 7806397051 5 7806397051 2 7806397051 2 9759091062 2 9759091062 1 9759091062 1 9759091062 5 9759091062 5	reviewText summary unit	IXReviewTime text label tex	e1 + 0 0 0 0 0 0 0 0 0 0
9788072216 5 9790790961 5 5 5 5 5 5 5 5 5	So I got this abo Lurrrrrrvv This product has Great Scent I'm very picky wh Spring Garden in +	ensUf", pattern="\\s+ [,.()\"]") Eysis loadDefaultStopWords("english"), input	0 0 +
<pre>#Logistic Booste lr = LogisticReg #Create a pipeli steps = [tokeni lr_pipeline = Pi #fit the trainin model = lr_pipel #Obtain the pred</pre>	classifier ssion(maxIter=100,regParam=0.02,elasticNetPara by combining all the functions we defined above, stopwords_remover, cv, idf, lr] line(stages=steps) dataset dataset into the pipeline e.fit(trainingData) tions from the model l.transform(testData)		v, idf, gbtc
evaluator = Bina areaUnderROC = e print('Test Area # COMMAND #Building the vo vocabulary = mod weights = model. weights = [float schema = StructT	Classification Evaluator function ClassificationEvaluator() Luator.evaluate(predictions) Inder ROC for Linear Regression model with Councilor Coulary to explore the coefficients Istages[2].vocabulary Index Rocefficients.toArray() Index Rocefficients Istages[-1].coefficients.toArray() Index Rocefficients Istages[-1].coefficients Index Rocefficients Index Roceffici	ot Vectorizer is ' , areaUnderROC)	
# COMMAND cdf.orderBy("wei	for Linear Regression model with Count Vector+ ight + 6475 2465 5672 0935 7106 7884 9578 8377	rizer is 0.9425358441407751	
highly 0.21 +	rows + rows + weight + 6692452 2958275 9348242 5194058 4936205 4399708 1882328 1715495 0842884 9916554 + rows		
<pre>lp.printSchema() lp.filter(~(F.co # COMMAND #model evaluatio lp = predictions counttotal = pre correct = lp.fil wrong = lp.filte ratioWrong = flo ratioCorrect=cor trueneg =(lp.fi truepos = (lp.fi falseneg = (lp.fi</pre>	elect("label", "prediction") ctions.count() r(F.col("label")== F.col("prediction")).count(r(F.col("label") == F.col("prediction"))).coun (wrong) / float(counttotal)	<pre>cit() f) == F.col("prediction")).count()) /count()) /count())/count())/count()) el") == F.col("prediction")).count()</pre>	unttotal)/counttotal
<pre>recall = truepo #fmeasure= 2 pr accuracy=(truepo # COMMAND print('counttota print('correct print('wrong print('ratioWron</pre>	<pre>:', counttotal :', correct :', wrong :', ratioWrong :', ratioCorrect :', trueneg :', truepos :', falseneg :', falsepos :', precision :', recall :', accuracy</pre>	falseneg)	
prediction: counttotal : 17 correct : 15 wrong : 21 ratioWrong : 0. ratioCorrect : 0. truen : 0. truep : 0. falsep : 0. precision : 0. recall : 0. accuracy : 0.	9 583506104584197	ove - tokenizer , stopwords_remover, c	v, idf, gbtc
<pre>bigram = NGram(i tfs = HashingTF #IDF model idf = IDF(inputCle) lr = LogisticReg steps = [tokeni</pre>	chased on the regex pattern ckenizer(inputCol="text",outputCol="reviewToke cutCol = "reviewTokensUf", outputCol = "bigrams nputCol="bigrams", outputCol="h_features") ="h_features",outputCol="features") ssion(maxIter=20) r, stopwords_remover, bigram, tfs, idf, lr] Pipeline(stages=steps)		
# COMMAND evaluator = Bina areaUnderROC = e print('Test Area # COMMAND #model evaluatio lp = bi_predicti counttotal = bi_	Deline.fit(trainingData) Dodel.transform(testData) ClassificationEvaluator() Luator.evaluate(bi_predictions) Doder ROC for Bigrams linear regression', areal		
<pre>wrong = lp.filte ratioWrong = flo ratioCorrect=cor trueneg =(lp.fi truepos = (lp.fi falseneg = (lp.fi falsepos = (lp.f) precision = true recall = truepo #fmeasure= 2 pr accuracy=(truepo print('counttota print('correct print('wrong print('ratioWrong)</pre>	<pre>c(F.col("label") == F.col("prediction"))).coun (wrong) / float(counttotal) er(F.col("label") == 0.0).filter(F.col("label") er(F.col("label") == 1.0).filter(F.col("label") ter(F.col("label") == 0.0).filter(~(F.col("label")) ter(F.col("label") == 1.0).filter(~(F.col("label")) ter(F.col("label") == 1.0).filter(~(F.col("label")) s / (truepos + falsepos) / (truepos + falseneg) ision recall / (precision + recall) trueneg) / (truepos + trueneg + falsepos + f :', counttotal :', correct :', wrong :', ratioWrong :', ratioCorrect)</pre>	<pre>it() f) == F.col("prediction")).count()) /count() /</pre>	unttotal)/counttotal
print('truep print('falsen print('falsep print('precision print('recall #print('fmeasure print('accuracy # COMMAND bi_predictions.s # COMMAND Test Area Under R counttotal : 17 correct : 15 wrong : 18	<pre>:', truepos :', falseneg :', falsepos :', precision :', recall :', fmeasure :', accuracy) cct(F.col("bigrams")).show(5) for Bigrams linear regression 0.94646558914524 9</pre>	2727	
truep : 0. falsen : 0. falsep : 0. precision : 0. recall : 0. accuracy : 0. +	ms + + rows	ls", pattern="\\W")	
<pre>stopwords_remove StopWordsRemover # bag of words c countVectors = C label_stringIdx rf = RandomFores rf = RandomFores</pre> steps = [regexT	topWords=StopWordsRemover.loadDefaultStopWords	s("english"),inputCol="words",outputCol="words",	l="filtered")
<pre># COMMAND rf_predictions.s # COMMAND evaluator = Bina areaUnderROC = e print('Test Area # COMMAND</pre>	sts_pipeline.fit(trainingData) odel.transform(testData) v(5) ClassificationEvaluator() luator.evaluate(rf_predictions) oder ROC for Random Forests Classification is	', areaUnderROC)	
<pre>counttotal = rf_ correct = lp.fil wrong = lp.filte ratioWrong = flo ratioCorrect=cor trueneg =(lp.fi truepos = (lp.fi falseneg = (lp.fi falsepos = (lp.fi precision = true recall = truepo #fmeasure= 2 pr</pre>	<pre>er(F.col("label") == 0.0).filter(F.col("label") er(F.col("label") == 1.0).filter(F.col("label") ter(F.col("label") == 0.0).filter(~(F.col("label")) ter(F.col("label") == 1.0).filter(~(F.col("label")) s / (truepos + falsepos) / (truepos + falseneg) ision recall / (precision + recall) + trueneg) / (truepos + trueneg + falsepos + falsepo</pre>	<pre>cit() f) == F.col("prediction")).count()) /count() /count()</pre>	unttotal)/counttotal
print('truen print('truep print('falsen print('falsep print('precision print('recall #print('fmeasure print('accuracy # COMMAND asin overafeatures new_labe	<pre>:', ratioWrong) t:', ratioCorrect) :', trueneg) :', truepos) :', falseneg) :', falsepos) :', precision) :', recall) :', fmeasure) :', accuracy)</pre>	ixReviewTime text labe	el words fi
(10000, [8, 22, 24, 2] B00004U9V2 5 (10000, [5, 24, 26, 2] B000052WYD 5 (10000, [7, 36, 55, 6] B000052WYD 5 (10000, [2, 3, 8, 15,] B000052WYD 5 (10000, [3, 30, 81, 8] +	for Random Forests Classification is 0.89567 4 96821008984105	37240 0.0 1399248000 A favorite My son 1 78709 0.0 1365033600 Really covers. I 1 91399 0.0 1309910400 Great little prod 1 93204 0.0 1391817600 MUST HAVE My must 1 16372 0.0	<pre>0 [a, favorite, my, [favorite, sor 0 [really, covers, [really, cover 0 [great, little, p [great, little 0 [must, have, my, [must, must, good</pre>
<pre>falsen : 0. falsep : 0. precision : 0. recall : 1. accuracy : 0. nb = NaiveBayes(steps = [regexT nb_pipeline = Pi # COMMAND model = nb_pipel</pre>	96821008984105 718562874251497 03178991015895 pothing=1) enizer, stopwords_remover_rf, countVectors, la line(stages=steps) e.fit(trainingData) odel.transform(testData)	dbel_stringIdx, nb]	
<pre>#gbtc = GBTClass #steps = [token #pipeline = Pipe # COMMAND #model = pipelin # COMMAND #predictions = m # COMMAND #evaluator = Bin #areaUnderROC =</pre>	er, stopwords_remover, cv, idf, gbtc] ne(stages=steps) fit(trainingData) el.transform(testData) yClassificationEvaluator() aluator.evaluate(predictions) Under ROC', areaUnderROC)	ove - tokenizer , stopwords_remover, c	v, idf, gbtc
<pre>#NaiveBayes Impl # COMMAND evaluator = Bina areaUnderROC = e print('Test Area # COMMAND #model evaluatio lp = nb_predicti counttotal = nb_ correct = lp.fil wrong = lp.filte</pre>	ClassificationEvaluator() Luator.evaluate(nb_predictions) nder ROC for Naive Bayes Classification is ', s.select("label", "prediction") edictions.count() r(F.col("label")== F.col("prediction")).count(r(F.col("label") == F.col("prediction"))).count(wrong) / float(counttotal)		
<pre>truepos = (lp.fi falseneg = (lp.fi falsepos = (lp.fi falsepos = (lp.fi precision = true recall = truepo #fmeasure= 2 pr accuracy=(truepo print('counttota print('correct print('wrong print('ratioWron print('ratioCorr print('truen print('truep print('falsen</pre>	<pre>:', correct :', wrong :', ratioWrong :', ratioCorrect :', trueneg :', truepos :', falseneg)</pre>	<pre>p == F.col("prediction")).count())/count() pel") == F.col("prediction"))).count() pel") == F.col("prediction"))).count()</pre>	unttotal)/counttotal
counttotal : 17 correct : 15 wrong : 19 ratioWrong : 0. truen : 0. truep : 0. falsep : 0. precision : 0. recall : 0.	<pre>:', precision) :', recall) :', fmeasure) :', accuracy) for Naive Bayes Classification is 0.54669678 4 7</pre> 040082930200415	347903211	
<pre>#Using XGBoost #Create a pipeli #Tokenize the se tokenizer = Rege word2Vec = Word2 steps = [tokeni xgboost_prepro = xg_boost_model = train_df = xg_boo test_df = xg_boo #Create a pipeli</pre>	by combining all the functions we defined above ence based on the regex pattern okenizer(inputCol="text",outputCol="reviewToker c(vectorSize=5, seed=42, inputCol="reviewToker	ensUf", pattern="\\s+ [,.()\"]") usUf", outputCol="features")	
<pre>#Remove Stop Wor stopwords_remove #converts word d cv = CountVector steps = [tokeni Pipeline_mo = Pi transformed_mode train_tr = trans test_tr = transf</pre>	chace based on the regex pattern okenizer(inputCol="text", outputCol="reviewToke") that do not contribute in any way to our analy stopWordsRemover(stopWords=StopWordsRemover. uments to vectors of token counts er(inputCol="reviewTokens", outputCol="cv", voca r, stopwords_remover, cv] line(stages=steps) Pipeline_mo.fit(trainingData) rmed_model.transform(trainingData) med_model.transform(testData) r(maxIter=5, maxDepth=2, labelCol="label", see	loadDefaultStopWords("english"),input	Col="reviewTokensUf",outputCol="revie
#Tokenize the se	_tr) radiennt Boosted Trees ence based on the regex pattern okenizer(inputCol="text",outputCol="reviewToke") that do not contribute in any way to our anal = StopWordsRemover(stopWords=StopWordsRemover. uments to vectors of token counts	ysis	Col="reviewTokensUf",outputCol="revie