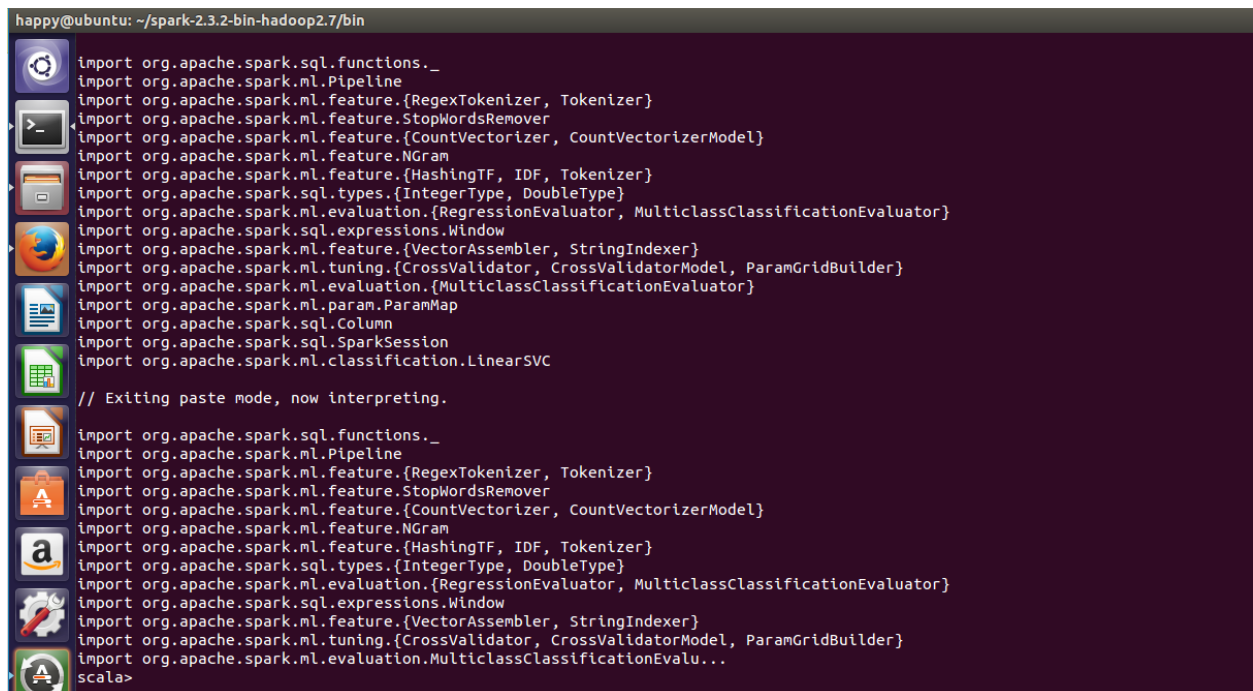


```

import org.apache.spark.sql.functions._
import org.apache.spark.ml.Pipeline
import org.apache.spark.ml.feature.{RegexTokenizer, Tokenizer}
import org.apache.spark.ml.feature.StopWordsRemover
import org.apache.spark.ml.feature.{CountVectorizer, CountVectorizerModel}
import org.apache.spark.ml.feature.NGram
import org.apache.spark.ml.feature.{HashingTF, IDF, Tokenizer}
import org.apache.spark.sql.types.{IntegerType, DoubleType}
import org.apache.spark.ml.evaluation.{RegressionEvaluator, MulticlassClassificationEvaluator}
import org.apache.spark.sql.expressions.Window
import org.apache.spark.ml.feature.{VectorAssembler, StringIndexer}
import org.apache.spark.ml.tuning.{CrossValidator, CrossValidatorModel, ParamGridBuilder}
import org.apache.spark.ml.evaluation.{MulticlassClassificationEvaluator}
import org.apache.spark.ml.param.ParamMap
import org.apache.spark.sql.Column
import org.apache.spark.sql.SparkSession
import org.apache.spark.ml.classification.LinearSVC

```



A terminal window with a dark purple background and a sidebar of application icons on the left. The terminal shows the same Spark imports as the previous block, followed by a comment and a second set of imports. The prompt 'scala>' is visible at the bottom.

```

happy@ubuntu: ~/spark-2.3.2-bin-hadoop2.7/bin
import org.apache.spark.sql.functions._
import org.apache.spark.ml.Pipeline
import org.apache.spark.ml.feature.{RegexTokenizer, Tokenizer}
import org.apache.spark.ml.feature.StopWordsRemover
import org.apache.spark.ml.feature.{CountVectorizer, CountVectorizerModel}
import org.apache.spark.ml.feature.NGram
import org.apache.spark.ml.feature.{HashingTF, IDF, Tokenizer}
import org.apache.spark.sql.types.{IntegerType, DoubleType}
import org.apache.spark.ml.evaluation.{RegressionEvaluator, MulticlassClassificationEvaluator}
import org.apache.spark.sql.expressions.Window
import org.apache.spark.ml.feature.{VectorAssembler, StringIndexer}
import org.apache.spark.ml.tuning.{CrossValidator, CrossValidatorModel, ParamGridBuilder}
import org.apache.spark.ml.evaluation.{MulticlassClassificationEvaluator}
import org.apache.spark.ml.param.ParamMap
import org.apache.spark.sql.Column
import org.apache.spark.sql.SparkSession
import org.apache.spark.ml.classification.LinearSVC

// Exiting paste mode, now interpreting.

import org.apache.spark.sql.functions._
import org.apache.spark.ml.Pipeline
import org.apache.spark.ml.feature.{RegexTokenizer, Tokenizer}
import org.apache.spark.ml.feature.StopWordsRemover
import org.apache.spark.ml.feature.{CountVectorizer, CountVectorizerModel}
import org.apache.spark.ml.feature.NGram
import org.apache.spark.ml.feature.{HashingTF, IDF, Tokenizer}
import org.apache.spark.sql.types.{IntegerType, DoubleType}
import org.apache.spark.ml.evaluation.{RegressionEvaluator, MulticlassClassificationEvaluator}
import org.apache.spark.sql.expressions.Window
import org.apache.spark.ml.feature.{VectorAssembler, StringIndexer}
import org.apache.spark.ml.tuning.{CrossValidator, CrossValidatorModel, ParamGridBuilder}
import org.apache.spark.ml.evaluation.MulticlassClassificationEvalu...
scala>

```

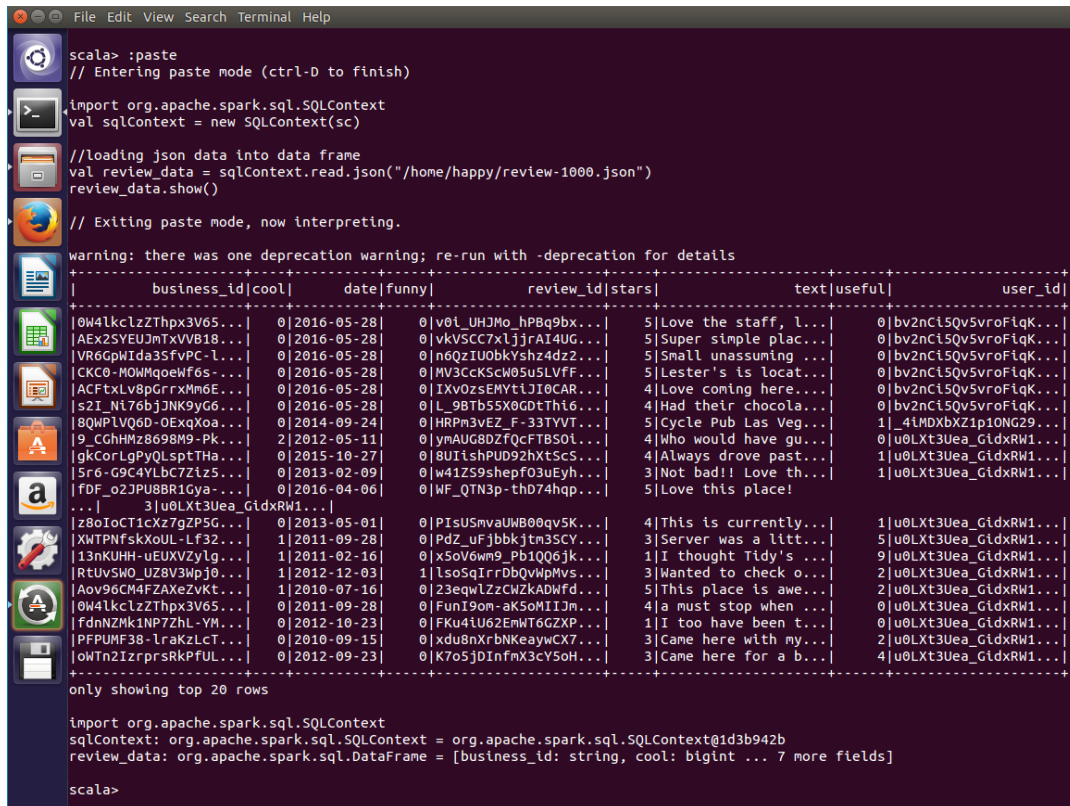
```
import org.apache.spark.sql.SQLContext
```

```
val sqlContext = new SQLContext(sc)
```

```
//loading json data into data frame
```

```
val review_data = sqlContext.read.json("/home/happy/review-1000.json")
```

```
review_data.show()
```



```
scala> :paste
// Entering paste mode (ctrl-D to finish)

import org.apache.spark.sql.SQLContext
val sqlContext = new SQLContext(sc)

//loading json data into data frame
val review_data = sqlContext.read.json("/home/happy/review-1000.json")
review_data.show()

// Exiting paste mode, now interpreting.

warning: there was one deprecation warning; re-run with -deprecation for details
-----
|      business_id|cool|      date|funny|      review_id|stars|      text|useful|      user_id|
-----+-----+-----+-----+-----+-----+-----+-----+-----+
|0W4lkc1zZThpx3V65...|0|2016-05-28|0|v0l_UHJMo_hPBq9bx...|5|Love the staff, l...|0|bv2nC15Qv5vroFlqK...|
|AEx2SYEUJmTxVVB18...|0|2016-05-28|0|vkV5CC7xLjJrAI4UG...|5|Super simple plac...|0|bv2nC15Qv5vroFlqK...|
|VR6GpWIda35fVPC-L...|0|2016-05-28|0|n6QzIU0bkYshz4dz2...|5|Small unassuming ...|0|bv2nC15Qv5vroFlqK...|
|CKC0-M0WMqoeWf6s-...|0|2016-05-28|0|MV3CcKScW05u5LVFF...|5|Lester's is locat...|0|bv2nC15Qv5vroFlqK...|
|ACFTxLv8pGrxMn6E...|0|2016-05-28|0|IXv0zSEMYtIjI0CAR...|4|Love coming here...|0|bv2nC15Qv5vroFlqK...|
|s2I_Ni76bjJNK9yG6...|0|2016-05-28|0|L_9BTb55X0GDTThi6...|4|Had their chocola...|0|bv2nC15Qv5vroFlqK...|
|8QWPLVQ6D-0ExqXoa...|0|2014-09-24|0|HRPm3VEZ_F-33TVVT...|5|Cycle Pub Las Veg...|1|4lMDXbXZ1p10NG29...|
|9_CghHMz8698M9-Pk...|2|2012-05-11|0|ymAUG8DZF0cFTBS0L...|4|Who would have gu...|0|u0LXt3Uea_GldxRW1...|
|gkCorLgPyQLsptTha...|0|2015-10-27|0|8UIiShPUD92hXtScS...|4|Always drove past...|1|u0LXt3Uea_GldxRW1...|
|5r6-G9C4Ylbc7ZLz5...|0|2013-02-09|0|w41Z59shepF03uEyh...|3|Not bad!! Love th...|1|u0LXt3Uea_GldxRW1...|
|fDF_o2JPU8BR1Gya-...|0|2016-04-06|0|WF_QTN3p-thD74hqp...|5|Love this place!...|
|...|3|u0LXt3Uea_GldxRW1...|
|z8oToCT1cX27gZP5G...|0|2013-05-01|0|PIsUSmvaUMB00qv5K...|4|This is currently...|1|u0LXt3Uea_GldxRW1...|
|XWTPNfSkXoUL-LF32...|1|2011-09-28|0|PdZ_uFjbbkjtM3SCY...|3|Server was a litt...|5|u0LXt3Uea_GldxRW1...|
|13nKUHH-uEUxVZylg...|1|2011-02-16|0|x5oV6wm9_Pb1Q06jk...|1|I thought Tidy's ...|9|u0LXt3Uea_GldxRW1...|
|RTUySW0_UZ8V3Wpj0...|1|2012-12-03|1|lsoSqIrrDbQvWpMwS...|3|Wanted to check o...|2|u0LXt3Uea_GldxRW1...|
|Aov96CM4fZAXeZVKt...|1|2010-07-16|0|23eqwLZzCHZkADWfd...|5|This place is awe...|2|u0LXt3Uea_GldxRW1...|
|0W4lkc1zZThpx3V65...|0|2011-09-28|0|FunI9om-aK5oMII3m...|4|a must stop when ...|0|u0LXt3Uea_GldxRW1...|
|fdnNZmk1NP7ZHL-YM...|0|2012-10-23|0|FKu4iU62EnWt6GZXP...|1|I too have been t...|0|u0LXt3Uea_GldxRW1...|
|PFPUF38-lrakZLct...|0|2010-09-15|0|xdu8nXrbNKcaywCX7...|3|Came here with my...|2|u0LXt3Uea_GldxRW1...|
|oWtn2IzrprSRKPFUL...|0|2012-09-23|0|K7o5jDInfmx3cY5oH...|3|Came here for a b...|4|u0LXt3Uea_GldxRW1...|
-----
only showing top 20 rows

import org.apache.spark.sql.SQLContext
sqlContext: org.apache.spark.sql.SQLContext = org.apache.spark.sql.SQLContext@d3b942b
review_data: org.apache.spark.sql.DataFrame = [business_id: string, cool: bigint ... 7 more fields]

scala>
```

```
//removing all review with 3 stars
```

```
var review_filter_3star=review_data.filter(col("stars") != "3")
```

```
review_filter_3star.show()
```

```

scala> :paste
// Entering paste mode (ctrl-D to finish)

//removing all review with 3 stars
var review_filter_3star=review_data.filter(col("stars") != "3")
review_filter_3star.show()

// Exiting paste mode, now interpreting.

warning: there was one deprecation warning; re-run with -deprecation for details
+-----+-----+-----+-----+-----+-----+
| business_id|cool| date|funny| review_id|stars| text|useful| user_id|
+-----+-----+-----+-----+-----+-----+
|0W41kclzZThpx3V65...|0|2016-05-28|0|v0l_UHJMo_hPBq9bx...|5|Love the staff, l...|0|bv2nC15Qv5vroF1qK...|
|AEx2SYEUJmTxVVB18...|0|2016-05-28|0|vkVSCC7x1jJrA14UG...|5|Super simple plac...|0|bv2nC15Qv5vroF1qK...|
|VR6GpWIda3SfvPC-L...|0|2016-05-28|0|n6QzIUObkYshz4dz2...|5|Small unassuming ...|0|bv2nC15Qv5vroF1qK...|
|CKC0-MOWMqoeWf6s-...|0|2016-05-28|0|MV3CcKScW05u5LVff...|5|Lester's is locat...|0|bv2nC15Qv5vroF1qK...|
|ACFTxLv8pGrrxMm6E...|0|2016-05-28|0|IXvOzsEMVtiJI0CAR...|4|Love coming here...|0|bv2nC15Qv5vroF1qK...|
|s2I_Ni76bjJNK9yG6...|0|2016-05-28|0|L_9BTb55X0GdtThi6...|4|Had their chocola...|0|bv2nC15Qv5vroF1qK...|
|8QWPLVQ6D-0ExqXoa...|0|2014-09-24|0|HRPm3vEZ_F-33TYVT...|5|Cycle Pub Las Veg...|1|_4lMDxbXZ1p10NG29...|
|9_CGHMz8698M9-Pk...|2|2012-05-11|0|ymAUG8DZfQcFTBS0l...|4|Who would have gu...|0|u0LXt3Uea_GldxRW1...|
|gkCorLgPyQLsptTHa...|0|2015-10-27|0|8UItshPUD92hXtScS...|4|Always drove past...|1|u0LXt3Uea_GldxRW1...|
|fDF_o2JPU8BR1Gya-...|0|2016-04-06|0|WF_QTN3p-thD74hq...|5|Love this place!|
|...|3|u0LXt3Uea_GldxRW1...|
|z8oIoCT1cXz7gZP5G...|0|2013-05-01|0|PiSUSmvaUMB00qv5K...|4|This is currently...|1|u0LXt3Uea_GldxRW1...|
|13nKUHH-uEUVZy1g...|1|2011-02-16|0|x5oV6wm9_Pb1Q06jk...|1|I thought Tidy's ...|9|u0LXt3Uea_GldxRW1...|
|Aov96CM4FZAXeZvkt...|1|2010-07-16|0|23eqwLzZCmZkADWfd...|5|This place is awe...|2|u0LXt3Uea_GldxRW1...|
|0W41kclzZThpx3V65...|0|2011-09-28|0|FunI9om-ak5oMIIJm...|4|a must stop when ...|0|u0LXt3Uea_GldxRW1...|
|fdnNZMk1NP7ZHL-YM...|0|2012-10-23|0|FKu4iU62EmWtG6ZXP...|1|I too have been t...|0|u0LXt3Uea_GldxRW1...|
|zgQhtqX0gqMwInLBZ...|1|2012-10-30|2|WYDFJ0BOL7cydc7gN...|1|really excited to...|9|u0LXt3Uea_GldxRW1...|
|RWGI8u00x5GghYCEZ...|1|2011-10-14|0|Kki2nwtP8U2qmWrv...|4|This place remind...|0|u0LXt3Uea_GldxRW1...|
|hjk3ox7wiakEu0gT...|0|2012-05-10|2|ypjTMQLKdAwKGRS-K...|1|Food is very blan...|4|u0LXt3Uea_GldxRW1...|
|zxJlg4XCHNofY78WZ...|0|2011-09-28|0|y2iFom8a_SdAyC6I0...|2|a few years ago, ...|0|u0LXt3Uea_GldxRW1...|
|toHRQgUpLyJd8JVQ...|1|2012-09-23|1|ku1sDwmQo2wIgWA...|5|OMG - Definitely ...|1|u0LXt3Uea_GldxRW1...|
+-----+-----+-----+-----+-----+-----+
only showing top 20 rows

review_filter_3star: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [business_id: string, cool: bigint ... 7 more fields]

scala>

```

//combining 1 & 2 stars as BAD & 4 & 5 stars as GOOD(giving 1&2 star rating 0 and 4&5 Good)

```
var review_category = review_filter_3star.withColumn("Rating_labels",when(col("Stars") === "4" || col("stars") === "5", "Good").otherwise(0))
```

review_category.show()

```

scala> :paste
// Entering paste mode (ctrl-D to finish)

//combining 1 & 2 stars as BAD & 4 & 5 stars as GOOD(giving 1&2 star rating 0 and 4&5 Good)
var review_category = review_filter_3star.withColumn("Rating_labels",when(col("Stars") === "4" || col("stars") === "5", "Good").otherwise(0))
review_category.show()

// Exiting paste mode, now interpreting.

+-----+-----+-----+-----+-----+-----+
| business_id|cool| date|funny| review_id|stars| text|useful| user_id|Rating_labels|
+-----+-----+-----+-----+-----+-----+
|0W41kclzZThpx3V65...|0|2016-05-28|0|v0l_UHJMo_hPBq9bx...|5|Love the staff, l...|0|bv2nC15Qv5vroF1qK...|Good|
|AEx2SYEUJmTxVVB18...|0|2016-05-28|0|vkVSCC7x1jJrA14UG...|5|Super simple plac...|0|bv2nC15Qv5vroF1qK...|Good|
|VR6GpWIda3SfvPC-L...|0|2016-05-28|0|n6QzIUObkYshz4dz2...|5|Small unassuming ...|0|bv2nC15Qv5vroF1qK...|Good|
|CKC0-MOWMqoeWf6s-...|0|2016-05-28|0|MV3CcKScW05u5LVff...|5|Lester's is locat...|0|bv2nC15Qv5vroF1qK...|Good|
|ACFTxLv8pGrrxMm6E...|0|2016-05-28|0|IXvOzsEMVtiJI0CAR...|4|Love coming here...|0|bv2nC15Qv5vroF1qK...|Good|
|s2I_Ni76bjJNK9yG6...|0|2016-05-28|0|L_9BTb55X0GdtThi6...|4|Had their chocola...|0|bv2nC15Qv5vroF1qK...|Good|
|8QWPLVQ6D-0ExqXoa...|0|2014-09-24|0|HRPm3vEZ_F-33TYVT...|5|Cycle Pub Las Veg...|1|_4lMDxbXZ1p10NG29...|Good|
|9_CGHMz8698M9-Pk...|2|2012-05-11|0|ymAUG8DZfQcFTBS0l...|4|Who would have gu...|0|u0LXt3Uea_GldxRW1...|Good|
|gkCorLgPyQLsptTHa...|0|2015-10-27|0|8UItshPUD92hXtScS...|4|Always drove past...|1|u0LXt3Uea_GldxRW1...|Good|
|fDF_o2JPU8BR1Gya-...|0|2016-04-06|0|WF_QTN3p-thD74hq...|5|Love this place!|
|...|3|u0LXt3Uea_GldxRW1...|Good|
|z8oIoCT1cXz7gZP5G...|0|2013-05-01|0|PiSUSmvaUMB00qv5K...|4|This is currently...|1|u0LXt3Uea_GldxRW1...|Good|
|13nKUHH-uEUVZy1g...|1|2011-02-16|0|x5oV6wm9_Pb1Q06jk...|1|I thought Tidy's ...|9|u0LXt3Uea_GldxRW1...|0|
|Aov96CM4FZAXeZvkt...|1|2010-07-16|0|23eqwLzZCmZkADWfd...|5|This place is awe...|2|u0LXt3Uea_GldxRW1...|Good|
|0W41kclzZThpx3V65...|0|2011-09-28|0|FunI9om-ak5oMIIJm...|4|a must stop when ...|0|u0LXt3Uea_GldxRW1...|Good|
|fdnNZMk1NP7ZHL-YM...|0|2012-10-23|0|FKu4iU62EmWtG6ZXP...|1|I too have been t...|0|u0LXt3Uea_GldxRW1...|0|
|zgQhtqX0gqMwInLBZ...|1|2012-10-30|2|WYDFJ0BOL7cydc7gN...|1|really excited to...|9|u0LXt3Uea_GldxRW1...|0|
|RWGI8u00x5GghYCEZ...|1|2011-10-14|0|Kki2nwtP8U2qmWrv...|4|This place remind...|0|u0LXt3Uea_GldxRW1...|Good|
|hjk3ox7wiakEu0gT...|0|2012-05-10|2|ypjTMQLKdAwKGRS-K...|1|Food is very blan...|4|u0LXt3Uea_GldxRW1...|0|
|zxJlg4XCHNofY78WZ...|0|2011-09-28|0|y2iFom8a_SdAyC6I0...|2|a few years ago, ...|0|u0LXt3Uea_GldxRW1...|0|
|toHRQgUpLyJd8JVQ...|1|2012-09-23|1|ku1sDwmQo2wIgWA...|5|OMG - Definitely ...|1|u0LXt3Uea_GldxRW1...|Good|
+-----+-----+-----+-----+-----+-----+
only showing top 20 rows

review_category: org.apache.spark.sql.DataFrame = [business_id: string, cool: bigint ... 8 more fields]

scala>

```

//Giving 1 rating to Good

```
val review_rating = review_category.withColumn("number_Rating", when((col("Rating_labels") === "Good"),1).otherwise(0))
```

```
review_rating.show()
```

```
//creating list of words based on text field
```

```
val review_tokenizer = new Tokenizer().setInputCol("text").setOutputCol("tokenized_review")
```

```
val review_tokenized_regexp = new
RegexTokenizer().setInputCol("text").setOutputCol("tokenized_review").setPattern("\\W") //
alternatively .setPattern("\\w+").setGaps(false)
```

```
val review_token_counts = udf { (tokenized_review: Seq[String]) => tokenized_review.length }
```

```
val review_tokenized = review_tokenizer.transform(review_category)
```

```
review_tokenized.select("text", "tokenized_review").withColumn("tokens",
review_token_counts(col("tokenized_review"))).show(false)
```

```
val review_regexp_tokenized = review_tokenized_regexp.transform(review_category)
```

```
review_regexp_tokenized.select("text", "tokenized_review").withColumn("tokens",
review_token_counts(col("tokenized_review"))).show(false)
```

```
//removing all stop words from list
```

```
val stopword_remover = new StopWordsRemover()
```

```
.setInputCol("tokenized_review")
```

```
.setOutputCol("review_without_stopword")
```

```
val stopword_removed_review = stopword_remover.transform(review_regexp_tokenized)
```

```
stopword_removed_review.show()
```

business_id	cool	date	funny	review_id	stars	text	useful	user_id	Rating_labels	tokenized_review	review_without_stopword
104141c12ZThp3V6S...	0	2016-05-28	0	lv0t_UH3Mo_hBqgBx...	5	I love the staff, l...	0	bv2nc15qvsvrofiqk...	Good	[love, the, staff...	[love, staff, lov...
1AE25VEUJmT4VB18...	0	2016-05-28	0	vkV5CC7x1jfrA14UG...	5	Super simple plac...	0	bv2nc15qvsvrofiqk...	Good	[super, simple, p...	[super, simple, p...
1VR6opWIda35fVpC-...	0	2016-05-28	0	neQ2IU0bkYshz4d2...	5	Small unassuming ...	0	bv2nc15qvsvrofiqk...	Good	[small, unassum...	[small, unassum...
1CKC8-MQWQoeWf6S...	0	2016-05-28	0	WV3CcK5Cw05uSLVFF...	5	Lester's is locat...	0	bv2nc15qvsvrofiqk...	Good	[lester, s, is, l...	[lester, located...
1ACFLxLV8GrrXm0e...	0	2016-05-28	0	IXv0ZsEMytl1f0CAR...	4	Love coming here...	0	bv2nc15qvsvrofiqk...	Good	[love, coming, he...	[love, coming, ye...
1s21_NlT6bJmK9Yd6...	0	2016-05-28	0	1_9B7b5SX0Gt1t1e...	4	Had their chocola...	0	bv2nc15qvsvrofiqk...	Good	[had, their, choc...	[chocolate, alon...
18QWPLV6D-0E9xQoa...	0	2014-09-24	0	HRPm3VEZ_F-33YVT...	5	Cycle Pub Las Veg...	1	41MDXbXZ1p10NG29...	Good	[cycle, pub, las...	[cycle, pub, las...
19_CohHMz8698M9-Pk...	2	2012-05-11	0	ymAUG8D2FcFTBS0L...	4	Who would have gu...	0	u0LXT3Uea_GldxRW1...	Good	[who, would, have...	[guess, able, get...
1gKorLgPyQLspTTha...	0	2015-10-27	0	BU1IshPUD92hXT5cS...	4	Always drove past...	1	u0LXT3Uea_GldxRW1...	Good	[always, drove, p...	[always, drove, p...
1f0F_o2JPU8BR1Gya...	0	2016-04-06	0	IMF_QTNsp-thD74hqp...	5	I love this place!	0				
...	3	u0LXT3Uea_GldxRW1...				[love, this, plac...					
1z8o1oC1cx7g2P5G...	0	2013-05-01	0	PisUsmvaUH809q5K...	4	This is currently...	1	u0LXT3Uea_GldxRW1...	Good	[this, is, curren...	[currently, paren...
13nKUHH-uEUXV2yLg...	1	2011-02-16	0	x5oV6wm9_Pb1Q06jK...	1	I thought Tidy's ...	9	u0LXT3Uea_GldxRW1...	0	[i, thought, tidy...	[thought, tidy, f...
1Aov96CH4FZAXeZvk...	1	2010-07-16	0	23eqwLz2CWZKADWfd...	5	This place is awe...	2	u0LXT3Uea_GldxRW1...	Good	[this, place, is...	[place, awesome, ...]
104141c12ZThp3V6S...	0	2011-09-28	0	Fun1Pon-ak5oM1Jm...	4	a must stop when ...	0	u0LXT3Uea_GldxRW1...	Good	[a, must, stop, w...	[must, stop, mont...
1f0nZMk1NP72hL-YH...	0	2012-10-23	0	Fku41Uq2EmWt6GZXP...	1	I too have been t...	0	u0LXT3Uea_GldxRW1...	0	[i, too, have, be...	[trying, book, ap...
1z9QhtqX0g0Wm1nB2...	1	2012-10-30	2	WYDFJ0B0L7cydc7g...	1	Really excited to...	9	u0LXT3Uea_GldxRW1...	0	[really, excited...	[really, excited...
1RWG18u08xSGghVCE...	1	2011-10-14	0	Kk12nwTP8U2gmWwR...	4	This place remind...	0	u0LXT3Uea_GldxRW1...	Good	[this, place, rem...	[place, reminds, ...]
1hJ3ox7Wak0Eu0g1...	0	2012-05-10	2	ypj1MQLKdKwGR5-K...	1	Food is very blan...	4	u0LXT3Uea_GldxRW1...	0	[food, is, very, ...]	[food, bland, aut...
1z219xGCh0rY78M2...	0	2011-09-28	0	2if0n8a_2d4yca1o...	2	a few years ago, ...	0	u0LXT3Uea_GldxRW1...	0	[a, few, years, a...	[years, ago, used...
1toH9qLupLy3dB3VQ...	1	2012-09-23	1	kuisDokmq02w1gMA...	5	OMG - Definitely ...	1	u0LXT3Uea_GldxRW1...	Good	[omg, definitely...	[omg, definitely...

```
only showing top 20 rows

review_tokenizer: org.apache.spark.ml.feature.Tokenizer = tok_fe740aa199d0
review_tokenized_regexp: org.apache.spark.ml.feature.RegexTokenizer = regexpTok_46f5262063ic
review_token_counts: org.apache.spark.sql.expressions.UserDefinedFunction = UserDefinedFunction(<function1>, IntegerType, Some(List(ArrayType(StringType, true))))
review_tokenized: org.apache.spark.sql.DataFrame = [business_id: string, cool: bigint ... 9 more fields]
review_regexp_tokenized: org.apache.spark.sql.DataFrame = [business_id: string, cool: bigint ... 9 more fields]
stopword_remover: org.apache.spark.ml.feature.StopWordsRemover = stopwords_d21535f3823
stopword_removed_review: org.apache.spark.sql.DataFrame = [business_id: string, cool: bigint ... 10 more fields]

scala>
```

```
//making pair of words from the list
```

```

val bigram = new
NGram().setN(2).setInputCol("review_without_stopword").setOutputCol("bigram_review")

val review_bigram = bigram.transform(stopword_removed_review)

review_bigram.select("bigram_review").show()

val review_hashingTF = new
HashingTF().setInputCol("review_without_stopword").setOutputCol("RawReviewFeature").setNumFeatures(25)

val featurized_review = review_hashingTF.transform(review_bigram)

val inverse_review = new
IDF().setInputCol("RawReviewFeature").setOutputCol("supressed_review_feature")

val inv_review_model = inverse_review.fit(featurized_review)

val review_rescale = inv_review_model.transform(featurized_review)

val review_rescaled = review_rescale.select("supressed_review_feature")

review_rescaled.show()

```

```

scala> :paste
// Entering paste mode (ctrl-D to finish)

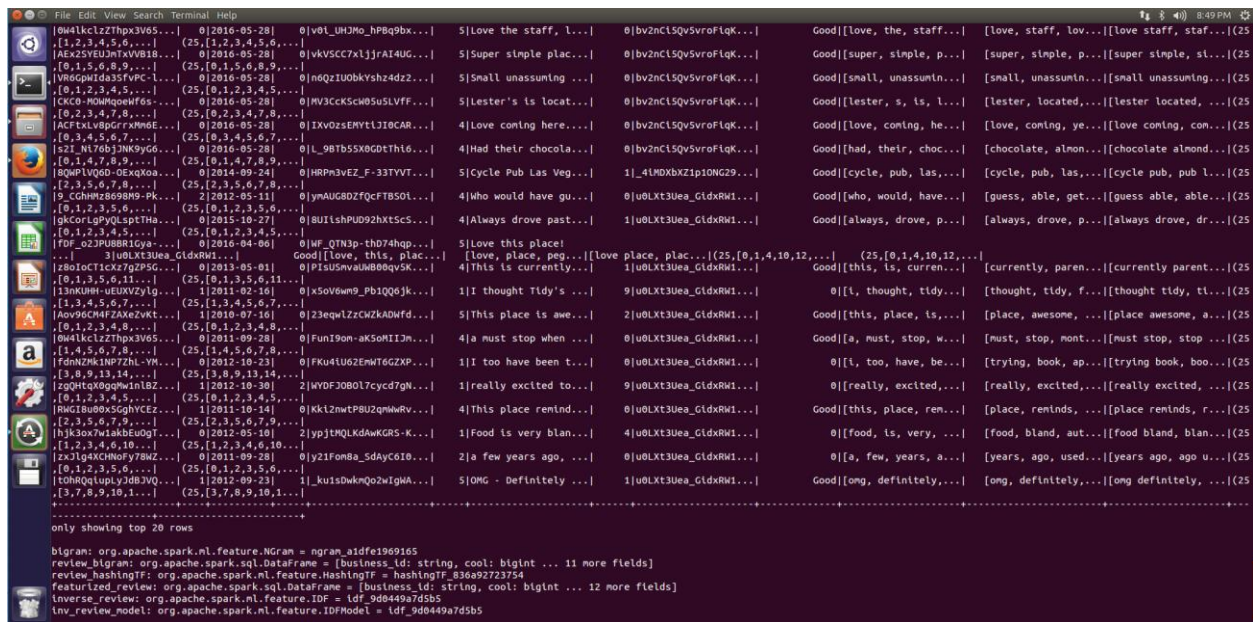
//making pair of words from the list
val bigram = new Ngram().setN(2).setInputCol("review_without_stopword").setOutputCol("bigram_review")
val review_bigram = bigram.transform(stopword_removed_review)
review_bigram.select("bigram_review").show()
val review_hashingTF = new HashingTF().setInputCol("review_without_stopword").setOutputCol("RawReviewFeature").setNumFeatures(25)
val featurized_review = review_hashingTF.transform(review_bigram)
val inverse_review = new IDF().setInputCol("RawReviewFeature").setOutputCol("supressed_review_feature")
val inv_review_model = inverse_review.fit(featurized_review)
val review_rescale = inv_review_model.transform(featurized_review)
val review_rescaled = review_rescale.select("supressed_review_feature")
review_rescaled.show()

// Exiting paste mode, now interpreting.

-----
+-----+
| bigram_review |
+-----+
|[love staff, staf...|
|[super single, st...|
|[small unassuming...|
|[lester located, ...|
|[love coming, con...|
|[chocolate almond...|
|[cycle pub, pub l...|
|[guess able, able...|
|[always drove, dr...|
|[love place, plac...|
|[currently parent...|
|[thought tidy, tl...|
|[place awesome, a...|
|[must stop, stop ...|
|[trying book, boo...|
|[really excited, ...|
|[place reminds, r...|
|[food bland, blan...|
|[years ago, ago u...|
|[ong definitely, ...|
+-----+
only showing top 20 rows

-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| business_id|cool| date|funny| review_id|stars| text|useful| user_id|Rating_labels| tokenized_review|review_without_stopword| bigram_review|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|RawReviewFeature|supressed_review_feature|
+-----+-----+-----+-----+-----+-----+-----+-----+
|[0M4tkclz27Hpx3V65...| 0|2016-05-28| 0|lv0t_UHMo_hP8q9bx...| 5|[Love the staff, l...| 0|bv2nc15QvSvr9figK...| Good|[love, the, staff...| [love, staff, lov...|[love staff, staf...|(25

```

```
val review_lsvc = new
LinearSVC().setFeaturesCol("supressed_review_feature").setLabelCol("number_Rating")

val review_pipeline = new
Pipeline().setStages(Array(review_tokenized_regexp,stopword_remover,bigram,review_hashingTF,inver
se_review,review_lsvc))

val review_evaluator = new MulticlassClassificationEvaluator()

.setLabelCol("number_Rating")

.setPredictionCol("prediction")

.setMetricName("accuracy")

val review_cross_validator = new CrossValidator()

.setEstimator(review_pipeline)

.setEvaluator(review_evaluator)

.setEstimatorParamMaps(new ParamGridBuilder().build)

.setNumFolds(4)

val Array(review_trainingData, review_testData) = review_rating.randomSplit(Array(0.85, 0.15),3296)

val review_train_model = review_cross_validator.fit(review_trainingData)

//Predicting reviews with test data

val review_predictions = review_train_model.transform(review_testData)
```

review_predictions

```
.select(col("number_Rating"), col("review_id"), col("business_id"), col("prediction"))
```

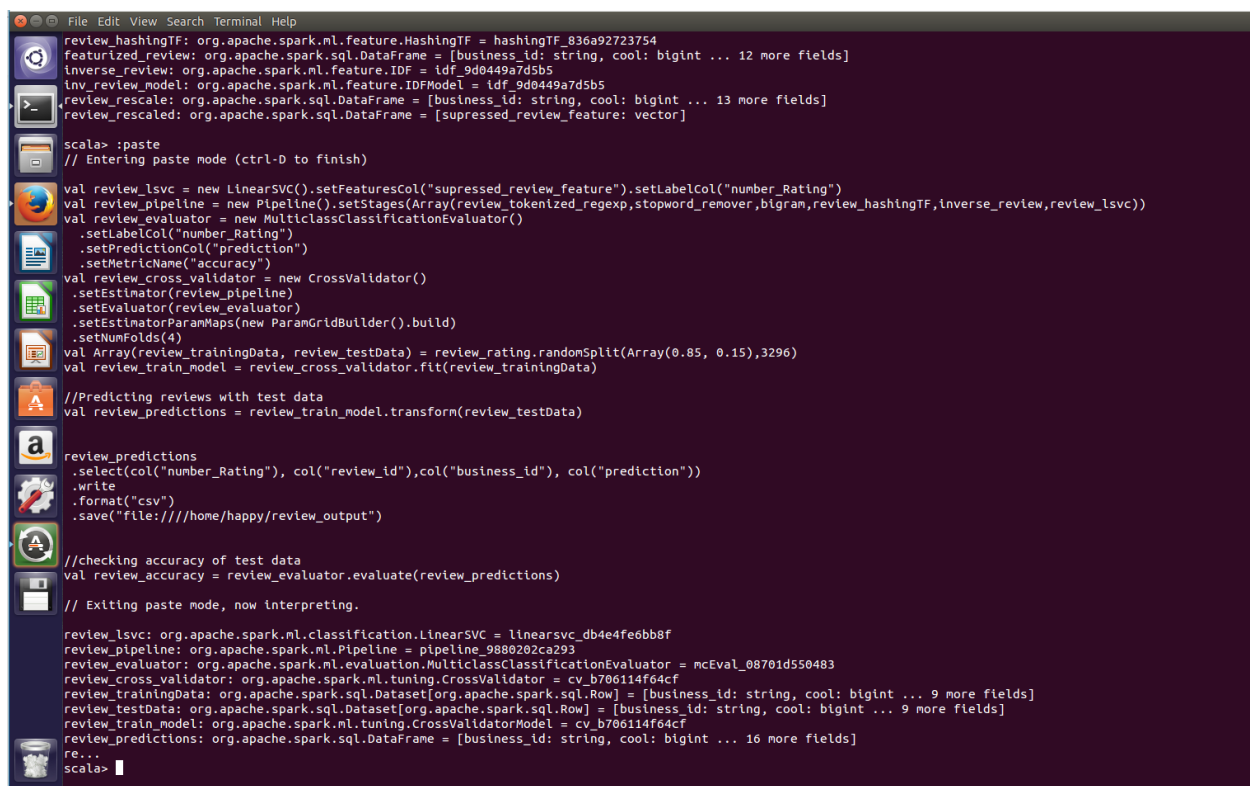
```
.write
```

```
.format("csv")
```

```
.save("file:///home/happy/review_output")
```

//checking accuracy of test data

```
val review_accuracy = review_evaluator.evaluate(review_predictions)
```



```
File Edit View Search Terminal Help
review_hashingTF: org.apache.spark.ml.feature.HashingTF = hashingTF_836a92723754
featurized_review: org.apache.spark.sql.DataFrame = [business_id: string, cool: bigint ... 12 more fields]
inverse_review: org.apache.spark.ml.feature.IDF = idf_9d0449a7d5b5
inv_review_model: org.apache.spark.ml.feature.IDFModel = idf_9d0449a7d5b5
review_rescale: org.apache.spark.sql.DataFrame = [business_id: string, cool: bigint ... 13 more fields]
review_rescaled: org.apache.spark.sql.DataFrame = [supressed_review_feature: vector]

scala> :paste
// Entering paste mode (ctrl-D to finish)

val review_lsvc = new LinearSVC().setFeaturesCol("supressed_review_feature").setLabelCol("number_Rating")
val review_pipeline = new Pipeline().setStages(Array(review_tokenized_regex, stopword_remover, bigram, review_hashingTF, inverse_review, review_lsvc))
val review_evaluator = new MulticlassClassificationEvaluator()
  .setLabelCol("number_Rating")
  .setPredictionCol("prediction")
  .setMetricName("accuracy")
val review_cross_validator = new CrossValidator()
  .setEstimator(review_pipeline)
  .setEvaluator(review_evaluator)
  .setEstimatorParamMaps(new ParamGridBuilder().build)
  .setNumFolds(4)
val Array(review_trainingData, review_testData) = review_rating.randomSplit(Array(0.85, 0.15), 3296)
val review_train_model = review_cross_validator.fit(review_trainingData)

//Predicting reviews with test data
val review_predictions = review_train_model.transform(review_testData)

review_predictions
  .select(col("number_Rating"), col("review_id"), col("business_id"), col("prediction"))
  .write
  .format("csv")
  .save("file:///home/happy/review_output")

//checking accuracy of test data
val review_accuracy = review_evaluator.evaluate(review_predictions)

// Exiting paste mode, now interpreting.

review_lsvc: org.apache.spark.ml.classification.LinearSVC = llinearsvc_db4e4fe0bb8f
review_pipeline: org.apache.spark.ml.Pipeline = pipeline_9880202ca293
review_evaluator: org.apache.spark.ml.evaluation.MulticlassClassificationEvaluator = mcEval_08701d550483
review_cross_validator: org.apache.spark.ml.tuning.CrossValidator = cv_b706114f64cf
review_trainingData: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [business_id: string, cool: bigint ... 9 more fields]
review_testData: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [business_id: string, cool: bigint ... 9 more fields]
review_train_model: org.apache.spark.ml.tuning.CrossValidatorModel = cv_b706114f64cf
review_predictions: org.apache.spark.sql.DataFrame = [business_id: string, cool: bigint ... 16 more fields]
re...
scala>
```