

db-emotion-analysis

Load Data

FINISHED

Took 0 sec. Last updated by yc7093_nyu_edu at December 10 2024, 2:06:48 PM.

```
val basePath = "/user/yc7093_nyu_edu/imdb-reviews-w-emotion/part"
val fileSuffixes = List("-01-all") //, "-02-all", "-03-all", "-04-all")
```

```
val initialPath = s"$basePath${fileSuffixes.head}"
var rawDF = spark.read.parquet(initialPath)
```

```
for (suffix <- fileSuffixes.tail) {
  val fullPath = s"$basePath$suffix"
  val part_df = spark.read.parquet(fullPath)
  rawDF = rawDF.union(part_df)
}
```

```
rawDF.show(5)
```

potter_tag,	reviewer_name,	helpful,	predicted_emotion,	badness,
joy	love	anger	fear	surprise emo
tion				
+-----+	+-----+	+-----+	+-----+	+-----+
-----+	-----+	-----+	-----+	-----+
-----+	-----+	-----+	-----+	-----+
----+-----+				
lrw5552176	FeastModelIt	Chapter Two (2...	2.0	bad and BORING
0 I was enjoying it...	[0, 2]	{score -> 0.1004...	0.10041969269514084	0.67338609695434
57 0.017813226208090782	0.1923339068889618	0.011774818412959576	0.004272174555808306	
joy				
lrw6455111	rapadgettral	Perry Mason: The ...	4.0	Not feasible
0 It's hard when gu...	[0, 0]	{score -> 0.0441...	0.04410260170698166	0.68059456348419
19 0.013725311495363712	0.24673616886138916	0.010360205546021461	0.004481049720197916	
joy				
lrw5178379	MehnaJain2	Fixerr (2019-)	10.0	Fixerr
0 Everyone plays th...	[0, 0]	{score -> 0.0035...	0.003539941739290...	0.98156654834747
31 0.007583207450807095	0.002556430874392...	0.002453002380207181	0.002300753956660...	

Took 1 sec. Last updated by anonymous at December 08 2024, 6:17:03 PM. (outdated)

Rating Distribution

FINISHED

Took 0 sec. Last updated by yc7093_nyu_edu at December 10 2024, 2:06:25 PM.

```
val ratingDistribution = rawDF.groupBy("rating")
  .count()
  .orderBy("rating")

ratingDistribution.show()
```

```
+-----+-----+
|rating|count|
+-----+-----+
|  1.0|35253|
|  2.0|13427|
|  3.0|13669|
|  4.0|14233|
|  5.0|18553|
|  6.0|23809|
|  7.0|31730|
|  8.0|38719|
|  9.0|35389|
| 10.0|79125|
+-----+-----+
```

ratingDistribution: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [rating: double, count: bigint]

Took 3 sec. Last updated by anonymous at December 08 2024, 6:17:09 PM. (outdated)

```
val outputPath = "/user/yc7093_nyu_edu/imdb-emotion-analysis/rating-distribution"

ratingDistribution.coalesce(1)
  .write
  .mode("overwrite")
  .option("header", "true")
  .csv(outputPath)
```

outputPath: String = /user/yc7093_nyu_edu/imdb-emotion-analysis/rating-distribution

Took 2 sec. Last updated by anonymous at December 08 2024, 6:59:05 PM. (outdated)

Emotion Distribution

FINISHED

Took 0 sec. Last updated by yc7093_nyu_edu at December 10 2024, 2:06:14 PM.

```
val emotionDistribution = rawDF.groupBy("emotion")
  .count()
  .orderBy("emotion")

emotionDistribution.show()
```

```

+-----+-----+
| emotion| count|
+-----+-----+
|   anger| 81480|
|   fear| 11217|
|   joy|158506|
|   love| 22475|
| sadness| 17854|
| surprise| 12375|
+-----+-----+

```

emotionDistribution: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [emotion: string, count: bigint]

Took 3 sec. Last updated by anonymous at December 08 2024, 6:17:12 PM. (outdated)

```

val outputPath = "/user/yc7093_nyu_edu/imdb-emotion-analysis/emotion-distribution"
emotionDistribution.coalesce(1)
  .write
  .mode("overwrite")
  .option("header", "true")
  .csv(outputPath)

```

outputPath: String = /user/yc7093_nyu_edu/imdb-emotion-analysis/emotion-distribution

Took 2 sec. Last updated by anonymous at December 08 2024, 6:58:48 PM. (outdated)

Ratings Distribution Within Emotion

FINISHED

Took 0 sec. Last updated by yc7093_nyu_edu at December 10 2024, 2:07:40 PM.

```

import org.apache.spark.sql.expressions.Window
import org.apache.spark.sql.functions._

val emotionWindow = Window.partitionBy("emotion")

val ratingDistributionWithinEmotion = rawDF.groupBy("emotion", "rating")
  .count()
  .withColumn("total_count", sum("count").over(emotionWindow))
  .withColumn("percentage", (col("count") / col("total_count")) * 100)
  .orderBy("emotion", "rating")

val rowCount = ratingDistributionWithinEmotion.count()
ratingDistributionWithinEmotion.show(rowCount.toInt, truncate = false)

```

SPARK JOB FINISHED

```

|sadness| 17.0| 158506| 17854| 17.779709453580152|
|sadness| 18.0| 11521| 17854| 18.519099361487621|

```

```
import org.apache.spark.sql.expressions.Window
import org.apache.spark.sql.functions._

val ratingWindow = Window.partitionBy("rating")

val emotionDistributionWithinRating = rawDF.groupBy("rating", "emotion")
  .count()
  .withColumn("total_count", sum("count").over(ratingWindow))
  .withColumn("percentage", (col("count") / col("total_count")) * 100)
  .orderBy("rating", "emotion")

val rowCount = emotionDistributionWithinRating.count()
```

```
emotionDistributionWithinRating.show(rowCount.toInt, truncate = false)
```

```
17.0  |love      |2420 |31730      |7.626851560037819 |
17.0  |sadness   |1389 |31730      |4.377560668137409 |
17.0  |surprise  |1308 |31730      |4.122281752284904 |
18.0  |anger     |8391 |38719      |21.6715307730055  |
18.0  |fear      |1617 |38719      |4.176244221183398 |
18.0  |joy       |22050|38719      |56.94878483431906 |
18.0  |love      |3274 |38719      |8.455796895580981 |

18.0  |sadness   |1521 |38719      |3.928303933469356 |
18.0  |surprise  |1866 |38719      |4.819339342441696 |
19.0  |anger     |6157 |35389      |17.39806154454774 |
19.0  |fear      |1185 |35389      |3.3484981208850204|
19.0  |joy       |21489|35389      |60.72225832885925 |
19.0  |love      |3362 |35389      |9.50012715815649  |
19.0  |sadness   |1292 |35389      |3.6508519596484783|
19.0  |surprise  |1904 |35389      |5.38020288790302  |
110.0 |anger     |11075|79125      |13.996840442338074|
110.0 |fear      |1848 |79125      |2.3355450236966826|
```

Took 5 sec. Last updated by anonymous at December 08 2024, 6:25:24 PM. (outdated)

```
val outputPath = "/user/yc7093_nyu_edu/imdb-emotion-analysis/emotion-distribution-within-rating"
emotionDistributionWithinRating.coalesce(1)
  .write
  .mode("overwrite")
  .option("header", "true")
  .csv(outputPath)
```

outputPath: String = /user/yc7093_nyu_edu/imdb-emotion-analysis/emotion-distribution-w-rating

Took 2 sec. Last updated by anonymous at December 08 2024, 6:59:47 PM. (outdated)

Keyword Distribution by Emotions

FINISHED

Took 0 sec. Last updated by yc7093_nyu_edu at December 10 2024, 2:08:27 PM.

```
import org.apache.spark.sql.functions._ SPARK JOB FINISHED

val basePath = "/user/yc7093_nyu_edu/imdb_partitioned_by_emotion_rating_word"
val specificEmotion = "sadness"

val emotionDf = spark.read.parquet(s"$basePath/emotion=$specificEmotion")
```

```
val keywordDistribution = emotionDf.groupBy("word")
  .count()
  .orderBy(desc("count"))

val rowCount = keywordDistribution.count()
keywordDistribution.show(rowCount.toInt, truncate = false)
```

```
+-----+-----+
|word      |count|
+-----+-----+
|good      |4971 |
|great     |3333 |
|bad       |3303 |
|love      |3143 |
|boring     |2861 |
|best       |2286 |
|fun        |1753 |
|interesting|1165 |
|worst      |1156 |
|family     |1149 |
|original   |990  |
|amazing    |973  |
|terrible   |965  |
|perfect    |902  |
|beautiful |897  |
```

Took 43 sec. Last updated by anonymous at December 08 2024, 7:03:52 PM. (outdated)

```
val outputPath = s"/user/yc7093_nyu_edu/imdb-emotion-analysis/$specificEmotion-keyword-distribution.csv"
keywordDistribution.coalesce(1)
  .write
  .mode("overwrite")
  .option("header", "true")
  .csv(outputPath)
```

outputPath: String = /user/yc7093_nyu_edu/imdb-emotion-analysis/sadness-keyword-distribution

Took 19 sec. Last updated by anonymous at December 08 2024, 7:04:11 PM. (outdated)

```
import org.apache.spark.sql.functions._
val basePath = "/user/yc7093_nyu_edu/imdb_partitioned_by_emotion_rating_word"
val specificEmotion = "love"

val emotionDf = spark.read.parquet(s"$basePath/emotion=$specificEmotion")

val keywordDistribution = emotionDf.groupBy("word")
```

```
.count()
.orderBy(desc("count"))
```

```
val rowCount = keywordDistribution.count()
keywordDistribution.show(rowCount.toInt, truncate = false)
```

```
+-----+-----+
|word      |count|
+-----+-----+
|love      |12167|
|good      |7470 |
|great     |6234 |
|best      |4045 |
|fun       |3664 |
|bad       |3074 |
|amazing   |2278 |
|beautiful |1974 |
|perfect   |1952 |
|funny     |1886 |
|family    |1723 |
|interesting|1646 |
|original  |1596 |
|comedy    |1495 |
|excellent |1201 |
```

Took 51 sec. Last updated by anonymous at December 08 2024, 7:05:42 PM. (outdated)

```
val outputPath = s"/user/yc7093_nyu_edu/imdb-emotion-analysis/$specificEmotion-keyword-distribution"
keywordDistribution.coalesce(1)
  .write
  .mode("overwrite")
  .option("header", "true")
  .csv(outputPath)
```

outputPath: String = /user/yc7093_nyu_edu/imdb-emotion-analysis/love-keyword-distribution

Took 21 sec. Last updated by anonymous at December 08 2024, 7:06:03 PM. (outdated)

```
import org.apache.spark.sql.functions._
```

SPARK JOB FINISHED

```
val basePath = "/user/yc7093_nyu_edu/imdb_partitioned_by_emotion_rating_word"
val specificEmotion = "joy"
```

```
val emotionDf = spark.read.parquet(s"$basePath/emotion=$specificEmotion")
```

```
val keywordDistribution = emotionDf.groupBy("word")
  .count()
  .orderBy(desc("count"))
```

```
val rowCount = keywordDistribution.count()
```

lvisually_stunning	l392	l
laction_flick	l388	l
ldie_hard	l382	l
labsolutely_amazing	l358	l
lgood_fun	l315	l
lcult_classic	l303	l
lhorror_flick	l295	l
lrom_com	l260	l
lbest_action	l250	l
lhidden_gem	l226	l
lbest_horror	l198	l
lrip_off	l177	l
lgreat_family	l158	l
lwow_wow	l45	l
lwaste_money	l21	l
lhear_warming	l1	l

+-----+-----+

Took 1 min 5 sec. Last updated by anonymous at December 08 2024, 7:07:08 PM. (outdated)

```
val outputPath = s"/user/yc7093_nyu_edu/imdb-emotion-analysis/$specificEmotion-keyword-distribution"
keywordDistribution.coalesce(1)
  .write
  .mode("overwrite")
  .option("header", "true")
  .csv(outputPath)
```

outputPath: String = /user/yc7093_nyu_edu/imdb-emotion-analysis/joy-keyword-distribution

Took 29 sec. Last updated by anonymous at December 08 2024, 7:07:37 PM. (outdated)

```
import org.apache.spark.sql.functions._
val basePath = "/user/yc7093_nyu_edu/imdb_partitioned_by_emotion_rating_word"
val specificEmotion = "surprise"

val emotionDf = spark.read.parquet(s"$basePath/emotion=$specificEmotion")

val keywordDistribution = emotionDf.groupBy("word")
  .count()
  .orderBy(desc("count"))

val rowCount = keywordDistribution.count()
keywordDistribution.show(rowCount.toInt, truncate = false)
```

lword	lcountl
lfun	l4469

+-----+-----+

lfunny	3913
lgood	3836
lgreat	3282
llove	2655
lbest	2139
linteresting	1911
lamazing	1751
lbad	1480
lcomedy	1000
lperfect	923
loriginal	763
lbeautiful	732
lbrilliant	683

Took 39 sec. Last updated by anonymous at December 08 2024, 7:08:16 PM. (outdated)

```
val outputPath = s"/user/yc7093_nyu_edu/imdb-emotion-analysis/$specificEmotion-keyword-distribution"
keywordDistribution.coalesce(1)
  .write
  .mode("overwrite")
  .option("header", "true")
  .csv(outputPath)
```

outputPath: String = /user/yc7093_nyu_edu/imdb-emotion-analysis/surprise-keyword-distribution

Took 17 sec. Last updated by anonymous at December 08 2024, 7:08:33 PM. (outdated)

```
import org.apache.spark.sql.functions._
```

SPARK JOB FINISHED

```
val basePath = "/user/yc7093_nyu_edu/imdb_partitioned_by_emotion_rating_word"
val specificEmotion = "anger"
```

```
val emotionDf = spark.read.parquet(s"$basePath/emotion=$specificEmotion")
```

```
val keywordDistribution = emotionDf.groupBy("word")
  .count()
  .orderBy(desc("count"))
```

```
val rowCount = keywordDistribution.count()
keywordDistribution.show(rowCount.toInt, truncate = false)
```

word	count
lgood	23585
lbad	17137
lgreat	14953
llove	12776

lbest	11210
lfun	10915
lworst	6221
linteresting	5674
lfamily	5450
lfunny	5449
loriginal	5293
lboring	4993
lhorror	4899
lcomedy	4709

Took 35 sec. Last updated by yc7093_nyu_edu at December 10 2024, 2:28:13 PM.

```
val outputPath = s"/user/yc7093_nyu_edu/imdb-emotion-analysis/$specificEmotion-keyword-distribution"
keywordDistribution.coalesce(1)
  .write
  .mode("overwrite")
  .option("header", "true")
  .csv(outputPath)
```

outputPath: String = /user/yc7093_nyu_edu/imdb-emotion-analysis/anger-keyword-distribution

Took 15 sec. Last updated by yc7093_nyu_edu at December 10 2024, 2:28:28 PM.

```
import org.apache.spark.sql.functions._
```

SPARK JOB FINISHED

```
val basePath = "/user/yc7093_nyu_edu/imdb_partitioned_by_emotion_rating_word"
val specificEmotion = "fear"
```

```
val emotionDf = spark.read.parquet(s"$basePath/emotion=$specificEmotion")
```

```
val keywordDistribution = emotionDf.groupBy("word")
  .count() // Count occurrences of each word within the emotion
  .orderBy(desc("count"))
```

```
val rowCount = keywordDistribution.count()
keywordDistribution.show(rowCount.toInt, truncate = false)
```

word	count
good	3629
great	2370
love	1947
bad	1868
fun	1831
best	1657
horror	1618
interesting	1191
funny	1096
family	846

original	1813	1
perfect	1766	1
comedy	1677	1
amazing	1668	1
incredible	1500	1

Took 36 sec. Last updated by yc7093_nyu_edu at December 10 2024, 2:29:04 PM.

SPARK JOB FINISHED 99%

```
keywordDistribution.coalesce(1)
  .write
  .mode("overwrite")
  .option("header", "true")
  .csv(outputPath)
```

outputPath: String = /user/yc7093_nyu_edu/imdb-emotion-analysis/fear-keyword-distribution

Took 9 sec. Last updated by yc7093_nyu_edu at December 10 2024, 2:29:13 PM.

Started a few seconds ago.

Rating Distribution Within Keywords

FINISHED

Took 0 sec. Last updated by yc7093_nyu_edu at December 10 2024, 2:10:04 PM.

```
import org.apache.spark.sql.functions._
```

SPARK JOB FINISHED

```
val basePath = "/user/yc7093_nyu_edu/imdb-partitioned_by_emotion_rating_word"
```

```
val partitionedDf = spark.read.parquet(basePath)
```

```
val ratingDistributionByKeyword = partitionedDf.groupBy("word", "rating")
  .count()
  .orderBy("word", "rating")
```

```
val pivotedDistribution = ratingDistributionByKeyword.groupBy("word")
  .pivot("rating")
  .sum("count")
```

```
pivotedDistribution.show(truncate = false)
```

awesome	1238	1137	1151	1163	1277	1408	1752	11374	11864	15086
good	16946	13921	14578	15372	17583	110765	115179	115780	111516	117553
wonderful	1205	1131	1190	1205	1336	1573	11256	12010	12155	14388
funny	12148	11078	11141	11228	11469	12038	12840	13546	12786	14629
disappointing	1740	1507	1659	1764	1824	1659	1466	1318	1192	1172
absolutely_amazing	18	14	12	12	19	110	125	151	176	1396
comedy	11387	1735	1842	1922	11281	11959	12568	12832	12178	13734
romantic_comedy	133	127	140	155	155	1111	1149	1167	191	1126
well_worth	17	13	13	110	132	1148	1380	1610	1422	1371
pleasantly_surprised	13	17	14	113	127	1106	1210	1223	1206	1215

pleasantly_surprised	17	14	115	127	1100	1219	1323	1200	1219	1
horror_flick	133	131	151	139	162	189	1109	1114	153	148

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+

only showing top 20 rows

```
import org.apache.spark.sql.functions._
basePath: String = /user/yc7093_nyu_edu/imdb_partitioned_by_emotion_rating_word
partitionedDf: org.apache.spark.sql.DataFrame = [review_id: string, movie: string ... 6 more columns]
```

Took 4 min 34 sec. Last updated by anonymous at December 08 2024, 6:48:21 PM. (outdated)

```
val outputPath = "/user/yc7093_nyu_edu/imdb-emotion-analysis/rating-distribution-w-word"
pivotedDistribution.coalesce(1)
  .write
  .mode("overwrite")
  .option("header", "true")
  .csv(outputPath)
```

outputPath: String = /user/yc7093_nyu_edu/imdb-emotion-analysis/rating-distribution-w-word

Took 2 min 8 sec. Last updated by anonymous at December 08 2024, 7:03:09 PM. (outdated)