CS167 Project D-1 Twitter Analysis

Group D1 Minsoo Kim - mkim410 Rajeev Thundiyil - rthun001 Christian Boroff - cboro002

Task 1 - Data Preparation

- Prepare the data for analysis later in the project.
- To give a brief summary, I first took the given data set and kept only the some of the attributes and created a new output file with it. Then I extracted the top 20 keywords from the newly created dataset.

Loading the file

Relevant Attributes

```
// Keep only the following attributes {id, text, entities.hashtags.txt, user.description,
retweet_count, reply_count, and quoted_status_id}
    val clean_tweets_df: DataFrame = tweetsDF.selectExpr("id", "text",
"transform(entities.hashtags, x -> x.text) AS hashtags", "user.description AS user_description",
"retweet_count", "reply_count", "quoted_status_id")
```

The relevant attributes mentioned above included the following,

- {id, text, entities.hashtags.txt, user.description, retweet_count, reply_count, and quoted_status_id}

Storing output to a new JSON file to tweets_clean

Results

Tweets_1K

oted_status_id			user_description retwe	hashtags	text	id
null			I'm not here anymore	111	saya tahu dia ter	921633443934433280
188	86574364	0 92133	やってます 🩋 甲子園 日	[] 野球ガ	ふみとおおおお 😭 🈭 💗 💗	921633444219596800
33444131680256	9		ا]] اا			
			• TxSU • 🤘 HTX	[1]	find a boo &	921633443976568838
			constantly grumpy	111	He referred to hi	921633443708096512
			transduce		@BarRefaeli my fi	921633444400115713
			Sierra High Schoo	[1]	Idk bout y'all bu	921633444437770241
53313819081523			David killed Goli		<i>ඎඎඎ</i> gwym mane∣	921633444223954944
			畑68→東京某学校/元I 0	D1 :	I'm at 北京餃子 in 仙台	921633444773400576
			Zpace for share w	[]]	I'm at G Tower in	921633444953772032
37211901820928	8 8		they hurt u and t	[1]	beach https://t.c	921633445326942211
			10 Années de dans	aquet, TeamSilv	@christophelicat [921633445360623617
			#HealthTips News	aturalCure, Hea	#NaturalCure #Hea [921633444693663745
			Cheif operating o	[]]	@ShirazHassan my	921633445788258304
91726585487366	8 9		Area Manager for	[]]		921633445499097088
				111	@Auisirilak1 maul	921633445729538848
			รักสันโดษ เจ็บแต่ 0		เป็นคนเชื่อว่าตัว	921633444915904512
0119140162355	8 9		∰:michellerosello	[1]	5mentarios https:	921633445301964800
			hopeless romantic	[1]	Mas okay talaga k	921633446518079488
null			tremeran le fogli	[1]	Paesaggiando http	921633444496576512

only showing top 20 rows

Frequent Hashtags

Then, I ran a top-k SQL query to select the top 20 most frequent hashtags. The below shows the SQL and how I collected the result in an array of keywords.

```
hashtag|count|
      ALDUBXEBLoveis
        FurkanPalal1
                no3091
                Laloni
              sbhawksl
|DoktorlarDenkliki...|
        Benimisteğim|
     احتاج بالوقت مذا
 CNIextravaganza2017
                 lovel
                happyl
      السعودية
           nowplaying
           beautiful|
             türkiyel
             vegaltal
           KittyLivel
                   鳂
           tossademar
```

Tweets_10K

	hashtag c	ount
ALDUBXE	BLoveis	84
Furka	nPalalı	51
	no309	51
	Lal0n	51
	chien	30
	job	28
	Hiring	22
	sbhawks	16
1	op3Apps	16
	perdu	15
	trouvé	15
[Ca	reerArc	14
	Job	12
trum	prussia	12
	trndnl	12
	Jobs	11
ShowtimeLets	Celebr8	91
	hiring	9
impeachtru	mppence	91
	music	8
+		+

```
user description|retweet count|reply count| quoted status id|
[921633443934433280]
                           saya tahu dia ter...|
|921633444219596800|ふみとおおおお 😭 😭 🐭 🐭 🕶 ...|
                                                               []|野球ガールの頭やってます|| 甲子園...|
                                                                                                                      0|921338657436459008|
                                                                                                                                   [921633444131680256]
                                                                                    • TxSU • ⊌ HTX...|
19216334439765688381
                           find a boo &amp: ...|
19216334437080965121
                           He referred to hi...!
                                                                                       constantly grumpy|
                                                                                                transducel
|921633444400115713|
                           @BarRefaeli mv fi...|
19216334444377702411
                           Idk bout y'all bu...|
                                                                                    Sierra High Schoo...|
[921633444223954944]

⊗⊗⊗⊗gwym mane...|

                                                                                     David killed Goli ... |
                                                                                                                                   0 | 921633138190815232 |
|921633444773400576| I'm at 北京餃子 in 仙台...|
                                                                     五橋→元茶畑68→東京某学校/元I...|
19216334449537720321
                           I'm at G Tower in...
                                                                                    Zpace for share w...|
19216334453269422111
                           beach https://t.c...|
                                                                                    they hurt u and t...!
19216334453606236171
                           @christophelicat ...|[Taquet, TeamSilv...|
                                                                                    10 Années de dans...|
19216334446936637451
                           #NaturalCure #Hea...|[NaturalCure, Hea...|
                                                                                    #HealthTips News ...|
[921633445788258304]
                           @ShirazHassan my ...|
                                                                                    Cheif operating o...|
[921633445499097088]
                           So true! https://...|
                                                                                    Area Manager for ...|
                                                                                                                                  0|921491726585487360|
|921633445729538048|
                           @Auisirilak1 maul...|
|921633444915904512|
                           เป็นคนเชื่อว่าตัว...|
                                                                                รักสันโดษ เจ็บแต่...|
                           5mentarios https:...|
                                                                                       #:michellerosello|
                                                                                                                                  0|921501191401623552|
|921633445301964800|
19216334465180794881
                           Mas okay talaga k...|
                                                                                       hopeless romantic|
19216334444965765121
                           Paesaggiando http...|
                                                                                     tremeran le fogli...|
only showing top 20 rows
```

Task 2

My task was to show the topic of the tweet by doing an array intersection of the most frequent hashtags from task 1 with the hashtags within the tweet itself. If it's a match, a new column is created with the topic shown, found through the intersection.

```
clean_tweets_df.createOrReplaceTempView( viewName = "tweets_clean")
  //convert keywords to an array separated with , so it can be used for array intersect in a query
  val topics: String = "'"+ keywords.mkString("','") + "'"

  //dataframe
```

Task 2 (cont)

```
val topics_df: DataFrame = sparkSession.sql(
    sqlText = s"""
    SELECT id, text,element_at(t1.tweet_topic,1), user_description, retweet_count, reply_count, quoted_status_id
    FROM ( SELECT *, array_intersect(hashtags, array($topics)) AS tweet_topic FROM tweets_clean) AS t1 WHERE size(tweet_topic) > 0;
    """)

//write to json
topics_df.write.json( path = "tweets_topic.json")
topics_df.show()
val t4 = System.nanoTime()
```

The query picks up all of the data types in order, as given in the project specifications.

```
root
|-- id: long (nullable = true)
|-- text: string (nullable = true)
|-- topic: string (nullable = true)
|-- user_description: string (nullable = true)
|-- retweet_count: long (nullable = true)
|-- reply_count: long (nullable = true)
|-- quoted_status_id: long (nullable = true)
```

Task 2 Results

Operations on file 'Tweets_1k.json' took 7.5779234 seconds						
id	text eleme	ent_at(tweet_topic, 1)	user_description retweet	t_count repl	y_count quote	ed_status_id
921633446644080641	#negramaroofficia	love	Negramanteinside			null
921633445045866497	#CNIextravaganza2	CNIextravaganza2017	Hebat Produknya H			null
921633449882128384	#DoktorlarDenklik	DoktorlarDenkliki	emin ben			null
921633451773648896	Na miss ko mag tw	ALDUBxEBLoveis	Resilient. Object			null
921633452289642497	#FurkanPalalı Değ	FurkanPalalı	null			null
921633451714920448	#KittyLive penuh	KittyLive	Semangat			null
null 0 0			#احتاج_بالوقت_هذا احمت		9216334537	15738624
921633455766728704	Don`t run from yo	ALDUBxEBLoveis	BEHIND THE THICK			null
921633464276996097	Künefe Ocağı Rezi	türkiye	Bym Isı Rezistans			null
921633465921028096 ブルーマウンテンさ		sbhawks 野球垢です! ホー			null	
921633468634632194	Pastinya dong, ka	CNIextravaganza2017	BLOGGER BUZZER		0 921627	555261702144
921633470413262848	#FurkanPalalı Değ	FurkanPalalı	null			null
921633476964646912	Start by doing wh	ALDUBxEBLoveis	BEHIND THE THICK			null
921633490482946048	Stay by Rihanna F	nowplaying	A live stream of			null
921633493569998848	#FurkanPalalı Değ	FurkanPalalı	null			null
921633494530392064	Saturday let's go	nowplaying	Student of Life!			null
null 0 0					9216334955118	96065
921633502453518336	There are many th	ALDUBxEBLoveis	BEHIND THE THICK			null
921633507558002688	#happy #Saturday	happy	Pastry chef, pers			null
921633511752257536	#FurkanPalalı Değ	FurkanPalalı	null			null

Task 3 - Machine Learning Model

• For this task we built a machine learning model pipeline using the code below

```
val tokenizer = new Tokenizer().setInputCol("text").setOutputCol("words")

val hashingTF = new HashingTF().setInputCol("words").setOutputCol("features")

val stringIndexer = new StringIndexer().setInputCol("element_at(tweet_topic, 1)").setOutputCol("label").setHandleInvalid("skip")

val logisticRegression = new LogisticRegression()

val pipeline = new Pipeline().setStages(Array(tokenizer, hashingTF, stringIndexer, logisticRegression))

val Array(trainingData, testData) = topics_df.randomSplit(Array(0.7, 0.3))

val logisticModel = pipeline.fit(trainingData)

val predictions = logisticModel.transform(testData)

predictions.select([col="id", [col="text", "element_at(tweet_topic, 1)", "user_description", "label", "prediction").show(numRows= 10)
```

• This produced the following output

		野球垢です!ホークスファンです! 6.8		
⊕⊕\nla ako masa				

Task 3 - Continued

 We then computed the precision and recall of our machine learning model on the 10k tweets dataset using the following code.

```
// Compute the number of true positives, false positives, and false negatives for each class

val tp = (0 ≤ to ≤ 10).map(c => predictions.filter(col( colName = "label") === c && col( colName = "prediction") === c).count()).sum

val fp = (0 ≤ to ≤ 10).map(c => predictions.filter(col( colName = "label") =!= c && col( colName = "prediction") === c).count()).sum

val fn = (0 ≤ to ≤ 10).map(c => predictions.filter(col( colName = "label") === c && col( colName = "prediction") =!= c).count()).sum

// Compute overall precision and recall

val overallPrecision = tp.toDouble / (tp + fp)

val overallRecall = tp.toDouble / (tp + fn)
```

Our machine learning model using logistic regression to predict the topic of a tweet had a
precision and recall of 0.94382. This means 94.38% of the time, our model correctly
identified true positive cases.

Question

Which classifier did we use for our machine learning model?

- a. logistic regression classifier
- b. multi-layer perceptron classifier
- c. decision tree classifier