

Mongo_requeteSparkSQL

September 10, 2020

1 5. Résultats: Requêtes Spark SQL

```
In [1]: # import des librairies
        from pyspark import SparkContext
        from pyspark.sql import SparkSession, Row
        from pyspark.sql.functions import explode

In [2]: # Initialisation de spark
        sc = SparkContext()
        spark = SparkSession.builder.config("spark.sql.broadcastTimeout", "36000").getOrCreate()
```

1.0.1 5.1. Observations de batchView:

```
In [3]: # Chargement des données dans une dataframe spark
        df_batchView = spark.read.format("mongo").option("uri", "mongodb://127.0.0.1/twitter.batchView").load()
        df_batchView.persist()
        df_batchView.printSchema()
        df_batchView.createOrReplaceTempView("dfbatch")

root
|-- _id: struct (nullable = true)
|   |-- oid: string (nullable = true)
|-- count: long (nullable = true)
|-- date: string (nullable = true)
|-- hashtags: string (nullable = true)
|-- heureDebut: string (nullable = true)
|-- heureFin: string (nullable = true)
```

```
In [4]: # C'est plus pratique d'utiliser une fonction qui affichera les résultats une par une
        def afficherTopTenHtagFromBatch(batchDataframe, date, heure_debut):
            """
            Fonction affichant le classement des hashtags les plus populaire
            à partir du dataframeBatch
            La fonction prend en paramètre:
            - date: String comme par exemple : 2020-09-03
```

```

- heure_debut: int entre 0 et 23
"""
#date = "2020-09-03"
#heure_debut = 16
df2 = spark.sql("""SELECT hashtags, count FROM dfbatch df1 WHERE df1.date == '{0}' AND
                  .format(date, heure_debut )).orderBy("count", ascending=False)

print("*****")
print("Le {0} entre {1}:00:00 et {2}:00:00, ci-dessous les hashtags les plus utilisés")
df2.show()

```

```

In [5]: date = "2020-09-05"
        # Afficher les résultat par heure du 5 septembre
        for x in range(25):
            afficherTopTenHtagFromBatch(df_batchView, date, x)

```

```

*****
Le 2020-09-05 entre 0:00:00 et 1:00:00, ci-dessous les hashtags les plus utilisés:

```

```

+-----+-----+
|          hashtags|count|
+-----+-----+
| ArtistoftheSummer| 764|
|      WeRespectVets| 125|
|           HiKai|   81|
|          BBNaija|   72|
|           BTS|    65|
|__...|    56|
|TrumpHatesOurMili...|   52|
|          IceCream|   36|
|          EP4|    34|
| SayangAdminShopee|   32|
+-----+-----+

```

```

*****
Le 2020-09-05 entre 1:00:00 et 2:00:00, ci-dessous les hashtags les plus utilisés:

```

```

+-----+-----+
|          hashtags|count|
+-----+-----+
| ArtistoftheSummer| 537|
| DraftDodgerDon| 132|
|      WeRespectVets| 108|
|          GoStars|   86|
|           BTS|   58|
|          BBNaija|   52|
|          losers|   49|
|          IceCream|   45|
|           HiKai|   42|
|          Mulan|   38|

```

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

Le 2020-09-05 entre 2:00:00 et 3:00:00, ci-dessous les hashtags les plus utilisés:

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

hashtags	count
----------	-------

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

ArtistoftheSummer	317
-------------------	-----

ROSE	122
------	-----

BLACKPINK	113
-----------	-----

SmackDown	109
-----------	-----

EXO	106
-----	-----

LISA	105
------	-----

ExolsGotPower	104
---------------	-----

ExoPower3rdAnnive...	103
----------------------	-----

PowerExo	101
----------	-----

DraftDodgerDon	96
----------------	----

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

Le 2020-09-05 entre 3:00:00 et 4:00:00, ci-dessous les hashtags les plus utilisés:

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

hashtags	count
----------	-------

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

BTS	757
-----	-----

500	
-----	--

BREAKTHESILENCE_T...	469
----------------------	-----

ArtistoftheSummer	325
-------------------	-----

SmackDown	167
-----------	-----

5Baje5Minutes	165
---------------	-----

WeRespectVets	128
---------------	-----

DraftDodgerDon	108
----------------	-----

EXO	96
-----	----

ExoPower3rdAnnive...	88
----------------------	----

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

Le 2020-09-05 entre 4:00:00 et 5:00:00, ci-dessous les hashtags les plus utilisés:

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

hashtags	count
----------	-------

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

5Baje5Minutes	350
---------------	-----

ArtistoftheSummer	343
-------------------	-----

BTS	289
-----	-----

3_...	198
-------	-----

speakup	158
---------	-----

RRBExamDates	156
--------------	-----

152	
-----	--

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

BREAKTHESILENCE_T...	130
DraftDodgerDon	103
rrbexamdates	76

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

Le 2020-09-05 entre 5:00:00 et 6:00:00, ci-dessous les hashtags les plus utilisés:

	hashtags count
--	----------------

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

5Baje5Minutes	604
RRBExamDates	332
speakup	319
ArtistoftheSummer	247
BTS	246
rrbexamdates	168
3_...	146
	119
BREAKTHESILENCE_T...	95
SpeakUpForSSCRail...	59

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

Le 2020-09-05 entre 6:00:00 et 7:00:00, ci-dessous les hashtags les plus utilisés:

	hashtags count
--	----------------

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

5Baje5Minutes	712
speakup	581
RRBExamDates	506
BTS	239
rrbexamdates	219
ArtistoftheSummer	194
SpeakUpForSSCRail...	124
3_...	123
	92
SpeakUp	87

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

Le 2020-09-05 entre 7:00:00 et 8:00:00, ci-dessous les hashtags les plus utilisés:

	hashtags count
--	----------------

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

RRBExamDates	4117
speakup	3611
5Baje5Minutes	1875
rrbexamdates	1390

SpeakUpForSSCRail...	806
SpeckUpForSSCRail...	525
SpeakUp	292
BiharBole_RozgarDo	264
SpeakUpForSSCRail...	193
rrbexamedates	192
+-----+-----+	

 Le 2020-09-05 entre 8:00:00 et 9:00:00, ci-dessous les hashtags les plus utilisés:

+-----+-----+	
hashtags count	
+-----+-----+	
RRBExamDates	6641
speakup	3207
rrbexamdates	1119
SpeakUpForSSCRail...	638
SpeckUpForSSCRail...	517
5Baje5Minutes	414
SpeakUp	261
rrbexamedates	234
BiharBole_RozgarDo	208
SpeakUpForSSCRail...	205
+-----+-----+	

 Le 2020-09-05 entre 9:00:00 et 10:00:00, ci-dessous les hashtags les plus utilisés:

+-----+-----+	
hashtags count	
+-----+-----+	
RRBExamDates	2319
speakup	747
SpeakUpForSSCRail...	236
rrbexamdates	149
SpeakUpFor69000Te...	134
55 117	
SpeakUp	104
rrbexamdate	93
SpeakUpForSSCRail...	91
minecraftmanhunt	79
+-----+-----+	

 Le 2020-09-05 entre 10:00:00 et 11:00:00, ci-dessous les hashtags les plus utilisés:

+-----+-----+	
hashtags count	
+-----+-----+	
RRBExamDates	5197

	speakup	2020
	SpeakUpFor69000Te...	489
	Dynamite300M	300
	SpeakUpForSSCRail...	289
	5Baje5Minutes	278
	rrbexamdates	259
	BTS	225
	BiharBole_RozgarDo	211
	rrbexamdate	187
+-----+-----+		

Le 2020-09-05 entre 11:00:00 et 12:00:00, ci-dessous les hashtags les plus utilisés:

+-----+-----+		
	hashtags	count
+-----+-----+		
	RRBExamDates	5201
	speakup	1916
	SpeakUpFor69000Te...	642
	rrbexamdates	297
	ExolsAreAlwaysWit...	244
	ExolsAreAlwaysWit...	240
		236
	CHANYEOL	225
		212
	SpeakUpForSSCRail...	209
+-----+-----+		

Le 2020-09-05 entre 12:00:00 et 13:00:00, ci-dessous les hashtags les plus utilisés:

+-----+-----+		
	hashtags	count
+-----+-----+		
	RRBExamDates	5190
	speakup	1751
	rrbexamdates	363
	SpeakUpFor69000Te...	334
	5Baje5Minutes	213
	DestinyClinicxMG	201
	SpeakUpForSSCRail...	185
	ArtistoftheSummer	183
	BTS	165
	SpeakUp	131
+-----+-----+		

Le 2020-09-05 entre 13:00:00 et 14:00:00, ci-dessous les hashtags les plus utilisés:

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

hashtags	count
RRBExamDates	4221
speakup	1943
55	1328
5Baje5Minutes	504
JENO	290
NCT	283
CHENLE	272
JUNGWOO	272
Awsaz	271
ItsAwkward_ButItsOk	271

 Le 2020-09-05 entre 14:00:00 et 15:00:00, ci-dessous les hashtags les plus utilisés:

hashtags	count
55	1676
RRBExamDates	1363
speakup	853
23YrsofSENSATiONA...	850
SooraraiPottru	832
5Baje5Minute	459
5Baje5Minutes	346
rrbexamdates	168
SarkaruVaariPaata	167
ArtistoftheSummer	158

 Le 2020-09-05 entre 15:00:00 et 16:00:00, ci-dessous les hashtags les plus utilisés:

hashtags	count
23YrsofSENSATiONA...	987
SooraraiPottru	906
55	568
RRBExamDates	422
speakup	260
5Baje5Minutes	188
SarkaruVaariPaata	188
5Baje5Minute	125
ArtistoftheSummer	111
RamCharan13YrsTre...	97

Le 2020-09-05 entre 16:00:00 et 17:00:00, ci-dessous les hashtags les plus utilisés:

```
+-----+-----+
|          hashtags|count|
+-----+-----+
|23YrsOfSENSATiONA...| 455|
|      SooraraiPottru| 372|
|          55| 334|
|      RRBExamDates| 255|
|  ArtistoftheSummer| 203|
|          speakup| 160|
|  PrabhasAdvBdayCDP| 120|
|      __| 117|
|  SarkaruVaariPaata| 111|
|ThreeDaysTo_Incar...| 111|
+-----+-----+
```

Le 2020-09-05 entre 17:00:00 et 18:00:00, ci-dessous les hashtags les plus utilisés:

```
+-----+-----+
|          hashtags|count|
+-----+-----+
|  ArtistoftheSummer| 1807|
|23YrsOfSENSATiONA...| 517|
|      SooraraiPottru| 465|
|          55| 336|
|          BTS| 323|
|      StrayKids| 291|
|      | 268|
|      RRBExamDates| 247|
|      MyDaySelcaDay| 171|
|          MDSD| 156|
+-----+-----+
```

Le 2020-09-05 entre 18:00:00 et 19:00:00, ci-dessous les hashtags les plus utilisés:

```
+-----+-----+
|          hashtags|count|
+-----+-----+
|  ArtistoftheSummer| 1568|
|23YrsOfSENSATiONA...| 651|
|      SooraraiPottru| 587|
|          55| 298|
|      RRBExamDates| 217|
|          BTS| 190|
|  SarkaruVaariPaata| 142|
|          5Baje5Minutes| 113|
|          speakup| 98|
```


|__| 96|

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

Le 2020-09-05 entre 19:00:00 et 20:00:00, ci-dessous les hashtags les plus utilisés:

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

| hashtags|count|

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

| ArtistoftheSummer| 1145|

|23YrsOfSENSATiONA...| 492|

| SooraraiPottru| 468|

| 55| 144|

| BTS| 134|

| RRBExamDates| 114|

| Shame_On_MahaGovt| 90|

| BBNaija| 85|

| SarkaruVaariPaata| 82|

| ShamelessTV9| 72|

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

Le 2020-09-05 entre 20:00:00 et 21:00:00, ci-dessous les hashtags les plus utilisés:

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

| hashtags|count|

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

| ArtistoftheSummer| 1281|

|23YrsOfSENSATiONA...| 368|

| SooraraiPottru| 332|

| BTS| 103|

| 55| 91|

| BBNaija| 72|

| charity| 59|

| charitywater| 58|

| Shame_On_MahaGovt| 55|

|__...| 54|

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

Le 2020-09-05 entre 21:00:00 et 22:00:00, ci-dessous les hashtags les plus utilisés:

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

| hashtags|count|

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

| ArtistoftheSummer| 1128|

|23YrsOfSENSATiONA...| 208|

| SooraraiPottru| 200|

| BTS| 99|

| BBNaija| 86|

|JusticeForBreonna...| 63|

	ShamelessTV9	42
	COVID19	37
	TrumpHatesOurVete...	37
	TrumpBoatParade	33
+-----+-----+		

Le 2020-09-05 entre 22:00:00 et 23:00:00, ci-dessous les hashtags les plus utilisés:

+-----+-----+		
	hashtags count	
+-----+-----+		
	ArtistoftheSummer	965
	23YrsOfSENSATiONA...	135
	SooraraiPottru	116
	BBNaija	107
	JusticeForBreonna...	79
	BTS	77
	BGT	39
	TurnUpWithLaycon	38
	TrumpHatesOurVete...	37
	AskBigSean	30
+-----+-----+		

Le 2020-09-05 entre 23:00:00 et 24:00:00, ci-dessous les hashtags les plus utilisés:

+-----+-----+		
	hashtags count	
+-----+-----+		
	ArtistoftheSummer	517
	BBNaija	115
	TurnUpWithLaycon	56
	SooraraiPottru	47
	23YrsOfSENSATiONA...	47
	BTS	40
	Dumbkirk	36
	BBNaijaShallWe	35
	JusticeForBreonna...	30
	bbnaija	30
+-----+-----+		

Le 2020-09-05 entre 24:00:00 et 25:00:00, ci-dessous les hashtags les plus utilisés:

+-----+-----+	
	hashtags count
+-----+-----+	
+-----+-----+	

```
In [6]: # Libérer la mémoire
        df_batchView.unpersist()
```

```
Out[6]: DataFrame[_id: struct<oid:string>, count: bigint, date: string, hashtags: string, heureD
```

1.0.2 5.2. Observations de speedView:

```
In [7]: df_speedView = spark.read.format("mongo").option("uri","mongodb://127.0.0.1/twitter.spee
        df_speedView.persist()
        df_speedView.printSchema()
```

```
root
|-- _id: struct (nullable = true)
|   |-- oid: string (nullable = true)
|-- date_debut: timestamp (nullable = true)
|-- date_fin: timestamp (nullable = true)
|-- hashtags: array (nullable = true)
|   |-- element: struct (containsNull = true)
|       |-- hashtagUsed: string (nullable = true)
|       |-- quantite: integer (nullable = true)
|-- nbTotalMessagesPosted: integer (nullable = true)
```

```
In [8]: # C'est plus pratique d'utiliser une fonction qui affichera les résultats une par une
        def afficherTopTenInTime(dataframe, date, heure_debut_int, heure_fin_int):
            """
                Cette fonction filtre les heures dans le dataframe provenant de mongoDB
                et affiche les classement des HASHTAGS dans l'ordre
                Cette fonction prend en paramètre:
                - df : dataframe provenant de mongoDB
                - date est au format 2020-08-24
                -heure de début : int entre 0, et 23
                -heure de début : int entre 1, et 23
            """
            df1 = dataframe.\
                filter("date_debut > timestamp'{0}' {1}:00:00' AND date_fin < timestamp'{0}' {2}:00:00'")\
                .select("hashtags")

            df1.persist()
            #print(df1.count())

            df_exploded = df1.select(explode(df1.hashtags))
            df_exploded2 = df_exploded\
                .select(df_exploded.col.hashtagUsed, df_exploded.col.quantite)\
                .withColumnRenamed('col.hashtagUsed', 'hashtagUsed')\
                .withColumnRenamed('col.quantite', 'quantitee')

            result_df = df_exploded2.groupBy('hashtagUsed').sum().orderBy("sum(quantitee)", asce
```

```

print("*****")
print("Le {0} entre {1}:00:00 et {2}:00:00, ci-dessous les hashtags les plus utilisés")

result_df.show()
df1.unpersist()

In [9]: # On parcourt toutes les heures depuis minuit à midi
        date = "2020-09-06"
        for x in range(12):
            afficherTopTenInTime(df_speedView, date, x, x+1)

*****
Le 2020-09-06 entre 0:00:00 et 1:00:00, ci-dessous les hashtags les plus utilisés:
+-----+-----+
|      hashtagUsed|sum(quantitee)|
+-----+-----+
| ArtistoftheSummer|      782|
|      BBNaija|      246|
|      BBNaija|      126|
| TurnUpWithLaycon|      88|
|      BTS|      74|
|      bbnaija|      69|
| SayangAdminShopee|      31|
| 23YrsOfSENSATiONAL...|      29|
|      BBNaijaShallWe|      27|
|      Dumbkirk|      25|
+-----+-----+

*****
Le 2020-09-06 entre 1:00:00 et 2:00:00, ci-dessous les hashtags les plus utilisés:
+-----+-----+
|      hashtagUsed|sum(quantitee)|
+-----+-----+
| ArtistoftheSummer|      670|
|      BBNaija|      255|
|      BBNaija|      239|
|      BTS|      89|
|      bbnaija|      60|
|      IceCream|      53|
| MostRequestedLive|      39|
| TurnUpWithLaycon|      38|
|      |      31|
| SayangAdminShopee|      28|
+-----+-----+

*****
Le 2020-09-06 entre 2:00:00 et 3:00:00, ci-dessous les hashtags les plus utilisés:
+-----+-----+

```

hashtagUsed	sum(quantitee)
ArtistoftheSummer	960
BTSARMY	331
BBNaija	306
BTS_Dynamite	296
BBNajia	181
AEWAllOut	179
23YrsofSENSATiONA...	80
BTS	75
bbnaija	73
SooraraiPottru	71

Le 2020-09-06 entre 3:00:00 et 4:00:00, ci-dessous les hashtags les plus utilisés:

hashtagUsed	sum(quantitee)
ArtistoftheSummer	893
23YrsofSENSATiONA...	302
SooraraiPottru	242
AEWAllOut	211
BBNaija	158
BTSARMY	111
BTS_Dynamite	107
BTS	107
MostRequestedLive	48
UFCVegas9	37

Le 2020-09-06 entre 4:00:00 et 5:00:00, ci-dessous les hashtags les plus utilisés:

hashtagUsed	sum(quantitee)
ArtistoftheSummer	1196
23YrsofSENSATiONA...	392
SooraraiPottru	316
2_...	180
BTSARMY	110
AEWAllOut	108
BBNaija	82
BTS	82
IceCream	68
BTS_Dynamite	62

Le 2020-09-06 entre 5:00:00 et 6:00:00, ci-dessous les hashtags les plus utilisés:

```
+-----+-----+
|      hashtagUsed|sum(quantitee)|
+-----+-----+
|  ArtistoftheSummer|      1277|
| 23YrsOfSENSATiONA...|      438|
|    SooraraiPottru|      387|
| 2_...|      137|
|      AEWAllOut|      118|
|      IceCream|      102|
|    ShamelessTV9|       80|
|      BTS|       67|
|  _...|       66|
|      BTSARMY|       59|
+-----+-----+
```

Le 2020-09-06 entre 6:00:00 et 7:00:00, ci-dessous les hashtags les plus utilisés:

```
+-----+-----+
|      hashtagUsed|sum(quantitee)|
+-----+-----+
|  ArtistoftheSummer|     1419|
| 23YrsOfSENSATiONA...|     551|
|    SooraraiPottru|     474|
|    ShamelessTV9|     192|
|  7YrsOfEvergreenVVS|      95|
| 2_...|       78|
|      |       71|
|      BTS|       60|
|      bbrightvc|       57|
| SpeakUpForUPSSSCS...|       47|
+-----+-----+
```

Le 2020-09-06 entre 7:00:00 et 8:00:00, ci-dessous les hashtags les plus utilisés:

```
+-----+-----+
|      hashtagUsed|sum(quantitee)|
+-----+-----+
|    DestinedSidNaaz|     1965|
|  ArtistoftheSummer|     1269|
| 23YrsOfSENSATiONA...|     559|
|    SooraraiPottru|     510|
|    ShamelessTV9|     180|
| HappyBreakfastDay...|     135|
|  7YrsOfEvergreenVVS|      91|
|      SB19|       68|
|      SidNaaz|       43|
```

SpeakUpForPunjabA...	42
----------------------	----

+-----+

Le 2020-09-06 entre 8:00:00 et 9:00:00, ci-dessous les hashtags les plus utilisés:

hashtagUsed	sum(quantitee)
DestinedSidNaaz	2069
ArtistoftheSummer	1110
23YrsOfSENSATiONA...	481
SooraraiPottru	440
ShamelessTV9	192
BTS	92
HappyBreakfastDay...	76
SpeakUpForPunjabA...	69
SB19	43
7YrsOfEvergreenVVS	43

+-----+

Le 2020-09-06 entre 9:00:00 et 10:00:00, ci-dessous les hashtags les plus utilisés:

hashtagUsed	sum(quantitee)
DestinedSidNaaz	1436
ArtistoftheSummer	804
23YrsOfSENSATiONA...	466
SooraraiPottru	370
ShamelessTV9	181
Dynamite6thWin	105
BTS	86
NCT	49
9Baje9Minute	46
RheaChakraborty	46

+-----+

Le 2020-09-06 entre 10:00:00 et 11:00:00, ci-dessous les hashtags les plus utilisés:

hashtagUsed	sum(quantitee)
DestinedSidNaaz	1300
ArtistoftheSummer	568
23YrsOfSENSATiONA...	405
SooraraiPottru	354
ShamelessTV9	221
9Baje9Minute	138

	SEHUN	93
	Dynamite6thWin	83
	RheaChakraborty	74
	BTS	44
+-----+-----+		

Le 2020-09-06 entre 11:00:00 et 12:00:00, ci-dessous les hashtags les plus utilisés:

+-----+-----+	
	hashtagUsed sum(quantitee)
+-----+-----+	
	DestinedSidNaaz 1334
	23YrsOfSENSATiONa... 503
	ArtistoftheSummer 452
	Sooraraipottru 451
	ShamelessTV9 205
	Valimai 160
	9Baje9Minute 146
	SidNaaz 96
	RheaChakraborty 61
	PILIkulaIvyLivesEp1 48
+-----+-----+	

```
In [10]: # Libérer la mémoire
         df_speedView.unpersist()
```

```
Out[10]: DataFrame[_id: struct<oid:string>, date_debut: timestamp, date_fin: timestamp, hashtags
```