

Assignment 2 - Report

Yuzhe Ma & Xin Jin & Hongyi Wang

October 30, 2017

1 Part A: Structured Streaming

We use spark streaming programming model to process the twitter dataset, which contains records for the interactions between twitter users. The dataset are formed in rows, each row corresponding to an event between two users: User_A and User_B, followed by the event time and the interaction type. The interaction type can be one of the threes: retweets (RT), mention (MT) and reply (RE). For instance, if the type is MT, then it means User_A mentioned User_B at a particular time as specified in the row.

1.1 Question 1

In question 1, we wrote an application to count the number of events corresponding to each type of interaction, and output the counts for the events that happen in a window of 1 hour. The starting time of the window alters with step size 30 minutes and thus we will see a sliding window in the output. Fig. 1 shows three batches printed to the console.

Batch: 0	Batch: 1	Batch: 2
+-----+ window[Interaction count] +-----+	+-----+ window[Interaction count] +-----+	+-----+ window[Interaction count] +-----+
2012-07-02 15:30... RT 2 2012-07-06 06:30... RE 1 2012-07-05 06:30... RT 5 2012-07-05 07:00... MT 3 2012-07-05 13:30... RE 1 2012-07-04 17:30... RT 2 2012-07-05 23:30... RE 1 2012-07-05 17:30... RT 1 2012-07-03 09:30... MT 1 2012-07-04 09:30... MT 1 2012-07-01 11:00... MT 9 2012-07-07 04:00... MT 1 2012-07-02 10:00... RT 1 2012-07-04 15:00... RE 2 2012-07-04 03:00... RT 11 2012-07-05 21:00... MT 1 2012-07-04 05:30... RT 5 2012-07-04 18:30... RT 7 2012-07-07 06:30... RT 1 +-----+ only showing top 20 rows	2012-07-02 15:30... RT 3 2012-07-02 13:30... RT 2 2012-07-06 01:30... RE 1 2012-07-06 06:30... RE 1 2012-07-05 06:30... RT 7 2012-07-05 07:00... MT 4 2012-07-05 13:30... RE 1 2012-07-04 17:30... RT 4 2012-07-05 23:30... RE 1 2012-07-05 17:30... RT 2 2012-07-03 09:30... RT 3 2012-07-03 03:00... MT 1 2012-07-04 03:00... MT 12 2012-07-01 11:00... MT 1 2012-07-07 04:00... MT 1 2012-07-02 10:00... RT 1 2012-07-04 03:00... RT 26 2012-07-04 15:00... RE 2 2012-07-05 21:00... MT 1 2012-07-04 05:30... RT 13 +-----+ only showing top 20 rows	2012-07-02 15:30... RT 3 2012-07-02 13:30... RT 2 2012-07-06 01:30... RE 1 2012-07-06 06:30... RE 1 2012-07-05 06:30... RT 8 2012-07-05 07:00... MT 7 2012-07-05 13:30... RE 1 2012-07-04 17:30... RT 5 2012-07-05 23:30... RE 1 2012-07-05 17:30... RT 2 2012-07-03 09:30... RT 3 2012-07-03 03:00... MT 1 2012-07-04 03:00... MT 19 2012-07-01 11:00... MT 1 2012-07-07 04:00... MT 1 2012-07-02 05:00... RT 1 2012-07-07 04:00... MT 1 2012-07-04 18:00... RT 1 2012-07-04 15:00... RE 3 2012-07-04 03:00... RT 43 +-----+ only showing top 20 rows
(a) Batch 0	(b) Batch 1	(c) Batch 2

Figure 1: Counts of the events for three different types of interactions. The first three batches are shown.

1.2 Question 2

In question 2, we wrote an application to emit the twitter IDs of users that have been mentioned by other users every 10 seconds. The output is required to be on HDFS, however, to show our result intuitively, we direct the output to console in this document and keep the output address as is required in the source code. Fig. 2 shows the output printed to the console.

Batch: 0		Batch: 1		Batch: 2	
ID		ID		ID	
383		58574		206726	
148		10574		116499	
38834		88		610	
29600		88		9181	
88		88		60120	
4022		105601		19913	
261475		8984		426630	
33476		147925		8732	
146171		80477		20487	
591952		6373		25184	
20882		6730		531	
14376		314965		7230	
88		677		13797	
509901		51444		32064	
31957		28233		88	
16783		35376		867	
88		14584		10	
406395		28076		205	
68285		2689		235	
309781		1062		519	

(a) Batch 0

(b) Batch 1

(c) Batch 2

Figure 2: Counts of the events for three different types of interactions. The first three batches are shown.

1.3 Question 3

In question 3, we wrote an application takes as input a list of target twitter user IDs and every 5 seconds, emits the number of tweet actions of a user if it is present in the input list. Fig. 3 shows the output printed to the console.

User	Batch: 0	Batch: 10	Batch: 20
51640	User count	User count	User count
194667	194667 1	194667 1	194667 1
227235	408352 1	408352 1	408352 1
408352	119701 1	119701 2	119701 2
119701	51640 1	51640 1	51640 2
234	227235 1	227235 1	227235 1
34546456			

(a) User list

(b) Batch 0

(c) Batch 10

(d) Batch 20

Figure 3: Counts of the events for three different types of interactions. User list, Batch 0, 10 and 20 are shown.

2 Part B - Apache Storm

2.1 Question 1

To answer this question, you need to extend the PrintSampleStream.java topology to collect tweets which are in English and match certain keywords (at least 10 distinct keywords) at your desire. You should collect at least 200'000 tweets which match the keywords and store them either in your local filesystem or HDFS. Your topology should be runnable both in local mode and cluster mode.

Solution:

Source code:

Main function: PartBQ1.java

Spout:

PartBQ1Spout.java

Use Twitter4j library to stream tweets in English language from twitter under the constriction of given keywords.

Bolt:

PartBQ1WriteLocalBolt.java(not used)

Write to local file

Bolt:

use HdfsBolt

Write to hdfs

Compile application:

```
cd /home/ubuntu/software/apache-storm-1.0.2/examples/storm-starter  
mvn clean install -DskipTests=true; mvn package -DskipTests=true
```

Run application:

```
cd /home/ubuntu/software/apache-storm-1.0.2/examples/storm-starter  
storm jar target/storm-starter-*.jar org.apache.storm.starter.PartBQ1 [lo-  
cal/cluster; ]keywords list;
```

e.g.:

```
storm jar target/storm-starter-*.jar org.apache.storm.starter.PartBQ1 clus-  
ter halloween love peace Running win sex girls boys man woman father mother
```

2.2 Question 2

To answer this question, you need to create a new topology which provides the following functionality: every 30 seconds, you should collect all the tweets which are in English, have certain hashtags and have the friendsCount field satisfying the condition mentioned below. For all the collected tweets in a time interval, you need to find the top 50% most common words which are not stop words (<http://www.ranks.nl/stopwords>). For every time interval, you should store both the tweets and the top common words in separate files either in your local filesystem or HDFS. Your topology should be runnable both in local mode and cluster mode.

Solution:

Main function: PartBQ2.java

Spout:
PartBQ2Spout.java
 Use Twitter4j library to stream tweets in English language from twitter under the constriction of hashtags.

SampleFriendCountSpout.java
 In every 30 seconds, it will random sample a friendcount from the given friend count lists.

SampleHashTagsSpout.java
 In every 30 seconds, it will random sample a list of hashtags based on the given hashtag lists.

Bolt:
TweetFilterBolt.java
 Based on the stream received from PartBQ2Spout, SampleFriendCountSpout, SampleHashTagsSpout, use friendcount and hashtags to filter the tweets.

TweetSplitWordBolt.java
 Based on the stream received from TweetFilterBolt, split the tweets into words, remove stopwords and remove non-character words. **TweetWordRankBolt.java**

Based on the stream received from TweetSplitWordBolt, in every 30 seconds, it will sort the words received in current time interval and keep the words with rankings above 50%.

Helper class:
MyWord.java
 To store word and its counter.

Compile application:
`cd /home/ubuntu/software/apache-storm-1.0.2/examples/storm-starter
mvn clean install -DskipTests=true; mvn package -DskipTests=true`

Run application:
`cd /home/ubuntu/software/apache-storm-1.0.2/examples/storm-starter
storm jar target/storm-starter-*.jar org.apache.storm.starter.PartBQ2 local/cluster; friendcount list; hashtags list;`

e.g.:
`storm jar target/storm-starter-*.jar org.apache.storm.starter.PartBQ2 cluster 5 6 7 8 9 10 20 40 60 80 100 1000 10000 100000 halloween love peace running win sex girls boys man woman father mother design sexy dress fashion free shop style sales today price pumpkin star iphone google airbnb`

3 Part C: Flink

3.1 Question 1

In Question 1, we wrote application to find both disjoint 1-min time intervals using *TumblingEventTimeWindows* in *Flink* while application for non-disjoint 1-min time interval was implemented using *SlidingEventTimeWindows*. Since we're using Event time here in Flink, we extracted timestamp provided in sorted

higgs dataset. Noticed that, the timestamp provided in the dataset are in second, we converted them into millisecond for compatibility with Flink (implementation details can be check in our code *PartCQuestion_1_1.java* is the implement of Tumbling Window while *PartCQuestion_1_2.java* implements Sliding Window), part of results can be found in Fig. 4. For full results, please run our applications and the full results will print out to the console.

partC_Q1_dtsJoint	
<i>partC_Q1 sliding output:</i>	
TimeWindow[start=1341377560000, end=1341377640000)	Count: 113 Type: RT
TimeWindow[start=1341377640000, end=1341377710000)	Count: 135 Type: RT
TimeWindow[start=1341377710000, end=1341377780000)	Count: 136 Type: RT
TimeWindow[start=1341377780000, end=1341377820000)	Count: 152 Type: RT
TimeWindow[start=1341377820000, end=1341377880000)	Count: 197 Type: MT
TimeWindow[start=1341377880000, end=1341377940000)	Count: 197 Type: MT
TimeWindow[start=1341377940000, end=1341377980000)	Count: 105 Type: RT
TimeWindow[start=1341377980000, end=1341377990000)	Count: 109 Type: RT
TimeWindow[start=1341377990000, end=1341378000000)	Count: 114 Type: MT
TimeWindow[start=1341378000000, end=1341378060000)	Count: 160 Type: RT
TimeWindow[start=1341378060000, end=1341378120000)	Count: 140 Type: RT
TimeWindow[start=1341378120000, end=1341378180000)	Count: 169 Type: RT
TimeWindow[start=1341378180000, end=1341378240000)	Count: 135 Type: MT
TimeWindow[start=1341378240000, end=1341378280000)	Count: 187 Type: MT
TimeWindow[start=1341378280000, end=1341378320000)	Count: 25 Type: MT
TimeWindow[start=1341378320000, end=1341378360000)	Count: 187 Type: RT
TimeWindow[start=1341378360000, end=1341378400000)	Count: 133 Type: MT
TimeWindow[start=1341378400000, end=1341378460000)	Count: 190 Type: RT
TimeWindow[start=1341378460000, end=1341378520000)	Count: 186 Type: RT
TimeWindow[start=1341378520000, end=1341378580000)	Count: 117 Type: MT
TimeWindow[start=1341378580000, end=1341378640000)	Count: 113 Type: MT
TimeWindow[start=1341378640000, end=1341378700000)	Count: 137 Type: RT
TimeWindow[start=1341378700000, end=1341378760000)	Count: 167 Type: RT
TimeWindow[start=1341378760000, end=1341378840000)	Count: 121 Type: MT
TimeWindow[start=1341378840000, end=1341378900000)	Count: 128 Type: RT
TimeWindow[start=1341378900000, end=1341378960000)	Count: 194 Type: MT
TimeWindow[start=1341378960000, end=1341379020000)	Count: 108 Type: MT
TimeWindow[start=1341379020000, end=1341379720000)	Count: 190 Type: RT
TimeWindow[start=1341379720000, end=1341379860000)	Count: 112 Type: MT
TimeWindow[start=1341379860000, end=1341379920000)	Count: 147 Type: MT
TimeWindow[start=1341379920000, end=1341379980000)	Count: 209 Type: RT
TimeWindow[start=1341379980000, end=1341379980000)	Count: 159 Type: RT
TimeWindow[start=1341379980000, end=1341379980000)	Count: 226 Type: RT
TimeWindow[start=1341379980000, end=1341379980000)	Count: 243 Type: RT
TimeWindow[start=1341379980000, end=1341379980000)	Count: 211 Type: RT
TimeWindow[start=1341379980000, end=1341379980000)	Count: 217 Type: MT
TimeWindow[start=1341379980000, end=1341379980000)	Count: 239 Type: RT
<i>partC_Q1 sliding output:</i>	
TimeWindow[start=1341377539000, end=1341377540000)	Count: 101 Type: RT
TimeWindow[start=1341377540000, end=1341377541000)	Count: 106 Type: RT
TimeWindow[start=1341377541000, end=1341377542000)	Count: 108 Type: RT
TimeWindow[start=1341377542000, end=1341377543000)	Count: 107 Type: RT
TimeWindow[start=1341377543000, end=1341377544000)	Count: 104 Type: RT
TimeWindow[start=1341377544000, end=1341377545000)	Count: 105 Type: RT
TimeWindow[start=1341377545000, end=1341377546000)	Count: 109 Type: RT
TimeWindow[start=1341377546000, end=1341377547000)	Count: 107 Type: RT
TimeWindow[start=1341377547000, end=1341377548000)	Count: 108 Type: RT
TimeWindow[start=1341377548000, end=1341377549000)	Count: 108 Type: RT
TimeWindow[start=1341377549000, end=1341377550000)	Count: 109 Type: RT
TimeWindow[start=1341377550000, end=1341377551000)	Count: 108 Type: RT
TimeWindow[start=1341377551000, end=1341377552000)	Count: 107 Type: RT
TimeWindow[start=1341377552000, end=1341377553000)	Count: 108 Type: RT
TimeWindow[start=1341377553000, end=1341377554000)	Count: 104 Type: RT
TimeWindow[start=1341377554000, end=1341377555000)	Count: 102 Type: RT
TimeWindow[start=1341377555000, end=1341377556000)	Count: 102 Type: RT
TimeWindow[start=1341377556000, end=1341377557000)	Count: 102 Type: RT
TimeWindow[start=1341377557000, end=1341377558000)	Count: 108 Type: RT
TimeWindow[start=1341377558000, end=1341377559000)	Count: 104 Type: RT
TimeWindow[start=1341377559000, end=1341377560000)	Count: 100 Type: RT
TimeWindow[start=1341377560000, end=1341377561000)	Count: 104 Type: RT
TimeWindow[start=1341377561000, end=1341377562000)	Count: 104 Type: RT
TimeWindow[start=1341377562000, end=1341377563000)	Count: 103 Type: RT
TimeWindow[start=1341377563000, end=1341377564000)	Count: 100 Type: RT
TimeWindow[start=1341377564000, end=1341377565000)	Count: 100 Type: RT
TimeWindow[start=1341377565000, end=1341377566000)	Count: 105 Type: RT
TimeWindow[start=1341377566000, end=1341377567000)	Count: 100 Type: RT
TimeWindow[start=1341377567000, end=1341377568000)	Count: 100 Type: RT
TimeWindow[start=1341377568000, end=1341377569000)	Count: 100 Type: RT
TimeWindow[start=1341377569000, end=1341377570000)	Count: 100 Type: RT
TimeWindow[start=1341377570000, end=1341377571000)	Count: 107 Type: RT
TimeWindow[start=1341377571000, end=1341377572000)	Count: 106 Type: RT
TimeWindow[start=1341377572000, end=1341377573000)	Count: 103 Type: RT
TimeWindow[start=1341377573000, end=1341377574000)	Count: 100 Type: RT
TimeWindow[start=1341377574000, end=1341377575000)	Count: 101 Type: RT
TimeWindow[start=1341377575000, end=1341377576000)	Count: 100 Type: RT
TimeWindow[start=1341377576000, end=1341377577000)	Count: 100 Type: RT
TimeWindow[start=1341377577000, end=1341377578000)	Count: 104 Type: RT
TimeWindow[start=1341377578000, end=1341377579000)	Count: 103 Type: RT
TimeWindow[start=1341377579000, end=1341377580000)	Count: 103 Type: RT
TimeWindow[start=1341377580000, end=1341377581000)	Count: 103 Type: RT
TimeWindow[start=1341377581000, end=1341377582000)	Count: 103 Type: RT
TimeWindow[start=1341377582000, end=1341377583000)	Count: 103 Type: RT
TimeWindow[start=1341377583000, end=1341377584000)	Count: 103 Type: RT
TimeWindow[start=1341377584000, end=1341377585000)	Count: 103 Type: RT
TimeWindow[start=1341377585000, end=1341377586000)	Count: 103 Type: RT
TimeWindow[start=1341377586000, end=1341377587000)	Count: 103 Type: RT
TimeWindow[start=1341377587000, end=1341377588000)	Count: 103 Type: RT
TimeWindow[start=1341377588000, end=1341377589000)	Count: 103 Type: RT
TimeWindow[start=1341377589000, end=1341377590000)	Count: 103 Type: RT
TimeWindow[start=1341377590000, end=1341377591000)	Count: 103 Type: RT
TimeWindow[start=1341377591000, end=1341377592000)	Count: 103 Type: RT
TimeWindow[start=1341377592000, end=1341377593000)	Count: 103 Type: RT
TimeWindow[start=1341377593000, end=1341377594000)	Count: 103 Type: RT
TimeWindow[start=1341377594000, end=1341377595000)	Count: 103 Type: RT
TimeWindow[start=1341377595000, end=1341377596000)	Count: 103 Type: RT
TimeWindow[start=1341377596000, end=1341377597000)	Count: 103 Type: RT
TimeWindow[start=1341377597000, end=1341377598000)	Count: 103 Type: RT
TimeWindow[start=1341377598000, end=1341377599000)	Count: 103 Type: RT
TimeWindow[start=1341377599000, end=1341377600000)	Count: 103 Type: RT
TimeWindow[start=1341377600000, end=1341377601000)	Count: 103 Type: RT
TimeWindow[start=1341377601000, end=1341377602000)	Count: 103 Type: RT
TimeWindow[start=1341377602000, end=1341377603000)	Count: 103 Type: RT
TimeWindow[start=1341377603000, end=1341377604000)	Count: 103 Type: RT
TimeWindow[start=1341377604000, end=1341377605000)	Count: 103 Type: RT
TimeWindow[start=1341377605000, end=1341377606000)	Count: 103 Type: RT
TimeWindow[start=1341377606000, end=1341377607000)	Count: 103 Type: RT
TimeWindow[start=1341377607000, end=1341377608000)	Count: 103 Type: RT
TimeWindow[start=1341377608000, end=1341377609000)	Count: 103 Type: RT
TimeWindow[start=1341377609000, end=1341377610000)	Count: 103 Type: RT
TimeWindow[start=1341377610000, end=1341377611000)	Count: 103 Type: RT
TimeWindow[start=1341377611000, end=1341377612000)	Count: 103 Type: RT
TimeWindow[start=1341377612000, end=1341377613000)	Count: 103 Type: RT
TimeWindow[start=1341377613000, end=1341377614000)	Count: 103 Type: RT
TimeWindow[start=1341377614000, end=1341377615000)	Count: 103 Type: RT
TimeWindow[start=1341377615000, end=1341377616000)	Count: 103 Type: RT
TimeWindow[start=1341377616000, end=1341377617000)	Count: 103 Type: RT
TimeWindow[start=1341377617000, end=1341377618000)	Count: 103 Type: RT
TimeWindow[start=1341377618000, end=1341377619000)	Count: 103 Type: RT
TimeWindow[start=1341377619000, end=1341377620000)	Count: 103 Type: RT
TimeWindow[start=1341377620000, end=1341377621000)	Count: 103 Type: RT
TimeWindow[start=1341377621000, end=1341377622000)	Count: 103 Type: RT
TimeWindow[start=1341377622000, end=1341377623000)	Count: 103 Type: RT
TimeWindow[start=1341377623000, end=1341377624000)	Count: 103 Type: RT
TimeWindow[start=1341377624000, end=1341377625000)	Count: 103 Type: RT
TimeWindow[start=1341377625000, end=1341377626000)	Count: 103 Type: RT
TimeWindow[start=1341377626000, end=1341377627000)	Count: 103 Type: RT
TimeWindow[start=1341377627000, end=1341377628000)	Count: 103 Type: RT
TimeWindow[start=1341377628000, end=1341377629000)	Count: 103 Type: RT
TimeWindow[start=1341377629000, end=1341377630000)	Count: 103 Type: RT
TimeWindow[start=1341377630000, end=1341377631000)	Count: 103 Type: RT

(a) Tumbling Window

(b) Sliding Window

Figure 4: Counts of the events (≥ 100) for tumbling and sliding windows

3.2 Question 2

In Question 2 we wrote application to handle out-of-order data. Different from Question 1, we use allow lateness interface in Flink this time. By tuning the late allowance we get the result in Fig.5. From this results we can tell that, allow lateness provides a good strategy for Flink to handle out-of-order data. More late allowance gives closer results to in order data.

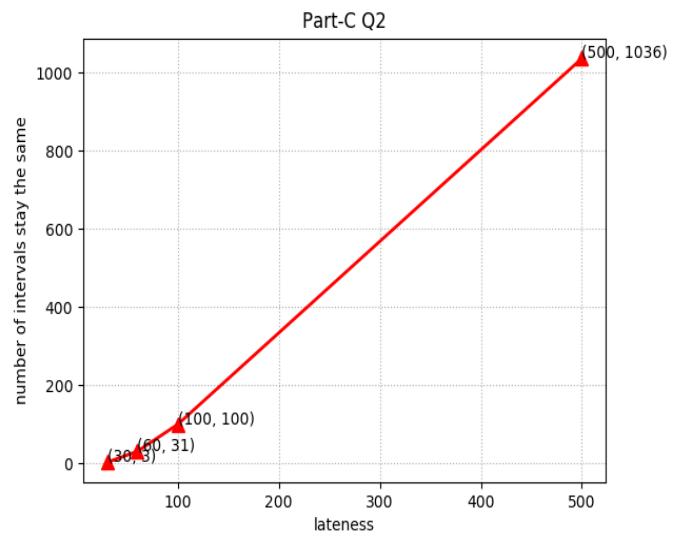


Figure 5: number of intervals that stay the same under various late allowances