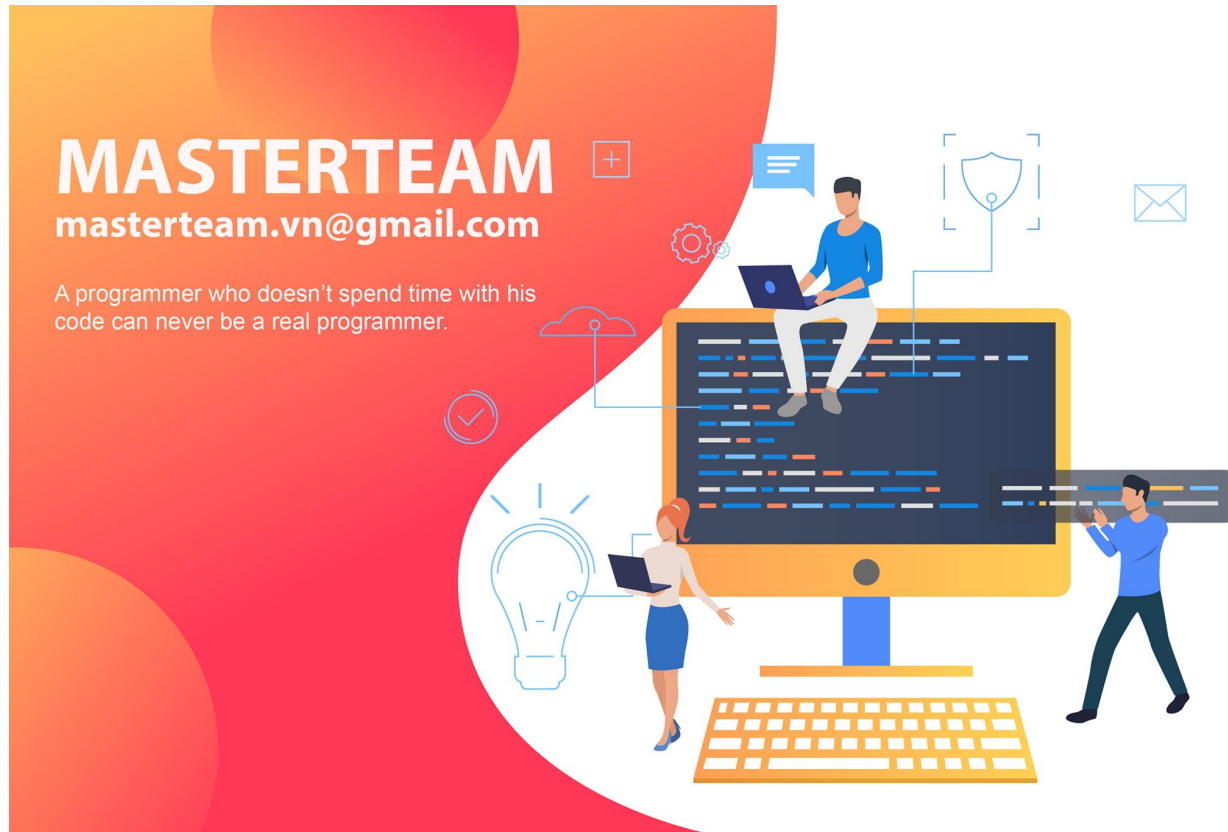


Viết NOTES:

Lưu dưới 2 dạng:

1. File docx dùng để chỉnh sửa
2. File .pdf dùng để preview trên Git



Demo một đoạn text dài

## Xử lý dữ liệu Twitter trong thời gian thực

Trong bài tutorial này, chúng ta sẽ viết chương trình sử dụng TwitterStreaming API để thu thập dữ liệu trong thời gian thực và truyền về Kafka. Sau đó ta sẽ sử dụng Spark như một Consumer của Kafka để tiến hành xử dữ liệu đổ về. Để hiểu rõ thêm về bài tutorial này, các bạn có thể tham khảo hai bài tutorial trước là: [Sử dụng Kafka với Twitter](#) và [Tích hợp Kafka với Spark sử dụng Structured Streaming](#).

Trong bài này, chúng ta sẽ sử dụng kết hợp cả hai ngôn ngữ Java và Scala để viết hai chương trình:

- KafkaTwitterStreaming.java: sử dụng Java để viết chương trình thu thập dữ liệu Twitter trong thời gian thực (Twitter Streaming Data) và truyền về Kafka
- KafkaSparkSS.scala: sử dụng Scala để viết chương trình Spark lấy dữ liệu từ Kafka broker.

Tương tự như trong bài tutorial [Sử dụng Kafka với Twitter](#), chúng ta sẽ sử dụng thư viện twitter4j nhưng là twitter4j-stream thay vì twitter4j-core. Do đó ta thêm dependency sau vào file pom.xml:

```

1      <!-- https://mvnrepository.com/artifact/org.twitter4j/twitter4j-stream -->
2      <dependency>
3          <groupId>org.twitter4j</groupId>
4          <artifactId>twitter4j-stream</artifactId>
5          <version>4.0.7</version>
6      </dependency>

```

Ta tiến hành cấu hình thuộc tính của Kafka và thực hiện kết nối với tài khoản Twitter tương tự như bài [Sử dụng Kafka với Twitter](#) với hai hàm sau (ta thêm .setJSONStoreEnabled(true) để truyền dữ liệu về Kafka dưới dạng Json):

```

1  private static Properties getKafkaProp()
2  {
3      // create instance for properties to access producer configs
4      Properties props = new Properties();
5      //Assign localhost id
6      props.put("bootstrap.servers", "localhost:9092");
7      //Set acknowledgements for producer requests.
8      props.put("acks", "all");
9      //If the request fails, the producer can automatically retry,
10     props.put("retries", 0);
11     //Specify buffer size in config
12     props.put("batch.size", 16384);
13     //Reduce the no of requests less than 0
14     props.put("linger.ms", 1);
15     //The buffer.memory controls the total amount of memory available to the producer
16     for buffering.
17     props.put("buffer.memory", 33554432);
18     props.put("key.serializer",
19         "org.apache.kafka.common.serialization.StringSerializer");
20     props.put("value.serializer",

```

```

1      "org.apache.kafka.common.serialization.StringSerializer");
9
2
0      return props;
2  }
1
2
2  private static Configuration getTwitterConf()
2  {
3      //Config Twitter API key to access Twitter API
4      //The String keys here are only examples and not valid.
2      //You need to use your own keys
2      ConfigurationBuilder cb = new ConfigurationBuilder();
6      cb.setDebugEnabled(true)
2      .setJSONStoreEnabled(true)
2      .setOAuthConsumerKey("Fljmu9Wp1YVNXhqfmDHDyEAz9")
8      .setOAuthConsumerSecret("7CZDMiqhaeV7FOsUTYLgi9utt4eYEVaxqVuKZj5VGHLyq
9 00mLU")
3      .setOAuthAccessToken("1060702756430729216-
0 1L9lL05TdEbanhGDFETkKMknmbw70w")
3      .setOAuthAccessTokenSecret("Qu41ydcAzTxClfVW4BMU6UjziS6Lv9Kkwz1zBXKh3J
1  Wrx");
3
2      return cb.build();
3
4  }
3
5
3
6
3
7
3
8
3
9

```

Tiếp theo ta viết hàm `getStreamTweets()` để thu thập dữ liệu Twitter trong thời gian thực như sau:

```

1  public static void getStreamTweets(final Producer<String, String> producer,String
2  topicName) {
3      TwitterStream twitterStream = new
4      TwitterStreamFactory(getTwitterConf()).getInstance();
5      StatusListener listener = new StatusListener(){
6
7          public void onStatus(Status status) {
8              //Status To JSON String
9              String statusJson = DataObjectFactory.getRawJSON(status);
10             ProducerRecord data = new ProducerRecord("TwitterStreaming", statusJson);
11
12
13
14

```

```

        System.out.println(statusJson);

        //Send data
        producer.send(data);

    }

    public void onException(Exception ex) {
        ex.printStackTrace();
    }

    public void onDeleteNotice(StatusDeletionNotice statusDeletionNotice) {

    }

    public void onTrackLimitationNotice(int numberOfLimitedStatuses) {

    }

    public void onScrubGeo(long userId, long upToStatusId) {

    }

    public void onStallWarning(StallWarning warning) {

    }
};

twitterStream.addListener(listener);
twitterStream.filter(topicName);
}

```

Chạy chương trình KafkaTwitterStreaming với chủ đề tìm kiếm là “Big Data” ta được kết quả sau (các bạn tham khảo thêm về Kafka cũng như cách khởi chạy Kafka trong [loạt bài tutorials về Kafka](#)):

```

1  public static void main(String[] args) {
2      //The Kafka Topic
3      String topicName="Big Data";
4
5      //Define a Kafka Producer
6      Producer<String, String> producer = new KafkaProducer<String, String>(getKafkaProp());
7      getStreamTweets(producer,topicName);
8
9  }

```

Kết quả chạy trong IntelliJ (dữ liệu hiển thị là dữ liệu dưới dạng Json, chứa đầy đủ thông tin về bài đăng như ID, nội dung(text), thời gian post(created\_at), ...):

```
{\"created_at\":\"Sat Jan 05 09:02:38 +0000 2019\",\"id\":\"1081476051077754882\",\"id_str\":\"1081476051077754882\",\"text\":\"10 Effective Big Data Strategies For Marketing Online H...\",\"source\":\"\",\"created_at\":\"Sat Jan 05 09:02:38 +0000 2019\",\"id\":\"1081476051967115264\",\"id_str\":\"1081476051967115264\",\"text\":\"Opinion | Big data is way hotter than you think - Liven...\",\"source\":\"\",\"created_at\":\"Sat Jan 05 09:02:39 +0000 2019\",\"id\":\"1081476056664535042\",\"id_str\":\"1081476056664535042\",\"text\":\"10 Effective Big Data Strategies For Marketing Online H...\",\"source\":\"\",\"created_at\":\"Sat Jan 05 09:02:44 +0000 2019\",\"id\":\"1081476078546370560\",\"id_str\":\"1081476078546370560\",\"text\":\"RT @EinsteinMillan: Proyección de producción @pdvsa 201...\",\"source\":\"\",\"created_at\":\"Sat Jan 05 09:03:15 +0000 2019\",\"id\":\"1081476208255217664\",\"id_str\":\"1081476208255217664\",\"text\":\"Apache Spark 2.0 with Scala - Hands On with Big Data!\\n...\",\"source\":\"\",\"created_at\":\"Sat Jan 05 09:03:22 +0000 2019\",\"id\":\"1081476235518042112\",\"id_str\":\"1081476235518042112\",\"text\":\"RT @ShawnaG NDP: Mr Kenney only mentions the Stats Can...\",\"source\":\"\",\"created_at\":\"Sat Jan 05 09:03:24 +0000 2019\",\"id\":\"1081476244732952576\",\"id_str\":\"1081476244732952576\",\"text\":\"RT @jennycohn1: \"Newly surfaced e-mails show that the f...\",\"source\":\"\",\"created_at\":\"Sat Jan 05 09:03:25 +0000 2019\",\"id\":\"1081476250642788352\",\"id_str\":\"1081476250642788352\",\"text\":\"RT @CharterKubiya: Yako that's my tweet. https://t.c...\",\"source\":\"\",\"created_at\":\"Sat Jan 05 09:03:27 +0000 2019\",\"id\":\"1081476257769054209\",\"id_str\":\"1081476257769054209\",\"text\":\"RT @davidsirota: Whether or not you think this story is...\",\"source\":\"\",\"created_at\":\"Sat Jan 05 09:03:37 +0000 2019\",\"id\":\"1081476299653349376\",\"id_str\":\"1081476299653349376\",\"text\":\"RT @drbobgill: Please note complete absence of privat...\",\"source\":\"\",\"created_at\":\"Sat Jan 05 09:04:05 +0000 2019\",\"id\":\"1081476419190849536\",\"id_str\":\"1081476419190849536\",\"text\":\"Unlock the keys to higher Facebook engagement with the...\",\"source\":\"\",\"created_at\":\"Sat Jan 05 09:04:08 +0000 2019\",\"id\":\"1081476428548489216\",\"id_str\":\"1081476428548489216\",\"text\":\"RT @jennycohn1: \"Newly surfaced e-mails show that the f...\",\"source\":\"\",\"created_at\":\"Sat Jan 05 09:04:12 +0000 2019\",\"id\":\"1081476447309580608\",\"id_str\":\"1081476447309580608\",\"text\":\"RT @jennycohn1: \"Newly surfaced e-mails show that the f...\",\"source\":\"\"}
```

Kết quả chạy Kafka Consumer (nhận kết quả tương tự như kết quả hiển thị trong IntelliJ, kết quả dưới đây chính là bài đăng thứ 4 có số ID là 1081476078546370560 ở trên):

```
{\"created_at\":\"Sat Jan 05 09:02:44 +0000 2019\",\"id\":\"1081476078546370560\",\"id_str\":\"1081476078546370560\",\"text\":\"RT @EinsteinMillan: P...\",\"source\":\"\",\"created_at\":\"Sat Jan 05 09:02:44 +0000 2019\",\"id\":\"1081476078546370560\",\"id_str\":\"1081476078546370560\",\"text\":\"RT @EinsteinMillan: P...\",\"source\":\"\"}
```

Như vậy ta đã hoàn thành việc viết chương trình KafkaTwitterStreaming.java để thu thập dữ liệu Twitter trong thời gian thực và truyền về Kafka (số lượng dữ liệu sẽ tăng theo thời gian do chương trình sẽ liên tục thu thập thông tin liên quan đến chủ đề cần tìm kiếm)

Bước tiếp theo ta viết chương trình KafkaSparkSS.scala để lấy dữ liệu về từ Kafka broker. Ta thực hiện việc kết nối Spark với Kafka tương tự như trong bài tutorial [Tích hợp Kafka với Spark sử dụng Structured Streaming](#)

```
1 //Define a Spark session
2 val spark=SparkSession.builder().appName(\"Spark Kafka Integration using Structured Streaming\")
3   .master(\"local\")
4   .getOrCreate()
5
6 //Set the Log file level
7 spark.sparkContext.setLogLevel(\"WARN\")
```

```

8
9 //Implicit methods available in Scala for converting common Scala objects into DataFrames
10 import spark.implicits._
11
12 //Subscribe Spark to topic 'TwitterStreaming'
13 val df=spark.readStream.format("kafka")
14   .option("kafka.bootstrap.servers","localhost:9092")
15   .option("subscribe","TwitterStreaming")
16   .load()

```

Như đã đề cập ở trên, dữ liệu truyền về là dưới dạng Json và Spark Structured Streaming lại không hỗ trợ việc tự động trích xuất Schema của dữ liệu Kafka. Do đó ta phải tiến hành định nghĩa Schema này trong chương trình KafkaSparkSS.scala. Chúng ta có thể định nghĩa Schema theo cách thủ công như trong bài [Phân tích và thống kê dữ liệu sử dụng Spark DataFrame](#), tuy nhiên Schema của Twiter tương đối phức tạp nên ta sẽ sử dụng phương pháp dùng Spark để tự động trích xuất Schema của một mẫu dữ liệu Twitter cho trước (Ở đây chúng ta copy một mẫu dữ liệu trong phần chạy KafkaTwitterStreaming.java và lưu thành file twitter.json rồi sử dụng đoạn code sau để trích xuất Schema)

```

1 //Extract the schema from a sample of Twitter Data
2 val twitterData=spark.read.json("src/main/resources/data_source/twitter.json").toDF()
3 val twitterDataScheme=twitterData.schema

```

```

root
|-- contributors: string (nullable = true)
|-- coordinates: string (nullable = true)
|-- created_at: string (nullable = true)
|-- entities: struct (nullable = true)
|   |-- hashtags: array (nullable = true)
|   |   |-- element: string (containsNull = true)
|   |-- symbols: array (nullable = true)
|   |   |-- element: string (containsNull = true)
|   |-- urls: array (nullable = true)
|   |   |-- element: string (containsNull = true)
|   |-- user_mentions: array (nullable = true)
|   |   |-- element: struct (containsNull = true)
|   |       |-- id: long (nullable = true)
|   |       |-- id_str: string (nullable = true)
|   |       |-- indices: array (nullable = true)
|   |       |   |-- element: long (containsNull = true)
|   |       |-- name: string (nullable = true)
|   |       |-- screen_name: string (nullable = true)
|-- favorite_count: long (nullable = true)
|-- favorited: boolean (nullable = true)
|-- filter_level: string (nullable = true)
|-- geo: string (nullable = true)
|-- id: long (nullable = true)
|-- id_str: string (nullable = true)
|-- in_reply_to_screen_name: string (nullable = true)
|-- in_reply_to_status_id: string (nullable = true)
|-- in_reply_to_status_id_str: string (nullable = true)
|-- in_reply_to_user_id: string (nullable = true)
|-- in_reply_to_user_id_str: string (nullable = true)
|-- is_quote_status: boolean (nullable = true)
|-- lang: string (nullable = true)
|-- place: string (nullable = true)
|-- quote_count: long (nullable = true)
|-- reply_count: long (nullable = true)
|-- retweet_count: long (nullable = true)
|-- retweeted: boolean (nullable = true)
|-- retweeted_status: struct (nullable = true)

```

Sau khi đã xác định được Schema của dữ liệu Twitter, ta tiến hành đọc dữ liệu nhận về với Schema này như sau:

```

1 //Reading the streaming json data with its schema
2 val twitterStreamData=df.selectExpr( "CAST(value AS STRING) as jsonData")

```



```

3 .select(from_json($"jsonData",schema = twitterDataScheme).as("data"))
4 .select("data.*")

```

Sau đó ta tiến hành hiển thị dữ liệu Spark lấy về từ Kafka với đoạn code sau:

```

1 // Display output (all columns)
2 val query = twitterStreamData
3   .writeStream
4   .outputMode("append")
5   .format("console")
6   .start()
7
8 // Display output (only few columns)
9 val query2 = twitterStreamData.select("created_at","user.name","text","user.lang")
10  .writeStream
11  .outputMode("append")
12  .format("console")
13  .start()
14
15 query.awaitTermination()
16 query2.awaitTermination()

```

Chạy chương trình KafkaSparkSS.scala ta được kết quả sau:

```

-----
Batch: 1
-----
+-----+-----+-----+-----+
|      created_at      |      name      |      text      |lang|
+-----+-----+-----+-----+
|Sat Jan 05 09:18:...|LAXMIKANTA BJP SORO|RT @RajeshkTripat...| en|
+-----+-----+-----+-----+

-----
Batch: 1
-----
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|contributors|coordinates|      created_at      |      entities      |favorite_count|favorited|filter_level|geo|      id      |      id_str      |in_reply_to_screen_name|in_reply_to_user_id_str|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|      null      |      null      |Sat Jan 05 09:18:...|[[]], [], [], [[12...|      0      |false|low|null|1081480130512343042|1081480130512343042|      null      |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+

-----
Batch: 2
-----
+-----+-----+-----+-----+
|      created_at      |      name      |      text      |lang|
+-----+-----+-----+-----+
|Sat Jan 05 09:18:...|Kim Amadril|RT @PreppyQ: Here...| en|
+-----+-----+-----+-----+

-----
Batch: 2
-----
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|contributors|coordinates|      created_at      |      entities      |favorite_count|favorited|filter_level|geo|      id      |      id_str      |in_reply_to_screen_name|in_reply_to_user_id_str|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|      null      |      null      |Sat Jan 05 09:18:...|[[]], [], [], [[37...|      0      |false|low|null|1081480149437112321|1081480149437112321|      null      |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+

```