

MIDS W205

Lab #	10	Lab Title	Spark Streaming Introduction
Related Module(s)	9	Goal	Spark Streaming, filtering, counting, sliding window.
Last Updated	4/1/16	Expected duration	60-90 minutes

Introduction

Spark Streaming is an extension to Spark that enables processing of streaming data. By streaming data processing we mean that data arrives continuously to an analytics process, and that we have a need to process and analyze the data as it comes in. Spark Streaming can receive data from any sources. For the purposes of this lab we will feed data from a set up we will characterize as a “poor” mans streams since it is for demonstration purposes only.

Before we start with the Lab we will briefly introduce Spark Streaming and its main concepts. From other labs you are familiar with Storm. Spark Streaming differs from Storm in several respects. Firstly Spark Streaming divides an incoming stream into batches. Each batch is processed as one unit using the core Spark infrastructure. This is in contrast to Storm where data is processed as it comes in. Storm like computing is sometimes called *record-at-a-time* processing. Spark Streaming like computation is sometimes called a *micro batching*. Micro batching has both advantages and disadvantages compared to a *record-at-a-time* streaming solution. One advantage of Spark Streaming is that it builds on the core Spark RDD processing, and hence enables usage of multiple paradigms such as stream processing, batching and interactive queries on one platform.



There are a few basic concepts that you need to understand in Spark Streaming: DStream, Transformations and Output Operations.

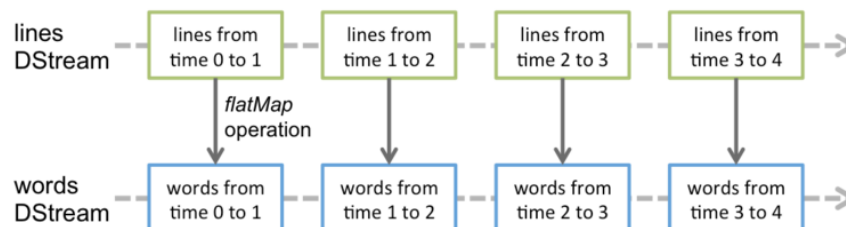
Discretized Streams, or DStreams for short, is an architectural concept that essentially captures the fact that an incoming continuous stream is chunked in to discrete RDDs for Spark processing. DStreams allows many of the transformations supported by RDD's. Below you will find a list of some of them. The full list is available in the programming guide.

- **map(func):** Return a new DStream by passing each element of the source DStream through a function `func`.

- **flatMap(func):** Similar to map, but each input item can be mapped to 0 or more output items.
- **filter(func):** Return a new DStream by selecting only the records of the source DStream on which func returns true.
- **repartition(numPartitions):** Changes the level of parallelism in this DStream by creating more or fewer partitions.
- **union(otherStream):** Return a new DStream that contains the union of the elements in the source DStream and other DStream.
- **count():** Return a new DStream of single-element RDDs by counting the number of elements in each RDD of the source DStream.
- **reduce(func):** Return a new DStream of single-element RDDs by aggregating the elements in each RDD of the source DStream using a function func (which takes two arguments and returns one). The function should be associative so that it can be computed in parallel.
- **countByKey():** When called on a DStream of elements of type K, return a new DStream of (K, Long) pairs where the value of each key is its frequency in each RDD of the source DStream.
- **join(otherStream, [numTasks]):** When called on two DStreams of (K, V) and (K, W) pairs, return a new DStream of (K, (V, W)) pairs with all pairs of elements for each key.

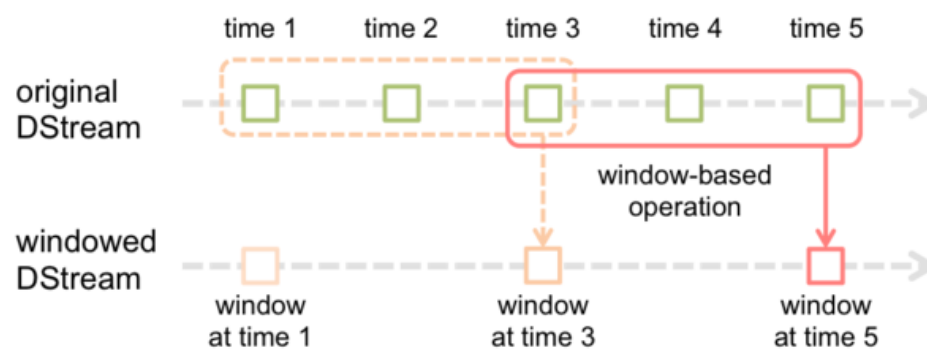
DStream provides a number of output operations such as `saveAsTextFiles` and `pprint`. You can find the complete list in the programming guide.

The figure below is taken from the programming guide and illustrates DStreams and transformations. Each DStream is partitioned into a set of RDD's also called *batches*. When you apply a stateless transformation you will create a new DStream of RDD's with the data that resulted from each execution of the transformation.



The above transformations are all state-less. This means that they create a new DStream and that they do not carry over any information from previous *batches* or computations. Unless specifically managed, there is no state information maintained that can be accessed or updated by subsequent transformations or other operations. Spark Streaming allows you to define a function that makes updates to a new defined state based on the current state and values from an incoming data stream. Such transformations are called *state-ful* transformations. We consider state-ful transformations outside of this introductory lab.

Spark Streaming provides a concept of *windowing*. This means that you can define a window that spans more than one RDD in a DStream. Spark Streaming would combine all the RDD's in the window and allow you to apply a transformation on the data in the window. It can be useful in many types of real-time computations. Lets say you receive new Stock quotes every ten seconds. These quotes are captured from an incoming stream socket as Distinct RDDs. Your objective is to calculate some value based on the last 30-second window but with a maximum delay of 10 seconds. Windowing would allow you to combine the ten-second DStream RDD's into 30-second windows and slide the window . This concept allows you to more easily compute the value you are seeking based on the last 30 seconds of real-time data. The figure below from the programming guide depicts the concept of sliding windows.



Below we provide a sample of operations available on windows.

- **window(windowLength, slideInterval)** : Return a new DStream which is computed based on windowed batches of the source DStream.
- **countByWindow(windowLength, slideInterval)** : Return a sliding window count of elements in the stream.
- **reduceByWindow(func, windowLength, slideInterval)** : Return a new single-element stream, created by aggregating elements in the stream over a sliding interval using func. The function should be associative so that it can be computed correctly in parallel.
- **reduceByKeyAndWindow(func, windowLength, slideInterval, [numTasks])**: When called on a DStream of (K, V) pairs, returns a new DStream of (K, V) pairs where the values for each key are aggregated using the given reduce function func over batches in a sliding window. Note: By default, this uses Spark's default number of parallel tasks (2 for local mode, and in cluster mode the number is determined by the config property spark.default.parallelism) to do the grouping. You can pass an optional numTasks argument to set a different number of tasks.

Finally, we like to discuss check pointing. Lets assume you are doing calculations cross multiple batches using windowing. And lets say that your system crashed after a calculation is done, and the oldest batch used was discarded. This would make it tricky for the system to recover since it

could not re-calculate the results from the last window. If the batches were persisted the application could recalculate the results during a recovery phase by using persisted batches. Check pointing allows you to save batches to enable or speed up recovery.

This problem would be even more exacerbated if we use a state-ful transformation. This transformation may use a compounded result from many batches and possibly windows. To recalculate it would need to access to the now obsoleted batches. Even if we had access to these batches it could take a long time to recalculate the results. In this case you also want to check point the calculated state. Spark streaming provides checkpoint functionality for the purpose of dealing with these situations.

Instructions, Resources, and Prerequisites

To run this lab you can use you own laptop with spark 1.5.x installed or one of the course AMI's. Spark Streaming was released with Spark 1.3.x, but we recommend using latest possible release. This lab has been tested with version 1.5.1.

We will mainly use the `pyspark` shell but also use `spark-submit`. In the below table you can find links to useful and necessary resources.

Resource	What
http://spark.apache.org/docs/latest/streaming-programming-guide.html#transformations-on-dstreams	Spark Streaming Programming Guide.
http://www.w3schools.com/json/json_syntax.asp	Introduction to JSON data format.
http://www.meetup.com/meetup_api/	Description of Meetup data API.
http://www.meetup.com/meetup_api/docs/stream/2/rsvps/#websockets	Meetup websocket API.
http://www.slideshare.net/spark-project/deep-divewithsparkstreaming-tathagatadassparkmeetup20130617	Technical presentation on Spark Streaming.
https://github.com/UC-Berkeley-I-School/w205-labs-exercises/blob/master/docs/Installing%20Cloudera%20Hadoop%20on%20UCB%20W205%20Base%20Image.pdf	Instructions for getting Hadoop set up on AMI.
http://www.eecs.berkeley.edu/Pubs/TechRpts/2012/EECS-2012-259.pdf	Spark Streaming Paper.
https://spark.apache.org/docs/0.9.0/python-programming-guide.html	Python programming guide.

Step-0. Pre-Requirements

Update the repository from git hub. The Lab directory will contain a sub-directory called “somedata”. We will be using this data to illustrate a few things.

First, if an example starts with the “\$” prompt, it is run in the terminal shell. If it starts with the “>>>” prompt, it is run in `pyspark` shell . If code is within a “box” you should consider it a code fragment for illustration purposes and not a command you are asked to execute.

If you are running this lab on an AMI, it is likely you have HDFS running. Refer to the hadoop set up guide to check how to start, stop and test HDFS. . In Step-1 we will be simulating a data

stream by copying data to a directory. If you use HDFS Step-1 access to files will look slightly different

On a laptop or other Linux system without HDFS create a directory using this command:

```
$mkdir /tmp/datastreams
```

If you are running on an AML with HDFS running, create a directory using this command:

```
$sudo -u hdfs hdfs dfs -mkdir /tmp/datastreams
```

Step-1. Getting Started

Start `pyspark` using (at least) 2 cores, and run the necessary import statements.

```
$MASTER=local[2] pyspark
>>>from pyspark import SparkContext
>>>from pyspark.streaming import StreamingContext
>>>sc
```

Pyspark already has a Spark context so we will not need to create one. But we want to make sure that pyspark uses multiple cores for parallel processing, so we provide a `MASTER` environment variable. If you were to implement this as a program you would need to add a statement as show below. You will do this later in this lab when we use `spark-submit`.

```
sc = SparkContext("local[2]", "MyApp")
```

For the time being we are using the interactive shell to make it easier for you to experiment and retry. We will start by doing a very simple streaming operation by converting all incoming words to upper case. We will first show how you do it from a file to simulate a stream. Secondly we will show how to have the Streaming application listen to a socket streaming data.

You will be submitting files to this directory for processing by the application. For the purpose of this lab we previously named it `/tmp/datastreams`. You can use a name of your choosing, but make sure to use your path for the commands below.

Create a streaming context using the following statement.

```
>>>ssc = StreamingContext(sc, 1)
```

Tell the context to read data from files in a directory. This means it will be monitoring the directory for new files and read them as they arrive. All files in the directory must have the same format, and must be created atomically. You can create them by moving or copying them there.

You cannot open a file and incrementally add to it and expect that these updates will be read. If you are using the local file system type the following statement.

```
>>>lines= ssc.textFileStream("file:///tmp/datastreams")
```

If you are using HDFS use the following statement. The difference is in the URI.

```
>>>lines= ssc.textFileStream("hdfs:///tmp/datastreams")
```

Now lets do a simple transformation that converts all the words in a RDD to upper case.

```
>>>uclines = lines.map(lambda word: word.upper())
```

We will output the words by printing the top results.

```
>>>uclines.pprint()
```

Finally we start the process. After this command you should see an output confirming the processing of a new RDD. We have not moved any data to the directory yet, so each result will be empty.

```
>>>ssc.start()
```

Open a separate terminal window. Lets assume you have a simple file called `words` with the following content.

```
hej
kalle
kula
nisse
Hello
coloraDO
DOrado
eldorade
gatorade
```

If you are on Spark only system copy the file to the `/tmp/datastreams` local directory.

```
$cp words /tmp/datastreams/
```

If you are using a system with HDSF running, copy using the following command:

```
$sudo -u hdfs hdfs dfs -put words /tmp/datastreams/
```

Note: If you have Spark Streaming reading from HDFS you may or may not encounter an exception saying something along the lines of “file not found” referring to a file ending with “_COPYING_”. It is because put is not atomic. Put will create a temporary file while copying from the local file system which will later be removed. For the purpose of this lab Spark Streaming recovers and eventually picks up the right file. If you want to fix this issue, copy to other location in HDFS and then use “hdfs -mv” command to move the file into /tmp/datastreams. The mv command is in a atomic operation.

If you look in the streaming window you will see the result of the processing. You will notice that the words were converted to upper case words after the processing. Try copying the file to the directory again. What happens? If you instead copy to a new file you will see a different result.

```
$ cp words /tmp/datastreams/w1
```

You can continue copying to a different file and the will see that the spark process will pick the data up as it arrives in the directory.

```
$ cp words /tmp/datastreams/w2
$ cp words /tmp/datastreams/w3
```

SUBMISSION 1: Provide a screenshot of the output from the Spark Streaming process.

You can stop the Spark Streaming application by typing the following stop command in the pyspark shell.

```
>>>ssc.stop()
```

It is very useful to understand how to process files from during development of your processing logic. This way you can test different logic based on some specific data and which eases debugging.

Step-2. Connecting to Socket and batch time

Now lets do the same transformation but have the process listen to a streaming socket. You will need to restart pyspark. It will not allow you to create a new StreamingContext associated with the same Spark Context. After restart execute the import statements and the creation of the StreamingContext. Remember to provide the MASTER variable so that we use more than one core. On some platforms not doing this may block the streaming process from receiving data and processing at the same time.

```
$MASTER=local[2] pyspark

>>>from pyspark import SparkContext
>>>from pyspark.streaming import StreamingContext
>>>ssc = StreamingContext(sc, 1)
```

Instead of reading from files in a directory we will listen to a socket. The process terminating the socket will be running on you local host and we will use the port number 9999.

```
>>>lines = ssc.socketTextStream("localhost", 9999)
```

Define the transformation and output statement as previously and start the process.

```
>>>uclines= lines.map(lambda word: word.upper())
>>>uclines.pprint()
>>>ssc.start()
```

You will see some errors from Spark since there is no active port to connect to. But the process will continue to try to connect to the port. To create a port to which the streaming process can connect we will use the Unix `nc` command. The name `nc` stands for `netcat`. You can think about it as the Unix `cat` but for pushing data to a socket rather than a file. If you try this Lab on Windows, you can use the Windows `netcat` command in place of `nc`.

In a terminal window type the command below. The `l` option tells `nc` to listen for incoming connections. This is the right behavior as we expect Spark Streaming to connect to the socket. The `k` command tells `nc` to continue listening even if a connection is completed. This way you can have `nc` running and restart you spark streaming application without needing to restart `nc`.

```
$nc -lk 9999
```

In the terminal window were you started `nc`, type some words. What happens on the streaming application side?

The batch duration for our simple streaming application is one second. This is a short time when a human is providing the input. Lets change that to 30 seconds and see what happens. Restart `pyspark` and run the following commands.

```
>>>from pyspark import SparkContext
>>>from pyspark.streaming import StreamingContext
>>>ssc = StreamingContext(sc, 30)
>>>lines = ssc.socketTextStream("localhost", 9999)
>>>uclines= lines.map(lambda word: word.upper())
>>>uclines.pprint()
>>>ssc.start()
```

Type some words in the `nc` terminal window. As you can see, spark is now batching things up in RDD representing 30 seconds of incoming data.

Step-3. Parsing JSON data

We will now parse a more complicated data structure. We will parse data feed from the [Meetup API](#). Specifically we will be looking at RSVP events. Meetup provides several streaming API's. For the purpose of this Lab we will be creating a simple test stream. In real application you may need to implement a Spark Streaming end point using the Custom Receiver framework.

JSON is a format for a human readable representation of data that is transferred between distributed components. The following is a very simple example. An object consists of name-value pairs or sub objects. A colon separates the components in a name value pair. An object is enclosed by curly brackets. A value can also be an array that contains multiple objects.

```
{
  "firstName": "John",
  "lastName": "Smith",    "isAlive": true,    "age": 25,
  "address": {
    "streetAddress": "21 2nd Street",    "city": "New York",
    "state": "NY",    "postalCode": "10021-3100"    },
  "children": [],
  "spouse": null
}
```

The following is an example of an rsvp JSON object. As you can see it has sub-objects such as venue. And values such as visibility.

```
{
  "venue": {
    "venue_name": "Couchbase San Francisco",
    "lon": -122.397354,
    "lat": 37.790005,
    "venue_id": 21741962,
    "visibility": "public",
    "response": "yes",
    "guests": 0,
    "member": {
      "member_id": 8301128,
      "photo": "http://photos3.meetupstatic.com/photos/member/9/0/a/8/thumb_29557032.jpeg",
      "member_name": "Che Hsu"
    },
    "rsvp_id": 1577033972,
    "mtime": 1446516373612,
    "event": {
      "event_name": "Full Stack Development with NoSQL & Node.js or GoLang",
      "event_id": "226259580",
      "time": 1446602400000,
      "event_url": "http://www.meetup.com/The-San-Francisco-Couchbase-Meetup-Group/events/226259580/"
    },
    "group": {
      "group_topics": [
        {
          "urlkey": "java",
          "topic_name": "Java"
        },
        {
          "urlkey": "php",
          "topic_name": "PHP"
        },
        {
          "urlkey": "newtech",
          "topic_name": "New Technology"
        },
        {
          "urlkey": "ria",
          "topic_name": "Rich Internet Applications"
        },
        {
          "urlkey": "mobile-development",
          "topic_name": "Mobile Development"
        },
        {
          "urlkey": "nosql",
          "topic_name": "NoSQL"
        },
        {
          "urlkey": "databasepro",
          "topic_name": "Database Professionals"
        },
        {
          "urlkey": "database-development",
          "topic_name": "Database Development"
        }
      ]
    }
  }
}
```

```
Development"}, {"urlkey": "softwaredev", "topic_name": "Software
Development"}, {"urlkey": "web-development", "topic_name": "Web
Development"}, {"urlkey": "ia", "topic_name": "Information
Architecture"}], "group_city": "San
Francisco", "group_country": "us", "group_id": 1693125, "group_name": "
The San Francisco Couchbase Group", "group_lon": -
122.42, "group_urlname": "The-San-Francisco-Couchbase-Meetup-
Group", "group_state": "CA", "group_lat": 37.75}}
```

Before we start receiving using a stream we want to make sure we can parse the data properly. Initially we just like to extract the venue from each incoming record. You can understand the `rsvp` stream call and response formats by looking at the API definition of [rsvps](#).

Connect to a socket directly using the curl command below. You should see `rsvp` events being printed in your terminal window as they arrive from the stream.

```
$curl -i http://stream.meetup.com/2/rsvps
```

Before we parse the stream, let's read data from the file system and try out a few transformations. In the `somedata` directory you will have a few files called: `meetup.data.1`, `meetup.data.2` and so forth. Each containing 200 previously received and saved `rsvp` events, except the last one that contains 162 `rsvp` events. The combined data in these files represents more than 1000 `rsvp`'s that we retrieved from the meetup api. We will use the by copying them to `/tmp/datastreams` so that a streaming application will pick them up.

The statement below parses the incoming JSON objects. It takes each row (which is a JSON `rsvp` object) and tests if it has a `venue` object. If it does it sends the object to the `json.loads` function for parsing. Once the parsing is done it returns just the `venue` object.

```
slines = lines.flatMap(lambda x: [ j['venue'] for j in
json.loads('['+x+']') if 'venue' in j] )
```

Load the following spark streaming program:

```
$MASTER=local[2] pyspark

>>>from pyspark import SparkContext
>>>from pyspark.streaming import StreamingContext
>>>import json
>>>ssc = StreamingContext(sc, 10)
>>>lines = ssc.textFileStream("file:///tmp/datastreams")
>>>slines = lines.flatMap(lambda x: [ j['venue'] for j in
json.loads('['+x+']') if 'venue' in j] )
>>>cnt=slines.count()
```

```
>>>cnt.pprint()
>>>slines.pprint()
```

Finally, start the process.

```
>>>ssc.start()
```

We are importing a `json` library. We are using this library to parse the json structure for an `rsvp` using the `json.loads` function. We then extract the `venue` element out of the structure and push that into a new RDD. We will print the total amount of `rsvp` events received in the batch, the number of events with a `venue`, and the top venues with `pprint`.

We can test this by just copying the `meetup.data.*` files into the `/tmp/datastreams` directory.

```
$cp meetup.data.1 /tmp/datastreams/
$cp meetup.data.2 /tmp/datastreams/
$cp meetup.data.3 /tmp/datastreams/
```

and so forth...Every time you copy a file it will be detected by Spark Streaming, picked up and included in the current batch. Stop the streaming process by typing:

```
>>>ssc.stop()
```

[SUBMISSION 2: Change the code so that you save the venue components to a text file. Submit you code.](#)

Step-4. Hooking up to a simple stream.

In this step we will analyze incoming streaming data. The way we integrate is not production grade, but it provides is with an easy way of actually analyzing a real stream.

First we will need to change our script so that it reads from a socket.

```
from pyspark import SparkContext
from pyspark.streaming import StreamingContext
import json
ssc = StreamingContext(sc, 10)
lines = ssc.socketTextStream("localhost", 9999)
jslines = lines.flatMap(lambda x: [ j['venue'] for j in
json.loads('['+x+']') if 'venue' in j] )
lcnt=lines.count()
lcnt.pprint()
c=jslines.count()
```

```
c.pprint()
jslines.pprint()
ssc.start()
```

In a separate terminal window run the command below.

```
curl -i http://stream.meetup.com/2/rsvps | nc -lk 9999
```

Note: you may encounter a parsing error when you first start parsing the stream. This is because a websocket stream sends a HTTP header before it send the data. This is obviously not JSON and the Spark Streaming JSON parsing will not like that.

This command uses `curl` to grab the content from the Meetup websocket API and feed it into `nc`. The `nc` process will make the data available on the socket on port 9999 on `localhost`. The behavior of this stream can become artificially bursty. This is due to that the Unix pipe command (`|`) buffer data before forwarding it. To more closely mimic the incoming rate of rsvp's you can type the following on OSX.

```
script -q /dev/null curl -
i http://stream.meetup.com/2/rsvps | nc -lk 9999
```

On Linux it would be the following command.

```
script -q -c/dev/null curl
-i http://stream.meetup.com/2/rsvps | nc -lk 9999
```

You would probably not see that much difference on the spark streaming end since it is batching the incoming events up anyways. But you can see in the logging that the spark streaming process receives a more continuous stream of data. This is because the sending process forwards them directly as they arrive, rather than letting the pipe program buffer them.

Step-5. Running with spark-submit

Until now we have been running the command in a pyspark shell. Clearly one also needs to be able to submit and execute streaming jobs and standalone programs. Put the following commands into a file called `venuecounter.py`. For the purpose of this lab I will assume you put the script in `$HOME`. If the location differs from that, change the code accordingly.

```
from pyspark import SparkContext
from pyspark.streaming import StreamingContext
import json
sc = SparkContext("local[2]", "MyApp")
ssc = StreamingContext(sc, 10)
lines = ssc.socketTextStream("localhost", 9999)
```

```
jslines = lines.flatMap(lambda x: [ j['venue'] for j in
json.loads('['+x+']') if 'venue' in j] )
lcnt=lines.count()
lcnt.pprint()
c=jslines.count()
c.pprint()
jslines.pprint()
ssc.start()
ssc.awaitTermination()
```

Start the “poor mans stream” using either of the following command.

```
$curl -i http://stream.meetup.com/2/rsvps | nc -lk 9999
```

Or if you want to avoid the pipe buffering (OSX command) you can use the following.

```
$script -q /dev/null curl -
i http://stream.meetup.com/2/rsvps | nc -lk 9999
```

Run the script with the submit command:

```
$spark-submit $HOME/venuecounter.py localhost 9999
```

You should see output similar to what is show in the screenshot below.

```
-----
Time: 2015-11-08 10:42:50
-----
15/11/08 10:42:52 WARN BlockManager: Block input-0-1447008172200 replicated to only 0 peer(s) instead of 1 peers
15/11/08 10:42:52 WARN BlockManager: Block input-0-1447008172400 replicated to only 0 peer(s) instead of 1 peers
15/11/08 10:42:53 WARN BlockManager: Block input-0-1447008172800 replicated to only 0 peer(s) instead of 1 peers
15/11/08 10:42:54 WARN BlockManager: Block input-0-1447008174200 replicated to only 0 peer(s) instead of 1 peers
15/11/08 10:42:56 WARN BlockManager: Block input-0-1447008175800 replicated to only 0 peer(s) instead of 1 peers
15/11/08 10:42:56 WARN BlockManager: Block input-0-1447008176200 replicated to only 0 peer(s) instead of 1 peers
15/11/08 10:42:56 WARN BlockManager: Block input-0-1447008176600 replicated to only 0 peer(s) instead of 1 peers
15/11/08 10:42:57 WARN BlockManager: Block input-0-1447008177000 replicated to only 0 peer(s) instead of 1 peers
15/11/08 10:42:57 WARN BlockManager: Block input-0-1447008177200 replicated to only 0 peer(s) instead of 1 peers
15/11/08 10:42:58 WARN BlockManager: Block input-0-1447008177800 replicated to only 0 peer(s) instead of 1 peers
15/11/08 10:42:58 WARN BlockManager: Block input-0-1447008178000 replicated to only 0 peer(s) instead of 1 peers
15/11/08 10:42:58 WARN BlockManager: Block input-0-1447008178200 replicated to only 0 peer(s) instead of 1 peers
15/11/08 10:42:58 WARN BlockManager: Block input-0-1447008178600 replicated to only 0 peer(s) instead of 1 peers
15/11/08 10:42:59 WARN BlockManager: Block input-0-1447008179000 replicated to only 0 peer(s) instead of 1 peers
15/11/08 10:42:59 WARN BlockManager: Block input-0-1447008179400 replicated to only 0 peer(s) instead of 1 peers
15/11/08 10:43:00 WARN BlockManager: Block input-0-1447008180400 replicated to only 0 peer(s) instead of 1 peers
-----
Time: 2015-11-08 10:43:00
-----
603
-----
15/11/08 10:43:01 WARN BlockManager: Block input-0-1447008181000 replicated to only 0 peer(s) instead of 1 peers
-----
Time: 2015-11-08 10:43:00
-----
465
-----
Time: 2015-11-08 10:43:00
-----
{u'lat': 33.391479, u'venue_id': 17582512, u'lon': -111.788322, u'venue_name': u'India Oven'}
{u'lat': 43.251106, u'venue_id': 1556769, u'lon': -79.870781, u'venue_name': u'The Undermount'}
{u'lat': 39.570477, u'venue_id': 23912315, u'lon': 2.637716, u'venue_name': u'Bar Carrera23'}
{u'lat': 29.689041, u'venue_id': 23713095, u'lon': -95.827621, u'venue_name': u'Westheimer Lakes '}
{u'lat': -27.399618, u'venue_id': 24017706, u'lon': 153.440933, u'venue_name': u'Amity point Stradbroke'}
{u'lat': 51.420166, u'venue_id': 22443312, u'lon': -0.987627, u'venue_name': u'Grow @Green Park a ConnectIVT Innovation Hub'}
{u'lat': 37.975342, u'venue_id': 24137559, u'lon': 23.710804, u'venue_name': u'Foundation'}
{u'lat': 21.30485, u'venue_id': 22696332, u'lon': -157.857758, u'venue_name': u'Makua Beach (Pray for Sex Beach)'}
{u'lat': 36.095478, u'venue_id': 7666372, u'lon': -86.80249, u'venue_name': u'JT Moore Middle School'}
{u'lat': 39.570477, u'venue_id': 23912315, u'lon': 2.637716, u'venue_name': u'Bar Carrera23'}
...
-----
```

SUBMISSION 3: Provide a screenshot showing the running Spark Streaming application.

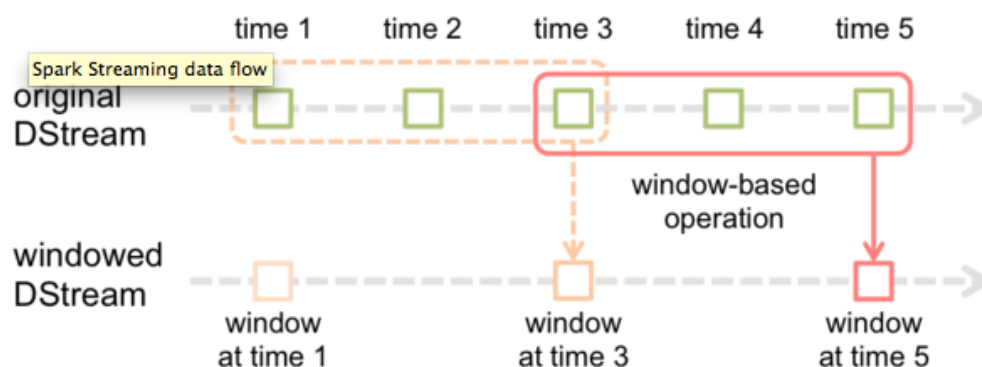
Step-6. Sliding Window

For the purpose of illustrating the concept we will now create a simple sliding window. With this application we like to know how many rsvp's we received in the last batch (10 seconds); the aggregated number of rsvp's the last 30 seconds and the number of rsvp's with a `venue` object in the last batch. To achieve this we will use windowing.

Windowing allows us to perform certain transformations for a sliding window over a number of batches. Our batch size is 10 seconds. We will be calling the DStream transform `countByWindow(windowLength, slideInterval)`. We want to perform the count calculation of last 30 seconds so the `windowLength` argument will be 30. We want to perform the calculation every 10 seconds, so the `slideInterval` will be set to 10.

```
wlcnt=lines.countByWindow(30,10)
```

Using the below illustration each green box represents a 10 second batch, the red box represents the sliding window of 30 seconds.



In order to enable windowing we will need to turn on checkpointing. This is done by defining a checkpoint directory. I will assume you have created the checkpoint directory in `/tmp/checkpointing`, but you can change that as you prefer. The statement looks as follows.

```
ssc.checkpoint("file:///tmp/checkpointing")
```

You could of course choose to store the check pointing on HDFS by changing the statement to

```
ssc.checkpoint("hdfs:///tmp/checkpointing")
```

The final program looks as follows, save this in `$HOME/venuecounter2.py`:

```
from pyspark import SparkContext
from pyspark.streaming import StreamingContext
import json
sc = SparkContext("local[2]", "MyApp")
ssc = StreamingContext(sc, 10)
ssc.checkpoint("file:///tmp/checkpointing")
lines = ssc.socketTextStream("localhost", 9999)
wlcnt=lines.countByWindow(30,10)
jslines = lines.flatMap(lambda x: [ j['venue'] for j in
json.loads('['+x+']') if 'venue' in j] )
lcnt=lines.count()
wlcnt.pprint()
lcnt.pprint()
c=jslines.count()
c.pprint()
jslines.pprint()
```

```
ssc.start()  
ssc.awaitTermination()
```

Now start the stream using the command below.

```
$curl -i http://stream.meetup.com/2/rsvps | nc -lk 9999
```

Or if you want to avoid the pipe buffering (OSX command) you can use the following.

```
$script -q /dev/null curl -  
i http://stream.meetup.com/2/rsvps | nc -lk 9999
```

Then start the stream processing application using the command:

```
$spark-submit $HOME/venuecounter2.py localhost 9999
```

SUBMISSION 4a: Provide a screenshot of the running Spark Streaming application that shows that the CountByWindow indeed provides an sum of the counts from the 3 latest batches. See example screenshot below.

SUBMISSION 4b: Also explain what the difference is between having 10 sec batches with a 30 sec sliding window and a 30 second batch length.

Troubleshooting

Hadoop and Spark on AMI

If you are running on a course AMI, make sure you have hadoop and spark set up properly. You can find instructions here: <https://github.com/UC-Berkeley-I-School/w205-labs-exercises/blob/master/docs/Installing%20Cloudera%20Hadoop%20on%20UCB%20W205%20Base%20Image.pdf>

If you like to install a new Spark 1.5 here are some instructions.

- Download Spark 1.5 from <http://spark.apache.org/downloads.html>

```
$gunzip spark-1.5.1-bin-hadoop2.6.tgz
$tar xvf spark-1.5.1-bin-hadoop2.6.tar
$ln -s ./spark-1.5.1-bin-hadoop2.6 /root/spark15
$export SPARK_HOME=/root/spark15
$export PATH=$SPARK_HOME/bin:$PATH
$# which pyspark
$/root/spark15/bin/pyspark
```