

Adam Reilly  
W205  
Lab 10

#### SUBMISSION 1

```
from pyspark import SparkContext
from pyspark.streaming import StreamingContext
ssc = StreamingContext(sc, 1)
lines= ssc.textFileStream("file:///tmp/datastreams")
uclines = lines.map(lambda word: word.upper())
overfive = uclines.filter(lambda word: len(word) > 5)
overfive.pprint()
ssc.start()
```

#### SUBMISSION 2.

```
$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()
prefix= ("file:///tmp/venues.txt")
slines.saveAsTextFiles(prefix)
ssc.start()
```

#### SUBMISSION 3

1. In this scenario, I would increase the batch processing time so that it would be better able to fit the flow of data.
2. One side effect of needing to update a streaming program is that you likely won't be able to access the data that is being produced while the program updates. You could handle this by writing an identical program before updating and run the data through that program (and you could even write the data to an outside file in case you would like to reimport it back into the updated program).

#### SUBMISSION 4a

```
-----
Time: 2017-04-17 02:43:10
-----
195
17/04/17 02:43:12 WARN BlockManager: Block input-0-1492396992000 replicated to only 0 peer(s) instead of 1 peers
-----
Time: 2017-04-17 02:43:10
-----
11
-----
Time: 2017-04-17 02:43:10
-----
6
-----
17/04/17 02:43:13 WARN BlockManager: Block input-0-1492396993400 replicated to only 0 peer(s) instead of 1 peers
-----
Time: 2017-04-17 02:43:10
-----
{u'lat': 47.666443000000001, u'venue_id': 22278342, u'lon': -122.371201, u'venue_name': u'Stoup Brewing'}
{u'lat': -32.034775000000003, u'venue_id': 25130966, u'lon': 115.753136, u'venue_name': u'Short Black Sheep'}
{u'lat': 38.694575999999998, u'venue_id': 25051820, u'lon': -77.299453999999997, u'venue_name': u'Kyle's House'}
{u'lat': 53.349803999999999, u'venue_id': 25033253, u'lon': -6.2603099999999996, u'venue_name': u'Dublin International airport'}
{u'lat': 43.644646000000002, u'venue_id': 24284662, u'lon': -79.394997000000004, u'venue_name': u'Lighthouse Labs'}
{u'lat': 43.768267999999999, u'venue_id': 1315192, u'lon': -79.412530000000004, u'venue_name': u'Toronto Dance Salsa '}
```

SUBMISSION 4B. The difference between a 10 sec batch with a 30 second sliding window and a 30 second batch is with the first option, the sorting is done in advance. Additionally, the sliding window can compare more data because the 30 second batch can only cover items added since the last batch was run, while the window has more flexibility.