# Homework – 3

## Fall 2024

## Spark Streaming

## Sahithya Arveti Nagaraju (Person number: 50559752, UBITName: sarvetin)

**Implement k-mer count program in Spark Streaming**

**• Create a local StreamingContext with two execution threads and a batch interval of 10 seconds.**

**Answer:**

Below is the code for creating StreamingContext with two execution threads and batch interval of 10 seconds:

```python
# Creating Spark context with 2 execution threads and naming App as "KMerCountApp"
conf=SparkConf().setAppName("KMerCountApp").setMaster("local[2]")
sc=SparkContext(conf=conf)
# Setting 10 second batch interval
ssc=StreamingContext(sc,10)
```

- SparkConf().setAppName("KMerCountApp") sets the name of the Spark application as 'KMerCountApp'.
- .setMaster("local[2]") sets the Spark context to run locally with 2 execution threads.
- StreamingContext(sc,10) initializes a StreamingContext with a batch interval of 10 seconds, which indicates that the system will handle and process the incoming data in chunks of 10 seconds.

DAG Visualization:

The above DAG shows a Spark Streaming job that continuously reads text data from a socket stream, processes it in batches, and writes the results to a destination. The windowing operation allows for real-time processing of the data stream, and the partitioning and mapping operations enable efficient parallel processing.

| 9 | Streaming job from [output operation 0, batch time 18:28:50]<br>call at C:\Users\siria\AppData\Local\Programs\Python\Python39\lib\site-packages\py4j\clientserver.py:617 | 2024/12/01 18:28:52 | 16 ms | 1/1 | 1/1 |
| 8 | Streaming job from [output operation 0, batch time 18:28:50]<br>call at C:\Users\siria\AppData\Local\Programs\Python\Python39\lib\site-packages\py4j\clientserver.py:617 | 2024/12/01 18:28:51 | 0.6 s | 1/1 (1 skipped) | 1/1 (1 skipped) |
| 7 | Streaming job from [output operation 0, batch time 18:28:50]<br>call at C:\Users\siria\AppData\Local\Programs\Python\Python39\lib\site-packages\py4j\clientserver.py:617 | 2024/12/01 18:28:50 | 2 s | 2/2 | 2/2 |
| 6 | Streaming job from [output operation 0, batch time 18:28:40]<br>call at C:\Users\siria\AppData\Local\Programs\Python\Python39\lib\site-packages\py4j\clientserver.py:617 | 2024/12/01 18:28:42 | 20 ms | 1/1 | 1/1 |
| 5 | Streaming job from [output operation 0, batch time 18:28:40]<br>call at C:\Users\siria\AppData\Local\Programs\Python\Python39\lib\site-packages\py4j\clientserver.py:617 | 2024/12/01 18:28:41 | 0.7 s | 1/1 (1 skipped) | 1/1 (1 skipped) |
| 4 | Streaming job from [output operation 0, batch time 18:28:40]<br>call at C:\Users\siria\AppData\Local\Programs\Python\Python39\lib\site-packages\py4j\clientserver.py:617 | 2024/12/01 18:28:40 | 2 s | 2/2 | 2/2 |
| 3 | Streaming job from [output operation 0, batch time 18:28:30]<br>call at C:\Users\siria\AppData\Local\Programs\Python\Python39\lib\site-packages\py4j\clientserver.py:617 | 2024/12/01 18:28:33 | 19 ms | 1/1 | 1/1 |
| 2 | Streaming job from [output operation 0, batch time 18:28:30]<br>call at C:\Users\siria\AppData\Local\Programs\Python\Python39\lib\site-packages\py4j\clientserver.py:617 | 2024/12/01 18:28:32 | 0.6 s | 1/1 (1 skipped) | 1/1 (1 skipped) |
| 1 | Streaming job from [output operation 0, batch time 18:28:30]<br>call at C:\Users\siria\AppData\Local\Programs\Python\Python39\lib\site-packages\py4j\clientserver.py:617 | 2024/12/01 18:28:30 | 2 s | 2/2 | 2/2 |

Above Spark UI shows that the each batch is processed every 10 seconds.

**• Create a DStream that represents streaming data from a TCP source (localhost:9999) [hint: use Netcat]**

**Answer:**

The DStream is created by connecting to a TCP source on localhost on port 9999. This is done by the following line in the code:
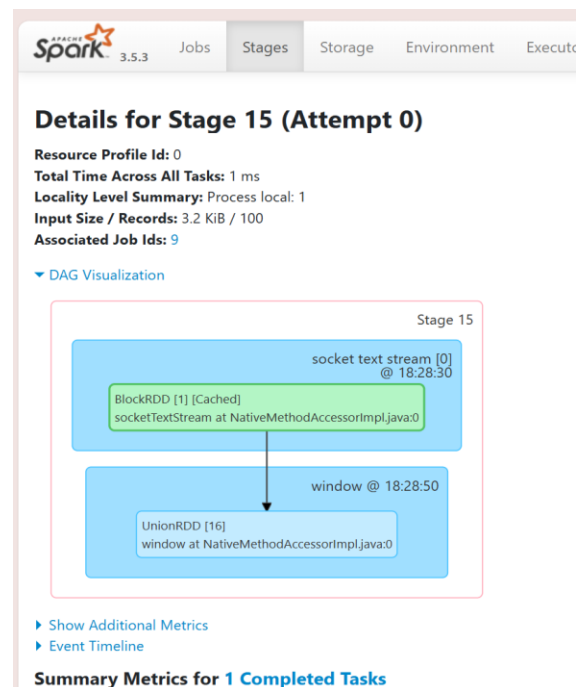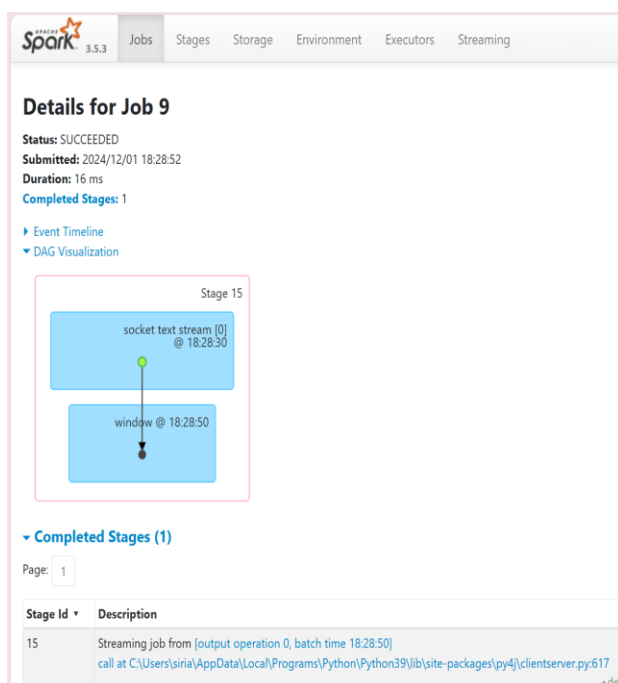
```
# DStream for recieving data from localhost (port 9999)
text_lines=ssc.socketTextStream("localhost",9999)
```

ssc.socketTextStream("localhost", 9999) sets up a DStream that listens for text data from a TCP socket on localhost at port 9999. The data from sentences.txt is streamed into this port using Netcat.

Netcat command for data streaming of sentences.txt into port 9999:

```
C:\Users\siria\DIC_Homework3>ncat -lk 9999 <sentences.txt
```

DAG Visualization:



Above DAG Visualization shows that a socket text stream (Dstream) created to stream data from TCP source (localhost:9999) with a window of 30 seconds with a sliding interval of 10 seconds.

**• Generate k-mers of length 3 from each line of text.**

**Answer:**

The create_kmers function in the code is responsible for generating k-mers of length 3 from each line of text:

```
5  # Creating k-mers of length 3
6  def create_kmers(text):
7      k = 3
8      return list(map(lambda i:text[i:i+k], range(len(text)- k+1)))
9
```

The function create_kmers takes a string text as input and generates a list of substrings (k-mers) of length 3. It uses the map function to iterate through all possible starting indices i of length-3 substrings within the text. The range range(len(text) - k + 1) ensures that the function doesn't attempt to create a k-mer starting too close to the end of the string. It generates a list of all the k-mers of length 3 which are possible in the text provided as shown in below screenshot:



Also, saved all the k-mers for each line in a text file named output_kmers.txt for reference.

**• Count the occurrences of each k-mer and print the k-mer counts. Take screenshots of your output.**

**Answer:**

The k-mer frequencies are counted by using the reduceByKey operation in the process_kmers function:

```
15      # Generating k-mers of length 3 for each line in the RDD and count k-mers
16      kmer_list = rdd.flatMap(lambda line:create_kmers(line))
17      kmer_count= kmer_list.map(lambda kmer:(kmer,1)).reduceByKey(lambda x,y:x+y)
```

- rdd.flatMap(lambda line: create_kmers(line)): Transforms every line of the RDD to a list of k-mers. The flatMap function will flatten all the lists of k-mers into one single RDD.

- reduceByKey(lambda x, y: x + y): reduces the RDD by key - i.e. k-mer - and sums up values, so counts a frequency of every k-mer.

The first top 10 k-mer counts are printed using the following code:

```
26      # Sorting the k-mers by frequency and taking the top 10 k-mers
27      top_kmers = kmer_count.takeOrdered(10, key=lambda x: -x[1])
28
29      print("\nTop 10 K-mer counts:")
30      for kmer, freq in top_kmers:
31          print(f"Count of '{kmer}': {freq}")
```

Explanation:

The kmer_count.collect() method collects all the k-mer counts as tuples of the form (kmer, frequency). A for loop iterates over these tuples, printing the first 10 k-mers and their frequencies.

```
Processing batch at: 2024-12-01 23:09:20

Top 10 K-mer counts:
Count of 'eha': 4
Count of 'xhn': 4
Count of 'byb': 3
Count of 'kqx': 3
Count of 'cjy': 3
Count of 'mkl': 3
Count of 'lpb': 3
Count of 'kfn': 3
Count of 'bhk': 3
Count of 'ykp': 2
```

```
Processing batch at: 2024-12-01 23:09:30

Top 10 K-mer counts:
Count of 'eha': 4
Count of 'xhn': 4
Count of 'byb': 3
Count of 'kqx': 3
Count of 'cjy': 3
Count of 'mkl': 3
Count of 'lpb': 3
Count of 'kfn': 3
Count of 'bhk': 3
Count of 'ykp': 2
```

```
Processing batch at: 2024-12-01 23:09:40

Top 10 K-mer counts:
Count of 'eha': 4
Count of 'xhn': 4
Count of 'byb': 3
Count of 'kqx': 3
Count of 'cjy': 3
Count of 'mkl': 3
Count of 'lpb': 3
Count of 'kfn': 3
Count of 'bhk': 3
Count of 'ykp': 2
```

The above is the output of kmer counts, here I'm only displaying only top 10 kmer counts since there are a lot of kmers. Additionally, it can be observed that Dstream is performing kmer counting for every 10 seconds interval.

Also, saved all the k-mers count in a text file named output_kmer_count.txt for reference.

File    Edit    View    Settings    Help

```
 1  mqn: 1
 2  qns: 1
 3  nsl: 1
 4  slv: 1
 5  lvv: 1
 6  vva: 1
 7  vay: 1
 8  ayk: 1
 9  ykp: 2
10  kpm: 1
11  pmw: 2
12  mwt: 1
13  wtv: 1
14  tvf: 1
15  vfi: 1
16  fic: 1
17  icf: 1
18  cft: 1
19  fty: 1
20  tyk: 2
21  ykn: 1
22  knj: 1
23  njb: 1
24  jbw: 1
25  bwd: 1
26  wdz: 1
27  dzk: 1
28  zkl: 1
29  nkp: 1
30  kpz: 2
31  pzo: 1
32  zov: 1
33  ova: 1
34  vac: 1
35  acr: 2
36  crm: 1
37  nmb: 1
```

## Alternative method:

The above method using netcat was getting whole data from sentences.txt at first time interval itself, so it was performing kmer operation at first time interval and was not getting any data streamed after that from netcat.

For alternative method, first run tcp.py file in one terminal like below screenshot and then run code.py in another terminal

So, I have used an alternative method for sending data using TCP socket program as shown below to stream data in localhost 9999,



Jupyter tcp.py Last Checkpoint: 51 minutes ago

File   Edit   View   Settings   Help

```python
1  import socket
2  import time
3
4  # Create a TCP/IP socket
5  server_socket = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
6  server_socket.bind(('localhost', 9999))
7  server_socket.listen(1)
8
9  print("Server is listening on port 9999...")
10
11 while True:
12     connection, client_address = server_socket.accept()
13     try:
14         print(f"Connection from {client_address}")
15         # Continuously send lines of data
16         with open('sentences.txt', 'r') as file:
17             for line in file:
18                 connection.sendall(line.encode())
19                 time.sleep(1)  # Simulate streaming with a delay
20     finally:
21         connection.close()
22
```

By using above method, data was sending continuously and able to get the kmer counts in batchs. Here are the few screenshots of kmer counts after using this method,

```
Processing batch at: 2024-12-01 22:42:50
24/12/01 22:42:51 WARN RandomBlockReplicationPolicy: Expecting 1 replicas with only 0 peer/s.
24/12/01 22:42:51 WARN BlockManager: Block input-0-1733110971000 replicated to only 0 peer(s) instead of 1 peers
24/12/01 22:42:52 WARN RandomBlockReplicationPolicy: Expecting 1 replicas with only 0 peer/s.
24/12/01 22:42:52 WARN BlockManager: Block input-0-1733110972000 replicated to only 0 peer(s) instead of 1 peers
24/12/01 22:42:53 WARN RandomBlockReplicationPolicy: Expecting 1 replicas with only 0 peer/s.
24/12/01 22:42:53 WARN BlockManager: Block input-0-1733110973000 replicated to only 0 peer(s) instead of 1 peers
24/12/01 22:42:54 WARN RandomBlockReplicationPolicy: Expecting 1 replicas with only 0 peer/s.
24/12/01 22:42:54 WARN BlockManager: Block input-0-1733110974000 replicated to only 0 peer(s) instead of 1 peers
24/12/01 22:42:55 WARN RandomBlockReplicationPolicy: Expecting 1 replicas with only 0 peer/s.
24/12/01 22:42:55 WARN BlockManager: Block input-0-1733110975000 replicated to only 0 peer(s) instead of 1 peers
24/12/01 22:42:56 WARN RandomBlockReplicationPolicy: Expecting 1 replicas with only 0 peer/s.
24/12/01 22:42:56 WARN BlockManager: Block input-0-1733110976000 replicated to only 0 peer(s) instead of 1 peers

Top 10 K-mer counts:
Count of 'lvv': 1
Count of 'ayk': 1
Count of 'tvf': 1
Count of 'zkl': 1
Count of 'nkp': 1
Count of 'rmb': 1
Count of 'bit': 1
Count of 'vkp': 1
Count of 'yjb': 1
Count of 'zly': 1
```

```
Processing batch at: 2024-12-01 22:43:00
```

```
Top 10 K-mer counts:
Count of 'cst': 2
Count of 'irw': 2
Count of 'ykp': 2
Count of 'kpx': 2
Count of 'ayk': 1
Count of 'tvf': 1
Count of 'nkp': 1
Count of 'rmb': 1
Count of 'bit': 1
Count of 'yjb': 1
```

```
Processing batch at: 2024-12-01 22:50:51
```

```
Top 10 K-mer counts:
Count of 'xhp': 2
Count of 'bhk': 2
Count of 'sjq': 2
Count of 'amw': 2
Count of 'scs': 2
Count of 'lqy': 2
Count of 'jcd': 2
Count of 'hno': 2
Count of 'dlq': 2
Count of 'amk': 2
```

```
Processing batch at: 2024-12-01 22:51:41
```

```
Top 10 K-mer counts:
Count of 'kpx': 2
Count of 'lqy': 2
Count of 'ogw': 2
Count of 'ykp': 2
Count of 'kna': 2
Count of 'cxb': 2
Count of 'hno': 2
Count of 'znh': 2
Count of 'scs': 2
Count of 'kqx': 2
```

```
Processing batch at: 2024-12-01 22:52:32
```

```
Top 10 K-mer counts:
Count of 'cst': 2
Count of 'irw': 2
Count of 'yjh': 2
Count of 'byb': 2
Count of 'bvn': 2
Count of 'nat': 2
Count of 'njj': 2
Count of 'jbx': 2
Count of 'cxb': 2
Count of 'yab': 2
24/12/01 22:53:07 WARN Random
```

```
Processing batch at: 2024-12-01 22:53:25
```

```
Top 10 K-mer counts:
Count of 'yda': 2
Count of 'fyd': 2
Count of 'kzd': 2
Count of 'amj': 2
Count of 'bvn': 2
Count of 'nat': 2
Count of 'cjy': 2
Count of 'njj': 2
Count of 'jbx': 2
Count of 'tzs': 2
```

```
Processing batch at: 2024-12-01 22:54:20
```

```
Top 10 K-mer counts:
Count of 'wwf': 2
Count of 'byb': 2
Count of 'amj': 2
Count of 'cjy': 2
Count of 'mkl': 2
Count of 'tzs': 2
Count of 'qwh': 2
Count of 'qcl': 2
Count of 'tlr': 2
Count of 'wyu': 2
```