

Principles of Computer Science II

Large Scale Computation with Apache Spark

Ioannis Chatzigiannakis

Sapienza University of Rome

Lecture 18



Profile most-frequent k-mer

```
CGGGGCTGGGTCGTCACATTCCCCTTTTCGATA
TTTGAGGGTGCCCAATAACCAAAGCGGACAAA
GGGATGCCGTTTGACGACCTAAATCAACGGCC
AAGGCCAGGAGCGCCTTTGCTGGTTCTACCTG
AATTTTCTAAAAAGATTATAATGTCGGTCCCTC
CTGCTGTACAACTGAGATCATGCTGCTTCAAC
TACATGATCTTTTGTGGATGAGGGAATGATGC
```

- ▶ Seven (7) 32-nucleotide DNA sequences
- ▶ A “secret” pattern $P=ATGCAACT$ of length $l=8$ implanted.
- ▶ Can you reconstruct P by analyzing the DNA sequences?



An example

```
CGGGGCTATGCAACTGGGTCGTCACATTCCCCTTTTCGATA
TTTGAGGGTGCCCAATAAATGCAACTCCAAAGCGGACAAA
GGATGCAACTGATGCCGTTTGACGACCTAAATCAACGGCC
AAGGATGCAACTCCAGGAGCGCCTTTGCTGGTTCTACCTG
AATTTTCTAAAAAGATTATAATGTCGGTCCATGCAACTTC
CTGCTGTACAACTGAGATCATGCTGCATGCAACTTTCAAC
TACATGATCTTTTGATGCAACTTGGATGAGGGAATGATGC
```

- ▶ The same DNA sequences with the implanted pattern $ATGCAACT$
- ▶ Can you spot the locations of the implanted pattern?



An example

```
CGGGGCTATGCAACTGGGTCGTCACATTCCCCTTTTCGATA
TTTGAGGGTGCCCAATAAATGCAACTCCAAAGCGGACAAA
GGATGCAACTGATGCCGTTTGACGACCTAAATCAACGGCC
AAGGATGCAACTCCAGGAGCGCCTTTGCTGGTTCTACCTG
AATTTTCTAAAAAGATTATAATGTCGGTCCATGCAACTTC
CTGCTGTACAACTGAGATCATGCTGCATGCAACTTTCAAC
TACATGATCTTTTGATGCAACTTGGATGAGGGAATGATGC
```

- ▶ Same as before but showing the implant locations.
- ▶ Devise an MapReduce algorithm to automatically identify the implanted pattern
- ▶ Length l is known.
- ▶ Sequence: <https://goo.gl/xN7WvE>



Profile most-frequent k-mer

```
1 from pyspark import SparkContext, SparkConf
2 from operator import add
3
4 conf = SparkConf().setAppName(" Profile").setMaster(" local")
5 sc = SparkContext(conf=conf)
6
7 val raw_data = sc.wholeTextFiles(" pattern.txt")
8
9 def splitLine(line):
10     pairs = []
11     if len(line) > 1:
12         for symbol in range(0, len(line)-8):
13             pairs.append((line[symbol:symbol+8], 1))
14
15     return pairs
16
17 pairs = raw_data.flatMap(splitLine) \
18                 .reduceByKey(add) \
19                 .sortBy(lambda a: -a[1])
20 print(pairs.take(5))
```



Profile most-Frequent first appearing k-mer

```
CGGGGCTGGGTCGTCACATTCCCCTTTTCGATA
TTTGAGGGTGCCCAATAACCAAAGCGGACAAA
GGGATGCCGTTTGACGACCTAAATCAACGGCC
AAGGCCAGGAGCGCCTTTGCTGGTTCTACCTG
AATTTTCTAAAAAGATTATAATGTCGGTCCTC
CTGCTGTACAACCTGAGATCATGCTGCTTCAAC
TACATGATCTTTTGTGGATGAGGGAATGATGC
```

- ▶ Identify most-frequent first-appearing k-mer
- ▶ In first line, 3-mer CGG appears before GGC, TGG, GTC, ...



Profile most-frequent k-mer

```
1 from pyspark import SparkContext, SparkConf
2 from operator import add
3
4 conf = SparkConf().setAppName(" Profile").setMaster(" local")
5 sc = SparkContext(conf=conf)
6
7 raw_data = sc.textFile(" yeast_chr1.txt")
8
9 def splitLine(line):
10     pairs = []
11     for symbol in range(0, len(line)-6, 3):
12         for second in range(symbol+3, len(line)-6, 3):
13             pairs.append(((line[symbol:symbol+3],
14                             line[second:second+3]), 1))
15
16     return pairs
17
18 pairs = raw_data.flatMap(splitLine) \
19                 .reduceByKey(add) \
20                 .sortBy(lambda a: -a[1])
21 print(pairs.take(10))
```

