



DANMARKS TEKNISKE UNIVERSITET

01415

Computational Tools for Big Data

Week 67 - Exercise

Anonymous

October 2016

order.orderID	product.productID	product.productName
11011	58	Escargots de Bourgogne
11011	71	Flotemysost
10952	6	Grandma's Boysenberry Spread
10952	28	Rössle Sauerkraut
10835	59	Raclette Courdavault
10835	77	Original Frankfurter grüne Soße
10692	63	Vegie-spread
10702	3	Aniseed Syrup
10702	76	Lakkalikööri
10643	39	Chartreuse verte
10643	46	Spegesild
10643	28	Rössle Sauerkraut

Exercise 5.3

Assignment :

The customer with customerID ALFKI has made a number of orders containing some products. Return orders made by ALFKI that contain at least 2 products. Also return the products.

Solution :

We selected relation between customers-orders and filtered entries from customer 'ALFKI'. Then we counted number of products in each order and selected just those which contain at least 2 products.

```

1 MATCH (customer:Customer)-[:PURCHASED]->(order:Order)
2 WHERE customer.customerID = "ALFKI"
3 MATCH (order)-[r:ORDERS]->(:Product)
4 WITH order, count(r) AS count
5 WHERE count => 2
6 MATCH (order)-[r:ORDERS]->(product:Product)
7 RETURN order.orderID, product.productID, product.productName, count;
```

Output:

order.orderID	product.productID	product.productName	count
10835	59	Raclette Courdavault	2
10835	77	Original Frankfurter grüne Soße	2
10643	39	Chartreuse verte	3
10643	46	Spegesild	3
10643	28	Rössle Sauerkraut	3
11011	58	Escargots de Bourgogne	2
11011	71	Flotemysost	2
10702	3	Aniseed Syrup	2
10702	76	Lakkalikööri	2
10952	6	Grandma's Boysenberry Spread	2
10952	28	Rössle Sauerkraut	2

Exercise 5.4

Assignment :

Determine how many and who has ordered "Uncle Bob's Organic Dried Pears" (productID 7).

Solution :

We selected relation between product-order-customer and then removed all entries which do not contain productID = 7. We sum quantity and return customers ids.

```

1 MATCH (product:Product)<-[r:ORDERS]-(order:Order)<-[[:PURCHASED]]-(customer:Customer)
2 WHERE product.productID = "7"
3 RETURN customer.customerID, SUM(r.quantity)
```

Output:
20 Customers

customer.customerID	SUM(r.quantity)
VAFFE	10
BOTTM	20
BONAP	58
SPLIR	10
GOURL	3
OCEAN	6
ERNSH	18
LILAS	15
SAVEA	45
QUICK	135
FOLKO	60
OTTIK	55
RATTC	76
LACOR	5
BSBEV	34
EASTC	100
SANTG	12
VICTE	26
FOLIG	35
REGGC	40

Exercise 5.5

Assignment :

How many different and which products have been ordered by customers who have also ordered “Uncle Bob’s Organic Dried Pears”?

Solution :

Similar to solution from Exercise 5.4, we selected customers which ordered productID = 7 and then on these customers we match all products orders they made and returned just distinct names.

```

1 MATCH (prod:Product)<-[r:ORDERS]-(Order)<-[:PURCHASED]-(customer:Customer)
2 WHERE prod.productID = "7"
3 MATCH (customer)-[:PURCHASED]->(ord:Order) -[:ORDERS] ->(product:Product)
4 RETURN DISTINCT(product.productName)

```

Output:

76 Products

Grandma’s Boysenberry Spread, Chef Anton’s Cajun Seasoning, Northwoods Cranberry Sauce, Uncle Bob’s Organic Dried Pears, Aniseed Syrup, Chang, Pavlova, Sir Rodney’s Marmalade , Tunnbröd , Mascarpone Fabioli , Ikura, Queso Manchego La Pastora, Konbu, Tofu , Pâté chinois , Camembert Pierrot, Wimmers gute Semmelknödel, Louisiana Hot Spiced Okra, Chartreuse verte , Jack’s New England Clam Chowder, Spegesild, Filo Mix , Original Frankfurter grüne Soße, Rhönbräu Klosterbier , Röd Kaviar , Guaraná Fantástica , Tarte au sucre , Queso Cabrales , Côte de Blaye, Gnocchi di nonna Alice , Alice Mutton , NuNuCa Nuß-Nougat-Creme, Schoggi Schokolade , Flotemysost, Gorgonzola Telino, Lakkalikööri , Perth Pasties, Raclette Courdavault , Louisiana Fiery Hot Pepper Sauce , Nord-Ost Matjeshering, Escargots de Bourgogne , Chai , Ipoh Coffee, Mozzarella di Giovanni , Chef Anton’s Gumbo Mix , Tourtière, Ravioli Angelo , Outback Lager, Maxilaku , Rössle Sauerkraut, Sir Rodney’s Scones, Gumbär Gummiäarchen, Longlife Tofu, Gudbrandsdalsost , Teatime Chocolate Biscuits , Genen Shouyu , Geitost, Steeleye Stout , Thüringer Rostbratwurst, Gustaf’s Knäckebröd, Manjimup Dried Apples, Vegie-spread , Zaanse koeken, Sirop d’érable , Gula Malacca , Scottish Longbreads, Singaporean Hokkien Fried Mee, Carnarvon Tigers , Boston Crab Meat , Rogede sild, Mishi Kobe Niku, Inlagd Sill, Valkoinen suklaa , Sasquatch Ale, Chocolate, Gravad lax

2 Week 7

Exercise 7.1

Assignment :

Define and implement a MapReduce job to count the occurrences of each word in a text file. Document that it works with a small example.

Solution :

We have implemented this function with MrJob using one single step, consisting of a mapper, a combiner and a reducer. Each line of the file goes to a mapper that divides the line into words and yields a (key,value) pair containing ("word",1). The combiners will receive a set of these tuples and for each key, they will add up their values. The reducer performs the same task, but in this case it receives all the (key,value) pair, thus finishing the job.

Then we output the ("word",N) pairs as values, with no key, so that the next job will receive them all under the same key. We have improved a little the system and created a next step in which a reducer orders the words in descending order according to the number of occurrences of the word.

The next image shows the partial output of the first step. The file used as input is the python code file itself.

```
montoya@montoya:~/Desktop/DTU Lec/1st Semester/5. Computational Tools for Big Data/2.Assingments/67$ python ./main_BoW.py ./main_BoW.py
No configs found; falling back on auto-configuration
Creating temp directory /tmp/main_BoW.montoya.20161014.175740.231398
Running step 1 of 1...
Streaming final output from /tmp/main_BoW.montoya.20161014.175740.231398/output.
..
null      [14, "#"]
null      [1, "'__main__':"]
null      [2, "(key,value)"]
null      [1, "(sum(values),"]
null      [1, "/desktop/dtu"]
null      [1, "1"]
null      [2, "="]
null      [1, "==" ]
null      [1, "[" ]
null      [1, "]" ]
null      [2, ", " ]
null      [1, "- name -"]
```

Figure 2: Step 1 Word Counter

The next image shows the partial output of the second step, with the ordered words. We can see that each word is a tuple with the number of occurrences and the the word itself.

```

montoya@montoya:~/Desktop/DTU Lec/1st Semester/5. Computational Tools for Big Data/2.Assingments/67$ python ./main_BoW.py ./main_BoW.py
No configs found; falling back on auto-configuration
Creating temp directory /tmp/main_BoW.montoya.20161014.180124.418347
Running step 1 of 2...
Running step 2 of 2...
Streaming final output from /tmp/main_BoW.montoya.20161014.180124.418347/output.
..
null    [[14, "#"], [9, "the"], [6, "def"], [6, "yield"], [5, "will"], [4, "key",
], [4, "of"], [4, "words"], [3, "all"], [3, "and"], [3, "from"], [3, "pairs"],
[3, "reducer"], [3, "with"], [2, "(key,value)"], [2, "="], [2, "_"], [2, "combi
ner"], [2, "each"], [2, "for"], [2, "get"], [2, "import"], [2, "in"], [2, "it"],
[2, "joining"], [2, "none,"], [2, "ordered"], [2, "print"], [2, "return"], [2,
"sum(values)"], [2, "this"], [2, "values):"], [1, "'__main__':"], [1, "(sum(valu
es),"), [1, "/desktop/dtu"], [1, "1"], [1, "!="], [1, "["], [1, ""]], [1, "__nam
e__"], [1, "a"], [1, "at"], [1, "can."], [1, "class"], [1, "combiner=self.combin
er_wc,"], [1, "combiner_wc(self,"], [1, "count."], [1, "counted"], [1, "create"]
, [1, "do"], [1, "every"], [1, "frequency"], [1, "get."], [1, "getKey(item):"],

```

Figure 3: Step 2 Word Counter

2.1 main_BoW.py

```

1  from mrjob.job import MRJob
2  from mrjob.step import MRStep
3
4  class MRWordWC(MRJob):
5
6      def mapper_WC(self, _, line):
7          # Each mapper will create the (key,value) pairs of words
8          # for every line they get.
9          words = line.split()
10         for word in words:
11             yield word.lower(), 1
12
13     def combiner_WC(self, key, values):
14         # Each combiner will get a set of (key,value) pairs
15         # And it will join the ones it can.
16         yield key, sum(values)
17
18     def reducer_WC(self, key, values):
19         # The reducer will do the same joining with all the joining
20         # from the combiner
21
22         # Instead of yielding key-value pairs with yield key, sum(values)
23         # We just yield values, with the words and their count.
24         # This way, the next reducer will get all values at once
25         # print list(key)
26         # print list(values)
27         yield None, (sum(values), key)
28
29     def reducer_orderedWords(self, _, word_count_pairs):
30         # This reducer gets all the counted words and yields
31         # then in order of frequency
32         def getKey(item):
33             return item[0]
34
35         ordered = sorted(word_count_pairs, key= getKey, reverse=True)
36         yield None, ordered
37
38     def steps(self):
39         return [

```

```

40         MRStep mapper=self.mapper_WC,
41             combiner=self.combiner_WC,
42             reducer=self.reducer_WC),
43         MRStep(reducer=self.reducer_orderedWords)
44     ]
45
46 if __name__ == '__main__':
47     MRWordWC.run()

```

Exercise 7.2

Assignment :

Define and implement a MapReduce job that determines if a graph has an Euler tour (all vertices have even degree) where you can assume that the graph you get is connected.

This file <https://www.dropbox.com/s/usdi0wpsqm3jb7f/eulerGraphs.txt?dl=0> has 5 graphs – for each graph, the first line tells the number of nodes N and the number of edges E. The next E lines tells which two nodes are connected by an edge. Two nodes can be connected by multiple edges.

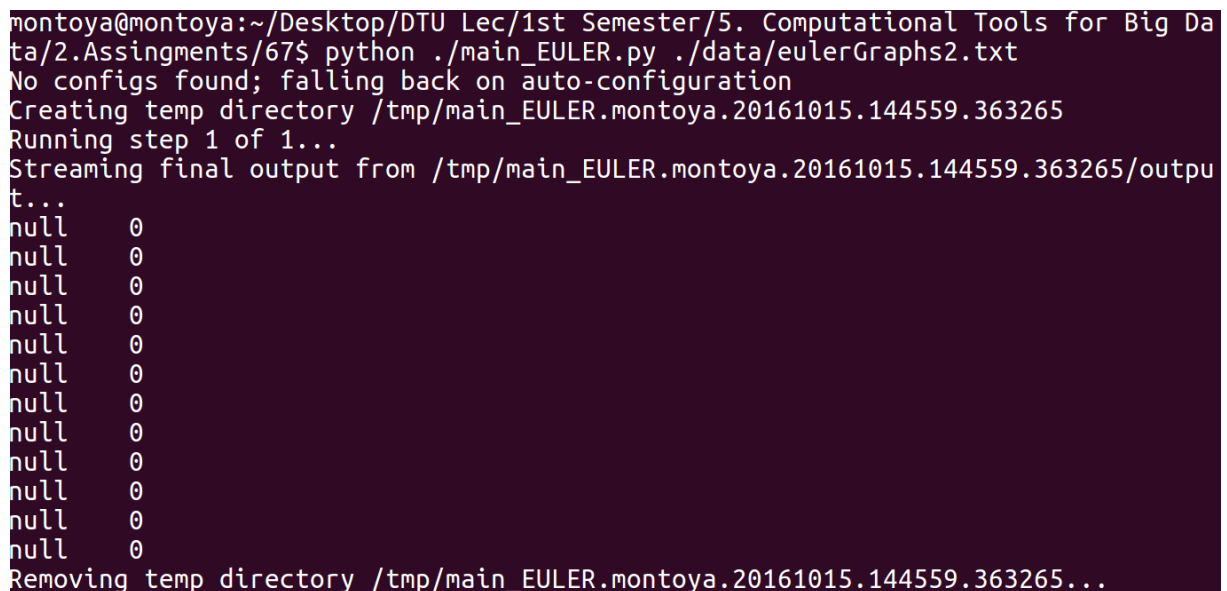
It is fine if you split the file into 5 different files. You do not need to keep the node and edge counts in the top of the file.

Solution :

We have splitted the file into 5 different files to make the code more reusable. From Graph Theory we know that there is an Euler Path if the number of edges of each node is even, or, there are exactly 2 nodes with an odd number of edges. To solve this exercise we have created 2 steps.

The **first step** calculates if the number of edges of each node is an even or odd number. For doing so, for each vertex, its mapper creates ("Node", 1) pairs, then the combiner and the reducer add the values for each node, and the reducer also applies modulus 2 operator. Then the **second step** is a reducer that checks if the number of odd degree nodes is 0 or 2 and outputs the solution.

The next image shows the partial output of the first step using the second graph given in the assignment. For every node, it yields 0 if the number of edges is even and 1 if it is odd. The key is the same for all nodes and it is "None".



```

montoya@montoya:~/Desktop/DTU Lec/1st Semester/5. Computational Tools for Big Data/2.Assingments/67$ python ./main_EULER.py ./data/eulerGraphs2.txt
No configs found; falling back on auto-configuration
Creating temp directory /tmp/main_EULER.montoya.20161015.144559.363265
Running step 1 of 1...
Streaming final output from /tmp/main_EULER.montoya.20161015.144559.363265/output...
null      0
null      0
null      0
null      0
null      0
null      0
null      0
null      0
null      0
null      0
null      0
null      0
null      0
null      0
null      0
Removing temp directory /tmp/main_EULER.montoya.20161015.144559.363265...

```

Figure 4: Euler path Step 1

The next image shows the output of the second step, it just the result of whether or not there exists an Euler path.

```

montoya@montoya:~/Desktop/DTU Lec/1st Semester/5. Computational Tools for Big Data/2.Assingments/67$ python ./main_EULER.py ./data/eulerGraphs2.txt
No configs found; falling back on auto-configuration
Creating temp directory /tmp/main_EULER.montoya.20161015.144624.463638
Running step 1 of 2...
Running step 2 of 2...
Streaming final output from /tmp/main_EULER.montoya.20161015.144624.463638/output...
null      "Number of odd vertex nodes 0"
1         "There is an Euler Path"

```

Figure 5: Euler path Step 2

2.2 main_EULER.py

```

1  from mrjob.job import MRJob
2  from mrjob.step import MRStep
3  # This Job is meant to determine if a graph has an Euler tour
4  # This happens if all vertices have even degree, or only 2 vertex
5  # have an even degree
6
7  class MRWordVC(MRJob):
8
9      def mapper_VC(self, _, line):
10         # This mapper assigns per every vertex, a 1 to every node in the vertex.
11         nodes_vertex = line.split()
12         if (len(nodes_vertex) == 2):
13             yield nodes_vertex[0], 1
14             yield nodes_vertex[1], 1
15
16     def combiner_VC(self, key, values):
17         # Each combiner will get a set of (key,value) pairs
18         # And it will join the ones it can.
19         yield key, sum(values)
20
21     def reducer_VC(self, key, values):
22         # The reducer will do the same joining with all the joining
23         # from the combiner. And it will also say if the number of
24         # vertex is odd or even for every node
25         yield None, sum(values)%2
26
27     def reducer_oddVertex(self, _, vertex_count_nodes):
28         # This reducer gets all the counted words and yields
29         # then in order of frequency
30         suma_odd = sum(vertex_count_nodes)
31         yield None, "Number of odd vertex nodes " + str(suma_odd)
32
33         if (suma_odd == 0 or suma_odd == 2):
34             yield 1, "There is an Euler Path"
35         else:
36             yield 0, "There is no Euler Path"
37
38     def steps(self):
39         return [
40             MRStep(mapper=self.mapper_VC,
41                   combiner=self.combiner_VC,
42                   reducer=self.reducer_VC),
43             MRStep(reducer=self.reducer_oddVertex)
44         ]
45
46     if __name__ == '__main__':

```


Exercise 7.3

Make a MapReduce job which counts the number of triangles in a graph. Use the following graph <http://www.cise.ufl.edu/research/sparse/matrices/SNAP/roadNet-CA.html>. The “Matrix Market” format is the preferred one as it has the same format as the Facebook data. (Remember to remove the header in the beginning of file)

Document that it works using a small example and run it on the large file to see that you get the right results.

Solution :

We have solved this problem using two different approaches, first we will show the most efficient one for this problem and then we will explain the second implementation.

First implementation :

In the first solution, we have used the Algorithm 3 MR-NodeIterator++(V,E) described in the document “Counting Triangles and the Curse of the Last Reducer” by Siddharth Suri Sergei and Vassilvitskii. In this implementation we have 3 steps:

- The first step calculates the adjacency matrix for all nodes. That is, for every node, it stores its neighbours in a sparse form. Each node is a row of the adjacency matrix and is stored as a dictionary, where each entry is a neighbour of that node and contain the value 1. The reducer of the method yields all the nodes (rows of the adjacency matrix) with the same key "A" and a tuple ("Node", DictionaryOfNeighbours) so that the next step has access to the whole adjacency matrix.
- The second step contains only a reducer that yields all the possible 2-node paths. For every node, and for each couple of neighbours of the node, it yields one of the nodes and the neighbours of the other node. It is optimized so that it does not yield each pair of neighbours twice. The yielded 2 node paths are composed by a key containing the first node and the two neighbours, and a value containing the first neighbour and the neighbours of the second neighbour.
- The third step contains a mapping that calculates for every 2-node paths, if there exists the triangle. For doing so, it just checks if the first node received is in the neighbourhood of the second node received. If it is, it yields a 1. Then the reducer sums all the yielded values, and divided by 3 to get the number of triangles. It is divided by 3 because we are gonna check the same path 3! times but since we optimize it so that for a given node, it does not output the same 2 neighbour nodes twice, we only have to divide by 3.

The next image shows the output for the first step using the third graph given in the assignment. We can appreciate that it yields all the nodes under the same key value "A". The nodes are rows of the adjacency matrix and we can see the tuples ("Node", DictionaryOfNeighbours).

```
montoya@montoya:~/Desktop/DTU Lec/1st Semester/5. Computational Tools for Big Data/2.Assingments/67$ python ./main_Triangle.py ./data/eulerGraphs3.txt
No configs found; falling back on auto-configuration
Creating temp directory /tmp/main_Triangle.montoya.20161015.145318.260322
Running step 1 of 1...
Streaming final output from /tmp/main_Triangle.montoya.20161015.145318.260322/output...
"A"      ["0", {"1": 1, "98": 1, "2": 1, "99": 1}]
"A"      ["1", {"99": 1, "0": 1, "3": 1, "2": 1}]
"A"      ["10", {"11": 1, "8": 1, "12": 1, "9": 1}]
"A"      ["11", {"9": 1, "10": 1, "13": 1, "12": 1}]
"A"      ["12", {"11": 1, "10": 1, "13": 1, "14": 1}]
"A"      ["13", {"11": 1, "12": 1, "15": 1, "14": 1}]
"A"      ["14", {"13": 1, "12": 1, "15": 1, "16": 1}]
"A"      ["15", {"13": 1, "14": 1, "17": 1, "16": 1}]
"A"      ["16", {"18": 1, "15": 1, "14": 1, "17": 1}]
"A"      ["17", {"19": 1, "18": 1, "15": 1, "16": 1}]
"A"      ["18", {"19": 1, "20": 1, "17": 1, "16": 1}]
"A"      ["19", {"18": 1, "20": 1, "21": 1, "17": 1}]
```

Figure 6: Counting Triangles step 1

The next image shows the output for the second step. We can see that the each yielded item contains in its value the first neighbour and the neighbours of the second neighbour.

```
Running step 2 of 2...
Streaming final output from /tmp/main_Triangle.montoya.20161015.150611.136503/output...
"24_25_26"      [{"25": {"24": 1, "25": 1, "27": 1, "28": 1}}]
"24_25_22"      [{"25": {"24": 1, "20": 1, "21": 1, "23": 1}}]
"24_25_23"      [{"25": {"24": 1, "25": 1, "21": 1, "22": 1}}]
"24_26_22"      [{"26": {"24": 1, "20": 1, "21": 1, "23": 1}}]
"24_26_23"      [{"26": {"24": 1, "25": 1, "21": 1, "22": 1}}]
"24_22_23"      [{"22": {"24": 1, "25": 1, "21": 1, "22": 1}}]
"25_24_26"      [{"24": {"24": 1, "25": 1, "27": 1, "28": 1}}]
"25_24_27"      [{"24": {"25": 1, "26": 1, "28": 1, "29": 1}}]
"25_24_23"      [{"24": {"24": 1, "25": 1, "21": 1, "22": 1}}]
"25_26_27"      [{"26": {"25": 1, "26": 1, "28": 1, "29": 1}}]
"25_26_23"      [{"26": {"24": 1, "25": 1, "21": 1, "22": 1}}]
"25_27_23"      [{"27": {"24": 1, "25": 1, "21": 1, "22": 1}}]
"26_24_25"      [{"24": {"24": 1, "26": 1, "27": 1, "23": 1}}]
"26_24_27"      [{"24": {"25": 1, "26": 1, "28": 1, "29": 1}}]
```

Figure 7: Counting Triangles step 2

The next image shows the output for the third step. It just consists of the number of triangles found.

```
Streaming final output from /tmp/main_Triangle.montoya.20161015.145201.833106/output...
"NT"      100
Removing temp directory /tmp/main_Triangle.montoya.20161015.145201.833106...
```

Figure 8: Counting Triangles step 3

For the big Graph roadNet problem, the obtained number of triangles is 120676 and in Intel(R) Core(TM) i7-4712MQ CPU @ 2.30GHz processor, along with other operating applications, the time required is lower than 8min30sec.

```
montoya@montoya:~/Desktop/DTU Lec/1st Semester/5. Computational Tools for Big Data/2.Assingments/67$ time python main_Triangle.py ./data/roadNet-CA.mtx
No configs found; falling back on auto-configuration
Creating temp directory /tmp/main_Triangle.montoya.20161015.160600.265625
Running step 1 of 3...
Running step 2 of 3...
Running step 3 of 3...
Streaming final output from /tmp/main_Triangle.montoya.20161015.160600.265625/output...
"NT"      120676
Removing temp directory /tmp/main_Triangle.montoya.20161015.160600.265625...

real      8m24.405s
user      8m0.068s
sys       0m12.776s
```

Figure 9: Counting Triangles roadNet problem

2.3 main_Traingle.py

```
1 from mrjob.job import MRJob
2 from mrjob.step import MRStep
3
4 # This Job determines the number of triangle sin a graph using the algorithm
5 # Algortihm 3 MR-NodeIterator++(V,E) described in the document
6 # "Counting Triangles and the Curse of the Last Reducer"
7 # by Siddharth Suri Sergei and Vassilvitskii
```

```

8
9 class MRNodesNT(MRJob):
10
11     #####
12     ## Job to create the adjacency matrix
13     #####
14     def mapper_AGeneration(self, _, line):
15         # This mapper assigns per every vertex, a 1 to every node in the vertex.
16         nodes_vertex = line.split()
17         if (len(nodes_vertex) == 2):
18             yield nodes_vertex[0], (nodes_vertex[1],1)
19             yield nodes_vertex[1], (nodes_vertex[0],1)
20
21     def combiner_AGeneration(self, key, values):
22         # Combines the adjacencies of the nodes.
23         yield key, list(values)
24
25     def reducer_AGeneration(self, key, values):
26         # Combines the adjacencies of the nodes again, all of them.
27         final_list = []
28         for value in values:
29             final_list.extend(list(value))
30         final_dict = dict(final_list)
31
32         row_elem = [key, final_dict]
33         # Each row contains on its data, the number of the row and the dictionary
34         # of values in sparse way. [Nrow, Dict of cols]
35
36         yield "A", row_elem # So that the next reducer get
37
38     #####
39     ## Generate the 2-length paths
40     #####
41     # We need to have access to all rows of the adjacency matrix to create
42     # the tasks so we implement a new reducer to join them
43     def reducer_inter(self, key, values):
44         val_list = list(values)
45         Matrix_A = dict(val_list) # We make another dictionary with the cols
46         A_keys = Matrix_A.keys()
47
48         # Output all the pairs !! And their neighbours
49         for k in A_keys:
50             neighs = Matrix_A[k]
51             neighs_keys = neighs.keys()
52             Nneigh = len(neighs_keys)
53             for i in range(Nneigh):
54                 for j in range(i, Nneigh):
55                     nk1 = neighs_keys[i]
56                     nk2 = neighs_keys[j]
57                     if (nk1 != nk2):
58                         yield str(k) + "_" + str(nk1) + "_" + str(nk2), [nk1, Matrix_A[nk2]]
59
60     #####
61     ## Check if you can form a triangle with every pair
62     #####
63     def mapper_2Path(self, key, values):
64         # This mapper outputs pairs of neighbours for each node.
65         # So that latter the mapping can
66         lista = list(values)
67         # print lista
68         if (lista[0] in lista[1].keys()):

```

```

69         yield None,1
70
71     def reducer_2Path(self, key, values):
72         yield "NT", sum(values)/3
73
74     def steps(self):
75         return [
76
77             MRStep(mapper=self.mapper_AGeneration,
78                   combiner=self.combiner_AGeneration,
79                   reducer=self.reducer_AGeneration),
80
81             MRStep(reducer=self.reducer_inter),
82
83             MRStep(mapper=self.mapper_2Path,
84                   reducer=self.reducer_2Path),
85         ]
86
87 if __name__ == '__main__':
88     MRNodesNT.run()

```

Second approach :

Before we implemented the above described algorithm, we found on the Internet that the number triangles in the Graph could be calculated as $tr(A^3/6)$ so we decided to implement a job that gets the diagonal of A^3 . It ended up being a little tricky to do with our limited knowledge but we will try to explain how we implemented it.

- As always, we first have a step to create the Adjacency Matrix. It is stored in a sparse way, as a dictionary, in which each row is an entry (referenced by the number of the node). Each row is described by another dictionary that contains for every non-zero elements, the entry (node, value). It is just a way of describing the non-zero elements of the Adjacency matrix in a sparse form.
- Next task is to perform AA multiplication in a sparse way so we created a another step consisting of only one reducer that yields a (key, value) pair for every element of the output matrix to calculate. The key consists of the position of the matrix to be calculated "Nodei_Nodej" and the values is the tuple [Row_i, Col_j], lucky enough, the A matrix is symmetric so the rows and the columns are the same.
- The next step calculates the elements of the AA matrix, first, its mapper performs the sparse multiplication in an efficient way and yields a (key, value) pair consisting on ("Row_i", "Col_j":Value) in order to construct the AA in the same sparse row-oriented way as matrix A . The combiner and the reducer collects all the outputs and finally create matrix. During all the stages, we had to propagate matrix A using a special key, in order to be able to use in the future to calculate AAA .
- The next step just combines both sparse matrices AA and A under the same key so that we can have access to both of them simultaneously.
- Step number 5 is a reducer similar to the number 2, it yields a (key, value) pair for every element of the output matrix to calculate AAA . But since we only need the diagonal, it only yields those values.
- Step number 6 performs the calculation of the elements of the output matrix in the same way as step 3, but its reducer just yields the value of the diagonal elements (if they are non-zero)
- The final step is a reducer that just adds all the values and divides by 6 to finally obtain the output.

2.4 main_Traingle2.py

```

1 from mrjob.job import MRJob
2 from mrjob.step import MRStep
3
4 # This Job determines the number of triangle sin a graph using

```

```

5  # the trace of the A^3 (A = Adjacency matrix)
6
7  class MRNodesNT(MRJob):
8      #####
9      ## Job to create the adjacency matrix
10     #####
11     def mapper_AGeneration(self, _, line):
12         # This mapper assigns per every vertex, a 1 to every node in the vertex.
13         nodes_vertex = line.split()
14         if (len(nodes_vertex) == 2):
15             yield nodes_vertex[0], (nodes_vertex[1],1)
16             yield nodes_vertex[1], (nodes_vertex[0],1)
17
18     def combiner_AGeneration(self, key, values):
19         # The combiner does need to do anything, the nodes will have all
20         # its nodes joined already
21         yield key, list(values)
22
23     def reducer_AGeneration(self, key, values):
24         # Same here
25         final_list = []
26         for value in values:
27             final_list.extend(list(value))
28         final_dict = dict(final_list)
29
30         row_elem = [key, final_dict]
31         # Each row contains the number of the row and the dictionary
32         # of values in sparse way. [Nrow, Dict of cols]
33
34         yield "Matrices", row_elem # So that the next reducer get
35
36     #####
37     ## Matrix multiplication
38     ## We are gonna multiply the matrix A*A in sparse way,
39     #####
40     def reducer_MatrixMult(self, key, values):
41         # This reducer output a par of (row_1, col_2) to multiply
42         # by the combiner so that it outputs one element of the
43         # output matrix
44         # SINCE THE MATRIX IS SYMETRIC
45         # We also want to obtain the number of nodes using the keys!
46         # For knowing how many multiplications
47
48         val_list = list(values)
49         Matrix_A = dict(val_list) # We make another dictionary with the cols
50         # First we yield the A matrix to perform A (A^2) later
51         yield "Matrices", val_list
52
53         keys_all = Matrix_A.keys()
54         nt = len(keys_all)
55         c = 1 # Just for having a sense of how much time it will take.
56         for i in keys_all:
57             if (c % 50 == 0):
58                 print str(c)+" / "+ str(nt) + " getting"
59                 c+= 1
60             for j in keys_all:
61                 key_aux = i + "_" + j
62                 yield key_aux, [Matrix_A[i], Matrix_A[j]]
63
64     #####
65     ## MAPPERS are gonna do multiplications

```

```

66 #####
67 def mapper_MatrixMult(self, key, values):
68     # This mapper expects the 2 vectors to multiply.
69     # The key contains the position
70     # Multiply the vectors and output the value only if it is nonzero
71     # The output key will be the row only.
72     if (key == "Matrices"):
73         yield key, list(values) # The first time we do not do [0] TODO
74
75     else:
76         val_list = list(values)
77         v1 = val_list[0]
78         v2 = val_list[1]
79         ##### Sparse multiplication of the vectors
80         # Get the intersection of elements
81         InterKeys = list(set(v1.keys()) & set(v2.keys()))
82
83         if (len(InterKeys) > 0): # If they have same keys in common
84             total = 0
85             for key_i in InterKeys:
86                 total += v1[key_i] * v2[key_i]
87
88             i = key.split("_")[0]
89             j = key.split("_")[1]
90
91             yield i, (j,total) # Yield the result
92
93 def combiner_AGeneration_2(self, key, values):
94     # Just combine the cells to create the output matrix
95     if (key == "Matrices"):
96         yield key,dict(list(values)[0])
97     else:
98         yield key, list(values)
99
100 def reducer_AGeneration_2(self, key, values):
101     # Same here, more combinations to create the matrix
102     if (key == "Matrices"):
103         yield key,list(values)[0]
104     else:
105         # Same here
106         final_list = []
107         for value in values:
108             final_list.extend(list(value))
109         final_dict = dict(final_list)
110         row_elem = [key, final_dict]
111         yield "AA", row_elem # So that the next reducer get
112
113 #####
114 ## JOIN MATRIX AA
115 #####
116 def reducer_AGeneration_3(self, key, values):
117     if (key == "Matrices"):
118         Matrix_A = list(values)[0]
119         yield "Matrices",Matrix_A
120
121     elif(key == "AA"):
122         val_list = list(values)
123         Matrix_AA = dict(val_list) # Make another dictionary with the cols
124         yield "Matrices",Matrix_AA
125
126 #####

```

```

127     ## MULTIPLY AGAIN
128     #####
129
130     def reducer_MatrixMult2(self, key, values):
131         final_list = list(values)
132         Matrix_A = final_list[0]
133         Matrix_AA = final_list[1]
134         # We are only interested in the diagonal elements so...
135         # we only calculate those. We do not need to calculate the whole matrix
136         print Matrix_AA.keys()
137         for i in Matrix_AA.keys():
138             key_aux = i + "_" + i
139             yield key_aux, [Matrix_AA[i], Matrix_A[i]]
140
141     #####
142     #### Finally get the trace of elements
143     #####
144     def reducer_trace_rows(self, key, values):
145         # Get the diagonal
146         lista = list(values)[0]
147         dictionary = dict(lista)
148
149         if key in dictionary: # Get the trace if it is non-zero
150             yield None, dictionary[key]
151         else:
152             yield None, 0
153
154     def reducer_trace_rows_2(self, key, values):
155         # Sum t
156         Trace = sum(values)
157         print "Trace " + str(Trace)
158         N_triangles = Trace/6
159
160         print "Number of triangles " + str(N_triangles)
161         yield "NT", N_triangles
162
163     def steps(self):
164         return [
165             # Get the Adjacency Matrix, the key is the row, and it is
166             # eliminated in the reduction so that we can handle all the rows
167             # in the next reduction
168             MRStep(mapper=self.mapper_AGeneration,
169                   combiner=self.combiner_AGeneration,
170                   reducer=self.reducer_AGeneration),
171
172             # Generate the vector multiplications
173             MRStep(reducer=self.reducer_MatrixMult),
174
175             # Multiply the matrices and obtain the result,
176             # While storing the initial matrix
177             MRStep(mapper=self.mapper_MatrixMult,
178                   combiner=self.combiner_AGeneration_2,
179                   reducer=self.reducer_AGeneration_2),
180
181             # We join the AA matrix and the A matrix
182             MRStep(reducer=self.reducer_AGeneration_3),
183             MRStep(reducer=self.reducer_MatrixMult2),
184
185             # We multiply them again
186             MRStep(mapper=self.mapper_MatrixMult,
187                   combiner=self.combiner_AGeneration_2,

```

```
188         reducer = self.reducer_trace_rows),
189
190         # Add all the traces and divide by 6
191         MRStep(reducer=self.reducer_trace_rows_2),
192     ]
193
194 if __name__ == '__main__':
195     MRNodesNT.run()
```