

Índice

1 Grafo multifichero:	2
1.1 Unión de rdds:	2
1.2 Triciclos	2
2 Múltiples grafos:	3

Listings

1 Grafo multifichero: mixed	2
2 Grafo multifichero: tag	2
3 Grafo multifichero: tricycles	2
4 Grafo multifichero: process data	2
5 Múltiples grafos: función para distinguir grafos	3
6 Múltiples grafos: get distinct edges	3
7 Múltiples grafos: nombres de los ficheros en los rdds	3
8 Múltiples grafos: tag	3
9 Múltiples grafos: construyendo los triciclos	3
10 Múltiples grafos: process data	3

1 Grafo multifichero:

En este ejercicio analizamos un grafo formado por multiples archivos y devolvemos todos sus triciclos.

1.1 Unión de rdds:

Primero unimos los rdd:

```
62 def mixed(sc, rdds):
63     for rdd in rdds:
64         rdd.cache()
65     data = sc.union(rdds)
```

1.2 Triciclos

Dado el rdd con los nodos y sus listas de adyacencia considerando solo nodos posteriores definimos la función tag que genera la lista de 'exist' y 'pending' para un nodo:

```
32 def tag(node_adj): # Función iterativa, no perezosa
33     node = node_adj[0]
34     adjs = list(node_adj[1])
35     adjs.sort()
36     result = [(node, x), 'exists'] for x in adjs
37     for i in range(len(adjs)):
38         for j in range(i, len(adjs)):
39             result.append(((adj[i], adj[j]), ('pending', node)))
40     return result
```

Generamos las listas de adyacencia y generamos los triciclos:

tricycles

```
43 def tricycles(tags):
44     return tags\
45         .groupByKey()\
46         .filter(lambda x: len(x[1]) > 1 and 'exists' in x[1])\
47         .flatMap(
48             lambda line: map(
49                 lambda x: (x[1], line[0][0], line[0][1]),
50                 filter(lambda x: not x == 'exists', line[1])
51             )
52     )
```

process data

```
55 def process_data(data):
56     edges = get_distict_edges(data)
57     node_adj = get_node_adj(edges)
58     tags = node_adj.flatMap(tag)
59     return tricycles(tags)
```

2 Multiples grafos:

Para analizar cada archivo como un grafo a parte cada arista contiene también el nombre del archivo al que pertenece:

Distinguiendo los grafos

```

20 | tupleMorph = lambda tup, x : (tup[0], (tup[1], x))
23 | def get_distict_edges(rdd):
24 |     return rdd\
25 |         .map(lambda x: tupleMorph(get_edges(x[0]),x[1]))\
26 |         .filter(lambda x: x[1] is not None)\
27 |         .distinct()
80 | def independent(sc, rdds, files):
81 |     assert len(rdds) == len(files)
82 |     tagged_rdds = [rdd.map(lambda x: (x,graph)).cache() for rdd, graph in zip(rdds, files)]
83 |     data = sc.union(tagged_rdds)
84 |     for tricycle in process_data(data).collect():
85 |         print(tricycle[0], tricycle[1])

```

Ahora realizamos un procedimiento análogo al del archivo anterior siempre teniendo en cuenta el archivo al que pertenece cada arista para que los triciclos esten compuestos por tres nodos del mismo grafo.

tag

```

34 | def tag(node_adj): # Función iterativa, no perezosa
35 |     node = node_adj[0]
36 |     adjs = list(node_adj[1])
37 |     adjs.sort()
38 |     result = [(node, x), (graph, 'exists')] for (x, graph) in adjs]
39 |     for i in range(len(adjs)):
40 |         for j in range(i, len(adjs)):
41 |             if adjs[j][1] == adjs[i][1] and not adjs[i][0] == adjs[j][0]: #Los nodos aparecen en el mismo grafo
42 |                 result.append(((adj[i][0], adj[j][0]), (adj[i][1], ('pending', node))))
43 |     return result

```

Construyendo los triciclos

```

52 | def findExists(tag_list):
53 |     for tag in tag_list:
54 |         if tag[1] == 'exists':
55 |             return True
56 |     return False
57 |
58 | def list2dict(lista):
59 |     result = {}
60 |     for key, value in lista:
61 |         if key not in result:
62 |             result[key] = [value]
63 |         else:
64 |             result[key].append(value)
65 |     return result
66 |
67 | def construct_tricycles(line):
68 |     edge, tags = line
69 |     return (edge, list2dict(tags))
70 |
71 |
72 | def tricycles(tags):
73 |     return tags\
74 |         .groupByKey()\
75 |         .filter(lambda x: len(x[1]) > 1 and findExists(x[1]))\
76 |         .map(construct_tricycles)

```

process data

```

77 | def process_data(data):
78 |     edges = get_distict_edges(data).cache()
79 |     node_adj = get_node_adj(edges).cache()
80 |     tags = node_adj.flatMap(tag).cache()
81 |     return tricycles(tags)

```