



11

BIG DATA AND GEOSPATIAL DATA MINING

3 DISTRIBUTED DATA PROCESSING

Tutorial: <https://media.readthedocs.org/pdf/mrjob/latest/mrjob.pdf>

mrjob

mrjob lets you write MapReduce jobs in Python 2.7/3.4+ and run them on several platforms. You can:

- Write multi-step MapReduce jobs in pure Python
- Test on your local machine
- Run on a Hadoop cluster
- Run in the cloud using Amazon Elastic MapReduce (EMR)
- Run in the cloud using Google Cloud Dataproc (Dataproc)
- Easily run *Spark* jobs on EMR or your own Hadoop cluster

mrjob is licensed under the Apache License, Version 2.0.

To get started, install with pip:

```
pip install mrjob
```

and begin reading the tutorial below.

3.2 MRJob

A. Martín

12

Overview

mrjob is the easiest route to writing Python programs that run on Hadoop. If you use mrjob, you'll be able to test your code locally without installing Hadoop or run it on a cluster of your choice.

Additionally, mrjob has extensive integration with Amazon Elastic MapReduce. Once you're set up, it's as easy to run your job in the cloud as it is to run it on your laptop.

Here are a number of features of mrjob that make writing MapReduce jobs easier:

- Keep all MapReduce code for one job in a single class
- Easily upload and install code and data dependencies at runtime
- Switch input and output formats with a single line of code
- Automatically download and parse error logs for Python tracebacks
- Put command line filters before or after your Python code

If you don't want to be a Hadoop expert but need the computing power of MapReduce, mrjob might be just the thing for you.

13

Classes in Python: Models that define and group the characteristics and properties that will have the objects that instantiate them. In python, a class is defined with the *class* statement followed by a generic name for the object.

```

File Edit Format Run Options Window Help
1 #!/usr/bin/env python2
2 # -*- coding: utf-8 -*-
3 """
4 Created on Wed Jul 18 19:08:33 2018
5
6 @author: angelmartinfurones
7 """
8 #CLASE QUE SIMULA UNA CALCULADORA
9 class Calculadora(object):
10     #EL CONSTRUCTOR INICIA "VALOR" A 0
11     def __init__(self): #SELF HACE REFERENCIA A SI MISMO
12         self.valor=0
13     #SUMA UN NUMERO 'N' AL VALOR
14     def suma(self,n):
15         self.valor +=n
16     def getValor(self):
17         return self.valor
18
19 calc = Calculadora() #INSTANCIANDO UN OBJETO EN LA VARIABLE CALC
20
21 calc.suma(2) #SUMA 2 AL VALOR DE LA CALCULADORA
22
23 print(calc.getValor()) #MOSTRARIA UN DOS EN LA SHELL
24
25 calc.suma(2) #SE SUMA 2 A LA VARIABLE VALOR
26
27 print(calc.getValor())
28

```

14

BIG DATA AND GEOSPATIAL DATA MINING

3 DISTRIBUTED DATA PROCESSING

```

1 from mrjob.job import MRJob
2
3 class Purchases(MRJob):
4
5     def mapper(self, key, line):
6         if len(line.strip().split(","))==6:
7             (date,time,store,item,cost,payment) = line.strip().split(",")
8             try:
9                 yield str(store),float(cost)
10            except:
11                pass
12
13     def reducer(self, store, cost):
14         yield str(store), sum(cost)
15
16 if __name__ == '__main__':
17     Purchases.run()

```

\$ Python3 SolMRJob.py purchases.txt | sort > salidaMRJob.txt

3.2 MRJob

A. Martín

15

BIG DATA AND GEOSPATIAL DATA MINING

3 DISTRIBUTED DATA PROCESSING

```

1 from mrjob.job import MRJob
2
3 class Purchases(MRJob):
4
5     def mapper(self, key, line):
6         if len(line.strip().split(","))==6:
7             (date,time,store,item,cost,payment) = line.strip().split(",")
8             try:
9                 yield str(store),float(cost)
10            except:
11                pass
12
13     def reducer(self, store, cost):
14         yield str(store), sum(cost)
15
16 if __name__ == '__main__':
17     Purchases.run()

```

MRJob class consisting of a Mapper function and Reducer function

Start, python is told to run the program

3.2 MRJob

A. Martín

16

BIG DATA AND GEOSPATIAL DATA MINING

3 DISTRIBUTED DATA PROCESSING

```

1 from mrjob.job import MRJob
2
3 class Purchases(MRJob):
4
5     def mapper(self, key, line):
6         if len(line.strip().split(","))==6:
7             (date,time,store,item,cost,payment) = line.strip().split(",")
8             try:
9                 yield str(store),float(cost)
10            except:
11                pass
12
13     def reducer(self, store, cost):
14         yield str(store), sum(cost)
15
16 if __name__ == '__main__':
17     Purchases.run()
18

```

• Line Refers to each line of raw input data
 • Key refers to the key of each of the input lines, in this case it is not necessary to assign any key so that key can be replaced by "_", so the key is ignored

3.2 MRJob

A. Martín

17

BIG DATA AND GEOSPATIAL DATA MINING

3 DISTRIBUTED DATA PROCESSING

```

1 from mrjob.job import MRJob
2
3 class Purchases(MRJob):
4
5     def mapper(self, key, line):
6         if len(line.strip().split(","))==6:
7             (date,time,store,item,cost,payment) = line.strip().split(",")
8             try:
9                 yield str(store),float(cost)
10            except:
11                pass
12
13     def reducer(self, store, cost):
14         yield str(store), sum(cost)
15
16 if __name__ == '__main__':
17     Purchases.run()
18

```

The output will be saved in the `yield` function as a key, value pair

3.2 MRJob

A. Martín

18

BIG DATA AND GEOSPATIAL DATA MINING

3 DISTRIBUTED DATA PROCESSING

```

1 from mrjob.job import MRJob
2
3 class Purchases(MRJob):
4
5     def mapper(self, key, line):
6         if len(line.strip().split(","))==6:
7             (date,time,store,item,cost,payment) = line.strip().split(",")
8             try:
9                 yield str(store),float(cost)
10            except:
11                pass
12
13     def reducer(self, store, cost):
14         yield str(store), sum(cost)
15
16 if __name__ == '__main__':
17     Purchases.run()
18

```

Sort and Group, MRJob does it for us, you don't have to program anything

3.2 MRJob A. Martín

19

BIG DATA AND GEOSPATIAL DATA MINING

3 DISTRIBUTED DATA PROCESSING

```

1 from mrjob.job import MRJob
2
3 class Purchases(MRJob):
4
5     def mapper(self, key, line):
6         if len(line.strip().split(","))==6:
7             (date,time,store,item,cost,payment) = line.strip().split(",")
8             try:
9                 yield str(store),float(cost)
10            except:
11                pass
12
13     def reducer(self, store, cost):
14         yield str(store), sum(cost)
15
16 if __name__ == '__main__':
17     Purchases.run()
18

```

The input is, in this case, the stores and the cost, but already ordered (we can assign another name to cost, for example occurrences)

3.2 MRJob A. Martín

20

```

1 from mrjob.job import MRJob
2
3 class Purchases(MRJob):
4
5     def mapper(self, key, line):
6         if len(line.strip().split(","))==6:
7             (date,time,store,item,cost,payment) = line.strip().split(",")
8             try:
9                 yield str(store),float(cost)
10            except:
11                pass
12
13     def reducer(self, store, cost):
14         yield str(store), sum(cost)
15
16 if __name__ == '__main__':
17     Purchases.run()
18

```

The output will be a key/value pair with the store and the sum of the costs

3.2 MRJob

21

22

```

Punto,Mes,Dia,Ano,Hora,Intensidad
106,4,1,2013,0,142
106,4,1,2013,1,82
106,4,1,2013,2,54
106,4,1,2013,3,68
106,4,1,2013,4,33
106,4,1,2013,5,41
106,4,1,2013,6,61
106,4,1,2013,7,71
106,4,1,2013,8,100
106,4,1,2013,9,134
106,4,1,2013,10,183
106,4,1,2013,11,227
106,4,1,2013,12,334
106,4,1,2013,13,342
106,4,1,2013,14,300
106,4,1,2013,15,214
106,4,1,2013,16,279
106,4,1,2013,17,345
106,4,1,2013,18,329
106,4,1,2013,19,429
106,4,1,2013,20,447
106,4,1,2013,21,387
106,4,1,2013,22,227
106,4,1,2013,23,171
106,4,2,2013,0,100
106,4,2,2013,1,59
106,4,2,2013,2,40
106,4,2,2013,3,26
106,4,2,2013,4,16

```

3.2 MRJob

A. Martín

23

```

File Edit Format Run Options Window Help
1 from mrjob.job import MRJob
2
3 class IntensidadPuntos(MRJob):
4
5     def mapper(self, key, line):
6         (punto,mes,dia,ano,hora,intensidad) = line.split(',')
7         if punto !='Punto':
8             yield str(punto), int(intensidad)
9
10    def reducer(self, punto, intensidad):
11        total = 0
12        numElements = 0
13        for x in intensidad:
14            total += x
15            numElements += 1
16
17        yield punto, int(total / (numElements))
18
19 if __name__ == '__main__':
20     IntensidadPuntos.run()
21

```

3.2 MRJob

A. Martín

24



25

```

1 #!/usr/bin/env python2
2 # -*- coding: utf-8 -*-
3 """Created on Mon Jul 23 08:58:29 2018
4 @author: angelmartinfurones
5 """
6 from mrjob.job import MRJob
7 from mrjob.step import MRStep
8
9 class IntensidadPuntos(MRJob):
10     def steps(self):
11         return [
12             MRStep(mapper=self.mapper_get_ratings,
13                   reducer_init=self.reducer_init,
14                   reducer=self.reducer_count_ratings)
15         ]
16
17     def mapper_get_ratings(self,key,line):
18         (punto,mes,dia,año,hora,intensidad) = line.split(',')
19         if punto != 'Punto':
20             yield key,(punto,int(intensidad))
21
22     def reducer_init(self):
23         self.Names = {}
24         file1=open("/home/angel/Documents/BigData/mapreduce/intensidad/TRA_ESPIRAS_P-2.CSV")
25         while True:
26             line=file1.readline()
27             if not line:break
28             fieldsline.split(';')
29             if str(fields[0]) != 'X' and fields[5] != None:
30                 coord=str(fields[0])+"," +str(fields[1])
31                 clave=str(fields[5])
32                 self.Names.update({clave:coord})
33
34     def reducer_count_ratings(self,punto,intensidad):
35         total = 0
36         numElements = 0
37         for x in intensidad:
38             total += x
39             numElements += 1
40         try:
41             yield self.Names[str(punto)], int(total / numElements)
42         except:
43             yield str(punto), total / numElements
44
45 if __name__ == '__main__':
46     IntensidadPuntos.run()

```

3.2 M

26

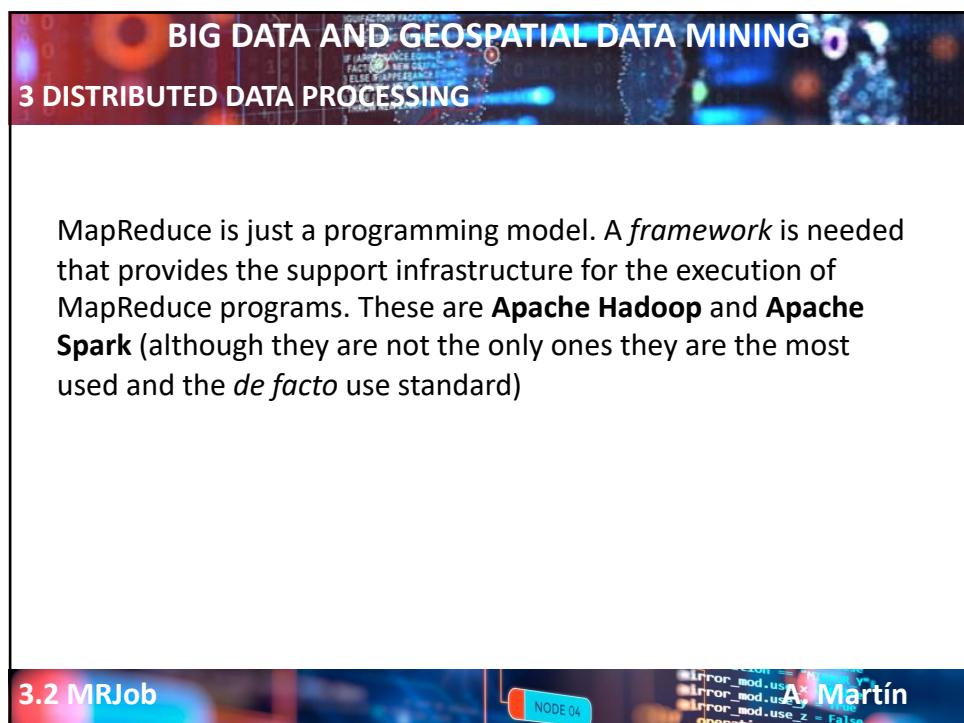
```

1 #!/usr/bin/env python2
2 # -*- coding: utf-8 -*-
3 # **Created on Mon Jul 23 08:58:29 2018
4 #@author: angelmartinforones
5 #
6 from mrjob.job import MRJob
7 from mrjob.step import MRStep
8
9 import utm
10
11 class IntensidadPuntos(MRJob):
12     def steps(self):
13         return [
14             MRStep(mapper=self.mapper_get_ratings,
15                   reducer_init=self.reducer_init,
16                   reducer=self.reducer_count_ratings)
17         ]
18     def mapper_get_ratings(self,key,line):
19         (punto,mes,dia,año,hora,intensidad) = line.split(',')
20         if punto != 'Punto':
21             yield str(punto),int(intensidad)
22     def reducer_init(self):
23         self.Names = {}
24         file1=open("/home/angel/Documents/BigData/mapreduce/intensidad/TRA_ESPIRAS_P-2.CSV")
25         while True:
26             line=file1.readline()
27             if not line:break
28             fields=line.split(',')
29             if str(fields[0]) == 'Y' and fields[5] != None:
30                 lat,lon=utm.toLatLon(float(fields[0]), float(fields[1]), 30, 'N')
31                 coord=str(str(lon)+"."+str(lat))
32                 clave=str(punto)+coord
33                 self.Names.update({clave:coord})
34
35     def reducer_count_ratings(self,punto,intensidad):
36         total = 0
37         numElements = 0
38         for x in intensidad:
39             total += x
40             numElements += 1
41         try:
42             yield self.Names[str(punto)], int(total / numElements)
43         except:
44             yield str(punto), int(total / numElements)
45
46 if __name__ == '__main__':
47     IntensidadPuntos.run()
48

```

Ln: 30 Col: 18
A. Martín

27



28