大规模分布式系统第二次实验报告
员司雨 17307110448


# PartI：统计其中各类文件的数量

# 一． 实验介绍

　　mrjob 是编写能够在 hadoop 上运行的 python 程序最简单的途径。如果使用 mrjob，可以在本地测试的代码，甚至不需要安装 hadoop 或者在选择的集群上运行。

　　另外，mrjob 可以和亚马逊的 EMR（Elastic MapReduce）服务无缝集成。只要设置完毕，就可以运行在 EMR 上，像在自己的笔记本上运行一样简单。


# 二． 实验环境

1. Ubuntu18.04
2. jdk 1.8.0_131
3. hadoop 2.7.7
4. Python3.6.2
5. mrjob 包


# 三． 实验过程

1. Python 版本查看



Python 版本为 3.6.2

2. mrjob 包安装



使用 pip 安装 mrjob 包

3. Python 代码编写

```
1.  from mrjob.job import MRJob
```

```
2.  import os
3.  class MRFILE_TYPE_Counter(MRJob):
4.
5.      def mapper(self, key, line):
6.          temp = line.split(' ')
7.          F = temp[-1]
8.          f = os.path.splitext(F)
9.          filename,ty = f
10.          yield ty, 1
11.
12.      def reducer(self, word, occurrences):
13.          yield word, sum(occurrences)
14.
15.  if __name__=='__main__':
16.      MRFILE_TYPE_Counter.run()
```

参见：type_count.py

4. 实验结果



```
ics@ubuntu:~/tools/Python-3.6.2/mycode$ python type_count.py -r local ./sample.t
xt
No configs found; falling back on auto-configuration
No configs specified for local runner
Creating temp directory /tmp/type_count.ics.20200330.112441.319208
Running step 1 of 1...
job output is in /tmp/type_count.ics.20200330.112441.319208/output
Streaming final output from /tmp/type_count.ics.20200330.112441.319208/output...
".docx"  1
".dwg"   32
".jpg"   9
".pdf"   43
Removing temp directory /tmp/type_count.ics.20200330.112441.319208...
ics@ubuntu:~/tools/Python-3.6.2/mycode$ python type_count.py -r local ./sample.t
xt >> output.txt
No configs found; falling back on auto-configuration
No configs specified for local runner
Creating temp directory /tmp/type_count.ics.20200330.112454.762896
Running step 1 of 1...
job output is in /tmp/type_count.ics.20200330.112454.762896/output
Streaming final output from /tmp/type_count.ics.20200330.112454.762896/output...
Removing temp directory /tmp/type_count.ics.20200330.112454.762896...
```

实验成功结果，具体参加：output.txt


# PartII：按文件的字节数大小降序排序输出文件名

## 一. 实验介绍

mrjob 是编写能够在 hadoop 上运行的 python 程序最简单的途径。如果使用 mrjob，可以在本地测试的代码，甚至不需要安装 hadoop 或者在选择的集群上运行。

另外，mrjob 可以和亚马逊的 EMR（Elastic MapReduce）服务无缝集成。只要设置完毕，就可以运行在 EMR 上，像在自己的笔记本上运行一样简单。

## 二. 实验环境

6.   Ubuntu18.04
7.   jdk 1.8.0_131
8.   hadoop 2.7.7
9.   Python3.6.2
10.  mrjob 包

# 三． 实验过程

1. Python 代码编写

```python
1.   from mrjob.job import MRJob
2.   from mrjob.step import MRStep
3.   import heapq
4.   import os
5.   class MRFILE_SIZE_Counter(MRJob):
6.
7.       def mapper(self, key, line):
8.           temp = line.split()
9.           Size = temp[2]
10.          if ',' in temp[-2]:
11.              F = temp[-1]
12.          else:
13.              F = temp[-2]+' '+temp[-1]
14.          yield (int(Size.replace(',','')),F),1
15.
16.      def reducer_1(self, key, value):
17.          yield None, key
18.
19.      def reducer_2(self, _, value):
20.          for s, f in heapq.nlargest(85, value):
21.              yield s,f
22.      def steps(self):
23.          return[MRStep(mapper = self.mapper,reducer = self.reducer_1),MRStep(reducer = self.reducer_2)]
24.
25.  if __name__=='__main__':
26.      MRFILE_SIZE_Counter.run()
```

参见：type_count.py

2. 实验结果

```
ics@ubuntu:~/tools/Python-3.6.2/mycode$ python size_count.py -r local ./sample.t
xt
No configs found; falling back on auto-configuration
No configs specified for local runner
Creating temp directory /tmp/size_count.ics.20200330.115333.878705
Running step 1 of 2...
Running step 2 of 2...
job output is in /tmp/size_count.ics.20200330.115333.878705/output
Streaming final output from /tmp/size_count.ics.20200330.115333.878705/output...
5272668 "\u603b\u4f53-\u5f31\u65bd-00-02.pdf"
5203305 "\u603b\u4f53-\u5f31\u65bd-00-01.pdf"
4520750 "10-09 \u5178\u578b\u5355\u5899\u8be6\u56fe.dwg"
3395145 "02-01 \u4e00\u5c42\u5e73\u9762\u56fe.dwg"
3389704 "10-08 \u6838\u5fc3\u7b52\u8be6\u56fe\uff08\u516b\uff09.dwg"
3389704 "10-07 \u6838\u5fc3\u7b52\u8be6\u56fe\uff08\u4e03\uff09.dwg"
3389704 "10-06 \u6838\u5fc3\u7b52\u8be6\u56fe\uff08\u516d\uff09.dwg"
3389704 "10-05 \u6838\u5fc3\u7b52\u8be6\u56fe\uff08\u4e94\uff09.dwg"
3389704 "10-04 \u6838\u5fc3\u7b52\u8be6\u56fe\uff08\u56db\uff09.dwg"
3389704 "10-03 \u6838\u5fc3\u7b52\u8be6\u56fe\uff08\u4e09\uff09.dwg"
3389704 "10-02 \u6838\u5fc3\u7b52\u8be6\u56fe\uff08\u4e8c\uff09.dwg"
3389704 "10-01 \u6838\u5fc3\u7b52\u8be6\u56fe\uff08\u4e00\uff09.dwg"
```

```
ics@ubuntu:~/tools/Python-3.6.2/mycode$ python size_count.py -r local ./sample.t
xt >> output2.txt
No configs found; falling back on auto-configuration
No configs specified for local runner
Creating temp directory /tmp/size_count.ics.20200330.115403.395286
Running step 1 of 2...
Running step 2 of 2...
job output is in /tmp/size_count.ics.20200330.115403.395286/output
Streaming final output from /tmp/size_count.ics.20200330.115403.395286/output...
Removing temp directory /tmp/size_count.ics.20200330.115403.395286...
```

实验成功结果，具体参加：output2.txt

3. 细节介绍

（1）由于 hadoop 默认是对 key 做升序排序输出，而我们的要求是倒序，所以引入包 heapq.

（2）由于输出是 unicode 编码，故而需要转化为中文。本地在线转换代码如下

```
1.  temp1 = []
2.
3.  address2 = "D:\output2.txt"
4.  with open(address2, 'r') as f1:
5.      for line in f1:
6.          line = line.strip('\n')
7.          temp1.append(line)
8.
9.  ls0 = []
10. for ele in temp1:
11.     a = ele.split('\t')
12.     a[1] = a[1].encode('utf-8').decode('unicode_escape')
13.     ls0.append(a)
14.
15. f = open("D:/output2_rewrite.txt", "w")
16. for element in ls0:
17.     f.write(str(element[0])+'\t'+element[1]+'\n')
```

18. f.close()

代码及输出结果参加 `file_rewrite.py` 和 `output2_rewrite.txt`。