# Bonus

## Bonus(1)



Run a large file (1G) from our HDFS, since there was an error when I ran the code in Jupyter notebook, so, I submit the task from the console.



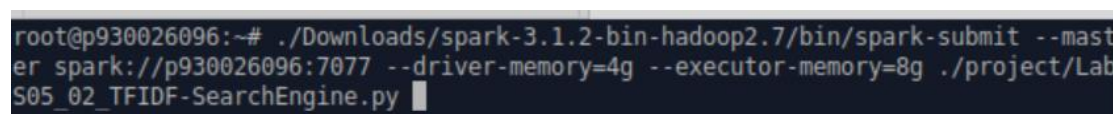Here is a screen shout of my file directory in HDFS which contains more than 1G.



Submit the file in the console by setting the driver memory and executor memory.

| Job Id ▾ | Description | Submitted | Duration | Stages: Succeeded/Total | Tasks (for all stages): Succeeded/Total |
|---|---|---|---|---|---|
| 4 | save at NativeMethodAccessorImpl.java:0<br>save at NativeMethodAccessorImpl.java:0 (kill) | 2021/12/02 16:19:38 | 2 s | 0/5 | 0/16 |

▾ Completed Jobs (4)

| Job Id ▾ | Description | Submitted | Duration | Stages: Succeeded/Total | Tasks (for all stages): Succeeded/Total |
|---|---|---|---|---|---|
| 3 | runJob at PythonRDD.scala:166<br>runJob at PythonRDD.scala:166 | 2021/12/02 16:18:58 | 37 s | 2/2 (3 skipped) | 5/5 (8 skipped) |
| 2 | sortByKey at /root/./project/LabS05_02_TFIDF-SearchEngine.py:33<br>sortByKey at /root/./project/LabS05_02_TFIDF-SearchEngine.py:33 | 2021/12/02 16:18:38 | 20 s | 1/1 (3 skipped) | 4/4 (8 skipped) |
| 1 | sortByKey at /root/./project/LabS05_02_TFIDF-SearchEngine.py:33<br>sortByKey at /root/./project/LabS05_02_TFIDF-SearchEngine.py:33 | 2021/12/02 16:03:59 | 15 min | 4/4 | 12/12 |
| 0 | count at /root/./project/LabS05_02_TFIDF-SearchEngine.py:24<br>count at /root/./project/LabS05_02_TFIDF-SearchEngine.py:24 | 2021/12/02 16:03:50 | 9 s | 1/1 | 2/2 |

The result of task. It takes approximately 16 mins.

Here is the output files of the codes.

Bonus（2）

```
In [2]: from pyspark import SparkConf, SparkContext
        from pyspark.sql import SparkSession, SQLContext
        import os, re
        import math
        import jieba

        sc = SparkContext.getOrCreate(SparkConf())
        sql = SQLContext(sc)

        data = sc.wholeTextFiles('./chinese_books')

        numFiles = data.count()

        wordcount = data.flatMap(lambda x: [((os.path.basename(x[0]) ,i) ,1) for i in jieba.cut(x[1])])\
                .reduceByKey(lambda a, b: a + b)
        wordcount.collect()

        tf = wordcount.map(lambda x: (x[0][1],(x[0][0],x[1])))

        idf = wordcount.map(lambda x: (x[0][1], (x[0][0], x[1], 1)))\
                .map(lambda x: (x[0], x[1][2]))\
                .reduceByKey(lambda x, y: x + y)\
                .map(lambda x: (x[0], math.log10(numFiles / x[1])))

        #Slightly modified map output as (doc, (term, tfidf))
        tfidf = tf.join(idf)\
                    .map(lambda x: (x[1][0][0], (x[0], x[1][0][1] * x[1][1])))\
                    .sortByKey()

        #Then we convert the TF-IDF to an DF, and save to the disk
        lines = tfidf.map(lambda x: (x[0], x[1][0], x[1][1])).toDF()
        lines.write.save("tfidf-index_chinese")
```

Run the code based on some Chinese books by using the package jieba.

```
lines.write.save( tfidf-index_chinese )

/usr/local/lib/python3.8/dist-packages/pyspark/sql/context.py:77: FutureWarning: Deprecated in 3.0.0. Use SparkSessi
on.builder.getOrCreate() instead.
  warnings.warn(
Building prefix dict from the default dictionary ...
Building prefix dict from the default dictionary ...
Dumping model to file cache /tmp/jieba.cache                    (0 + 2) / 2]
Loading model cost 0.750 seconds.
Prefix dict has been built successfully.
Dumping model to file cache /tmp/jieba.cache
Loading model cost 0.750 seconds.
Prefix dict has been built successfully.
```

The program is running.



The files that generate by the codes.