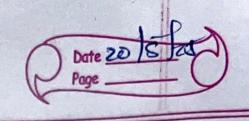
Lab - 10 Rach word appears in a file & write out a 11st of woods whose rount is strictly quester Them be using spark word court py from pysparn import SparkContext sc = SparkContext ("file1.txt") counts = text-file.flotmap (landa: line: line split() · map (landa word: (word, 1))
· rinture By Key Clanda a, b: a+b



output = counts.collect()

for (mord, count) in output:

print (f' Enerved 3: Ecount 3")

Output;

hello: 2 hello brusce,

brusce: 3 brusce how are you

how 3 = 1 brusee

are: 1

you; 1

```
word count.py
 Open ~
1 from pyspark import SparkContext
3 sc = SparkContext("local", "SimpleWordCount")
4 text file = sc.textFile("file1.txt")
5 counts = text file.flatMap(lambda line: line.split()) \
                     .map(lambda word: (word, 1)) \
                     .reduceByKey(lambda a, b: a + b)
8 output = counts.collect()
10 for (word, count) in output:
      print(f"{word}: {count}")
11
```

```
hello: 2
bmsce: 3
how: 1
are: 1
you: 1
```