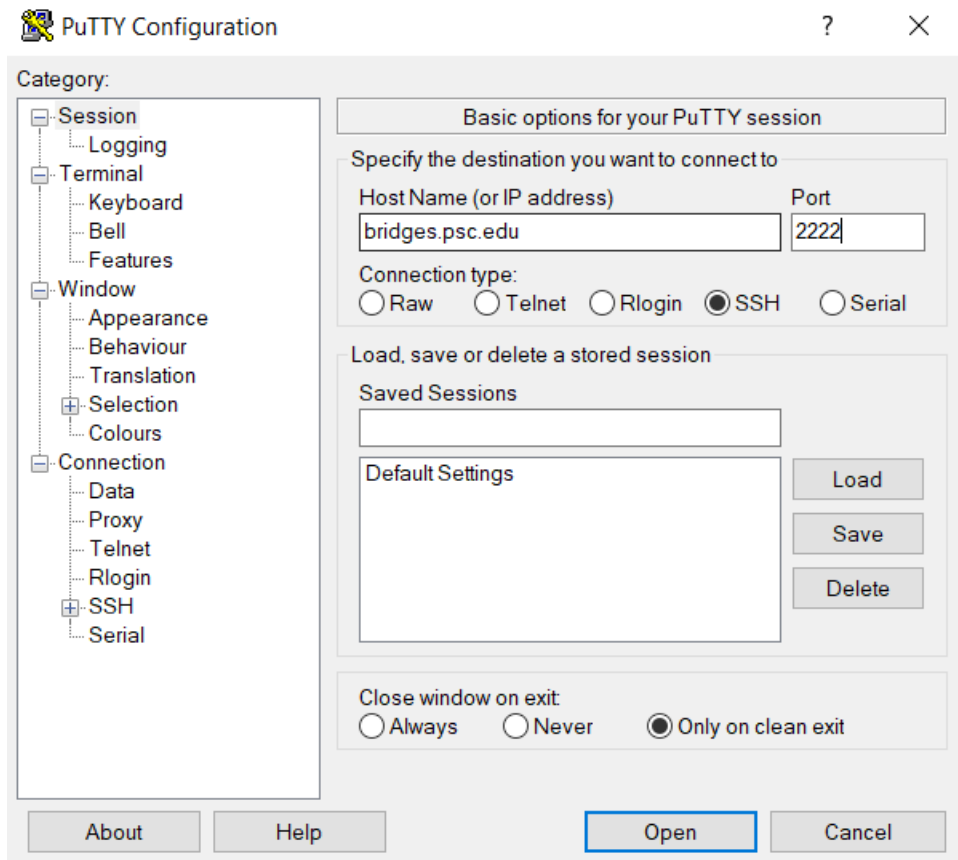Ocean Lu
CS 4080.02
Professor Yang
11.12.2019
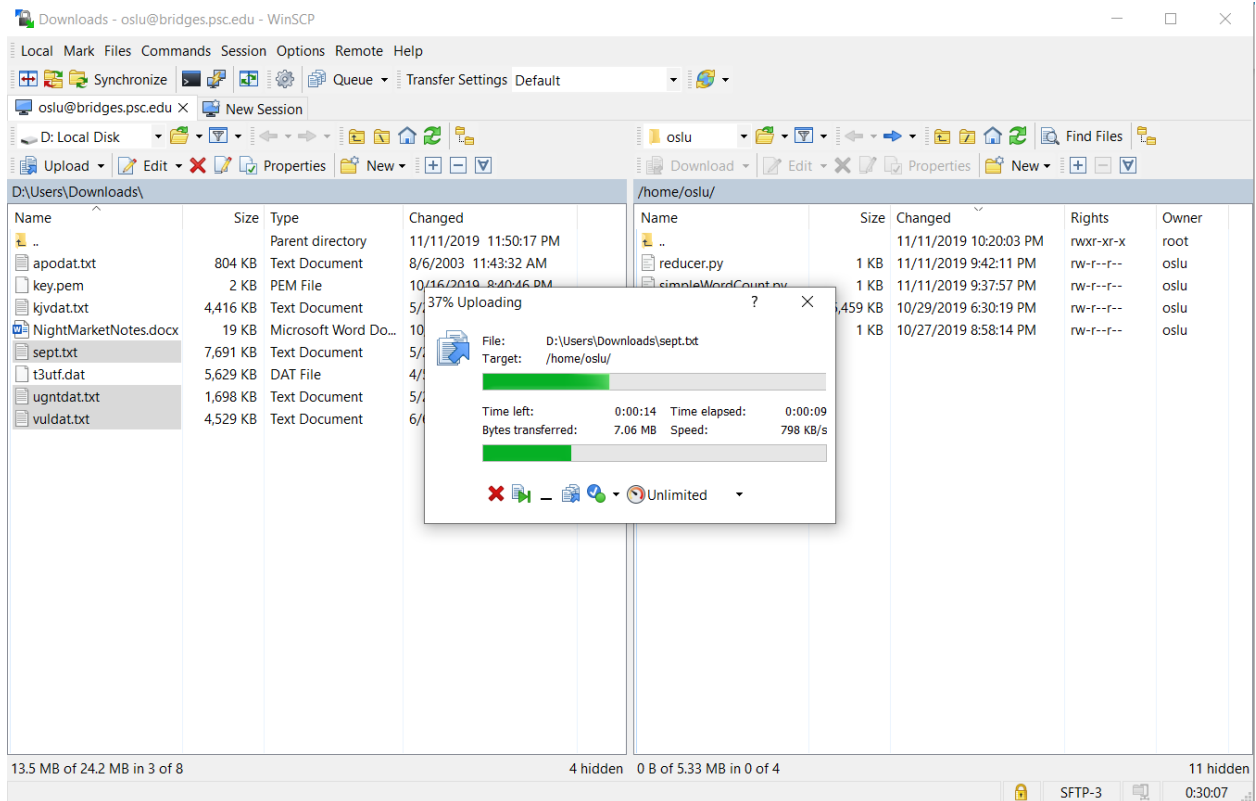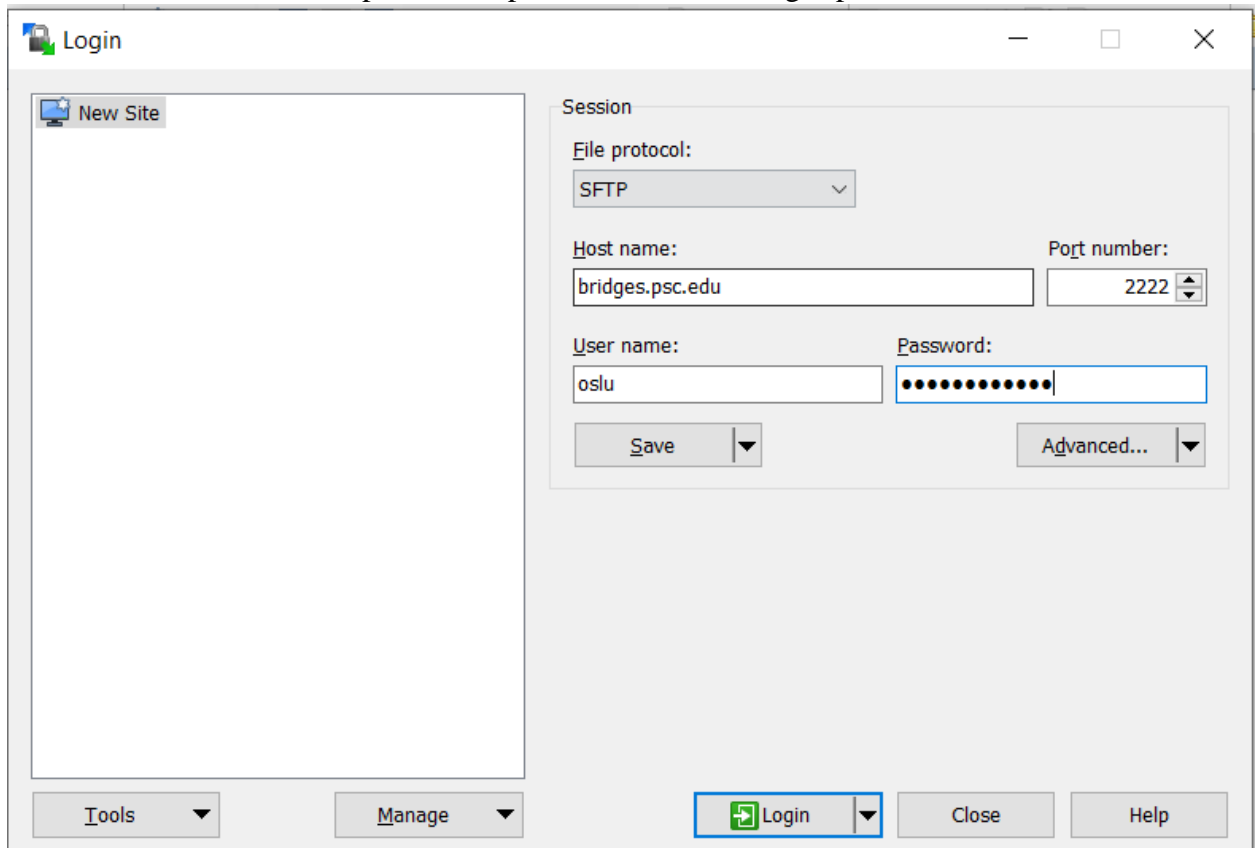
Project 2: MapReduce Programming on Hadoop and Spark

Task 1:

1. Log into bridges.psc.edu port 2222; used puTTY



2. Went to the bible site and downloaded the zip files

3. WinSCP the files of Shakespeare and apodat.txt into the bridges.psc.edu

Confirming that the file has been saved:

```
[oslu@r145 ~]$ ls
apodat.txt               kjvdat.txt   reducer.py   simpleWordCount.py   vuldat.txt
bridges.psc.xsede.org    mapper.py    sept.txt     ugntdat.txt
[oslu@r145 ~]$
```

4. I edited my reducer.py:

```
reducer.py ×
D: > Users > Downloads > reducer.py
 1   #!/usr/bin/env python
 2   # coding: utf-8
 3
 4   import sys
 5
 6   totalCount = 0
 7   prevWord = None
 8
 9   # "line" will iterate over each line of text read from standard input (Will receive intermediate key-value pairs from mapper after the sort and shuffle phase)
10   for line in sys.stdin:
11   # Saving key-value pair into array where the first element is the word and the second is 1
12     line = line.split()
13     #if we have encountered our first word yet and if the previoius word was different from the one currently being read in output the word and the total number of occurence
14     if prevWord and prevWord != line[0]:
15       print("%s\t%s"% (prevWord, totalCount))
16       totalCount = 0
17
18     prevWord = line[0]
19     totalCount += float(line[1])
20   # Ensures that we output the last line from standard input
21   if prevWord != None:
22     print("%s\t%s"% (prevWord, totalCount))
23
24   # CS 4080.02 Project 2 Part 1
25   for words in sorted(dictionary, key=dictionary.get, reverse=True):
26     print("%s\t%s"% (words,dictionary[words]))
```

5. I also scped the reducer.py over with winscp
6. Startup Hadoop (interact -N 4 -t 01:00:00, module load hadoop, start-hadoop.sh)
7. Load all bible text files and reducer.py to the HDFS storing it in the 'in' directory

```
[oslu@r437 ~]$ ls
apodat.txt               kjvdat.txt   reducer.py   simpleWordCount.py   vuldat.txt
bridges.psc.xsede.org    mapper.py    sept.txt     ugntdat.txt
[oslu@r437 ~]$ hadoop fs -put kjvdat.txt in
[oslu@r437 ~]$ hadoop fs -put vuldat.txt in
[oslu@r437 ~]$ hadoop fs -put sept.txt in
[oslu@r437 ~]$ hadoop fs -put ugntdat.txt in
[oslu@r437 ~]$ hadoop fs -ls in
Found 4 items
-rw-r--r--   2 oslu supergroup    4521345 2019-11-12 02:56 in/kjvdat.txt
-rw-r--r--   2 oslu supergroup    7875557 2019-11-12 02:56 in/sept.txt
-rw-r--r--   2 oslu supergroup    1738497 2019-11-12 02:56 in/ugntdat.txt
-rw-r--r--   2 oslu supergroup    4637519 2019-11-12 02:56 in/vuldat.txt
[oslu@r437 ~]$
```

8. Run commands such as:
   hadoop jar /opt/packages/hadoop-testing/hadoop/hadoop-2.7.3/share/hadoop/tools/lib/hadoop-streaming-2.7.3.jar -input in/apodat.txt -output out/ApodatOutput.txt -mapper mapper.py -file /home/oslu/mapper.py -reducer reducer.py -file /home/oslu/reducer.py
   to get output data for the bible files (same general commands)

oslu@br018:~

```
[oslu@r338 ~]$ hadoop jar /opt/packages/hadoop-testing/hadoop/hadoop-2.7.3/share
/hadoop/tools/lib/hadoop-streaming-2.7.3.jar -input in/apodat.txt -output out/Ap
odatOutput.txt -mapper mapper.py -file /home/oslu/mapper.py -reducer reducer.py
-file /home/oslu/reducer.py
19/11/12 01:14:47 WARN streaming.StreamJob: -file option is deprecated, please u
se generic option -files instead.
packageJobJar: [/home/oslu/mapper.py, /home/oslu/reducer.py, /tmp/hadoop-unjar79
57663381218299414/] [] /tmp/streamjob3042388271063423551.jar tmpDir=null
19/11/12 01:14:48 INFO client.RMProxy: Connecting to ResourceManager at r338.opa
.bridges.psc.edu/10.4.117.86:8032
19/11/12 01:14:48 INFO client.RMProxy: Connecting to ResourceManager at r338.opa
.bridges.psc.edu/10.4.117.86:8032
19/11/12 01:14:54 INFO mapred.FileInputFormat: Total input paths to process : 1
19/11/12 01:14:58 INFO mapreduce.JobSubmitter: number of splits:2
19/11/12 01:15:00 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_15
73539202748_0001
19/11/12 01:15:00 INFO impl.YarnClientImpl: Submitted application application_15
73539202748_0001
19/11/12 01:15:00 INFO mapreduce.Job: The url to track the job: http://r338.opa.
bridges.psc.edu:8088/proxy/application_1573539202748_0001/
19/11/12 01:15:00 INFO mapreduce.Job: Running job: job_1573539202748_0001
19/11/12 01:15:13 INFO mapreduce.Job: Job job_1573539202748_0001 running in uber
 mode : false
19/11/12 01:15:13 INFO mapreduce.Job:  map 0% reduce 0%
```

oslu@br018:~

```
                CPU time spent (ms)=4300
                Physical memory (bytes) snapshot=1203040256
                Virtual memory (bytes) snapshot=19310891008
                Total committed heap usage (bytes)=4506779648
        Shuffle Errors
                BAD_ID=0
                CONNECTION=0
                IO_ERROR=0
                WRONG_LENGTH=0
                WRONG_MAP=0
                WRONG_REDUCE=0
        File Input Format Counters
                Bytes Read=826377
        File Output Format Counters
                Bytes Written=86998
19/11/12 01:15:37 INFO streaming.StreamJob: Output directory: out/ApodatOutput.t
xt
[oslu@r338 ~]$ hdfs dfs -ls out/ApodatOutput.txt/
Found 2 items
-rw-r--r--   2 oslu supergroup          0 2019-11-12 01:15 out/ApodatOutput.txt/
_SUCCESS
-rw-r--r--   2 oslu supergroup      86998 2019-11-12 01:15 out/ApodatOutput.txt/
part-00000
[oslu@r338 ~]$
```
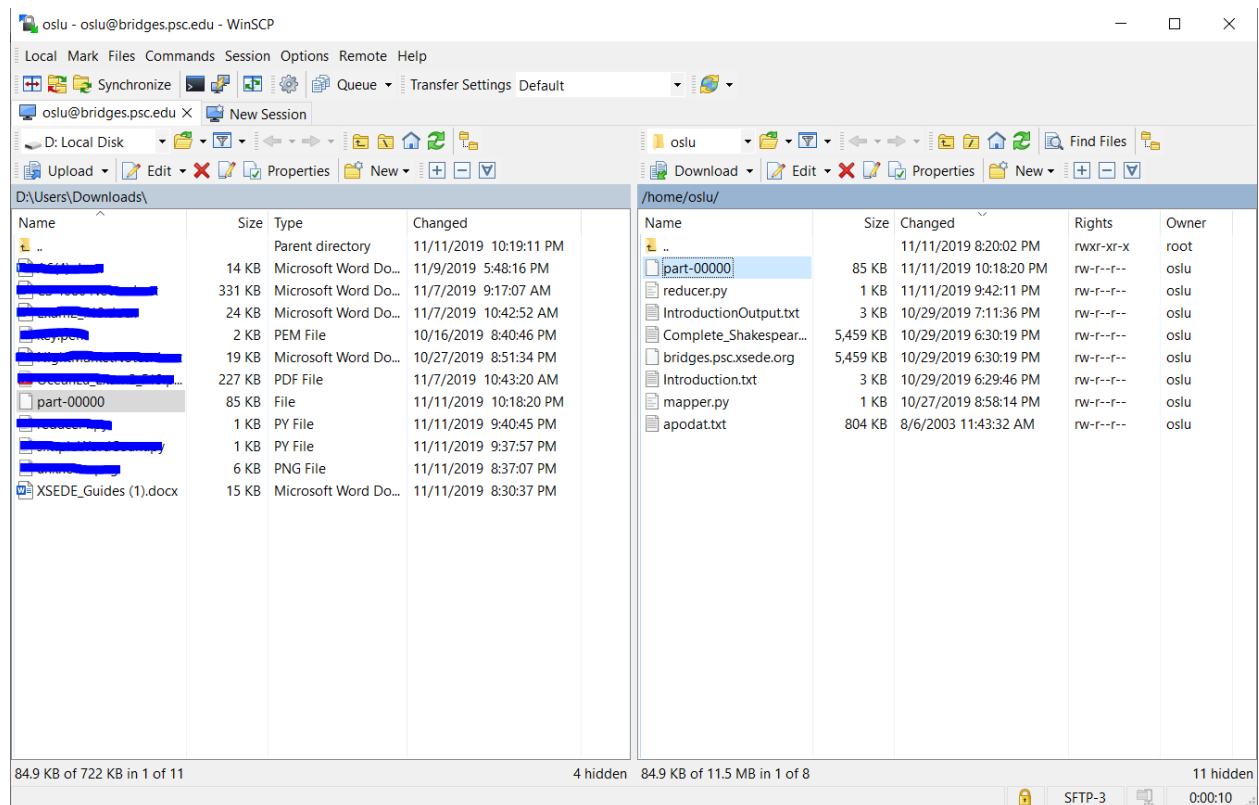
9.  Retrieve the file with the command: hadoop fs -get out/ApodatOutput.txt/part-00000
    /home/oslu
    *Make sure part-00000 other files are deleted, this does not override current files that

already exist



```
                    Physical memory (bytes) snapshot=1203040256
                    Virtual memory (bytes) snapshot=19310891008
                    Total committed heap usage (bytes)=4506779648
            Shuffle Errors
                    BAD_ID=0
                    CONNECTION=0
                    IO_ERROR=0
                    WRONG_LENGTH=0
                    WRONG_MAP=0
                    WRONG_REDUCE=0
            File Input Format Counters
                    Bytes Read=826377
            File Output Format Counters
                    Bytes Written=86998
19/11/12 01:15:37 INFO streaming.StreamJob: Output directory: out/ApodatOutput.t
xt
[oslu@r338 ~]$ hdfs dfs -ls out/ApodatOutput.txt/
Found 2 items
-rw-r--r--   2 oslu supergroup          0 2019-11-12 01:15 out/ApodatOutput.txt/
_SUCCESS
-rw-r--r--   2 oslu supergroup      86998 2019-11-12 01:15 out/ApodatOutput.txt/
part-00000
[oslu@r338 ~]$ hadoop fs -get out/ApodatOutput.txt/part-00000 /home/oslu
[oslu@r338 ~]$
```
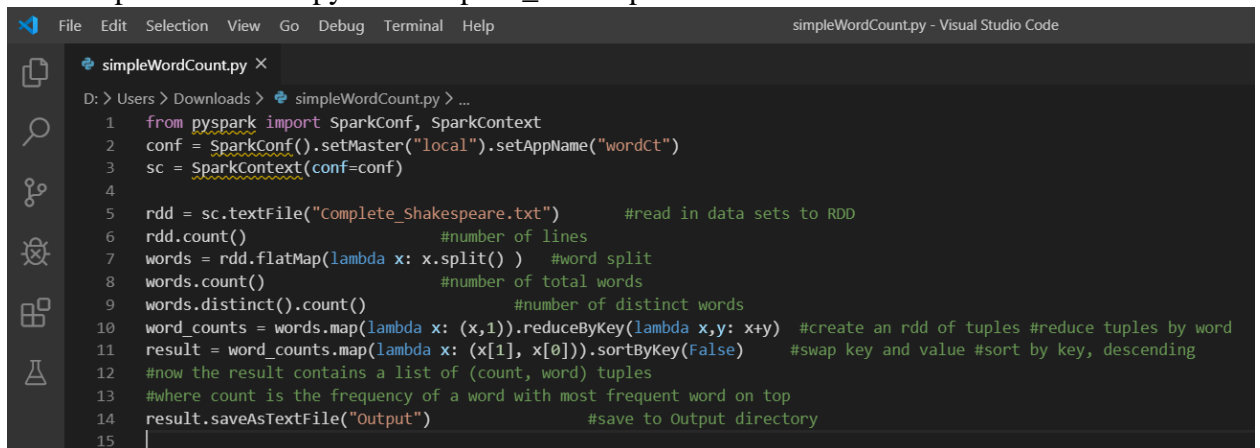
10. WinSCP over it to local host



Finished, now we have the output file of apodat.txt

Task 2:

1. Simple counting program given to us on in class slides with XSEDE & Spark notes, renamed "simpleWordCount.py"
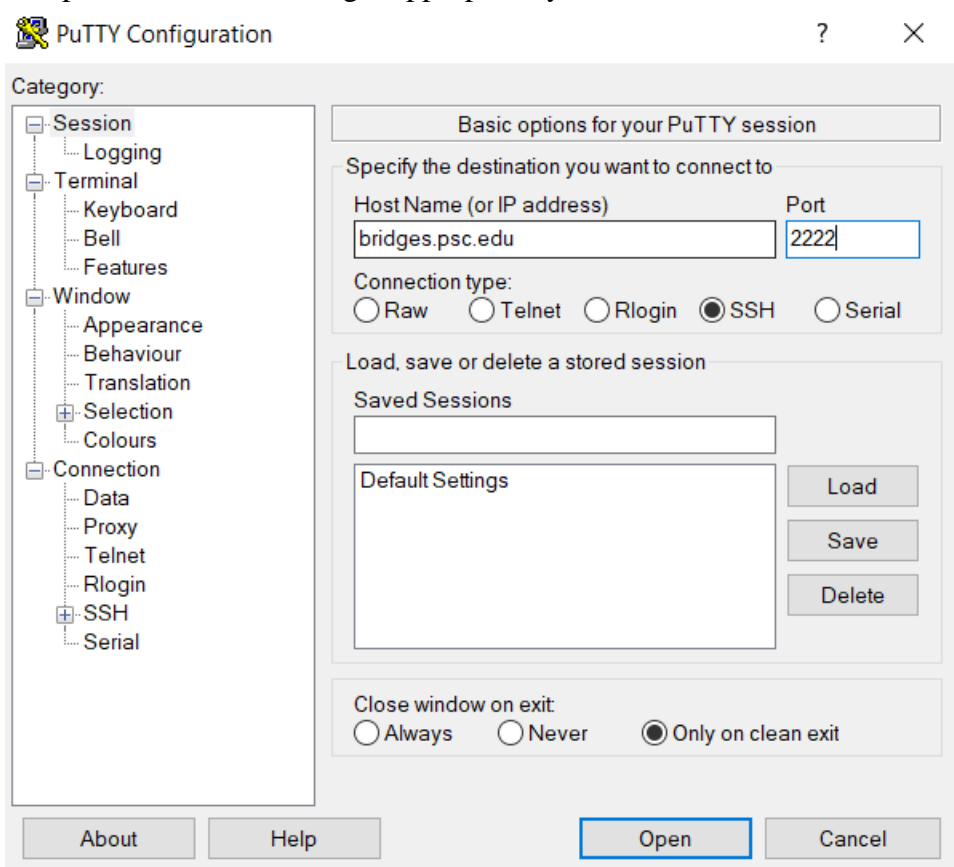


2. Take the SimpleWordCount.py and Complete_Shakespeare.txt and Bibles data files, and winscp it onto the drive

3. Edit simpleWordCount.py for Complete_Shakespeare.txt



Edit the sc.textfile("") appropriately when going through the bible data files
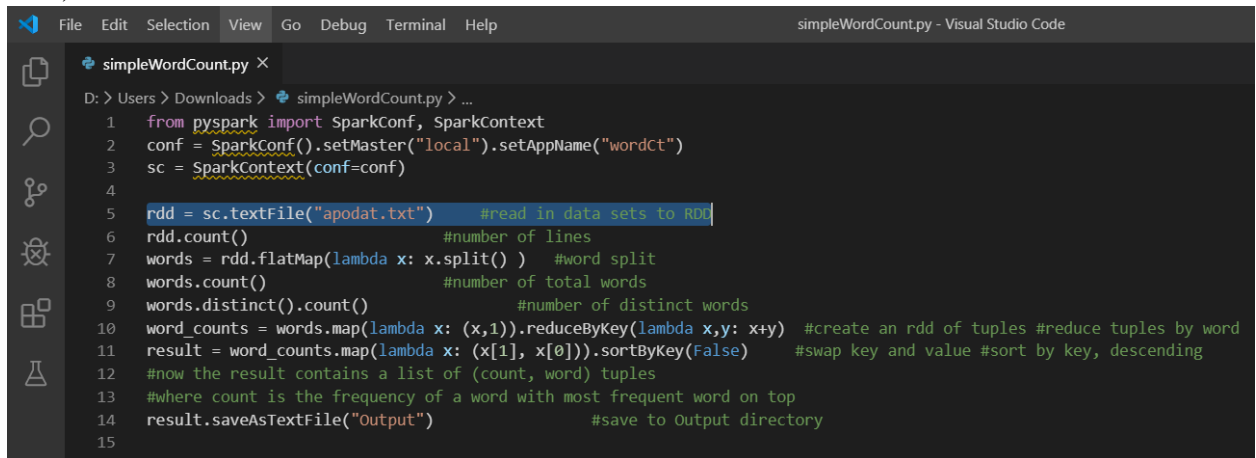
4. Use puTTY to SSH and login appropriately



5. Go into Hadoop and load spark and have the command: spark-submit simpleWordCount.py
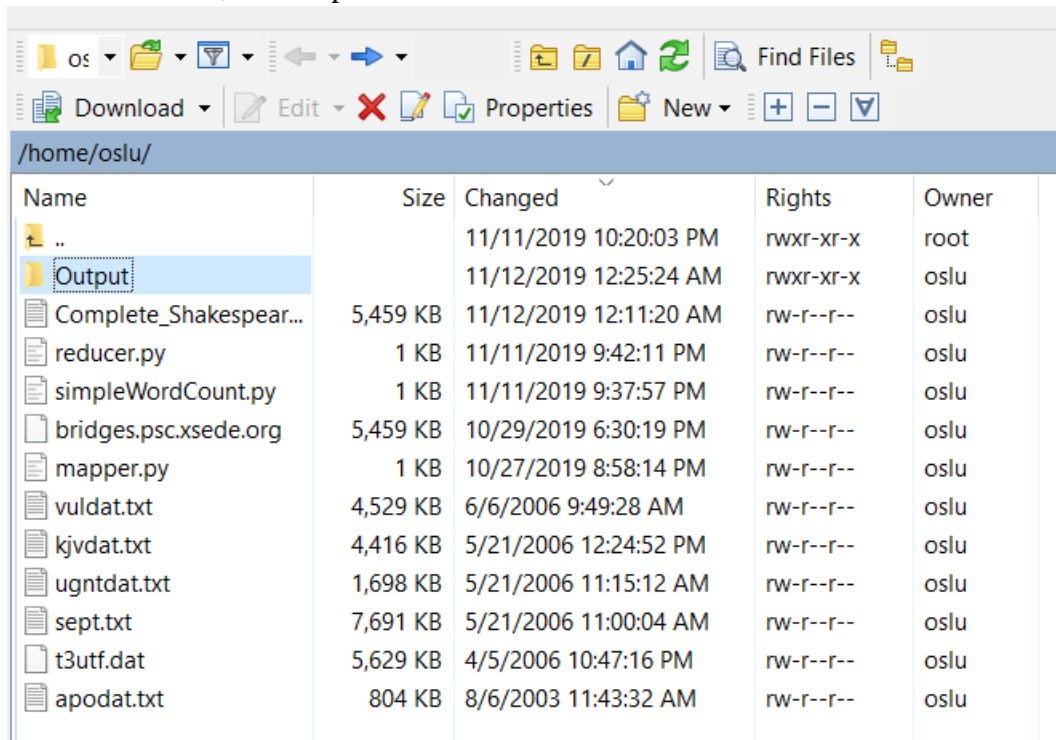
6. Edit simpleWordCount.py's sc.textfile approrpiately when changing files (into bible files)

```python
from pyspark import SparkConf, SparkContext
conf = SparkConf().setMaster("local").setAppName("wordCt")
sc = SparkContext(conf=conf)

rdd = sc.textFile("apodat.txt")      #read in data sets to RDD
rdd.count()                    #number of lines
words = rdd.flatMap(lambda x: x.split() )    #word split
words.count()                  #number of total words
words.distinct().count()           #number of distinct words
word_counts = words.map(lambda x: (x,1)).reduceByKey(lambda x,y: x+y)  #create an rdd of tuples #reduce tuples by word
result = word_counts.map(lambda x: (x[1], x[0])).sortByKey(False)     #swap key and value #sort by key, descending
#now the result contains a list of (count, word) tuples
#where count is the frequency of a word with most frequent word on top
result.saveAsTextFile("Output")            #save to Output directory
```

7. Do it for all files, and output should be in a folder form.

/home/oslu/

| Name | Size | Changed | Rights | Owner |
|------|------|---------|--------|-------|
| .. | | 11/11/2019 10:20:03 PM | rwxr-xr-x | root |
| Output | | 11/12/2019 12:25:24 AM | rwxr-xr-x | oslu |
| Complete_Shakespear... | 5,459 KB | 11/12/2019 12:11:20 AM | rw-r--r-- | oslu |
| reducer.py | 1 KB | 11/11/2019 9:42:11 PM | rw-r--r-- | oslu |
| simpleWordCount.py | 1 KB | 11/11/2019 9:37:57 PM | rw-r--r-- | oslu |
| bridges.psc.xsede.org | 5,459 KB | 10/29/2019 6:30:19 PM | rw-r--r-- | oslu |
| mapper.py | 1 KB | 10/27/2019 8:58:14 PM | rw-r--r-- | oslu |
| vuldat.txt | 4,529 KB | 6/6/2006 9:49:28 AM | rw-r--r-- | oslu |
| kjvdat.txt | 4,416 KB | 5/21/2006 12:24:52 PM | rw-r--r-- | oslu |
| ugntdat.txt | 1,698 KB | 5/21/2006 11:15:12 AM | rw-r--r-- | oslu |
| sept.txt | 7,691 KB | 5/21/2006 11:00:04 AM | rw-r--r-- | oslu |
| t3utf.dat | 5,629 KB | 4/5/2006 10:47:16 PM | rw-r--r-- | oslu |
| apodat.txt | 804 KB | 8/6/2003 11:43:32 AM | rw-r--r-- | oslu |