# Hadoop Python MapReduce

## Mapping Function

```python
#! usr/bin/env  python
import sy
for line in sys.stdin:
  line = line.strip()
  words = line.split()
  for word in words:
    print("%s\t%s" %(word, "1"))
```

And then we control + x exit

cat test.txt| python wc_mapper.py

## Reduce Function

```python
#read values
nano wc_reduce.py
#! usr/bin/env  python
import sys
word2count = {}
for line in sys.stdin:
  word,count = line.strip().split('\t')
  try:
    count = int(count)
  except ValueError:
    continue
  try:
    word2count[word]=word2count[word]+count
  except:
    word2count[word]=count
for word in word2count.keys():
  print("%s\t%s" %(word, word2count[word]))

```

cat test.txt | python wc_mapper.py|python wc_reduce.py
hdfs dfs - put test.txt
#Create the batch
nano runmr.sh
Write:

```bash
#! /bin/bash
hadoop jar blablabla
\" continue in the next line
– input /user/cloudera/text.txt
```

```
 5 - output /user/cloudera/wc_output
 6 - file wc_mapper.py
 7 - file wc_reducer.py
 8 - mapper "python wc_mapper.py"
 9 - reducer "python wc_reduce.py"
10 # if you run into error
11 chmod +x wc_mapper.py
```

```
1 sort -k 2rn
2 #This will deriectly go into the reduce file.
3 os.system('|sort -k 2rn')
4 #By default, 3 sections sent to the HDFS server.
5 #The sort funtion can only sort the chunk in that server
```