

Q1. Word Count Program

mapper.py

```
import sys
for line in sys.stdin: # input comes from STDIN (standard input)
    line = line.strip() # remove leading and trailing whitespace
    words = line.split() # split the line into words
    # increase counters
    for word in words:
        # write the results to STDOUT (standard output);
        # what we output here will be the input for the
        # Reduce step, i.e. the input for reducer.py #
        # tab-delimited; the trivial word count is 1
    print("%s\t%s" %(word, 1))
```

reducer.py

```
from operator import itemgetter
import sys
current_word = None
current_count = 0
word = None
# input comes from STDIN
for line in sys.stdin:
    # remove leading and trailing whitespace
    line = line.strip()
    # parse the input we got from mapper.py
    word, count = line.split('\t', 1)
    # convert count (currently a string) to int
    try:
        count = int(count)
    except ValueError:
        # count was not a number, so silently
        # ignore/discard this line
        continue
    # this IF-switch only works because Hadoop sorts map output
    # by key (here: word) before it is passed to the reducer
    if current_word == word:
        current_count += count
    else:
        if current_word:
            # write result to STDOUT
            print('%s\t%s' % (current_word, current_count) )
        current_count = count
        current_word = word
    # do not forget to output the last word if needed!
    if current_word == word:
        print('%s\t%s' % (current_word, current_count))
```



```
# output last word
if lastWord == word:
print( '%s\t%s' % (lastWord, sum ) )
```

freqmap2.py

```
from __future__ import print_function
import sys
# input comes from STDIN (standard input)
for line in sys.stdin:
    # print(line.strip().split())
    word, count = line.strip().split('\t')
    count = int(count)
    print( '%d,%s' % (count, word) )
```

freqred2.py

```
from __future__ import print_function
import sys
mostFreq = []
currentMax = -1
for line in sys.stdin:
    count, word = line.strip().split(', ', 1)
    count = int(count)
    if count > currentMax:
        currentMax = count
        mostFreq = [ word ]
    elif count == currentMax:
        mostFreq.append( word )
    # output mostFreq word(s)
    for word in mostFreq:
        print( '%s, %s' % ( word, currentMax ) )
```

Q3. Count and Summary using MapReduce

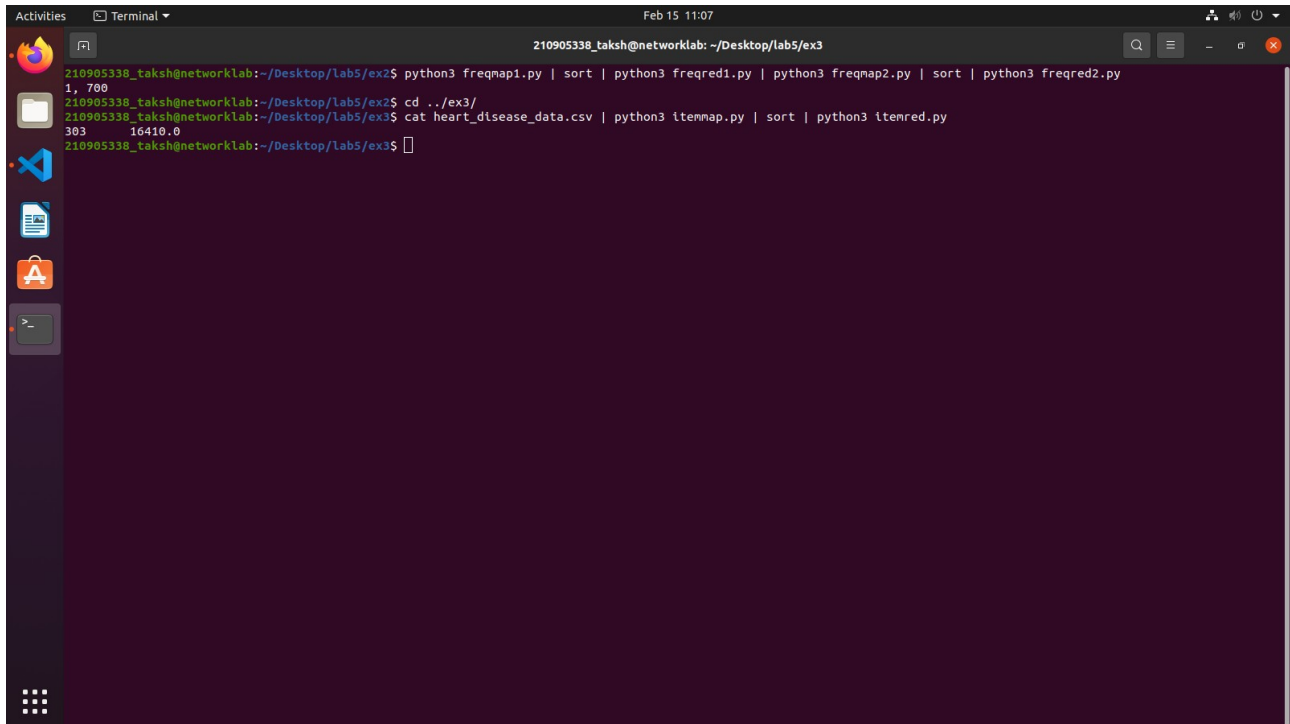
itemmap.py

```
import sys
for line in sys.stdin:
    data = line.strip().split(",")
    if len(data) == 14:
        age,sex,cp,trestbps,chol,fbs,restecg,thalach,exang,oldpeak,slope,ca,thal,target = data
        print (" {0}\t{1}".format(age, target))
```

itemred.py

```
import sys
transactions_count = 0
age_total = 0
for line in sys.stdin:
    data = line.strip().split("\t")
    if len(data) != 2:
```

```
# Something has gone wrong. Skip this line.  
continue  
age, target = data  
transactions_count += 1  
if age.isdigit():  
    age_total += float(age)  
print (transactions_count, "\t", age_total)
```



A terminal window titled '210905338_taksh@networklab: ~/Desktop/lab5/ex3' is shown. The terminal displays the following commands and their outputs:

```
210905338_taksh@networklab:~/Desktop/lab5/ex2$ python3 freqmap1.py | sort | python3 freqred1.py | python3 freqmap2.py | sort | python3 freqred2.py  
1, 700  
210905338_taksh@networklab:~/Desktop/lab5/ex2$ cd ../ex3/  
210905338_taksh@networklab:~/Desktop/lab5/ex3$ cat heart_disease_data.csv | python3 itemmap.py | sort | python3 itemred.py  
303      16410.0  
210905338_taksh@networklab:~/Desktop/lab5/ex3$
```

The terminal window has a dark purple background and a sidebar on the left with icons for various applications. The top of the window shows the date and time as 'Feb 15 11:07'.