# Lab #3 – Hadoop

1. Write Python program for the following three questions of project part 1.

    1) Find a sales breakdown by product category across all of our stores.

    2) Find the monetary value for the highest individual sale for each separate store.

    3) Find the total sales value across all the stores, and the total number of sales. (Assume there is only one reducer.

**2. Take a screen shot after running MapReduce code for question 1. Copy and paste the mapper and reducer code for question 1. Copy and paste the result for question 1.**



## Mapper.py:

```
1.  #!/usr/bin/python
2.
3.
4.  import sys
5.
6.  for line in sys.stdin:
7.      data = line.strip().split("\t")
8.      if len(data) == 6:
9.          date, time, store, item, cost, payment = data
10.         print "{0}\t{1}".format(item, cost)
```

## Reducer.py:

```
1.  #!/usr/bin/python
2.
3.  import sys
4.
5.  salesTotal = 0
6.  oldKey = None
7.
8.
9.  for line in sys.stdin:
10.     data_mapped = line.strip().split("\t")
11.     if len(data_mapped) != 2:
12.         # Something has gone wrong. Skip this line.
13.         continue
14.
15.     thisKey, thisSale = data_mapped
16.
17.     if oldKey and oldKey != thisKey:
18.         print oldKey, "\t", salesTotal
19.         oldKey = thisKey;
20.         salesTotal = 0
21.
22.     oldKey = thisKey
23.     salesTotal += float(thisSale)
24.
25. if oldKey != None:
26.     print oldKey, "\t", salesTotal
```

```
training@localhost:~/udacity_training/code2          _ □ ×

File  Edit  View  Search  Terminal  Help
Albuquerque     499.98
Anaheim         499.98
Anchorage       499.99
Arlington       499.95
Atlanta         499.96
Aurora   499.97
Austin   499.97
Bakersfield     499.97
Baltimore       499.99
Baton Rouge     499.98
Birmingham      499.99
Boise   499.98
Boston  499.99
Buffalo         499.99
Chandler        499.98
Charlotte       499.98
Chesapeake      499.98
Chicago         499.99
Chula Vista     499.99
Cincinnati      499.98
Cleveland       499.98
Colorado Springs        499.99
Columbus        499.98
Corpus Christi  499.96
Dallas   499.99
Denver   499.97
Detroit         499.99
:
```

**3. Take a screen shot after running MapReduce code for question 2. Copy and paste the mapper and reducer code for question 2. What are the values for the following store:**

Anchorage      **499.99**

Bakersfield        **499.97**

Colorado Springs   **499.99**



## **Mapper.py:**

```
1.  #!/usr/bin/python
2.
3.
4.  import sys
5.
6.  for line in sys.stdin:
7.      data = line.strip().split("\t")
8.      if len(data) == 6:
9.          date, time, store, item, cost, payment = data
10.         print "{0}\t{1}".format(store, cost)
```
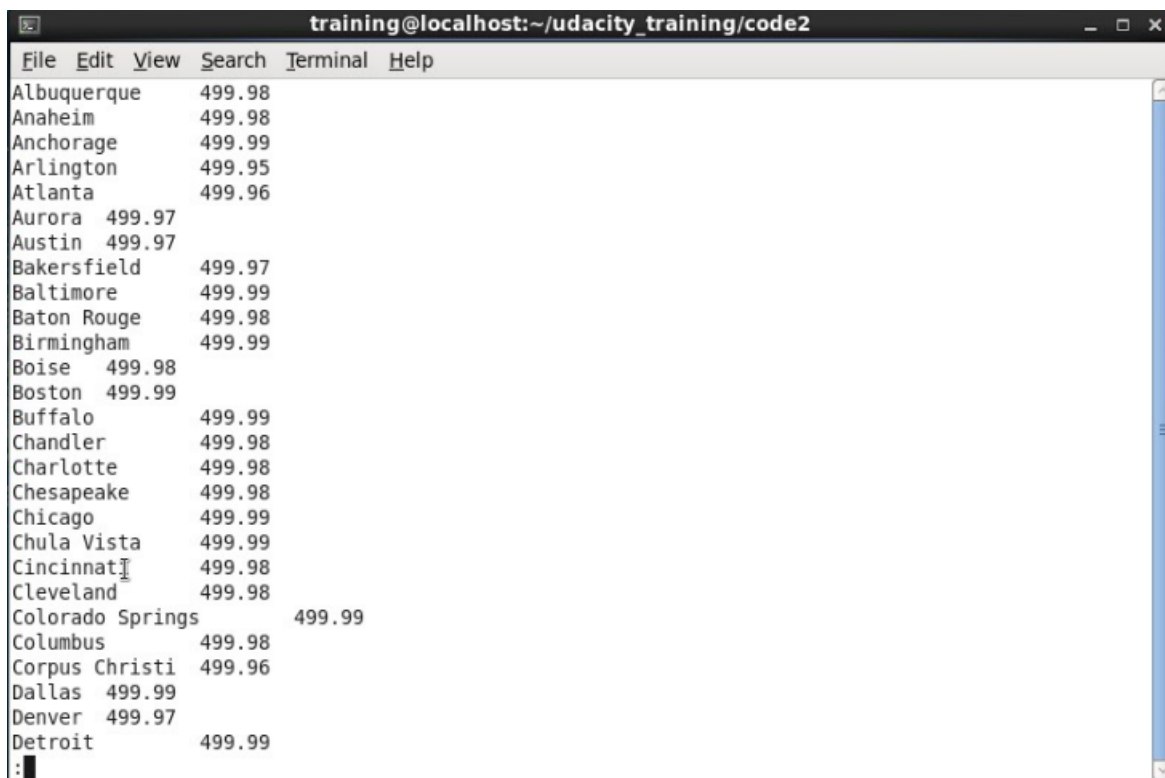
## **Reducer.py:**

```
1.  #!/usr/bin/python
2.
3.  import sys
4.
5.  sales = 0
6.  oldKey = None
7.
8.
9.  for line in sys.stdin:
10.     data_mapped = line.strip().split("\t")
11.     if len(data_mapped) != 2:
12.         continue
13.
```

```
14.    thisKeyStore, thisSale = data_mapped
15.
16.    if oldKey and oldKey != thisKeyStore:
17.        print oldKey, "\t", sales
18.        oldKey = thisKeyStore;
19.        sales = 0
20.
21.    oldKey = thisKeyStore
22.    #sales += float(thisSale)
23.    if sales < float(thisSale):
24.    sales = float(thisSale)
25.
26. if oldKey != None:
27.    print oldKey, "\t", sales
```

```
training@localhost:~/udacity_training/code2          _ □ ✕

File  Edit  View  Search  Terminal  Help
Albuquerque      499.98
Anaheim          499.98
Anchorage        499.99
Arlington        499.95
Atlanta          499.96
Aurora   499.97
Austin   499.97
Bakersfield      499.97
Baltimore        499.99
Baton Rouge      499.98
Birmingham       499.99
Boise    499.98
Boston   499.99
Buffalo          499.99
Chandler         499.98
Charlotte        499.98
Chesapeake       499.98
Chicago          499.99
Chula Vista      499.99
Cincinnat I      499.98
Cleveland        499.98
Colorado Springs         499.99
Columbus         499.98
Corpus Christi   499.96
Dallas   499.99
Denver   499.97
Detroit          499.99
:
```

4. Take a screen shot after running MapReduce code for question 3. Copy and paste the mapper and reducer code for question 3. What is the total number of sales and the total sales value from all the stores?

## Mapper.py:

```python
1.  #!/usr/bin/python
2.
3.
4.  import sys
5.
6.  for line in sys.stdin:
7.      data = line.strip().split("\t")
8.      if len(data) == 6:
9.          date, time, store, item, cost, payment = data
10.         print "{0}\t{1}".format(store, cost)
```

## Reducer.py:

```python
1.  #!/usr/bin/python
2.
3.  import sys
4.
5.  salesTotal = 0
6.  numOfSales = 0
7.
8.  for line in sys.stdin:
9.      data_mapped = line.strip().split("\t")
10.     if len(data_mapped) != 2:
11.         continue
12.
13.     thisKeyStore, thisCost = data_mapped
14.
15.     if thisKeyStore in line:
16.     numOfSales += 1
```

```
17.         salesTotal += float(thisCost)
18.
19.
20.
21. print salesTotal, "\t", numOfSales
```