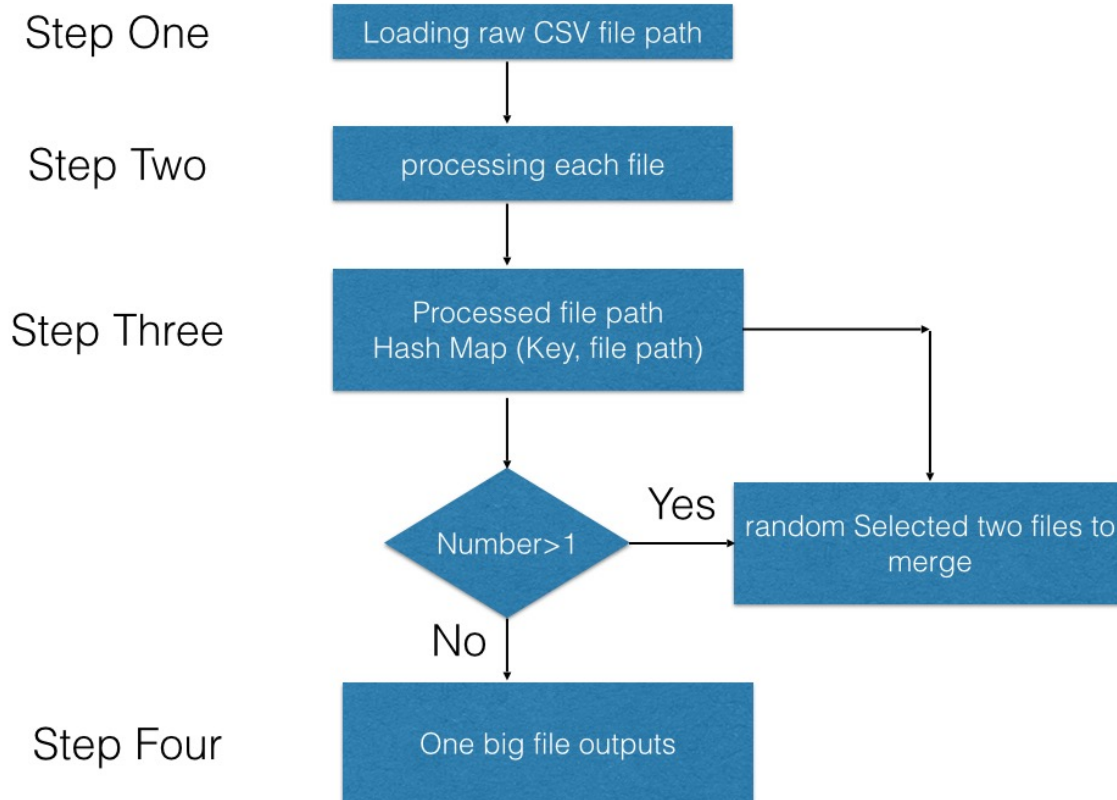Project 2:

Data aggregation:

The code flow chart is below ( No VIVO in my personal computer, sorry for that):

| Step One | Loading raw CSV file path |
|----------|---------------------------|

Step Two — processing each file

Step Three — Processed file path Hash Map (Key, file path)

Number>1 — Yes → random Selected two files to merge

No

Step Four — One big file outputs

## Step One: Loading raw CSV file path:

Input: raw CSV data file path, it may has multiply paths. In this code, we consider all CSV files are in one folder.

Outputs: lists of CSV data files

Code:

```python
1  # -*- coding: utf-8 -*-
2  """
3  Created on Sun Sep 20 11:17:41 2015
4
5  @author: weizhi
6  """
7
8
9  #%% loading each files
10 import glob, os
11 import numpy as np
12 import random
13 import pandas as pd
14
15 # https://docs.python.org/2/library/os.html
16 def findFilePath(path):
17     """
18     Input: raw CSV data path
19     Output: list of file path
20
21     """
22     os.chdir(path)
23     filePaths = []
24     for file in glob.glob("*.csv"):
25         filePaths.append(file)
26     return filePaths
27
```

```python
140 if __name__ == '__main__':
141     #%% raw csv data
142     path = '/Users/weizhi/Desktop/everStringProject'
143     filePaths = findFilePath(path)
```

## Step 2: Process each file

```python
29 # deal with each csv file
30 class logDataAnalysis(object):
31
32     """
33     From raw file to procssed outputs
34     """
35
36     def __init__(self,path,hour,fileName):
37         """
38         Path: csv files store
39         Hour: every # hours, like 2, 00-01 (00, 01), 02-03 (02,03),.....
40                                   like 3, 00-02, 03-05, .....
41         fileName: save file name after the processing
42         """
43
44         self.path = path
45         self.hour = hour
46         self.fileName = fileName
47
48     def readCSV(self,filePath):
49         data = pd.read_csv(filePath)
50         return data
51     def transfromTimeFormat(self,data):
52         timeColumn = data[data.keys()[-1]]
53         for i in range(len(timeColumn)):
54             item = timeColumn.loc[i]
55             T1 = item.split(' ')
56             T2 = T1[1].split(':')[0]
57             binTime = int(T2)/(self.hour)
58             leftBin = str(binTime*(self.hour)).zfill(2)
59             rightBin = str((binTime+1)*(self.hour)-1).zfill(2)
60             T3 = leftBin + '_' + rightBin
61             T4 = T1[0] + '_'+ T3
62             timeColumn.loc[i] = T4
63         data[data.keys()][-1] = timeColumn
64         return data
```

```python
66     def generateOutputs(self,path):
67         """
68         Input: one raw CSV file path
69
70         Output: processing outputs: import pandas as pd
71                 And, the file save path
72
73         """
74         data = self.transfromTimeFormat(self.readCSV(path))
75         keys = [key for key in data.keys()]
76         keys.reverse()
77         dataGroup = data.groupby(keys).groups
78         keysOuput = sorted(dataGroup.iterkeys())  # keep the keys sorted rather than hashing
79         outputs = pd.DataFrame(columns = ['period','content_id','uid','count'])
80         count = 0
81         for key in keysOuput:
82             curr = list(key)    # write to each columns to outputs.csv
83             curr.append(len(dataGroup[key])) # get the count of keys from groupby
84             outputs.loc[count] = curr
85             count +=1
86         saveFile = self.createSavePath() + '/'+'Output_'+ str(self.hour) + 'hours_'+ self.fileName
87     #   print savePath
88     #     saveFile = savePath
89         outputs.to_csv(saveFile,index=False)
90         return outputs,saveFile
91
92     def createSavePath(self):
93         savePath = self.path + '/' + 'Outputs'
94         try:
95             os.makedirs(savePath)
96         except OSError:
97             pass
98         return savePath
```

```python
144   #%% generate the raw data (CSV) to each CSV outputs formate
145   outputIndex = {}   # HashMap, file key: filePath. filePath has the processed file path
146   for index in range(len(filePaths)):
147
148       Obj = logDataAnalysis(path,2,filePaths[index])
149       data = Obj.readCSV(filePaths[0])
150
151       #data.to_csv(filePaths[0],index=False)
152       outputs,savePath = Obj.generateOutputs(filePaths[0])
153       outputIndex[str(index)] = savePath
154   # savePath = Obj.createSavePath()
```

# Step 3: Merge each processed files and generate the output

```python
99  #%% merge the files
100 class mergeFile():
101     def __init__(self,path):
102         """
103         Path, the raw CSV file path
104         """
105         self.path = path
106     def funct(self,df):
107         """
108         Input: dataframe
109         Output: dataframe, data aggregatoin update in count
110         """
111         df['count'] = df['count'].sum()
112         return df
113     def mergerResult(self,fileOne,fileTwo):
114         """
115         Input: fileone, fileTwo: two processing files
116         Output: merge outputs
117         """
118         fileCombine = pd.concat([fileOne,fileTwo])
119         fileCombine = fileCombine.reset_index()
120         column = list(fileOne)
121         result = fileCombine.groupby(column[:3]).apply(self.funct)
122         return result[column]   # get the merge files
123     def savePath(self,mergeResult,name):
124         """
125         input: merge result and name to save outputs
126         output: file save path
127
128         """
129         savePath = self.path + '/' + 'Outputs'
130         try:
131             os.makedirs(savePath)
132         except OSError:
133             pass
134         savePath = savePath + '/' + name
135         mergeResult.to_csv(savePath,index = False)
136         return savePath

155     #%% merge the output from the index, we assume the memory can hold all outputs
156     # merge each two files randomly, from n to n/2, then merge again, n/2 to n/4, .... until to get big files
157     merge = mergeFile(path)
158     while(len(outputIndex.keys())>1):
159         key1 = random.choice(outputIndex.keys()) # random select key
160         fileOne = pd.read_csv(outputIndex[key1])  # read the file
161         os.remove(outputIndex[key1])   # delete the file
162         outputIndex.pop(key1,None) # delete the key
163
164         key2 = random.choice(outputIndex.keys())  # random select key
165         fileTwo = pd.read_csv(outputIndex[key2])  # read the file
166         os.remove(outputIndex[key2])   # delete the file
167         outputIndex.pop(key2,None)  # delete the key
168
169         # merge
170         result = merge.mergerResult(fileOne,fileTwo)
171         # update the key
172         keyNew = key1 + '_' + key2
173         print keyNew
174         savePath = merge.savePath(result,keyNew)
175         outputIndex[keyNew] = savePath
176
177
178
179
180
181
```

Conclusion:

I have duplicated your attached seven times, and run the code above, I have got the the merge files.

Input:

| Name | Date Modified | Size | Kind |
|------|---------------|------|------|
| mapreduce_question_data copy 2.csv | Today, 4:37 PM | 3 KB | comm...values |
| mapreduce_question_data copy 3.csv | Yesterday, 2:56 PM | 3 KB | comm...values |
| mapreduce_question_data copy 4.csv | Yesterday, 2:56 PM | 3 KB | comm...values |
| mapreduce_question_data copy 5.csv | Yesterday, 2:56 PM | 3 KB | comm...values |
| mapreduce_question_data copy 6.csv | Yesterday, 2:56 PM | 3 KB | comm...values |
| mapreduce_question_data copy 7.csv | Yesterday, 2:56 PM | 3 KB | comm...values |
| mapreduce_question_data copy.csv | Yesterday, 2:56 PM | 3 KB | comm...values |

Outputs:

| Name | Date Modified | Size | Kind |
|------|---------------|------|------|
| 1_0_6_2_5_3_4.csv | Today, 10:12 PM | 8 KB | comm...values |

I have verified it one file (Right) and duplicate 7 times (Left). It show me correct answer:

One file                                    Duplicate 7 times

| | period | content_ | uid | count | |
|---|--------|----------|-----|-------|---|
| 1 | period | content_ | uid | count | |
| 2 | 2015-02-01_00_01 | content_ | uid_1 | 3 | |
| 3 | 2015-02-01_00_01 | content_ | uid_2 | 1 | |
| 4 | 2015-02-01_00_01 | content_ | uid_3 | 2 | |
| 5 | 2015-02-01_00_01 | content_ | uid_1 | 3 | |
| 6 | 2015-02-01_00_01 | content_ | uid_2 | 1 | |
| 7 | 2015-02-01_00_01 | content_ | uid_3 | 2 | |
| 8 | 2015-02-01_00_01 | content_ | uid_1 | 3 | |
| 9 | 2015-02-01_00_01 | content_ | uid_2 | 1 | |
| 10 | 2015-02-01_00_01 | content_ | uid_3 | 2 | |
| 11 | 2015-02-01_02_03 | content_ | uid_1 | 6 | |
| 12 | 2015-02-01_02_03 | content_ | uid_2 | 2 | |
| 13 | 2015-02-01_02_03 | content_ | uid_3 | 4 | |
| 14 | 2015-02-01_02_03 | content_ | uid_1 | 6 | |
| 15 | 2015-02-01_02_03 | content_ | uid_2 | 2 | |
| 16 | 2015-02-01_02_03 | content_ | uid_3 | 4 | |
| 17 | 2015-02-01_02_03 | content_ | uid_1 | 6 | |
| 18 | 2015-02-01_02_03 | content_ | uid_2 | 2 | |
| 19 | 2015-02-01_02_03 | content_ | uid_3 | 4 | |
| 20 | 2015-02-01_04_05 | content_ | uid_1 | 3 | |
| 21 | 2015-02-01_04_05 | content_ | uid_2 | 1 | |
| 22 | 2015-02-01_04_05 | content_ | uid_3 | 2 | |
| 23 | 2015-02-01_04_05 | content_ | uid_1 | 3 | |
| 24 | 2015-02-01_04_05 | content_ | uid_2 | 1 | |
| 25 | 2015-02-01_04_05 | content_ | uid_3 | 2 | |
| 26 | 2015-02-01_04_05 | content_ | uid_1 | 3 | |
| 27 | 2015-02-01_04_05 | content_ | uid_2 | 1 | |
| 28 | 2015-02-01_04_05 | content_ | uid_3 | 2 | |
| 29 | 2015-02-01_12_13 | content_ | uid_3 | 2 | |

| period | content_ | uid | count | |
|--------|----------|-----|-------|---|
| 2015-02-01_00_01 | content_ | uid_1 | 21 | |
| 2015-02-01_00_01 | content_ | uid_2 | 7 | |
| 2015-02-01_00_01 | content_ | uid_3 | 14 | |
| 2015-02-01_00_01 | content_ | uid_1 | 21 | |
| 2015-02-01_00_01 | content_ | uid_2 | 7 | |
| 2015-02-01_00_01 | content_ | uid_3 | 14 | |
| 2015-02-01_00_01 | content_ | uid_1 | 21 | |
| 2015-02-01_00_01 | content_ | uid_2 | 7 | |
| 2015-02-01_00_01 | content_ | uid_3 | 14 | |
| 2015-02-01_02_03 | content_ | uid_1 | 42 | |
| 2015-02-01_02_03 | content_ | uid_2 | 14 | |
| 2015-02-01_02_03 | content_ | uid_3 | 28 | |
| 2015-02-01_02_03 | content_ | uid_1 | 42 | |
| 2015-02-01_02_03 | content_ | uid_2 | 14 | |
| 2015-02-01_02_03 | content_ | uid_3 | 28 | |
| 2015-02-01_02_03 | content_ | uid_1 | 42 | |
| 2015-02-01_02_03 | content_ | uid_2 | 14 | |
| 2015-02-01_02_03 | content_ | uid_3 | 28 | |
| 2015-02-01_04_05 | content_ | uid_1 | 21 | |
| 2015-02-01_04_05 | content_ | uid_2 | 7 | |
| 2015-02-01_04_05 | content_ | uid_3 | 14 | |
| 2015-02-01_04_05 | content_ | uid_1 | 21 | |
| 2015-02-01_04_05 | content_ | uid_2 | 7 | |
| 2015-02-01_04_05 | content_ | uid_3 | 14 | |
| 2015-02-01_04_05 | content_ | uid_1 | 21 | |
| 2015-02-01_04_05 | content_ | uid_2 | 7 | |
| 2015-02-01_04_05 | content_ | uid_3 | 14 | |
| 2015-02-01_12_13 | content_ | uid_3 | 14 | |