

DATASCI W261: Machine Learning at Scale

David Rose
david.rose@berkeley.edu
W261-1
Week 04
2015.09.23

HW4.0

- **MRJob**: a Python API framework for accessing the Hadoop streaming capabilities. It differs from MapReduce in that it acts as a higher-level interface to MapReduce, yet utilizes the Hadoop MapReduce functionality. It provides a mechanism for creating data processing pipelines that MapReduce, on its own, cannot.
 - The ***_final() methods** are defined in the MRJob class, and as such can be overridden by classes that extend MRJob. The ***_final()** methods are executed when the input stream to the respective task is closed.
-

HW4.1

- **Serialization** converts an object into a bytestream that can be used for transporting the object over a network or to and from disk storage.
 - Within Hadoop and MRJob processes data must be **serialized**, at a minimum, when it is first submitted to a map task, when it is spilled to disk, when it is submitted to a reduce task, and again when the results are written to disk.
 - **Default modes** for MRJob serialization are:
 - INPUT_PROTOCOL = mrjob.protocol.RawValueProtocol
 - INTERNAL_PROTOCOL = mrjob.protocol.JSONProtocol
 - OUTPUT_PROTOCOL = mrjob.protocol.JSONProtocol
-

HW4.2

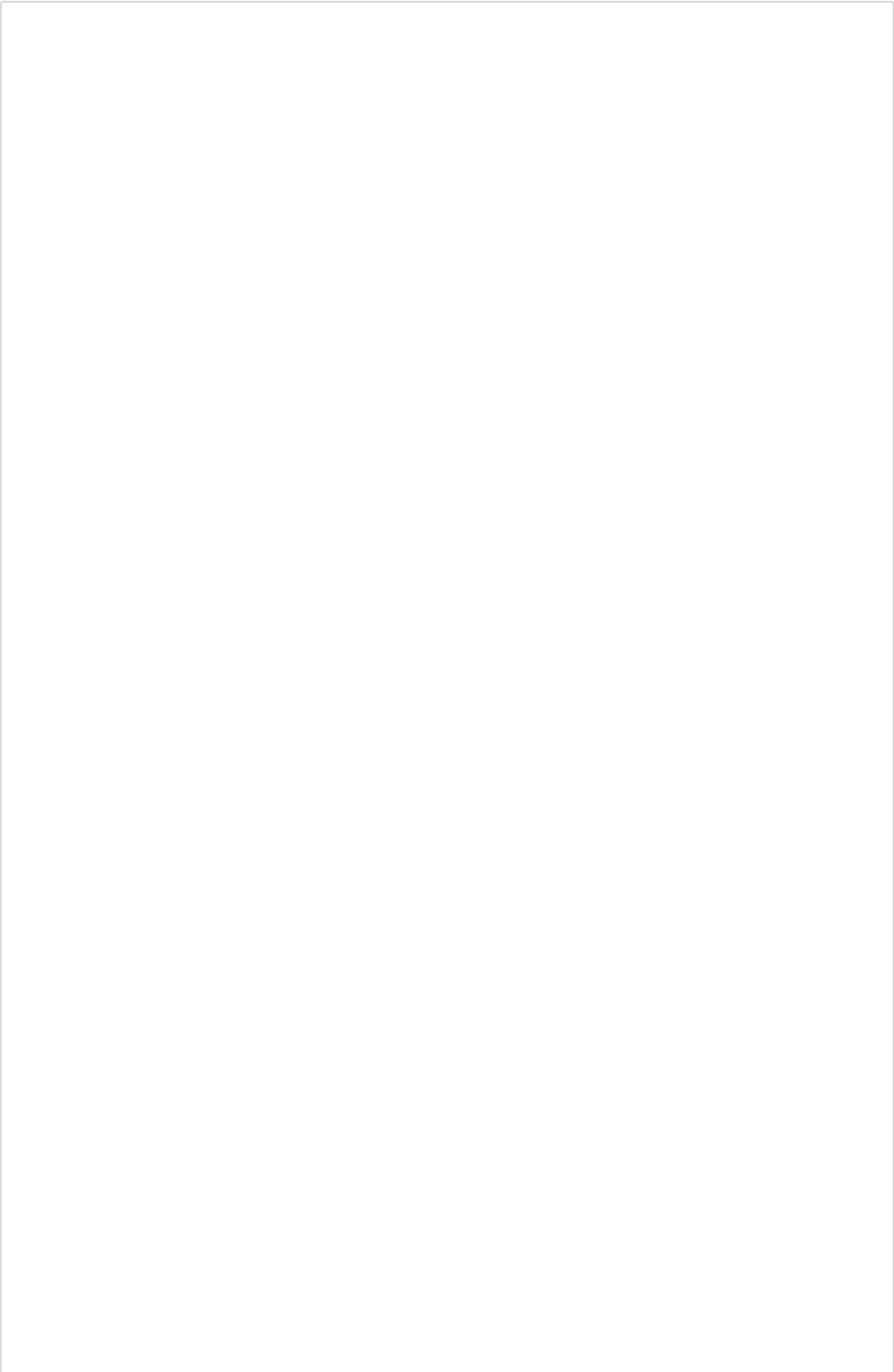
```
In [45]: 1 %%writefile hw_4_2_flatten.py
2 #!/usr/bin/python
3 '''
4 '''
5 from __future__ import print_function
6 import csv
7 import sys
8 with open(sys.argv[1], 'rb') as fin, open(sys.argv[2], 'w') as fout:
9     csvreader = csv.reader(fin, delimiter = ',', quotechar = '"')
10     visitorid = ''
11     for line in csvreader:
12         if line[0] == 'C':
13             visitorid = line[2]
14             continue
15         if line[0] == 'V':
16             line.append(visitorid)
17         print(','.join(line), file=fout)
```

Overwriting hw_4_2_flatten.py

```
In [46]: 1 # combine the page id and visitor id onto a single line for
2 # subsequent processing
3 !python hw_4_2_flatten.py anonymous-msweb.data flattened.data
```

HW4.3

In [47]:



```

1 %%writefile hw_4_3_mrjob.py
2 ''' count and list the top five most frequently visited pages
3 '''
4 from __future__ import print_function
5 from mrjob.job import MRJob
6 from mrjob.step import MRStep
7 import sys
8
9 class FrequentPages(MRJob):
10
11     def steps(self):
12         return [
13             MRStep(mapper=self.mapper,
14                   combiner=self.combiner,
15                   reducer=self.reducer,
16                   reducer_final=self.reducer_final)
17         ]
18
19     def mapper(self, _, line):
20         ''' enumerate page visits
21         '''
22         row = line.split(',')
23         if row[0] == 'V':
24             yield row[1], 1
25
26     def combiner(self, page, i):
27         ''' combine local results
28         '''
29         yield page, sum(i)
30
31     # track top five frequent pages in sorted order
32     topfive = [['',0]]
33
34     def inserttopfive(self, page, total):
35         ''' data structure and logic to maintain list of
36             top five most frequent pages.
37             This operation is performed here in the reducer
38             and again in the driver to capture output from
39             multiple reducers
40         '''
41         for j in range(0, len(self.topfive)):
42             if total > self.topfive[j][1]:
43                 self.topfive.insert(j, (page, total))
44                 if len(self.topfive) > 5:
45                     self.topfive.pop()
46                 break
47
48     def reducer(self, page, i):
49         ''' sum and sort page counts
50         '''
51         total = sum(i)
52         self.inserttopfive(page, total)
53
54     def reducer_final(self):

```

```

55         ''' emit results
56         '''
57         for i in range(0, len(self.topfive)):
58             yield(self.topfive[i][0], self.topfive[i][1])
59
60 if __name__ == '__main__':
61     FrequentPages.run()
62

```

Overwriting hw_4_3_mrjob.py

```

In [48]: 1 %%writefile hw_4_3_driver.py
2 from __future__ import print_function
3 from mrjob import util
4 import sys
5 from hw_4_3_mrjob import FrequentPages
6 util.log_to_null() # to suppress a 'no handler found' message
7
8 # list for storing most frequent pages
9 # we do this step here since multiple reducer tasks may run and th
10 # combined output needs to be processed
11 topfive = [['',0]]
12 def inserttopfive(page, total):
13     for j in range(0, len(topfive)):
14         if total > topfive[j][1]:
15             topfive.insert(j, (page, total))
16             if len(topfive) > 5:
17                 topfive.pop()
18             break
19
20 mr_job = FrequentPages(args=sys.argv[1:])
21 with mr_job.make_runner() as runner:
22     runner.run()
23     for line in runner.stream_output():
24         page, total = line.split()
25         inserttopfive(page, int(total))
26 for i in range(0, len(topfive)):
27     print('page: {}, visits: {}'.format(topfive[i][0],
28                                         topfive[i][1]), file=sys.s
29

```

Overwriting hw_4_3_driver.py

```

In [49]: 1 # HW 4.3: Find the 5 most frequently visited pages using mrjob fro
2 !python hw_4_3_driver.py flattened.data --strict-protocols -r loca

```

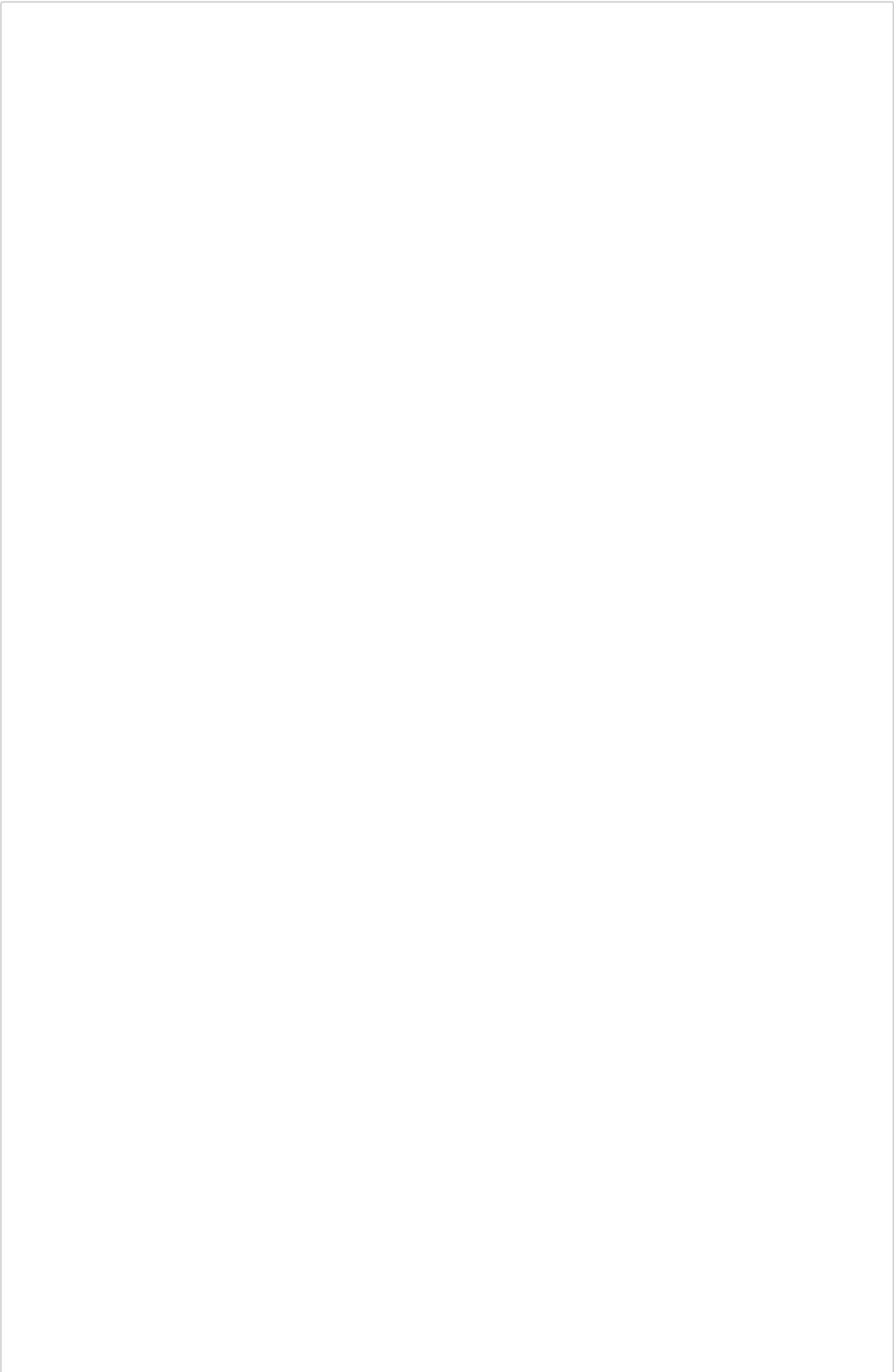
```

page: "1008", visits: 10836
page: "1034", visits: 9383
page: "1004", visits: 8463
page: "1018", visits: 5330
page: "1017", visits: 5108

```

HW4.4

In [50]:



```

1 %%writefile hw_4_4_mrjob.py
2 ''' count the number of page-visitor combinations and
3     for each page list the most frequent visitors
4
5     the data does not effectively support this operation since it
6     only lists unique page visits, therefore every visitor
7     will show up as having visited once, and therefore every visit
8     is the most frequent visitor
9 '''
10 from __future__ import print_function
11 from mrjob.job import MRJob
12 import sys
13
14 class FrequentVisitors(MRJob):
15
16     def mapper(self, _, line):
17         ''' enumerate page visitors
18         '''
19         row = line.split(',')
20         if row[0] == 'V':
21             # page ID, visitor ID
22             yield row[1], row[3]
23
24     # data structures to manage reducer logic
25     visitors = {}
26     currentpage = ''
27
28     def reducer(self, page, visitor):
29         ''' sum page visitor counts
30         '''
31         if not page == self.currentpage:
32             ''' page id has changed in the stream, so process and
33             the information for the current page
34             '''
35             if len(self.visitors) > 0:
36                 frequentv = []
37                 maxv = 0
38                 for v in self.visitors:
39                     if self.visitors[v] > maxv:
40                         frequentv = [v]
41                         maxv = self.visitors[v]
42                     elif self.visitors[v] == maxv:
43                         frequentv.append(v)
44                 # emit results
45                 for v in frequentv:
46                     yield self.currentpage, v
47             # reset counters
48             self.visitors = {}
49             self.currentpage = page
50         for v in visitor:
51             if not v in self.visitors:
52                 self.visitors[v] = 0
53             self.visitors[v] += 1
54

```



```

55     # process any remaining values after stream closes
56     def reducer_final(self):
57         if len(self.visitors) > 0:
58             frequentv = []
59             maxv = 0
60             for v in self.visitors:
61                 if self.visitors[v] > maxv:
62                     frequentv = [v]
63                 elif self.visitors[v] == maxv:
64                     frequentv.append(v)
65             for v in frequentv:
66                 yield self.currentpage, v
67
68 if __name__ == '__main__':
69     FrequentVisitors.run()
70

```

Overwriting hw_4_4_mrjob.py

```

In [51]: 1 %%writefile hw_4_4_driver.py
2 import csv
3 from mrjob import util
4 import sys
5 from hw_4_4_mrjob import FrequentVisitors
6
7 util.log_to_null() # to suppress a 'no handler found' message
8
9 # construct list of page ids and urls to satisfy the requirement
10 # this is essentially a join; it makes better sense to do this in
11 # the driver as doing so in hadoop offers no advantages and incurs
12 # additional network overhead; mrjob output is parsed and augmented
13 # the url information
14 pageinfo = {}
15 # read in the page attributes including the url; this could also have
16 # been preprocessed so that the page attributes were in their own
17 with open('flattened.data', 'rb') as fin:
18     csvreader = csv.reader(fin, delimiter = ',', quotechar = '"')
19     for line in csvreader:
20         if line[0] == 'A':
21             pageid = line[1]
22             url = line[4]
23             pageinfo[pageid] = url
24
25 mr_job = FrequentVisitors(args=sys.argv[1:])
26 with mr_job.make_runner() as runner:
27     runner.run()
28     for line in runner.stream_output():
29         page, visitor = line.replace('"', '').split()
30         print pageinfo[page], page, visitor
31

```

Overwriting hw_4_4_driver.py

```
In [52]: 1 # HW 4.4: Find the most frequent visitor of each page using
        2 # mrjob and the output of 4.2
        3 # In this output please include the webpage URL, webpageID
        4 # and Visitor ID.
        5 !python hw_4_4_driver.py flattened.data --strict-protocols -r loca
        6 !head -n50 mrjob_4_4_output
        7 !tail -n50 mrjob_4_4_output
```

```
/regwiz 1000 10001
/regwiz 1000 10010
/regwiz 1000 10039
/regwiz 1000 10073
/regwiz 1000 10087
/regwiz 1000 10101
/regwiz 1000 10132
/regwiz 1000 10141
/regwiz 1000 10154
/regwiz 1000 10162
/regwiz 1000 10166
/regwiz 1000 10201
/regwiz 1000 10218
/regwiz 1000 10220
/regwiz 1000 10324
/regwiz 1000 10348
/regwiz 1000 10376
/regwiz 1000 10384
/regwiz 1000 10409
/regwiz 1000 10420
```

HW4.5

```
In [53]: 1 ''' utility script to get some statistics on the data set
2         '''
3         with open('topUsers_Apr-Jul_2014_1000-words.txt', 'rb') as fin:
4             max = 0
5             total = 0
6             linecount = 0
7             classes = {}
8             for line in fin:
9                 row = line.split(',')
10                count = int(row[2])
11                clazz = row[1]
12                if count > max: max = count
13                total += count
14                linecount += 1
15                if not clazz in classes:
16                    classes[clazz] = 0
17                classes[clazz] += 1
18            print(linecount, max, total, str(classes))
19
(1000, 1724608, 61819567, "{ '1': 91, '0': 752, '3': 103, '2': 54 }")
```

```
In [54]: 1 ''' utility script to get some statistics on the data set
2         look for the maximum word ratio after the data is normalized
3         '''
4         with open('topUsers_Apr-Jul_2014_1000-words.txt', 'rb') as fin:
5             max = 0.0
6             for line in fin:
7                 row = line.split(',')
8                 count = int(row[2])
9                 for i in range(3, len(row)):
10                    p = float(row[i]) / count
11                    if p > max:
12                        max = p
13            print(max)

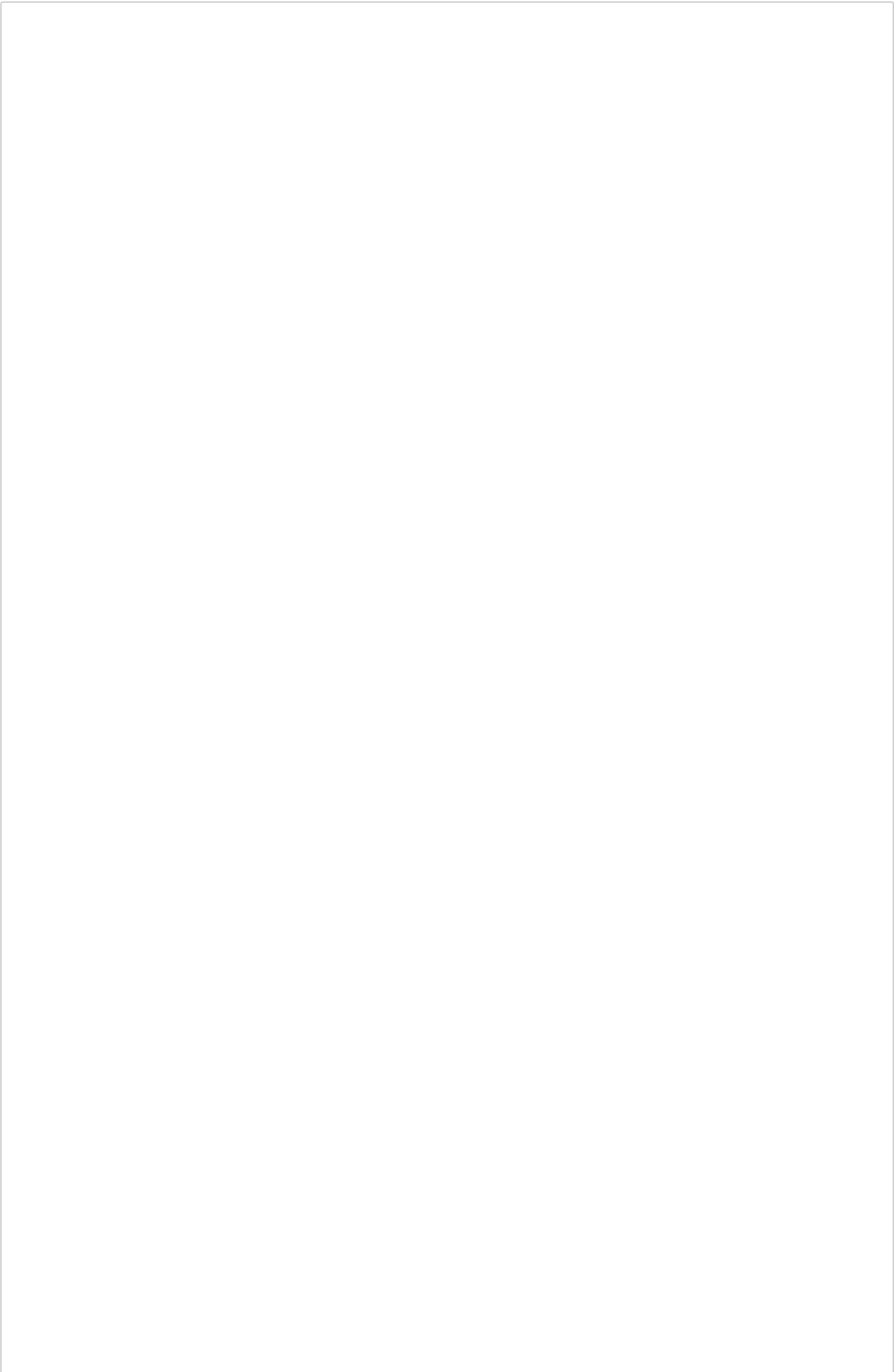
0.409909665665
```

```
In [55]: 1 %%writefile hw_4_5_preprocess.py
2 ''' preprocess the data set, normalizing the word counts
3 '''
4 from __future__ import print_function
5 import sys
6 with open(sys.argv[1], 'rb') as fin, open(sys.argv[2], 'w') as fout:
7     for line in fin:
8         row = map(int, line.split(','))
9         userid = row[0]
10        code = row[1]
11        total = row[2]
12        for j in range(2, len(row)):
13            row[j] = float(row[j]) / total
14        print('{}'.format(row[0]), file = fout, end = '')
15        for j in range(1, 3):
16            print(',{}'.format(row[j]), file = fout, end = '')
17        for j in range(3, len(row)):
18            print(',{:0.8f}'.format(row[j]), file = fout, end = '')
19        print('', file = fout)
20
```

Overwriting hw_4_5_preprocess.py

```
In [56]: 1 !python hw_4_5_preprocess.py topUsers_Apr-Jul_2014_1000-words.txt
```

In [57]:



```

1  %%writefile hw_4_5_mrjob.py
2  ''' map/reduce approach to determining stable centroids using
3      a K-means algorithm
4  '''
5  from __future__ import print_function
6  from numpy import argmin, array, random
7  from mrjob.job import MRJob
8  from mrjob.step import MRStep
9  from itertools import chain
10 import sys
11
12 # Calculate find the nearest centroid for data point
13 def MinDist(datapoint, centroid_points):
14     datapoint = array(datapoint)
15     centroid_points = array(centroid_points)
16     diff = datapoint - centroid_points
17     diffsq = diff*diff
18     # Get the nearest centroid for each instance
19     sumofsquares = list(diffsq.sum(axis = 1))
20     minindex = argmin(sumofsquares)
21     return minindex
22
23 # Check whether centroids converge
24 def stop_criterion(centroid_points_old, centroid_points_new, T):
25     oldvalue = list(chain(*centroid_points_old))
26     newvalue = list(chain(*centroid_points_new))
27     Diff = [abs(x-y) for x, y in zip(oldvalue, newvalue)]
28     Flag = True
29     for i in Diff:
30         if i > T:
31             Flag = False
32             break
33     return Flag
34
35 class MRKmeans(MRJob):
36
37     centroid_points=[]
38     k = -1
39     CENTROIDFILE = '/tmp/centroids.txt'
40
41     def steps(self):
42         return [
43             MRStep mapper_init = self.mapper_init,
44                   mapper=self.mapper,
45                   combiner = self.combiner,
46                   reducer=self.reducer)
47         ]
48     # load initial centroids
49     def mapper_init(self):
50         self.centroid_points=[]
51         with open(self.CENTROIDFILE, 'rb') as fin:
52             header = True
53             for line in fin:
54                 if header:

```

```

55         self.k = int(line.strip())
56         header = False
57     else:
58         self.centroid_points.append(map(float, line.strip().split(',')))
59     # print the value of k back to the cache file
60     with open(self.CENTROIDFILE, 'w') as fout:
61         print('{}'.format(self.k), file = fout)
62
63     #load data and output the nearest centroid index and data point
64     def mapper(self, _, line):
65         D = (map(float, line.split(',')[3:]))
66         centroid = MinDist(D, self.centroid_points)
67         yield centroid, (D, 1)
68
69     # aggregate data points locally
70     def combiner(self, centroid, inputdata):
71         count = 0
72         bucket = [0] * 1000
73         for data, n in inputdata:
74             count += n
75             data = map(float, data)
76             for j in range(0, len(data)):
77                 bucket[j] += data[j]
78         yield centroid, (bucket, count)
79
80
81     # aggregate values for each centroid, then recalculate centroid
82     def reducer(self, idx, inputdata):
83         centroids = []
84         with open(self.CENTROIDFILE, 'rb') as fin:
85             self.k = int(fin.readline().strip())
86         num = [0] * self.k
87         for i in range(self.k):
88             centroids.append([0.0]*1000)
89         for data, n in inputdata:
90             num[idx] += n
91             data = map(float, data)
92             for j in range(0, len(data)):
93                 centroids[idx][j] += data[j]
94         for j in range(0, len(centroids[idx])):
95             centroids[idx][j] = centroids[idx][j] / num[idx]
96         with open(self.CENTROIDFILE, 'a') as fout:
97             print(','.join(str(i) for i in centroids[idx]), file = fout)
98         yield idx, (centroids[idx])
99
100 if __name__ == '__main__':
101     MRKmeans.run()

```

Overwriting hw_4_5_mrjob.py

In [58]:




```

1 %%writefile hw_4_5_mrjob_labeler.py
2 ''' assigns centroid labels to data points
3
4     this step is done separately from the centroid calculations
5     to simplify the data processing
6 '''
7 from __future__ import print_function
8 from numpy import argmin, array, random
9 from mrjob.job import MRJob
10 from mrjob.step import MRStep
11 import sys
12
13 from hw_4_5_mrjob import MinDist
14
15 class MRLabeler(MRJob):
16
17     centroid_points=[]
18     k = -1
19     CENTROIDFILE = '/tmp/centroids.txt'
20
21     def steps(self):
22         return [
23             MRStep mapper_init = self.mapper_init,
24                   mapper=self.mapper,
25                   combiner = self.combiner,
26                   reducer=self.reducer,
27                   reducer_final=self.reducer_final)
28         ]
29     #load centroids info from file
30     def mapper_init(self):
31         self.centroid_points=[]
32         with open(self.CENTROIDFILE, 'rb') as fin:
33             header = True
34             for line in fin:
35                 if header:
36                     self.k = int(line.strip())
37                     header = False
38                 else:
39                     self.centroid_points.append(map(float,line.str
40
41     # determine the closest centroid for each data point
42     def mapper(self, _, line):
43         row = line.strip().split(',')
44         label = int(row[1])
45         D = (map(float,row[3:]))
46         centroid = MinDist(D, self.centroid_points)
47         yield centroid, (label, 1)
48
49     # combine sum of data points locally
50     def combiner(self, centroid, inputdata):
51         counts = [0, 0, 0, 0]
52         for label, n in inputdata:
53             counts[label] += n
54         for i in range(len(counts)):

```

```

55         yield centroid, (i, counts[i])
56
57
58     # sum the counts for centroids and classes
59     currentcentroid = ''
60     counts = []
61     def reducer(self, centroid, inputdata):
62         if not centroid == self.currentcentroid:
63             for i in range(len(self.counts)):
64                 yield self.currentcentroid, (i, self.counts[i])
65             self.counts = [0, 0, 0, 0]
66             self.currentcentroid = centroid
67         for label, n in inputdata:
68             self.counts[label] += n
69
70     def reducer_final(self):
71         for i in range(len(self.counts)):
72             yield self.currentcentroid, (i, self.counts[i])
73
74
75 if __name__ == '__main__':
76     MRKmeans.run()

```

Overwriting hw_4_5_mrjob_labeler.py

In [59]:



```

1 %%writefile hw_4_5_driver.py
2 ''' driver script for Kmeans job
3 '''
4 from __future__ import print_function
5 from numpy import random
6 from hw_4_5_mrjob import MRKmeans, stop_criterion
7 from hw_4_5_mrjob_labeler import MRLabeler
8 from mrjob import util
9 import sys
10
11 util.log_to_null() # to suppress a 'no handler found' message
12 CENTROIDFILE = '/tmp/centroids.txt'
13
14 def countlabels():
15     ''' count the number of each classification as labeled in the
16         original data set
17     '''
18     with open('topUsers_Apr-Jul_2014_1000-words.txt', 'rb') as f:
19         labels = {}
20         for line in f:
21             row = line.split(',')
22             label = row[1]
23             if not label in labels:
24                 labels[label] = 0
25             labels[label] += 1
26     return labels
27
28 def init_centroids_random_internal(k):
29     ''' select initial centroids by choosing data points
30         randomly from the data set
31     '''
32     randoms = sorted(random.randint(0, 1000, size = k))
33     centroids = []
34     with open('topUsers_Apr-Jul_2014_1000-words.txt', 'rb') as f:
35         lineno = 0
36         count = 0
37         for line in f:
38             if lineno in randoms:
39                 row = line.strip().split(',')
40                 data = map(float, row[3:])
41                 # normalize the values
42                 data = [i / float(row[2]) for i in data]
43                 centroids.append(data)
44                 count += 1
45                 if count == k:
46                     break
47             lineno += 1
48     with open(CENTROIDFILE, 'w') as f:
49         print('{}'.format(k), file = f)
50         for tuple in centroids:
51             print(','.join(str(i) for i in tuple), file = f)
52     return centroids
53
54 def init_centroids_random_external(k):

```

```

55     ''' create initial centroids by generating random values
56     '''
57     centroids = []
58     for i in range(k):
59         centroid = []
60         for j in range(1000):
61             centroid.append(random.uniform(0.0, .5))
62         centroids.append(centroid)
63     with open(CENTROIDFILE, 'w') as fout:
64         print('{}'.format(k), file = fout)
65         for tuple in centroids:
66             print(','.join(str(i) for i in tuple), file = fout)
67     return centroids
68
69 def init_centroids_perturbed(k):
70     ''' create initial centroids by using aggregated values and
71         perturbing them with random noise
72     '''
73     centroids = []
74     with open('topUsers_Apr-Jul_2014_1000-words_summaries.txt', 'r') as f:
75         for line in f:
76             row = line.strip().split(',')
77             if row[0] == 'ALL_CODES':
78                 data = map(float, row[3:])
79                 # normalize the values
80                 data = [i / float(row[2]) for i in data]
81                 for i in range(k):
82                     centroid = []
83                     for j in range(len(data)):
84                         # modify each value with random number in
85                         # range of +- the value; avoids a small number
86                         # by a large number
87                         centroid.append(data[j] + random.uniform(-1, 1))
88                     centroids.append(centroid)
89                 break
90     with open(CENTROIDFILE, 'w') as f:
91         print('{}'.format(k), file = f)
92         for tuple in centroids:
93             print(','.join(str(i) for i in tuple), file = f)
94     return centroids
95
96 def init_centroids_trained(k):
97     ''' create initial centroids by choosing the class-specific aggregated
98     '''
99     centroids = []
100     with open('topUsers_Apr-Jul_2014_1000-words_summaries.txt', 'r') as f:
101         for line in f:
102             row = line.strip().split(',')
103             if row[0] == 'CODE':
104                 code = row[1]
105                 total = int(row[2])
106                 data = map(float, row[3:])
107                 # normalize the values
108                 data = [i / total for i in data]

```

```

109         centroids.append(data)
110     with open(CENTROIDFILE, 'w') as f:
111         print('{}'.format(k), file = f)
112         for tuple in centroids:
113             print(','.join(str(i) for i in tuple), file = f)
114     return centroids
115
116 def getpurity(clusters):
117     majority = 0
118     for c in clusters:
119         counts = map(int, clusters[c]);
120         majority += max(counts)
121     return majority / float(1000)
122
123 def go(centroid_points):
124     ''' submit the centroid calculation job;
125         follow that with the data labeling job
126     '''
127     mr_job = MRKmeans(args = ['normalized.txt'])
128     iteration = 0
129     while(1):
130         # save previous centroids to check convergency
131         centroid_points_old = centroid_points[:]
132         with mr_job.make_runner() as runner:
133             runner.run()
134             # capture reducer output
135             for line in runner.stream_output():
136                 key,value = mr_job.parse_output_line(line)
137                 centroid_points[key] = value
138             iteration += 1
139             if stop_criterion(centroid_points_old, centroid_points, 0.001):
140                 break
141     print('centroids converged: {} iterations'.format(iteration),
142         mr_job = MRLabeler(args = ['normalized.txt'])
143     labels = countlabels()
144     clusters = {}
145     with mr_job.make_runner() as runner:
146         runner.run()
147         for line in runner.stream_output():
148             centroid, value = mr_job.parse_output_line(line)
149             label = int(value[0])
150             count = int(value[1])
151             if not centroid in clusters:
152                 clusters[centroid] = [0,0,0,0]
153             clusters[centroid][label] += count
154     for c in sorted(clusters):
155         results = map(float, clusters[c])
156         print('centroid {}: label 0: {:.4f} label 1: {:.4f} label 2: {:.4f} label 3: {:.4f}'.format(c, results[0] / labels['0'], results[1] / labels['1'], results[2] / labels['2'], results[3] / labels['3']),
157             file = sys.stdout)
158     # print the actual counts
159     print('centroid {}: counts: {}'.format(c, clusters[c]),
160         file = sys.stderr)

```

```
163         print('purity: {:.4f}'.format(getpurity(clusters)))
164         print('', file = sys.stderr)
165
166     print('k = 4, random initialization, v1', file = sys.stderr)
167     go(init_centroids_random_internal(4))
168     print('k = 4, random initialization, v2', file = sys.stderr)
169     go(init_centroids_random_external(4))
170     print('k = 2, perturbed initialization', file = sys.stderr)
171     go(init_centroids_perturbed(2))
172     print('k = 4, perturbed initialization', file = sys.stderr)
173     go(init_centroids_perturbed(4))
174     print('k = 4, trained initialization', file = sys.stderr)
175     go(init_centroids_trained(4))
176
177
```

Overwriting hw_4_5_driver.py

```
In [60]: 1 !python hw_4_5_driver.py
```

```
k = 4, random initialization, v1
centroids converged: 9 iterations
centroid 0: label 0: 0.0000 label 1: 0.1868 label 2: 0.2778 label 3:
0.0388
centroid 1: label 0: 0.0000 label 1: 0.5604 label 2: 0.0000 label 3:
0.0000
centroid 2: label 0: 0.9987 label 1: 0.0330 label 2: 0.0370 label 3:
0.9612
centroid 3: label 0: 0.0013 label 1: 0.2198 label 2: 0.6852 label 3:
0.0000
purity: 0.8560
```

```
k = 4, random initialization, v2
centroids converged: 3 iterations
centroid 0: label 0: 1.0000 label 1: 1.0000 label 2: 1.0000 label 3:
1.0000
purity: 0.7520
```

```
k = 2, perturbed initialization
centroids converged: 4 iterations
centroid 0: label 0: 0.9987 label 1: 0.0330 label 2: 0.2593 label 3:
0.9612
centroid 1: label 0: 0.0013 label 1: 0.9670 label 2: 0.7407 label 3:
0.0388
purity: 0.8390
```

```
k = 4, perturbed initialization
centroids converged: 4 iterations
centroid 0: label 0: 0.0000 label 1: 0.0000 label 2: 0.0370 label 3:
0.0000
centroid 1: label 0: 0.8816 label 1: 0.0110 label 2: 0.0000 label 3:
0.6214
centroid 2: label 0: 0.0013 label 1: 0.9670 label 2: 0.7037 label 3:
0.0388
centroid 3: label 0: 0.1170 label 1: 0.0220 label 2: 0.2593 label 3:
0.3398
purity: 0.8410
```

```
k = 4, trained initialization
centroids converged: 5 iterations
centroid 0: label 0: 0.9960 label 1: 0.0330 label 2: 0.2593 label 3:
0.3689
centroid 1: label 0: 0.0000 label 1: 0.5604 label 2: 0.0000 label 3:
0.0000
centroid 2: label 0: 0.0013 label 1: 0.4066 label 2: 0.7407 label 3:
0.0388
centroid 3: label 0: 0.0027 label 1: 0.0000 label 2: 0.0000 label 3:
0.5922
purity: 0.9010
```

Discussion: centroid initialization has a significant impact on the results. The purity scores of the five runs range from 0.75 to 0.90. Not surprisingly the best results came from the centroids initialized with the class-specific aggregate data, and the worst results came from using randomized synthetic data. Ideally the results would suggest a diagonal of populated cells with all others being zero. The final scenario above comes closest, but there is still room for improvement.

In []:

| | |
|---|--|
| 1 | |
|---|--|