

DATASCI W261: Machine Learning at Scale

David Rose

david.rose@berkeley.edu

W261-1

Week 03

2015.09.18

HW3.0

- Merge sort is an algorithm that combines two or more lists that are in sorted order individually and outputs a single list in sorted order. The algorithm compares the first elements of each input list to determine the element with the lowest (or highest) value. This value is removed from the input list and appended to the output list. This process repeats until there are no remaining input values. The algorithm runs in $O(m \log n)$ time where m is the total number of elements in the combined input lists and n is the number of input lists. Hadoop uses a merge sort when preparing data to be sent to a reduce task to provide the reducer with a list of input sorted by key.
- A combiner function aggregates or otherwise combines output from a mapper function before it is sent to a reducer. It functions similarly to a reducer but with some constraints. In part the constraints are imposed because the combiner is not guaranteed to execute. Whether it does or not is a runtime decision by the Hadoop process.
 - One constraint is that the output of the combiner must match the signature of the output of the mapper.
 - Another is that the combiner will only operate on the output of the mapper tasks of a single node so that in a multi-node environment one combiner will not see all of the data but rather only the data that has been assigned to the map tasks on that particular node.
 - Combiners can be very effective in reducing the amount of bandwidth required for transmitting data from a map task to a reduce task.
 - An example usage would be in an item counting operation, whether it be words in a document or items in a shopping basket. A combiner can take the canonical map output, $[\text{key}, 1]$, and produce for each key encountered a combined output, $[\text{key}, n]$, where n is the sum of counts for the key.
- The Hadoop shuffle is the core operation of the Hadoop framework. The shuffle is responsible for moving data from the map task output to the reduce task input.
 - Map side of shuffle:

- map output is written to a circular memory buffer
- when the buffer reaches a capacity threshold, data is spilled to disk
- before being written to spill files, data is partitioned based on the reducer targets and sorted in memory by key
- when map function completes, multiple spill files, if present, are merged into a single partitioned and sorted file.
- Reduce side of shuffle:
 - the reducer receives notice from the job tracker that a map task is complete
 - the reducer begins copying the sorted files from the identified map nodes
 - the files are merged together, potentially in several rounds, maintaining the sort order
 - the files are fed to the reduce function
- The Apriori algorithm is used for gaining efficiency when deriving association rules. The algorithm is based on the insight that if an itemset $I = \{i_1..i_n\}$ is frequent (i.e., occurs greater than a pre-defined threshold in the data set), then all subsets of I are also frequent. In practice this is implemented by determining the frequent itemsets of length $k = 1$, removing the itemsets that are not frequent, and using the remaining itemsets to construct a candidate set of $k + 1$ items. This process is repeated iteratively until the itemset size of interest is reached or when there are no frequent itemsets remaining.
 - Although I can find no references to the use of the apriori algorithm in public health outbreak investigations, it could potentially be used to identify potential association rules between behaviors, events, symptoms, and diagnoses. Maybe.
- Confidence: a measure of a rule $(X \rightarrow Y)$ calculated as the ratio of the support of an itemset to the support of a chosen subset of the itemset, i.e., the left hand side of the rule:

$$confidence(X \rightarrow Y) = \frac{support(X \cup Y)}{support(X)}$$

- Lift: a measure of a rule $(X \rightarrow Y)$ calculated as the ratio of the confidence of the rule to the expected support of each constituent of the rule if the constituents were statistically independent:

$$lift(X \rightarrow Y) = \frac{confidence(X \rightarrow Y)}{support(Y)} = \frac{support(X \cup Y)}{support(X)support(Y)}$$

Type *Markdown* and LaTeX: α^2

In [193]:

```
1 # load data file into hdfs
2 !hdfs dfs -put -f ProductPurchaseData.txt /in
```

HW3.1

Mapper for exploratory data analysis

```
In [194]: 1 %%writefile map-3.1.py
          2 #!/usr/bin/python
          3 ''' count some things
          4 '''
          5 from __future__ import print_function
          6 import random
          7 import sys
          8 import time
          9 random.seed()
         10 for line in sys.stdin:
         11     # split line into tokens of items
         12     tokens = line.split()
         13     basketsize = len(tokens)
         14     # generate a basket key
         15     # this is probably overkill
         16     basketkey = ''.join([str(random.random()), str(time.time())])
         17     for token in tokens:
         18         print('\t'.join([token, str(basketsize), str(basketkey)]),
         19
```

Overwriting map-3.1.py

Reducer for exploratory data analysis.

```

In [195]: 1 %%writefile reduce-3.1.py
          2 #!/usr/bin/python
          3 ''' reducer aggregates some statistics
          4 '''
          5 from __future__ import print_function
          6 import sys
          7 items = {}
          8 baskets = {}
          9 maxbasket = 0
         10 for line in sys.stdin:
         11     tokens = line.split()
         12     item = tokens[0]
         13     basketsize = int(tokens[1])
         14     basketkey = tokens[2]
         15     items[item] = 1
         16     # store size of each individual basket
         17     baskets[basketkey] = basketsize
         18     # keep running maximum of basket size
         19     maxbasket = max(maxbasket, basketsize)
         20 # calculate some statistics
         21 # mean basket size
         22 meanbasket = sum(baskets.values())/float(len(baskets))
         23 # standard deviation of basket size
         24 std = (sum([(x - meanbasket)**2 for x in baskets.values()])
         25         / float(len(baskets))**0.5)
         26 # emit results
         27 results = ('total baskets: {}, unique items: {}, largest basket si
         28           + 'mean basket size: {:.2f}, std deviation: {:.2f}')
         29         .format(len(baskets), len(items), maxbasket, meanbasket, std)
         30 print(results, file=sys.stdout)
         31

```

Overwriting reduce-3.1.py

```

In [196]: 1 !chmod +x *.py

```

```
In [197]: 1 # HW3.1. Do some exploratory data analysis of this dataset.
2 # Report your findings such as number of unique products; largest
3 # basket, etc. using Hadoop Map-Reduce.
4 !printf "preparing output directory\n"
5 !hdfs dfs -rm -r -f -skipTrash /out/out-3.1
6 !printf "executing yarn task\n"
7 !yarn jar /usr/local/Cellar/hadoop/2.7.0/libexec/share/hadoop/tool
8 -files ./map-3.1.py,./reduce-3.1.py \
9 -mapper ./map-3.1.py -reducer ./reduce-3.1.py \
10 -input /in/ProductPurchaseData.txt -output /out/out-3.1
11 !printf "checking output directory: "
12 !hdfs dfs -ls /out/out-3.1
13 !printf "displaying results\n"
14 !hdfs dfs -cat /out/out-3.1/part-00000
```

preparing output directory

Deleted /out/out-3.1

executing yarn task

checking output directory: Found 2 items

-rw-r--r-- 1 david supergroup 0 2015-09-21 17:50 /out/ou
t-3.1/_SUCCESS

-rw-r--r-- 1 david supergroup 115 2015-09-21 17:50 /out/ou
t-3.1/part-00000

displaying results

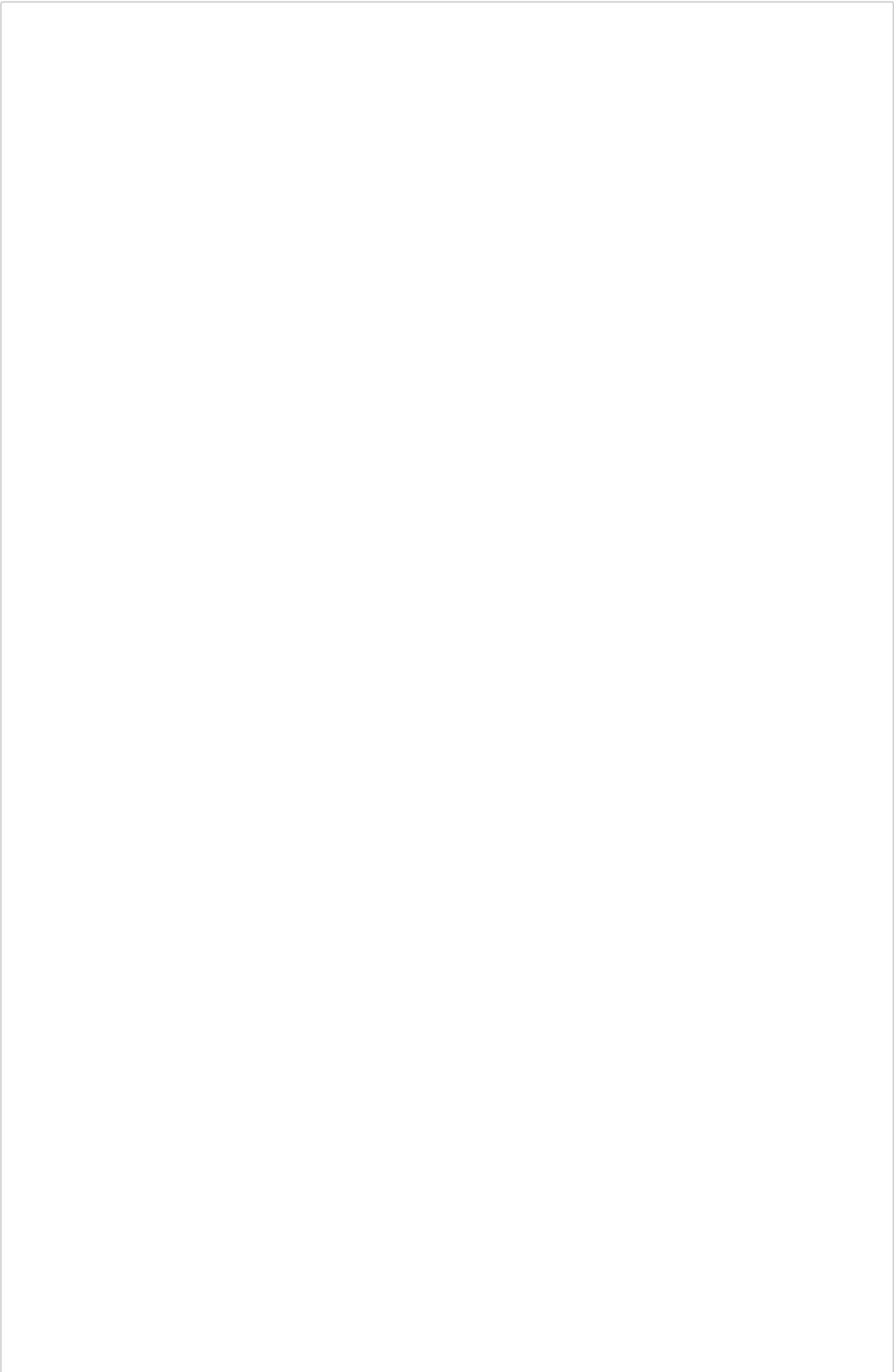
total baskets: 31101, unique items: 12592, largest basket size: 37,
mean basket size: 12.24, std deviation: 5.03

HW3.2

```
In [198]: 1 %%writefile map-3.2.py
2          #!/usr/bin/python
3          ''' map function to generate 2-itemsets
4
5             implements in-memory combiner to aggregate counts of each
6             2-itemset to reduce output file size
7             '''
8          from __future__ import print_function
9          import sys
10         itemsets = {}
11         magic = '*'
12         for line in sys.stdin:
13             # split line into tokens of items
14             items = line.split()
15             for i in range(0, len(items)):
16                 item1 = items[i]
17                 # create fake itemset to keep count of 1-itemset
18                 # for use in reduce task
19                 itemset = ' '.join([item1, magic])
20                 if itemset not in itemsets:
21                     itemsets[itemset] = 0
22                 itemsets[itemset] += 1
23                 # create all possible combinations of 2-itemsets
24                 # for this basket
25                 for j in range(i + 1, len(items)):
26                     item2 = items[j]
27                     itemset = ' '.join(sorted([item1, item2]))
28                     if itemset not in itemsets:
29                         itemsets[itemset] = 0
30                     itemsets[itemset] += 1
31         for itemset in itemsets:
32             results = '\t'.join([itemset, str(itemsets[itemset])])
33             print(results, file=sys.stdout)
34
```

Overwriting map-3.2.py

In [199]:



```

1 %%writefile reduce-3.2.py
2 #!/usr/bin/python
3 ''' reducer reads stream consisting of [2-itemset, count] pairs
4
5     utilizes order inversion pattern to enable efficient
6     stream processing
7
8 '''
9 from __future__ import print_function
10 import sys
11 # value for identifying left hand side of rule
12 magic = '*'
13 minsupport = 100
14 lhscurrent = ''
15 rhscurrent = ''
16 lhscount = 0
17 itemsetcount = 0
18 # loop through input
19 # when key changes, take action depending on which
20 # component of the key changes
21 for line in sys.stdin:
22     tokens = line.split('\t')
23     itemset = tokens[0]
24     lhs = itemset.split()[0]
25     rhs = itemset.split()[1]
26     count = int(tokens[1])
27
28     if not rhs == rhscurrent:
29         if itemsetcount > minsupport and not rhscurrent == magic:
30             # calculate the confidence for the current itemset
31             confidence = itemsetcount / float(lhscount)
32             print('{} -> {} \t{:0.8f}'.format(lhscurrent, rhscurrent, confidence),
33                   file=sys.stdout)
34             # reset loop values for right hand side
35             itemsetcount = 0
36             rhscurrent = rhs
37         if not lhs == lhscurrent:
38             # initialize the new lhs
39             # reset loop values for left hand side
40             lhscount = 0
41             lhscurrent = lhs
42         if rhs == magic:
43             # increment support count for left hand side
44             lhscount += count
45         else:
46             # increment support count for itemset
47             itemsetcount += count
48 # process the last line
49 if itemsetcount > minsupport and not rhscurrent == magic and lhscount > 0:
50     # calculate the confidence for the current itemset
51     confidence = itemsetcount / float(lhscount)
52     print('{} -> {} \t{:0.8f}'.format(lhscurrent, rhscurrent, confidence),
53         file=sys.stdout)

```



```
--  
55  
56
```

Overwriting reduce-3.2.py

```
In [200]: 1 !chmod +x *.py
```

```
In [201]: 1 # HW3.2.  
2  
3 !printf "preparing output directory\n"  
4 !hdfs dfs -rm -r -f -skipTrash /out/out-3.2  
5 !printf "executing yarn task\n"  
6 !yarn jar /usr/local/Cellar/hadoop/2.7.0/libexec/share/hadoop/tool  
7 -files ./map-3.2.py,./reduce-3.2.py \  
8 -mapper ./map-3.2.py -reducer ./reduce-3.2.py \  
9 -input /in/ProductPurchaseData.txt -output /out/out-3.2  
10 !printf "checking output directory: "  
11 !hdfs dfs -ls /out/out-3.2  
12 !printf "displaying results\n"  
13 !hdfs dfs -cat /out/out-3.2/part-00000 | sort -k 4nr -k 1 2>/dev/n
```

preparing output directory

Deleted /out/out-3.2

executing yarn task

checking output directory: Found 2 items

```
-rw-r--r--    1 david supergroup          0 2015-09-21 17:50 /out/ou  
t-3.2/_SUCCESS
```

```
-rw-r--r--    1 david supergroup      41952 2015-09-21 17:50 /out/ou  
t-3.2/part-00000
```

displaying results

```
DAI93865 -> FRO40251      1.00000000  
ELE12951 -> FRO40251      0.99056604  
DAI88079 -> FRO40251      0.98672566  
DAI43868 -> SNA82528      0.97297297  
DAI23334 -> DAI62779      0.95454545
```

HW3.3

NA

HW3.4

The approach outlined here relies on multiple invocations of hadoop, each with a single map and reduce function. Each invocation produces a file that will be made available through the hadoop shared cache architecture. This file will be used by the tasks in the next invocation.

- Round 1
 - Map1:
 - emit count (itemset, 1) for all 1-itemsets
 - Reduce1:
 - aggregate counts for all 1-itemsets
 - filter itemsets based on minimum support value
 - emit [itemset, count] for all frequent itemsets
 - file(s) produced by Reduce1 are combined, if more than 1, into file frequentitems1
- Round 2
 - Map2:
 - read in list of frequent 1-itemsets from file frequentitems1
 - emit count (itemset, 1) for all 2-itemsets iff the itemset contains at least one frequent 1-itemset
 - Reduce2:
 - aggregate counts for all 2-itemsets
 - filter itemsets based on minimum support value
 - emit [itemset, count] for all frequent itemsets
 - file(s) produced by Reduce2 are combined, if more than 1, into file frequentitems2
- Round 3
 - Map3:
 - read in list of frequent 2-itemsets from file frequentitems2
 - emit count (itemset, 1) for all 3-itemsets iff the itemset contains at least one frequent 2-itemset
 - Reduce3
 - aggregate counts for all 3-itemsets
 - filter itemsets based on minimum support value
 - read in list of frequent 2-itemsets from file frequentitems2 to make support levels for 2-itemsets available
 - for each frequent 3-itemset (X, Y, Z)
 - calculate confidence values for rules $((X, Y) \rightarrow Z), ((X, Z) \rightarrow Y), ((Y, Z) \rightarrow X)$ using the support values for the 3-itemsets and 2-itemsets, respectively.
 - e.g.,

$$confidence((X, Y) \rightarrow Z) = \frac{support((X, Y) \cup Z)}{support((X, Y))}$$
 - emit rule and confidence value

```

In [202]: 1 from __future__ import print_function
          2 def confidence(itemA, itemB):
          3     with open('./ProductPurchaseData.txt', 'r') as fin:
          4         atotal = 0
          5         btotal = 0
          6         aonlytotal = 0
          7         bonlytotal = 0
          8         bothtotal = 0
          9         for line in fin:
         10             items = line.split()
         11             if itemA in items:
         12                 atotal += 1
         13             if itemB in items:
         14                 btotal += 1
         15             if itemA in items and not itemB in items:
         16                 aonlytotal += 1
         17             if not itemA in items and itemB in items:
         18                 bonlytotal += 1
         19             if itemA in items and itemB in items:
         20                 bothtotal += 1
         21         print('LHS total: {}, LHS exclusive: {}, RHS total: {}, '
         22               + 'RHS exclusive: {}, both: {}'.format(atotal, aonlytotal, btotal, bonlytotal, bothtotal))
         23         print('{} -> {}, confidence: {}'.format(itemA, itemB, bothtotal / float(atotal)))
         24         #print('{} -> {}, confidence: {}'.format(itemB, itemA, bothtotal))
         25
         26
         27

```

```

In [203]: 1 confidence('DAI93865', 'FRO40251')
          2 confidence('SNA98488', 'SNA99873')

```

```

LHS total: {}, LHS exclusive: {}, RHS total: {}, RHS exclusive: 208,
both: 0
DAI93865 -> FRO40251, confidence: 1.0
LHS total: {}, LHS exclusive: {}, RHS total: {}, RHS exclusive: 1, b
oth: 0
SNA98488 -> SNA99873, confidence: 1.0

```