

Program Frequent Itemsets

Consider the following frequent itemset problem: we are given a set of items $A =$

$\{a_1; a_2; \dots; a_n\}$ and a set of baskets $B_1; B_2; \dots; B_m$ such that each basket is a subset

containing a certain number of items, $B_i \subseteq A$. Given a support s , frequent itemsets are

subsets $S \subseteq A$ such that S is a subset of t baskets (that is, the items in S all appear in at least s baskets).

For example, suppose that we have the following baskets:

$B_1 = \{\text{beer; screws; hammer}\}$

$B_2 = \{\text{saw; lugnut}\}$

$B_3 = \{\text{beer; screws; hammer; lugnut}\}$

$B_4 = \{\text{beer; screws; hammer; router}\}$

If our support $s = 2$ then we're looking for all subsets of items that appear in at least 2

of the baskets. In particular:

– $\{\text{lugnut}\}$ as it appears in $B_2; B_3$

– $\{\text{beer}\}$ as it appears in $B_1; B_3; B_4$

– $\{\text{hammer}\}$ as it appears in $B_1; B_3; B_4$

– $\{\text{screw}\}$ as it appears in $B_1; B_3; B_4$

– $\{\text{beer; screw}\}$ as it appears in $B_1; B_3; B_4$

– $\{\text{hammer; screw}\}$ as it appears in $B_1; B_3; B_4$

– $\{\text{hammer; beer}\}$ as it appears in $B_1; B_3; B_4$

This problem is widely seen in commerce, language analysis, search, and financial ap-

plications to learn association rules. For example, a customer analysis may show that

people often buy nails when they buy a hammer (the pair has high support), so run a

sale on hammers and jack up the price of nails!

For this exercise, we will simplify this problem by computing the support for all combinations of 1, 2, and 3 items. That is, for every single item, for every pair, and for every triple, you will compute the number of baskets in which the combination appears.

Your program will accept a file name as a command line argument which will include an instance of this problem with the following format: the first line will contain a single integer m representing the number of baskets. Each subsequent line contains a comma delimited list of items in the basket. An example:

beer, screws, hammer

saw, lugnuts

beer, screws, hammer, lugnuts

beer, screws, hammer, router

Your output will include the combination as well as the number of times it appears. For

brevity, you may exclude those combinations for which the count is zero.

For example:

Items: [beer, router, screws, saw, lugnuts, hammer]

Number of baskets: 4

3 => [hammer]

3 => [beer]

1 => [router]

3 => [screws]

1 => [saw]

2 => [lugnuts]

1 => [saw, lugnuts]

1 => [beer, lugnuts]

1 => [router, hammer]

1 => [lugnuts, hammer]

1 => [router, screws]

1 => [beer, router]

3 => [beer, screws]

3 => [screws, hammer]

1 => [screws, lugnuts]
3 => [beer, hammer]
3 => [beer, screws, hammer]
1 => [beer, lugnuts, hammer]
1 => [beer, router, screws]
1 => [router, screws, hammer]
1 => [screws, lugnuts, hammer]
1 => [beer, screws, lugnuts]
1 => [beer, router, hammer]

Using PHP, name your script itemSet.php ; it should be invocable from the command line as:

itemSet.php infile.txt and output the results to the standard output.