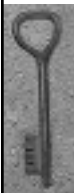



Online Association Rule Mining

By Christian Hidber
Department of Electrical Engineering
and Computer Science, UC Berkeley



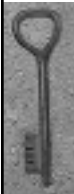
Background

- ◆ Mining for association rules is a form of data mining
- ◆ An example:
 - 65% of all customers who buy pasta and tomato sauce also buy parmesan cheese and red wine
- ◆ Useful for customer segmentation, cross-marketing, catalog design and product placement
- ◆ Online aggregation




Association Rule

- ◆ An association rule is an expression $X \Rightarrow Y$ where X and Y are disjoint itemsets
 - The *confidence* of this rule is the fraction of all transactions containing X that also contain Y
 - The *support* of this rule is the support of $X \cup Y$.
- ◆ A support must \geq a user-specified support threshold
- ◆ A confidence must \geq a user-specified confidence threshold



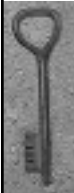
What's the problem?

- ◆ Finding association rules is a very CPU and memory intensive task
- ◆ All traditional algorithms operate offline
- ◆ User does not know the appropriate thresholds in advance



Traditional algorithms

- ◆ Apriori
- ◆ Partition
- ◆ Dynamic Itemset Counting
- ◆ OLAP-style
 - Support threshold s specified before the precomputation of the large itemsets
 - Large Itemset computation remains offline
 - Only rules with support $\geq s$ can be generated
- ◆ User fixed support threshold in advance
- ◆ No feedback to user
- ◆ May need more than two scans



Carma

- ◆ Continuous feedback
 - continuously produces association rules, while the list of purchases is scanned
- ◆ user controllable
 - During the first scan the user is free to change the support and confidence thresholds "on the fly"
- ◆ deterministic and accurate results
 - guarantees that it produces all association rules after at most 2 scans and for each rule its precise support and confidence value

Carma Algorithm: Phase I

- Count(v): number of occurrence of itemset v since v was inserted
- First Trans(v): index of the transaction at which v was inserted
- maxMissd(v): upper bound on the number of occurrences of v before v was inserted
- Minsupport: $\text{count}(v)/i$ lower bound
- Maxsupport: $(\text{maxMissd}(v) + \text{count}(v))/i$ upper bound

Carma Algorithm: Phase I

- Support lattice: a superset of all large itemsets
- Support sequence: a sequence of support threshold $\sigma = (\sigma_1, \sigma_2, \dots)$ σ_i denotes support threshold for the i-th transaction
- $\lceil \sigma \rceil_i$ denotes the least monotone decreasing sequence. It is called ceiling of σ up to i. A sharp lower bound relative to which V is a support lattice.
- $\text{Avg}(\sigma) = 1/i * \sum \sigma_j$ (where $1 \leq j \leq i$) It's the running average of σ up to i.

Carma Algorithm: Phase I

```

Function PhaseI( transaction sequence  $\{t_1, \dots, t_n\}$ , support sequence  $\sigma = (\sigma_1, \dots, \sigma_n)$ , support lattice V )
begin
   $V := \{ \emptyset \}$ ;  $\text{maxMissd}(v) := 0$ ;  $\text{firstTrans}(v) := 0$ ;  $\text{count}(v) := 0$ ;
  for  $i$  from 1 to n do
    // O Insertion
    for all  $v \subseteq V$  with  $v \subseteq t_i$  do  $\text{count}(v) := i$ ;  $cd$ ;
    // O Prune
    for all  $v \subseteq t_i$  with  $v \notin V$  do
      if  $\lceil \sigma \rceil_i(v) < \sigma_i$  and  $\text{maxSupport}(v) \geq \sigma_i$  then
         $V := V \cup \{v\}$ ;
         $\text{firstTrans}(v) := i$ ;
         $\text{count}(v) := 1$ ;
         $\text{maxMissd}(v) := \min \{ \lceil \sigma \rceil_i(v) - \text{firstTrans}(v) + 1, \text{count}(v) \} + i - 1$ ;
         $\text{maxMissd}(w) := \max \{ \text{maxMissd}(w) + \text{count}(w) - 1 \mid w \subseteq v \}$ ;
      else  $\text{count}(v) := 0$ ;
    end
  end
  // O Prune
   $E := \{ v \in V \mid \text{count}(v) = 0 \}$ ;
   $V := \{ v \in V \mid \text{maxSupport}(v) \geq \sigma_i \text{ or } |v| = 1 \}$ ;
end
return V;
end

```

Figure 3

Carma Algorithm: Phase I

- Heart of the algorithm: $\text{maxMissed}(v)$
- $\text{maxMissed}(v) \leq \text{maxMissed}(w) + \text{count}(w) - 1$
- Support $i(w) \geq \text{support } i(v)$ for all subsets w of v and w in t_i
- $\text{maxMissed}(v) \leq \lfloor (i-1) \text{avg } i-1(\lceil \sigma \rceil_{i-1}) \rfloor + |v| - 1$
- Support $(i-1)(v) \leq \text{avg } i-1(\lceil \sigma \rceil_{i-1}) + (|v| - 1)/(i-1)$
- $\text{maxMissed}(v)$ is defined as
 $\text{Min} \{ \lfloor (i-1) \text{avg } i-1(\lceil \sigma \rceil_{i-1}) \rfloor + |v| - 1, \text{maxMissed}(w) + \text{count}(w) - 1 \mid w \subseteq v \}$.
- $\text{maxMissed}(v) \leq i - 1$ for the current transaction index.

Carma: theorem 1

Theorem 1 Let V be the lattice returned by PhaseI(T, σ) for a transaction sequence T of length n and support sequence σ . Then V is a support lattice relative to the support threshold

$$\text{avg}_i(\sigma) = \frac{\sum_{j=1}^{i-1} \sigma_j}{i-1} \quad (1)$$

with v the maximal cardinality of a large itemset in T. For any itemset v

$$\text{support}_i(v) \geq \text{avg}_i(\sigma) + \frac{|v| - 1}{i-1} \quad \text{implies } v \in V$$

- The term $(c+1)/n$ is desirable
- The term $\text{avg}_i(\lceil \sigma \rceil_i)$ is a sharp lower bound relative to which V is a support lattice
- Support guarantee may not match the threshold specified by the user, but this guarantees *converges* to the user-specified threshold if the user keeps it constant for a large number of transactions

Carma: Phase I algorithm example

- Count(v): number of occurrence of itemset v since v was inserted
- First Trans(v): index of the transaction at which v was inserted
- maxMissd(v): upper bound on the number of occurrences of v before v was inserted
- Minsupport: $\text{count}(v)/i$ lower bound
- Maxsupport: $(\text{maxMissd}(v) + \text{count}(v))/i$ upper bound
- $\text{maxMissed}(v) = \text{Min} \{ \lfloor (i-1) \text{avg } i-1(\lceil \sigma \rceil_{i-1}) \rfloor + |v| - 1, \text{maxMissed}(w) + \text{count}(w) - 1 \mid w \subseteq v \}$.

Carma: Changing support thresholds

- ◆ To improve the speed of convergence
 - Run phase I with a lower threshold of $s*0.9$ instead of s .
 - Increase the threshold from $s*0.9$ to s , as the guaranteed threshold reaches s .

Carma: phase II

```

Function PhaseII (support lattice V, transaction sequence  $t_1, \dots, t_n$ )
  lattice  $f_1, f_2 = \emptyset$ 
  support sequence  $s = \{ \}$  support lattice
  begin
     $V := V \cup \{ t_1 \}$ ,  $maxSupport(t_1) < \alpha_1 \}$ 
    while  $\exists t \in V$ ,  $t < t$ ,  $f_1 \text{ firstTransaction}$  do
       $t := t + 1$ 
      for all  $v \in V$  do
         $f_1 := firstTransaction$ 
        if  $v \subseteq t_1$  and  $f_1 < t$  then  $insert(v) \cup v$ ,  $maxSupport(f_1) < \alpha_1$ 
        if  $f_1 = t$  then
           $maxSupport(f_1) := 0$ 
          for all  $w \in V$ ,  $w \subseteq v$  and  $maxSupport(w) > maxSupport(f_1)$  do
             $maxSupport(w) :=  $\maxSupport(w) -  $\maxSupport(f_1)$$ 
      end
    end
    if  $maxSupport(t) < \alpha_1$ , then  $V := V \cup \{ t \}$ 
  end while
  return V
  end$ 
```

Figure 7

Carma Implementation

- ◆ Dataset with 100k transactions of an average size of 10 items chosen from 10k items and an average large itemset size of 4
- ◆ All itemsets are stored in a single hashtable
 - Itemsets as keys: quickly access any subset

Performance: Carma, Apriori and DIC

- ◆ At thresholds 0.25% and below Carma outperform Apriori and DIC
 - Less number of scans
 - Smaller lattice maintained by Carma

Support Intervals

- ◆ Phase I maintains a superset of the large itemsets but not necessarily for the full transaction sequence
- ◆ Size of the support intervals (given by minsupport and maxsupport)
 - Average size 0.042% at threshold of 0.1% while 50% of all itemsets with an interval size below 0.004%

Conclusion

- ◆ Carma-compute large itemsets online
- ◆ Continuously produces large itemsets along with a shrinking support interval for each itemset
- ◆ Allow user to change the support threshold anytime during the first scan and always completes in at most 2 scan
- ◆ Carma's itemset lattice quickly approximates a superset of all large itemset while the sizes of the corresponding support intervals shrink rapidly
- ◆ Second scan is not needed when shrinking support intervals suffice so phase I can be used continuously