

# ExAnte: Anticipated Data Reduction in Constrained Pattern Mining

Francesco Bonchi<sup>†‡</sup>, Fosca Giannotti<sup>†</sup>, Alessio Mazzanti,<sup>‡</sup> and Dino Pedreschi<sup>‡</sup>

Pisa KDD Laboratory \*\*  
<http://www-kdd.cnuce.cnr.it>

<sup>†</sup>ISTI - CNR Area della Ricerca di Pisa, Via Giuseppe Moruzzi, 1 - 56124 Pisa, Italy  
e-mail [Giannotti@cnuce.cnr.it](mailto:Giannotti@cnuce.cnr.it)

<sup>‡</sup>Department of Computer Science, University of Pisa  
Via F. Buonarroti 2, 56127 Pisa, Italy  
e-mail {[bonchi](mailto:bonchi@di.unipi.it),[mazzanti](mailto:mazzanti@di.unipi.it),[pedre](mailto:pedre@di.unipi.it)}@di.unipi.it

**Abstract.** Constraint pushing techniques have been proven to be effective in reducing the search space in the frequent pattern mining task, and thus in improving efficiency. But while pushing anti-monotone constraints in a level-wise computation of frequent itemsets has been recognized to be always profitable, the case is different for monotone constraints. In fact, monotone constraints have been considered harder to push in the computation and less effective in pruning the search space. In this paper, we show that this prejudice is ill founded and introduce ExAnte, a pre-processing data reduction algorithm which reduces dramatically both the search space and the input dataset in constrained frequent pattern mining. Experimental results show a reduction of orders of magnitude, thus enabling a much easier mining task. ExAnte can be used as a pre-processor with any constrained pattern mining algorithm.

## 1 Introduction

Constrained itemset mining i.e., finding all itemsets included in a transaction database that satisfy a given set of constraints, is an active research theme in data mining [3, 6–12]. The most studied constraint is the frequency constraint, whose anti-monotonicity is used to reduce the exponential search space of the problem. Exploiting the anti-monotonicity of the frequency constraint is known as *a priori trick* [1, 2]: it dramatically reduces the search space making the computation feasible. Frequency is not only computationally effective, it is also semantically important since frequency provides “support” to any discovered knowledge. For these reasons frequency is the base constraint of what is generally referred to as *frequent itemset mining*. However, many other constraints can facilitate user-focussed exploration and control, as well as reduce the computation. For instance, a user could be interested in mining all frequently purchased itemsets

---

\*\* The present research is funded by “Fondazione Cassa di Risparmio di Pisa” under the “WebDigger Project”.

having a total price greater than a given threshold and containing at least two products of a given brand. Among these constraints, classes have been individuated which exhibit nice properties. The class of anti-monotone constraints is the most effective and easy to use in order to prune the search space. Since any conjunction of anti-monotone constraints is in turn anti-monotone, we can use the *a priori trick* to exploit completely the pruning power of the conjunction: the more anti-monotone constraints, the more selective the *a priori trick* will be.

The dual class, monotone constraints, has been considered more complicated to exploit and less effective in pruning the search space. As highlighted by Boulicaut and Jeudy in [3], pushing monotone constraints can lead to a reduction of anti-monotone pruning. Therefore, when dealing with a conjunction of monotone and anti-monotone constraints we face a tradeoff between anti-monotone and monotone pruning. Our observation is that the above consideration holds only if we focus completely on the search space of all itemsets, which is the approach followed by the work done so far.

In this paper we show that the most effective way of attacking the problem is to reason on both the itemsets search space and the transactions input database *together*. In this way, pushing monotone constraints does not reduce anti-monotone pruning opportunities, on the contrary, such opportunities are boosted. Dually, pushing anti-monotone constraints boosts monotone pruning opportunities: the two components strengthen each other recursively. We prove our previous statement by introducing ExAnte, a pre-processing data reduction algorithm which reduces dramatically both the search space and the input dataset in constrained frequent pattern mining.

ExAnte can exploit any constraint which has a monotone component, therefore also succinct monotone constraints [9] and convertible monotone constraints [10, 11] can be used to reduce the mining computation. Being a preprocessing algorithm, ExAnte can be coupled with any constrained pattern mining algorithm, and it is always profitable to start any constrained patterns computation with an ExAnte preprocess. The correctness of ExAnte is formally proven in this paper, by showing that the reduction of items and transaction database does not affect the set of constrained frequent patterns, which are solutions to the given problem, as well as their support. We discuss a thorough experimentation of the algorithm, which points out how effective the reduction is, and which potential benefits it offers to subsequent frequent pattern computation.

### **Our contributions:**

Summarizing, the data reduction algorithm proposed in this paper is characterized by the following:

- ExAnte uses, for the first time, the real synergy of monotone and anti-monotone constraints to prune the search space and the input dataset: the total benefit is greater than the sum of the two individual benefits.
- ExAnte can be used with any constraint which has a monotone component: therefore also succinct monotone constraints and convertible monotone constraints can be exploited.

- ExAnte maintains the exact support of each solution itemset: a necessary condition if we want to compute Association Rules.
- ExAnte can be used to make feasible the discovery of particular patterns which can be discovered only at very low support level, for which the computation is unfeasible for traditional algorithms.
- Being a pre-processing algorithm, ExAnte can be coupled with any constrained pattern mining algorithm, and it is always profitable to start any constrained pattern computation with an ExAnte preprocess.
- ExAnte is efficient and effective: even a very large input dataset can be reduced of an order of magnitude in a small computation.
- A thorough experimental study has been performed with different monotone constraints on various datasets (both real world and synthetic datasets), and the results are described in details.

## 2 Problem Definition

Let  $Items = \{x_1, \dots, x_n\}$  be a set of distinct literals, usually called **items**. An **itemset**  $X$  is a subset of  $Items$ . If  $|X| = k$  then  $X$  is called a **k-itemset**. A **transaction** is a couple  $\langle tID, X \rangle$  where  $tID$  is the unique transaction identifier and  $X$  is the content of the transaction (an itemset). A **transaction database**  $TDB$  is a finite set of transactions. An itemset  $X$  is contained in a transaction  $\langle tID, Y \rangle$  if  $X \subseteq Y$ . Given a transaction database  $TDB$  the subset of transactions which contain an itemset  $X$  is denoted  $TDB[X]$ . The **support** of an itemset  $X$ , written  $supp_{TDB}(X)$  is the cardinality of  $TDB[X]$ . Given a user-defined **minimum support**  $\delta$ , an itemset  $X$  is called **frequent** in  $TDB$  if  $supp_{TDB}(X) \geq \delta$ . This defines the frequency constraint  $\mathcal{C}_{freq}[TDB]$ : if  $X$  is frequent we write  $\mathcal{C}_{freq}[TDB](X)$  or simply  $\mathcal{C}_{freq}(X)$  when the dataset is clear from the context.

Let  $Th(\mathcal{C}) = \{X | \mathcal{C}(X)\}$  denotes the set all itemsets  $X$  that satisfy constraint  $\mathcal{C}$ . The *frequent itemset mining problem* requires to compute the set of all frequent itemsets  $Th(\mathcal{C}_{freq})$ . In general given a conjunction of constraints  $\mathcal{C}$  the *constrained itemset mining problem* requires to compute  $Th(\mathcal{C})$ ; the *constrained frequent itemsets mining problem* requires to compute  $Th(\mathcal{C}_{freq}) \cap Th(\mathcal{C})$ .

We now formally define the notion of anti-monotone and monotone constraints.

**Definition 1.** Given an itemset  $X$ , a constraint  $\mathcal{C}_{AM}$  is anti-monotone if

$$\forall Y \subseteq X : \mathcal{C}_{AM}(X) \Rightarrow \mathcal{C}_{AM}(Y)$$

If  $\mathcal{C}_{AM}$  holds for  $X$  then it holds for any subset of  $X$ . □

The frequency constraint is clearly anti-monotone. This property is used by the APRIORI algorithm with the following heuristic: if an itemset  $X$  does not satisfy  $\mathcal{C}_{freq}$ , then no superset of  $X$  can satisfy  $\mathcal{C}_{freq}$ , and hence they can be pruned. This pruning can affect a large part of the search space, since itemsets form a lattice. Therefore the APRIORI algorithm operates in a level-wise fashion moving bottom-up on the itemset lattice, and each time it finds an infrequent itemset it prunes away all its supersets.

**Definition 2.** Given an itemset  $X$ , a constraint  $\mathcal{C}_M$  is monotone if:

$$\forall Y \supseteq X : \mathcal{C}_M(X) \Rightarrow \mathcal{C}_M(Y)$$

independently from the given input transaction database. If  $\mathcal{C}_M$  holds for  $X$  then it holds for any superset of  $X$ .  $\square$

Note that in the last definition we have required a monotone constraint to be satisfied independently from the given input transaction database. This is necessary since we want to distinguish between simple monotone constraints and global constraints such as the “*infrequency constraint*”:

$$supp_{TDB}(X) \leq \delta.$$

This constraint is still monotone but has different properties since it is dataset dependent and it requires dataset scans in order to be computed. Obviously, since our pre-processing algorithm reduces the transaction dataset, we want to exclude the infrequency constraint from our study. Thus, our study focuses on “local” monotone constraints, in the sense that they depend exclusively on the properties of the itemset (as those ones in Table 1), and not on the underlying transaction database.

The general problem that we consider in this paper is the mining of itemsets which satisfy a conjunction of monotone and anti-monotone constraints:

$$Th(\mathcal{C}_{AM}) \cap Th(\mathcal{C}_M).$$

Since any conjunction of anti-monotone constraints is an anti-monotone constraint, and any conjunction of monotone constraints is a monotone constraint, in this paper we focus on the problem given by the conjunction of frequency anti-monotone constraint ( $\mathcal{C}_{AM} \equiv supp_{TDB}(X) \geq \delta$ ), with various simple monotone constraints (see Table 1).

$$Th(\mathcal{C}_{freq}) \cap Th(\mathcal{C}_M).$$

However, our algorithm can work with any conjunction of anti-monotone constraints, provided that the frequency constraint is included in the conjunction: the more anti-monotone constraints, the larger the data reduction will be.

Monotone constraint	$\mathcal{C}_M \equiv$
cardinality	$ X  \geq n$
sum of prices	$sum(X.prices) \geq n$
maximum price	$max(X.prices) \geq n$
minimum price	$min(X.prices) \leq n$
range of prices	$range(X.prices) \geq n$

**Table 1.** Monotone constraints considered in our analysis.

### 3 Search Space and Input Data Reduction

As already stated, if we focus only on the itemsets lattice, pushing monotone constraint can lead to a less effective anti-monotone pruning. Suppose that an itemset has been removed from the search space because it does not satisfy some monotone constraints  $\mathcal{C}_M$ . This pruning avoids checking support for it, but it may be that if we check support, the itemset could result to be infrequent, and thus all its supersets could be pruned away. By monotone pruning an itemset we risk to lose anti-monotone pruning opportunities given from the itemset itself. The tradeoff is clear [3]: pushing monotone constraint can save tests on anti-monotone constraints, however the results of these tests could have lead to more effective pruning. In order to obtain a real amalgam of the two opposite pruning strategies we have to consider the constrained frequent patterns problem in its whole: not focussing only on the itemsets lattice but considering it together with the input database of transactions. In fact, as proved by the theorems in the following section, monotone constraints can prune away transactions from the input dataset *without losing solutions*. This monotone pruning of transactions has got another positive effect: while reducing the number of transactions in input it reduces the support of items too, hence the total number of frequent 1-itemsets. In other words, the monotone pruning of transactions strengthens the anti-monotone pruning. Moreover, infrequent items can be deleted by the computation and hence pruned away from the transactions in the input dataset. This anti-monotone pruning has got another positive effect: reducing the size of a transaction which satisfies a monotone constraint can make the transaction violates the monotone constraint. Therefore a growing number of transactions which do not satisfy the monotone constraint can be found. We are clearly inside a loop where two different kinds of pruning cooperates to reduce the search space and the input dataset, strengthening each other step by step until no more pruning is possible (a fix-point has been reached). This is precisely the idea underlying ExAnte.

#### 3.1 ExAnte Properties

In this section we formalize the basic ideas of ExAnte. First we define the two kinds of reduction, then we prove the completeness of the method. In the next section we provide the pseudo-code of the algorithm.

**Definition 3.** [ $\mu$ -reduction] Given a transaction database  $TDB$  and a monotone constraint  $\mathcal{C}_M$ , we define the  $\mu$ -reduction of  $TDB$  as the dataset resulting from pruning the transactions that do not satisfy  $\mathcal{C}_M$ .

$$\mu[TDB]_{\mathcal{C}_M} = \{\langle tID, X \rangle \mid \langle tID, X \rangle \in TDB \wedge X \in Th(\mathcal{C}_M)\}$$

□

**Definition 4.** [ $\alpha$ -reduction] Given a transaction database  $TDB$ , a transaction  $\langle tID, X \rangle$  and a frequency constraint  $\mathcal{C}_{freq}[TDB]$ , we define the  $\alpha$ -reduction of

$\langle tID, X \rangle$  as the transaction resulting from pruning the items in  $X$  that do not satisfy  $\mathcal{C}_{freq}[TDB]$ .

$$\alpha[\langle tID, X \rangle]_{\mathcal{C}_{freq}[TDB]} = \langle tID, (F_1 \cap X) \rangle$$

Where:  $F_1 = \{i \in Items \mid \{i\} \in Th(\mathcal{C}_{freq}[TDB])\}$ . We define the  $\alpha$ -reduction of  $TDB$  as the dataset resulting from the  $\alpha$ -reduction of all transactions in  $TDB$ .  $\square$

The following two key theorems state that we can always  $\mu$ -reduce and  $\alpha$ -reduce a dataset without reducing the support of solution itemsets. Moreover, since satisfaction of  $\mathcal{C}_M$  is independent from the transaction dataset, all solution itemsets will still satisfy it. Therefore, we can always  $\mu$ -reduce and  $\alpha$ -reduce a dataset without losing solutions.

**Theorem 5 ( $\mu$ -reduction correctness).** *Given a transaction database  $TDB$ , a monotone constraint  $\mathcal{C}_M$ , and a frequency constraint  $\mathcal{C}_{freq}$ , we have that:*

$$\forall X \in Th(\mathcal{C}_{freq}[TDB]) \cap Th(\mathcal{C}_M) : supp_{TDB}(X) = supp_{\mu[TDB]_{\mathcal{C}_M}}(X).$$

*Proof.* Since  $X \in Th(\mathcal{C}_M)$ , all transactions containing  $X$  will also satisfy  $\mathcal{C}_M$  for the monotonicity property. Therefore no transaction containing  $X$  will be  $\mu$ -pruned (in other words:  $TDB[X] \subseteq \mu[TDB]_{\mathcal{C}_M}$ ). This, together with the definition of support, implies the thesis.  $\square$

**Theorem 6 ( $\alpha$ -reduction correctness).** *Given a transaction database  $TDB$ , a monotone constraint  $\mathcal{C}_M$ , and a frequency constraint  $\mathcal{C}_{freq}$ , we have that:*

$$\forall X \in Th(\mathcal{C}_{freq}[TDB]) \cap Th(\mathcal{C}_M) : supp_{TDB}(X) = supp_{\alpha[TDB]_{\mathcal{C}_{freq}}}(X).$$

*Proof.* Since  $X \in Th(\mathcal{C}_{freq})$ , all subsets of  $X$  will be frequent (by the anti-monotonicity of frequency). Therefore no 1-itemsets in  $X$  will be  $\alpha$ -pruned (in other words:  $TDB[X] \subseteq \alpha[TDB]_{\mathcal{C}_{freq}}$ ). This, together with the definition of support, implies the thesis.  $\square$

### 3.2 ExAnte Algorithm

The two theorems above suggest a fix-point computation. ExAnte starts the first iteration as any frequent pattern mining algorithm: counting the support of singleton items. Items that are not frequent are thrown away once and for all. But during this first count only transactions that satisfy  $\mathcal{C}_M$  are considered. The other transactions are signed to be pruned from the dataset ( $\mu$ -reduction). Doing so we reduce the number of interesting 1-itemsets. Even a small reduction of this number represents a huge pruning of the search space. At this point ExAnte deletes from alive transactions all infrequent items ( $\alpha$ -reduction). This pruning can reduce the monotone value (for instance, the total sum of prices) of some alive transactions, possibly resulting in a violation of the monotone

---

Procedure: **ExAnte**( $TDB, \mathcal{C}_M, min\_supp$ )

---

1.  $I = \emptyset$ ;
2. **forall** transactions  $t$  in  $TDB$  **do**
3.     **if**  $\mathcal{C}_M(t)$  **then forall** items  $i$  in  $t$  **do**
4.          $i.count++$ ; **if**  $i.count == min\_supp$  **then**  $I = I \cup \{i\}$ ;
5.  $old\_number\_interesting\_items = |Items|$ ;
6. **while**  $|I| < old\_number\_interesting\_items$  **do**
7.      $TDB = \alpha[TDB]_{\mathcal{C}_{freq}}$ ;
8.      $TDB = \mu[TDB]_{\mathcal{C}_M}$ ;
9.      $old\_number\_interesting\_items = |I|$ ;
10.      $I = \emptyset$ ;
11.     **forall** transactions  $t$  in  $TDB$  **do**
12.         **forall** items  $i$  in  $t$  **do**
13.              $i.count++$ ;
14.             **if**  $i.count == min\_supp$  **then**  $I = I \cup \{i\}$ ;
15.     **end while**

---

**Fig. 1.** The ExAnte algorithm pseudo-code.

constraints. Therefore we have another opportunity of  $\mu$ -reducing the dataset. But  $\mu$ -reducing the dataset we create new opportunities for  $\alpha$ -reduction, which can turn in new opportunities for  $\mu$ -reduction, and so on, until a fix-point is reached. The pseudo-code of ExAnte algorithm follows:

Clearly, a fix-point is eventually reached after a finite number of iterations, as at each step the number of alive items strictly decreases.

### 3.3 Run-through Example

Suppose that the transaction and price datasets in Table 2 are given. Suppose that we want to compute frequent itemsets ( $min\_supp = 4$ ) with a sum of prices  $\geq 45$ . During the first iteration the total price of each transaction is checked to avoid using transactions which do not satisfy the monotone constraint. All transaction with a sum of prices  $\geq 45$  are used to count the support for the singleton items. Only the fourth transaction is discarded. At the end of the count we find items  $a, e, f$  and  $h$  to be infrequent. Note that, if the fourth transaction had not been discarded, items  $a$  and  $e$  would have been counted as frequent. At this point we perform an  $\alpha$ -reduction of the dataset: this means removing  $a, e, f$  and  $h$  from all transactions in the dataset. After the  $\alpha$ -reduction we have more opportunities to  $\mu$ -reduce the dataset. In fact transaction 2, which at the beginning has a total price of 63, now has its total price reduced to 38 due to the pruning of  $a$  and  $e$ . This transaction can now be pruned away. The same reasoning holds for transactions number 7 and 9. At this point ExAnte counts once again the support of alive items with the reduced dataset. The item  $g$  which initially has got a support of 5 now has become infrequent (see Table 2 (c) for items support iteration by iteration). We can  $\alpha$ -reduce again the dataset, and

then  $\mu$ -reduce. After the two reductions transaction number 5 does not satisfy anymore the monotone constraint and it is pruned away. ExAnte counts again the support of items on the reduced datasets but no more items are found to have turned infrequent. The fix-point has been reached at the third iteration: the dataset has been reduced from 9 transactions to 4 transactions (number 1,3,6 and 8), and interesting itemsets have shrunk from 8 to 3 ( $b, c$  and  $d$ ). At this point any constrained frequent pattern mining algorithm would find very easily the unique solution to problem which is the 3-itemset  $\{b, c, d\}$ .

item	price
a	5
b	8
c	14
d	30
e	20
f	15
g	6
h	12

(a)

tID	Itemset	Total price
1	b,c,d,g	58
2	a,b,d,e	63
3	b,c,d,g,h	70
4	a,e,g	31
5	c,d,f,g	65
6	a,b,c,d,e	77
7	a,b,d,f,g,h	76
8	b,c,d	52
9	b,e,f,g	49

(b)

Supports			
Items	1 <sub>st</sub>	2 <sub>nd</sub>	3 <sub>rd</sub>
a	3	†	†
b	7	4	4
c	5	5	4
d	7	5	4
e	3	†	†
f	3	†	†
g	5	3	†
h	2	†	†

(c)

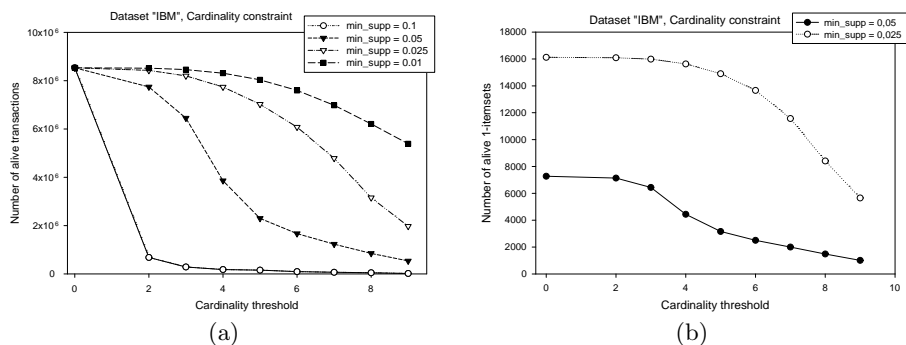
**Table 2.** Run-through Example: price table (a) and transaction database (b), items and their supports iteration by iteration (c).

## 4 Experimental Results

In this section we deeply describe the experimental study that we have conducted with different monotone constraints on various datasets. In particular, the monotone constraints used in the experimentation are in Table 1. In addition, we have experimented a harder to exploit constraint:  $avg(X.prices) \geq n$ . This constraint is clearly neither monotone nor anti-monotone, but can exhibit a monotone (or anti-monotone) behavior if items are ordered by ascending (or descending) price, and frequent patterns are computed following a prefix-tree approach. This class of constraints, named *convertible*, has been introduced in [10]. In our experiments the constraint  $avg(X.prices) \geq n$  is treated by inducing a weaker monotone constraint:  $max(X.prices) \geq n$ . Note that in every reported experiment we have chosen monotone constraints thresholds that are not very selective: there are always solutions to the given problem. In the experiments reported in this paper we have used two datasets. “IBM” is a synthetic dataset obtained with the most commonly adopted dataset generator, available from IBM Almaden<sup>1</sup>. We have generate a very large dataset since we have not been

<sup>1</sup> <http://www.almaden.ibm.com/software/quest/Resources/datasets/synndata.html#assocSynData>





**Fig. 2.** Transactions reduction (a), and interesting 1-itemsets reduction (b), on dataset “IBM”.

able to find a real-world dataset over one million transactions. “Italian” is a real-world dataset obtained from an Italian supermarket chain within a market-basket analysis project conducted by our research lab, few years ago (note that the prices are in the obsolete currency Italian Lira).

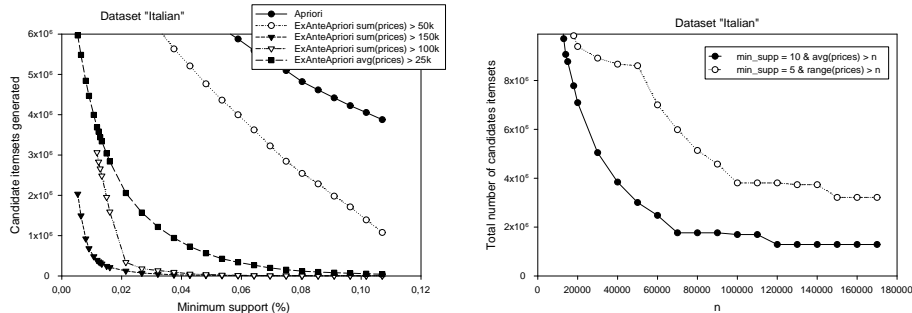
Dataset	Transactions	Items	Max Trans Size	Avg Trans Size
IBM	8,533,534	100,000	37	11.21
Italian	186,824	4800	31	10.42

Dataset	Min Price	Max Price	Avg Price
Italian	100	900,000	6454.87

**Table 3.** Characteristics of the datasets used in the experiments.

For a more detailed report of our experiments see [5]. In Figure 2 (a) the reduction of the number of transactions w.r.t the cardinality threshold is shown for four different support thresholds on the synthetic dataset. When the cardinality threshold is equal to zero the number of transactions equals the total number of transactions in the database, since there is no monotone pruning. Already for a low support threshold as 0.1% with a cardinality threshold equals to 2 the number of transactions decreases dramatically. Figure 2 (b) describes the reduction of number of interesting 1-itemsets on the same dataset.

As already stated, even a small reduction in the number of relevant 1-itemsets represents a very large pruning of the search space. In our experiments, as a measure of the search space explored, we have considered the number of candidate itemsets generated by a level-wise algorithm such as Apriori. In Figure 3 is



**Fig. 3.** Search space reduction on dataset "Italian".

reported a comparison of the number of candidate itemsets generated by Apriori and by *ExAnteApriori* (ExAnte pre-processing followed by Apriori) on the "Italian" dataset with various constraints. The dramatic search space reduction is evident, and it will be confirmed by computation time reported in the next section. How the number of candidate itemsets shrinks by increasing strength of the monotone constraint is also reported in Figure 3. This figure also highlights another interesting feature of ExAnte: even at very low support level ( $\text{min\_supp} = 5$  on a dataset of 186,824 transactions) the frequent patterns computation is feasible if coupled with a monotone constraint. Therefore, ExAnte can be used to make feasible the discovery of particular patterns which can be discovered only at very low support level, for instance:

- extreme purchasing behaviors (such as patterns with a very high average of prices);
- very long patterns (using the cardinality constraint coupled with a very low support threshold).

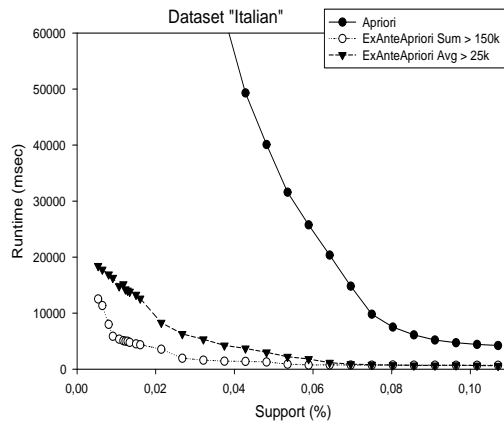
We report run-time comparison between Apriori and *ExAnteApriori*. We have chosen Apriori as the "standard" frequent pattern mining algorithm. Recall that every frequent pattern mining algorithm can be coupled with ExAnte pre-processing obtaining similar benefits. Execution time is reported in Figure 4. The large search space pruning reported in the previous section is here confirmed by the execution time.

## 5 Related Work

Being a pre-processing algorithm, ExAnte can not be directly compared with any previously proposed algorithm for constrained frequent pattern mining. However, it would be interesting to couple ExAnte data reduction with those algorithms and to measure the improve in efficiency. Among constrained frequent pattern mining algorithms, we would like to mention  $\mathcal{FIC}^M$  [11] and the recently proposed DualMiner [4].

Iteration	Transactions	1-itemsets
0	17306	2010
1	13167	1512
2	11295	1205
3	10173	1025
4	9454	901
5	9005	835
6	8730	785
7	8549	754
8	8431	741
9	8397	736
10	8385	734
11	8343	729
12	8316	726
13	8312	724
14	8307	722
15	8304	722

Execution time: 1.5 sec



**Fig. 4.** A typical execution of ExAnte: Dataset “Italian” with  $min\_sup = \%40$  and sum of prices  $\geq 100000$  (on the left); and a runtime comparison between Apriori and ExAnteApriori with two different constraints (on the right).

## 6 Conclusions and Future Work

In this paper we have introduced ExAnte, a pre-processing data reduction algorithm which reduces dramatically the search space the input dataset, and hence the execution time, in constrained frequent pattern mining. We have proved experimentally the effectiveness of our method, using different constraints on various datasets. Due to its capacity in focussing on any particular instance of the problem, ExAnte exhibits very good performance also when one of the two constraints (the anti-monotone or the monotone) is not very selective. This feature makes ExAnte useful to discover particular patterns which can be discovered only at very low support level, for which the computation is unfeasible for traditional algorithms.

We are actually developing a new algorithm for constrained frequent pattern mining, which will take full advantage of ExAnte pre-processing. We are also interested in studying in which other mining tasks ExAnte can be useful. We will investigate its applicability to the constrained mining of closed itemsets, sequential patterns, graphs structure, and other complex kinds of data and patterns. ExAnte executable can be downloaded by our web site:  
<http://www-kdd.cnuce.cnr.it/>

**Acknowledgements.** We are indebted with Laks V.S. Lakshmanan, who first suggested the problem to the first author.

## References

1. R. Agrawal, T. Imielinski, and A. N. Swami. Mining association rules between sets of items in large databases. In P. Buneman and S. Jajodia, editors, *Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data*, pages 207–216, Washington, D.C., 26–28 May 1993.
2. R. Agrawal and R. Srikant. Fast Algorithms for Mining Association Rules in Large Databases. In *Proceedings of the Twentieth International Conference on Very Large Databases*, pages 487–499, Santiago, Chile, 1994.
3. J.-F. Boulicaut and B. Jeudy. Using constraints during set mining: Should we prune or not? In *Actes des Seizime Journes Bases de Donnes Avances BDA'00, Blois (F)*, pages 221–237, 2000.
4. C. Bucila, J. Gehrke, D. Kifer, and W. White. Dualminer: A dual-pruning algorithm for itemsets with constraints. In *Proceedings of the 8th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2002.
5. F. Bonchi, F. Giannotti, A. Mazzanti, and D. Pedreschi. Exante: a preprocessing algorithm for constrained frequent pattern mining. Technical Report ISTI-B4-2003-07, ISTI, 2003.
6. G. Grahne, L. Lakshmanan, and X. Wang. Efficient mining of constrained correlated sets. In *16th International Conference on Data Engineering (ICDE' 00)*, pages 512–524. IEEE, 2000.
7. J. Han, L. V. S. Lakshmanan, and R. T. Ng. Constraint-based, multidimensional data mining. *Computer*, 32(8):46–50, 1999.
8. L. V. S. Lakshmanan, R. T. Ng, J. Han, and A. Pang. Optimization of constrained frequent set queries with 2-variable constraints. *SIGMOD Record (ACM Special Interest Group on Management of Data)*, 28(2), 1999.
9. R. T. Ng, L. V. S. Lakshmanan, J. Han, and A. Pang. Exploratory mining and pruning optimizations of constrained associations rules. In *Proceedings of the ACM SIGMOD International Conference on Management of Data (SIGMOD-98)*, volume 27,2 of *ACM SIGMOD Record*, pages 13–24, New York, June 1–4 1998. ACM Press.
10. J. Pei and J. Han. Can we push more constraints into frequent pattern mining? In R. Ramakrishnan, S. Stolfo, R. Bayardo, and I. Parsa, editors, *Proceedings of the 6th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'00)*, pages 350–354, N. Y., Aug. 20–23 2000. ACM Press.
11. J. Pei, J. Han, and L. V. S. Lakshmanan. Mining frequent item sets with convertible constraints. In *(ICDE'01)*, pages 433–442, 2001.
12. R. Srikant, Q. Vu, and R. Agrawal. Mining association rules with item constraints. In D. Heckerman, H. Mannila, D. Pregibon, and R. Uthurusamy, editors, *Proc. 3rd Int. Conf. Knowledge Discovery and Data Mining, KDD*, pages 67–73. AAAI Press, 14–17 Aug. 1997.