

PaMPa-HD: EXTENSION - TO MODIFY

Daniele Apiletti, Elena Baralis, Tania Cerquitelli,
Paolo Garza, Fabio Pulvirenti^a, and Pietro Michiardi^b

*^aDipartimento di Automatica e Informatica
Politecnico di Torino
Torino, Italy*

Email: name.surname@polito.it

*^bData Science Department
Eurecom
Sophia Antipolis, France
Email: pietro.michiardi@eurecom.fr*

Abstract

Frequent closed itemset mining is among the most complex exploratory techniques in data mining, and provides the ability to discover hidden correlations in transactional datasets. The explosion of Big Data is leading to new parallel and distributed approaches. Unfortunately, most of them are designed to cope with low-dimensional datasets, whereas no distributed high-dimensional frequent closed itemset mining algorithms exists. This work introduces PaMPa-HD, a parallel MapReduce-based frequent closed itemset mining algorithm for high-dimensional datasets, based on Carpenter. The experimental results, performed on both real and synthetic datasets, show the efficiency and scalability of PaMPa-HD.

Keywords:

1. Introduction

In the last years, the increasing capabilities of recent applications to produce and store huge amounts of information, the so called "Big Data", have changed dramatically the importance of the intelligent analysis of data. In both academic and industrial domains, the interest towards data mining, which focuses on extracting effective and usable knowledge from large collections of data, has risen. The need for efficient and highly scalable data mining tools increases with the size of the datasets, as well as their value for businesses and researchers aiming at extracting meaningful insights increases.

Frequent (closed) itemset mining is among the most complex exploratory techniques in data mining. It is used to discover frequently co-occurring items according to a user-provided frequency threshold, called minimum support. Existing mining algorithms revealed to be very efficient on basic datasets but very resource intensive in Big Data contexts. In general, applying data mining techniques to Big Data collections has often entailed to cope with computational costs that represent a critical bottleneck. For this reason, we are witnessing the explosion of parallel and distributed approaches, typically based on distributed frameworks, such as Apache Hadoop [1] and Spark [2]. Unfortunately, most of the scalable distributed techniques for frequent itemset mining have been designed to cope with datasets characterized by few items per transaction (low dimensionality, short transactions), focusing, on the contrary, on very large datasets in terms of number of transactions. Currently, only single-machine implementations exist to address very long transactions, such as Carpenter [3], and no distributed implementations at all.

Nevertheless, many researchers in scientific domains such as bioinformatics or networking, often require to deal with this type of data. For instance, most gene expression datasets are characterized by a huge number of items (related to tens of thousands of genes) and a few records (one transaction per patient or tissue). Many applications in computer vision deal with high-dimensional data, such as face recognition. Some smart-cities studies have built this type of large datasets measuring the occupancy of different car lanes: each transaction describes the occupancy rate in a captor location and in a given timestamp [4]. In the networking domain, instead, the heterogeneous environment provides many different datasets characterized by high-dimensional data, such as URL reputation, advertising, and social network datasets [4, 5].

This work introduces PaMPa-HD, a parallel MapReduce-based frequent closed itemset mining algorithm for high-dimensional datasets, based on the Carpenter algorithm. PaMPa-HD outperforms the single-machine Carpenter implementation and the best state-of-the-art distributed approaches, in both execution time and minimum support threshold. Furthermore, the implementation takes into account crucial design aspects, such as load balancing and robustness to memory-issues.

The paper is organized as follows: Section 2 introduces the frequent (closed) itemset mining problem, Section 3 briefly describes the centralized version of Carpenter, and Section 4 presents the proposed PaMPa-HD algorithm. Section 6 describes the experimental evaluations proving the effectiveness of the proposed technique, Section 7 provides a brief review of the state of the art, and Section 8 discusses possible applications of PaMPa-HD. Finally, Section 9 introduces future works and conclusions.

\mathcal{D}	
tid	items
1	a,b,c,l,o,s,v
2	a,d,e,h,l,p,r,v
3	a,c,e,h,o,q,t,v
4	a,f,v
5	a,b,d,f,g,l,q,s,t

(a) Horizontal representation of \mathcal{D}

TT	
item	tidlist
a	1,2,3,4,5
b	1,5
c	1,3
d	2,5
e	2,3
f	4,5
g	5
h	2,3
l	1,2,5
o	1,3
p	2
q	3,5
r	2
s	1,5
t	3,5
v	1,2,3,4

(b) Transposed representation of \mathcal{D}

$TT _{\{2,3\}}$	
item	tidlist
a	4,5
e	-
h	-
v	4

(c) $TT|_{\{2,3\}}$: example of conditional transposed table

Figure 1: Running example dataset \mathcal{D}

2. Frequent itemset mining background

Let \mathcal{I} be a set of items. A transactional dataset \mathcal{D} consists of a set of transactions $\{t_1, \dots, t_n\}$, where each transaction $t_i \in \mathcal{D}$ is a set of items (i.e.,

$t_i \subseteq \mathcal{I}$) and it is identified by a transaction identifier (tid_i). Figure 1a reports an example of a transactional dataset with 5 transactions. The dataset reported in Figure 1a is used as a running example through the paper.

An itemset I is defined as a set of items (i.e., $I \subseteq \mathcal{I}$) and it is characterized by a tidlist and a support value. The tidlist of an itemset I , denoted by $tidlist(I)$, is defined as the set of tids of the transactions in \mathcal{D} containing I , while the support of I in \mathcal{D} , denoted by $sup(I)$, is defined as the ratio between the number of transactions in \mathcal{D} containing I and the total number of transactions in \mathcal{D} (i.e., $|tidlist(I)|/|\mathcal{D}|$). For instance, the support of the itemset $\{aco\}$ in the running example dataset \mathcal{D} is $2/5$ and its tidlist is $\{1, 3\}$. An itemset I is considered frequent if its support is greater than a user-provided minimum support threshold $minsup$.

Given a transactional dataset \mathcal{D} and a minimum support threshold $minsup$, the Frequent Itemset Mining [6] problem consists in extracting the complete set of frequent itemsets from \mathcal{D} . In this paper, we focus on a valuable subset of frequent itemsets called frequent closed itemsets [3]. Closed itemsets allow representing the same information of traditional frequent itemsets in a more compact form.

A transactional dataset can also be represented in a vertical format, which is usually a more effective representation of the dataset when the average number of items per transactions is orders of magnitudes larger than the number of transactions. In this representation, also called transposed table TT , each row consists of an item i and its list of transactions, i.e., $tidlist(\{i\})$. Let r be an arbitrary row of TT , $r.tidlist$ denotes the tidlist of row r . Figure 6a reports the transposed representation of the running example reported

in Figure 1a.

Given a transposed table TT and a tidlist X , the conditional transposed table of TT on the tidlist X , denoted by $TT|_X$, is defined as a transposed table such that: (1) for each row $r_i \in TT$ such that $X \subseteq r_i.tidlist$ there exists one tuple $r'_i \in TT|_X$ and (2) r'_i contains all tids in $r_i.tidlist$ whose tid is higher than any tid in X .

For instance, consider the transposed table TT reported in Figure 6a. The projection of TT on the tidlist $\{2,3\}$ is the transposed table reported in Figure 1c.

Each transposed table $TT|_X$ is associated with an itemset composed by the items in $TT|_X$. For instance, the itemset associated with $TT|_{\{2,3\}}$ is $\{aehv\}$ (see Figure 1c).

3. The Carpenter algorithm

As discussed in section 7, the most popular techniques (e.g., Apriori and FP-growth) adopt the itemset enumeration approach to mine the frequent itemsets. However, itemset enumeration revealed to be ineffective with datasets with a high average number of items per transactions [3]. To tackle this problem, the Carpenter algorithm [3] was proposed. Specifically, Carpenter is a frequent itemset extraction algorithm devised to handle datasets characterized by a relatively small number of transactions but a huge number of items per transaction. To efficiently solve the itemset mining problem, Carpenter adopts an effective depth-first transaction enumeration approach based on the transposed representation of the input dataset. To illustrate the centralized version of Carpenter, we will use the running example dataset

\mathcal{D} reported in Figure 1a, and more specifically, its transposed version (see Figure 6a). As already described in Section 2, in the transposed representation each row of the table consists of an item i with its tidlist. For instance, the last row of Figure 6a points that item v appears in transactions 1, 2, 3, 4.

Carpenter builds a transaction enumeration tree where each node corresponds to a conditional transposed table $TT|_X$ and its related information (i.e., the tidlist X with respect to which the conditional transposed table is built and its associated itemset). The transaction enumeration tree, when pruning techniques are not applied, contains all the tid combinations (i.e., all the possible tidlists X). Figure 2 reports the transaction enumeration tree obtained by processing the running example dataset. To avoid the generation of duplicate tidlists, the transaction enumeration tree is built by exploring the tids in lexicographical order (e.g., $TT|_{\{1,2\}}$ is generated instead of $TT|_{\{2,1\}}$). Each node of the tree is associated with a conditional transposed table on a tidlist. For instance, the conditional transposed table $TT|_{\{2,3\}}$ in Figure 1c, matches the node $\{2, 3\}$ in Figure 2.

Carpenter performs a depth first search of the enumeration tree to mine the set of frequent closed itemsets. Referring to the tree in Figure 2, the depth first search would lead to the visit of the nodes in the following order: $\{1\}$, $\{1,2\}$, $\{1,2,3\}$, $\{1,2,3,4\}$, $\{1,2,3,4,5\}$, $\{1,2,3,5\}$, $\{\dots\}$. For each node, Carpenter applies a procedure that decides if the itemset associated with that node is a frequent closed itemset or not. Specifically, for each node, Carpenter decides if the itemset associated with the current node is a frequent closed itemset by considering: 1) the tidlist X associated with the node, 2)

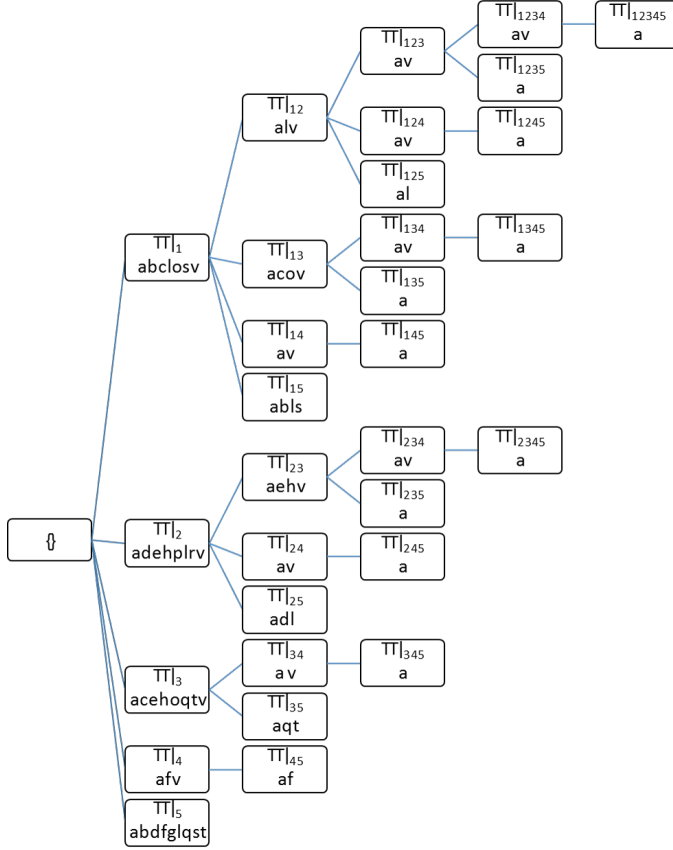


Figure 2: The transaction enumeration tree of the running example dataset in Figure 1a. For the sake of clarity, no pruning rules are applied to the tree.

the conditional transposed table $TT|_X, 3)$ the set of frequent closed itemsets found up to the current step of the tree search, and 4) the enforced minimum support threshold ($minsup$). Based on the theorems reported in [3], if the itemset I associated with the current node is a frequent closed itemset then

I is included in the frequent closed itemset set. Moreover, by exploiting the analysis performed on the current node, part of the remaining search space (i.e., part of the enumeration tree) can be pruned, to avoid the analysis of nodes that will never generate new closed itemsets. At this purpose, three pruning rules are applied on the enumeration tree, based on the evaluation performed on the current node and the associated transposed table $TT|_X$:

- **Pruning rule 1.** If the size of X , plus the number of distinct tids in the rows of $TT|_X$ does not reach the minimum support threshold, the subtree rooted in the current node is pruned.
- **Pruning rule 2.** If there is any tid tid_i that is present in all the tidlists of the rows of $TT|_X$, tid_i is deleted from $TT|_X$. The number of discarded tids is updated to compute the correct support of the itemset associated with the pruned version of $TT|_X$.
- **Pruning rule 3.** If the itemset associated with the current node has been already encountered during the depth first search, the subtree rooted in the current node is pruned because it can never generate new closed itemsets.

The tree search continues in a depth first fashion moving on the next node of the enumeration tree. More specifically, let tid_l be the lowest tid in the tidlists of the current $TT|_X$, the next node to explore is the one associated with $X' = X \cup \{tid_l\}$.

Among the three rules mentioned above, pruning rule 3 assumes a global knowledge of the enumeration tree explored in a depth first manner. This, as detailed in section 4, is very challenging in a distributed environment.

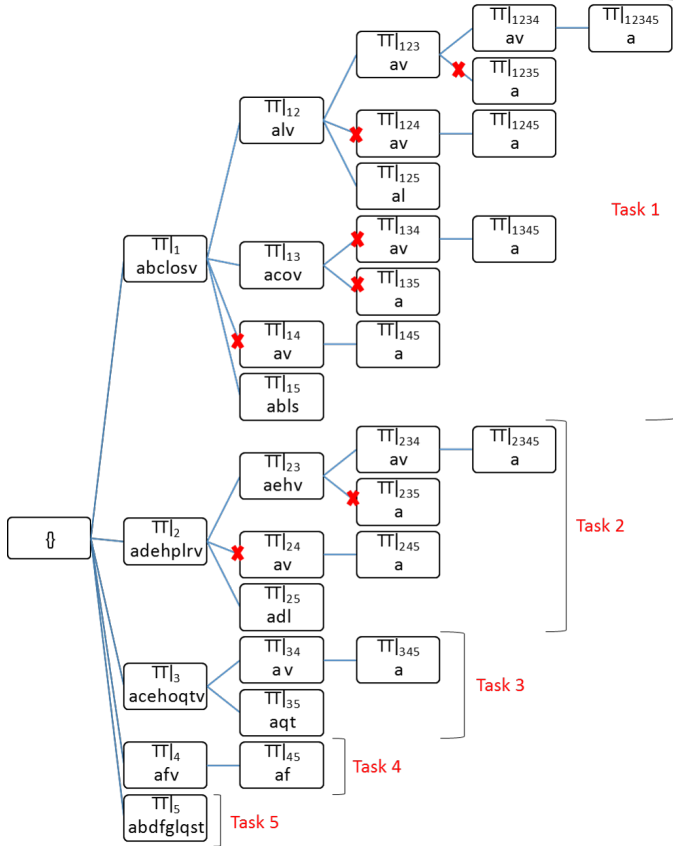


Figure 3: Running toy example: each node expands a branch of the tree independently. Pruning rule 1 and 2 are not applied. The pruning rule 3 is applied only within the same task: the red crosses on the edges represent pruned nodes due to local pruning rule 3, e.g. the one on node $\{2\ 4\}$ represents the pruning of node $\{2\ 4\}$.

4. The PaMPa-HD algorithm

Given the complete enumeration tree (see Figure 2), the centralized Carpenter algorithm extracts the whole set of closed itemsets by performing a

depth first search (DFS) of the tree. Carpenter also prunes part of the search space by applying the three pruning rules illustrated above. The PaMPa-HD algorithm proposed in this paper splits the depth first search process in a set of (partially) independent sub-processes, that autonomously evaluate sub-trees of the search space. Specifically, the whole problem can be split by assigning each subtree rooted in $TT|_X$, where X is a single transaction id in the initial dataset, to an independent sub-process. Each sub-process applies the centralized version of Carpenter on its conditional transposed table $TT|_X$ and extracts a subset of the final closed itemsets. The subsets of closed itemsets mined by each sub-process are merged to compute the whole closed itemset result. Since the sub-processes are independent, they can be executed in parallel by means of a distributed computing platform, e.g., Hadoop. Figure 3 shows the application of the proposed approach on the running example. Specifically, five independent sub-processes are executed in the case of the running example, one for each row (transaction) of the original dataset.

Partitioning the enumeration tree in sub-trees allows processing bigger enumeration trees with respect to the centralized version. However, this approach does not allow fully exploiting pruning rule 3 because each sub-process works independently and is not aware of the partial results (i.e., closed itemsets) already extracted by the other sub-processes. Hence, each sub-process can only prune part of its own search space by exploiting its “local” closed itemset list, while it cannot exploit the closed itemsets already mined by the other sub-processes. For instance, Task T2 in Figure 3 extracts the closed itemset av associated with node $TT|_{2,3,4}$. However, the same closed

itemset is also mined by T1 while evaluating node $TT|_{1,2,3}$. In the centralized version of Carpenter, the duplicate version of av associated with node $TT|_{1,2,4}$ is not generated because $TT|_{1,2,4}$ follows $TT|_{1,2,3}$ in the depth first search, i.e., the tasks are serialized and not parallel. **Fin qui vecchio con running example nuovo, rifrasare tutto?** Since pruning rule 3 has a high impact of the reduction of the search space, as detailed in Section 6, its inapplicability leads to a negative impact on the execution time of the distributed algorithm as described so far. To address this issue, we share partial results among the sub-processes. Each independent sub-process analyzes only a part of the search subspace, then, when memory is full or when a maximum number of visited node is reached, the partial results are synchronized through a synchronization phase. Of course, the exploration of the tree finishes also when the subspace has been completely explored. Specifically, the sync phase filters the partial results (i.e. nodes of the tree still to be analyzed and found closed itemsets) globally applying pruning rule 3. The pruning strategy consists of two phases. In the first one, all the transposed tables and the already found closed itemsets are analyzed. The transposed tables and the closed itemsets related to the same itemset are grouped together in a bucket. For instance, in our running example, each element of the bucket B_{av} can be:

- a frequent closed itemset av extracted during the subtree exploration of the node $TT_{3,4}$,
- a transposed table associated to the itemset av among the ones that still have to be expanded (nodes $TT_{1,2,3}$ and $TT_{2,3,4}$).

We remind the readers that, because of the independent nature of the Car-

penter subprocesses, the elements related to the same itemset can be numerous, because obtained in different subprocesses. Please note that all the extracted closed itemsets come together with the tidlist of the node in which they have been extracted.

In the second phase, in order to respect the depth-first pruning strategy of the rule 3, for each bucket it is kept only the oldest element (transposed table or closed itemset) based on a depth-first order. The depth-first sorting of the elements can be easily obtained comparing the tidlists of the elements of the bucket. Therefore, in our running example, as shown in Figure 5, from the bucket B_{av} , it is kept the node $TT_{1,2,3}$. **ho inserito questa cosa dei buckets sperando di migliorare la comprensibilita'.**

Afterwards, a new set of sub-processes is defined from the filtered results, starting a new iteration of the algorithm. In the new iteration, the Carpenter tasks ignore the frequent closed itemsets obtained in the previous iteration, which are just . They just process the remaining transposed tables, that are expanded, as before, until the maximum number of processed tables is reached or there is not enough memory to continue. The overall process is applied iteratively by instantiating new sub-processes and synchronizing their results, until there are no nodes left.

DA QUI IN POI VECCHIO con running example aggiornato, rfrasare? The application of this approach to our running example is represented in Figure 4. The table related to the itemset av associated with the tidlist/node $\{2, 3, 4\}$ is pruned because the synchronization job discovers a previous table with the same itemset, i.e. the node associated with the transaction ids combination $\{1, 2, 3\}$. The use of this approach allows the

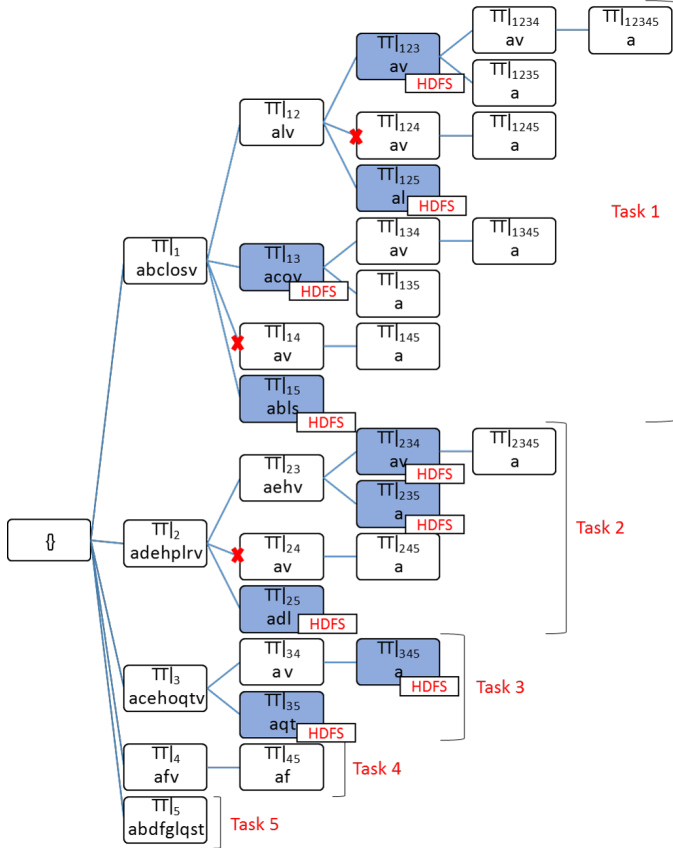


Figure 4: Execution of PaMPa-HD on the running example dataset. For sake of clarity, pruning rules 1 and 2 are not applied. The dark nodes represent the node that have been written to hdfs in order to apply the synchronization job

parallel execution of the mining process, providing at the same time a very high reliability dealing with heavy enumeration trees, which can be split and pruned according to pruning rule 3.

4.1. Implementation details

PaMPa-HD implementation exploits Hadoop MapReduce. The algorithm consists of three MapReduce jobs as shown in Figure 4.

The first job is developed to distribute the input dataset to the independent tasks, which will run a local version of the Carpenter algorithm. Each mapper is fed with a transaction of the input dataset, which is supposed to be in a vertical representation, together with the minsup parameter. As detailed in Algorithm 4.1, each transaction is in the form $item, tidlist$. For each transaction, the mapper performs the following steps. For each tid t_i of the input tidlist, given $TL_{greater}$ the set of tids $(t_{i+1}, t_{i+2}, \dots, t_n)$ greater than the considered tid t_i .

- If $|TL_{greater}| + 1 \leq minsup$, output a key-value pair $\langle key = t_i; value = TL_{greater}, item \rangle$, then analyze t_{i+1} of the tidlist.
- Else discard the tidlist.

For instance, if the input transaction is the tidlist of item b (b, 1 2 3) and minsup is 1, the mapper will output three pairs: $\langle key=1; value=2\ 3, b \rangle$, $\langle key=2; value=3, b \rangle$, $\langle key=3; value=b \rangle$.

After the mapper phase, the MapReduce shuffle and sort phase aggregates the $\langle key, value \rangle$ pairs and delivers to reducers the nodes of the first level of the tree, which represent the transposed tables projected on a single tid. The tables in Figure 6 resume the processing of a row of the initial Transposed representation of D . Reducers run a local Carpenter implementation from the input tables. Given that each key matches a single transposed table TT_X , each reducer builds the transposed tables with the tidlists contained in

the “value” fields. **TO DO INSERIRE ESEMPIO? o forse troppo? effettivamente forse gi abbastanza chiaro!**

From this table, a local Carpenter job is run. As already described in Section 3, Carpenter recursively processes a transposed table expanding it in a depth-first manner. At each iteration of the Carpenter subroutine, a counter is increased. When the count is over the given maximum expansion threshold or the memory of the cluster node is almost finished, the main routine is not invoked anymore. In this case, all the intermediate results are written to HDFS.

1. for each key value the transposed table is composed using the tidlists and a local Carpenter job is run
2. each recursion of the Carpenter subroutine increases a counter which is compared to the expansion threshold before each recursion
3. if the count is below the threshold and there is enough memory in the node to continue expanding the branch, another Carpenter recursion is scheduled
4. else, Carpenter main routine is not invoked anymore but all the intermediate results are written to HDFS

During the local Carpenter process, the found closed itemsets and the explored branches are stored in memory in order to apply a local pruning. Furthermore, closed itemsets are emitted as output as soon as they are mined, together with the tidlist of the node of the tree in which they have been found. This information is required by the synchronization phase in order to establish which element is the eldest in a depth first exploration

After this phase, the synchronization job is launched. It is a straightfor-

ward MapReduce job in which mappers input is the output of the previous job: it is composed of the closed frequent itemsets found in the previous Carpenter tasks and intermediate Transposed Tables that still have to be expanded. The itemsets are associated to their minsup and the tidlist related to the node of the tree in which they have been found; the Transposed Tables are associated to the table content, the corresponding itemset and the table tidlist. For each itemset, the mappers output a pair of the form $\langle \text{key}=\text{itemset}; \text{value}=\text{tidlist}, \text{minsup} \rangle$; for each tables the mappers out a pair of the form $\langle \text{key}=\text{itemset}; \text{value}=\text{tidlist}, \text{table_content} \rangle$. The shuffle and sort phase delivers to the reducers the pairs aggregated by keys. The reducers, which matches the buckets introduced in Section 4, compare the entries and emit, for the same key or itemset, only the eldest version in a depth first exploration. For instance, referring to our running example in Figure 5, in the bucket of the itemset av are collected the entries related to the nodes T_{123} and T_{234} . Since the tidlist 123 is previous than 234 in a depth-first exploration order, the reducer keeps and emits only the entry related to the node T_{123} . With this design, the redundant tables are discarded with a pruning very similar to the one related to a centralized memory at the cost of a very MapReduce-like job.

Finally, the last MapReduce job can be seen as a mixture of the two previous jobs. In the Map phase, in fact, all the remaining tables are expanded by a local Carpenter routine. In the Reduce phase, instead, is applied the same kind of synchronization that is run in the synchronization job. The job has two types of input: transposed tables and frequent closed itemsets. The former are expanded until it is reached the maximum expansion threshold

or the memory is not enough. At that moment, the tables are sent to the reducers, together with the frequent closed itemsets mined during the process. The format is the same as the synchronization job. The latter (i.e. the frequent closed itemsets of the previous PaMPa-HD job) are processed in the following way. If in memory there is already an oldest depth-first entry of the same itemset, the closed itemset is discarded. If there is not, it is saved into memory and used to improve the local pruning effectiveness. This job is iterated until all the Transposed Tables have been processed.

Thanks to the introduction of a global synchronization phase (job #2 and job#3), the proposed PaMPa-HD approach is able to apply pruning rule 3 and handle high-dimensional datasets, otherwise not manageable due to memory issues.

Job 1 Pseudo code

```

1: procedure MAPPER( $minsup; item_i; tidlist\ TL$ )
2:   for  $j = 0$  to  $|TL| - 1$  do
3:     tidlist  $TL_{greater}$  : set of tids greater than the considered
       tid  $t_j$ .
4:     if  $|TL_{greater}| + 1 \geq minsup$  then
5:       output  $\langle key = t_j; value = TL_{greater}, item \rangle$ 
6:     else Break
7:     end if
8:   end for
9: end procedure
10: procedure REDUCER( $key = tid\ X, value = list\ of\ tidlists\ TL[$ 
     $]$ )
11:   Create new transposed table  $TT|_X$ 
12:   for each tidlist  $TL_i$  of  $TL[ ]$  do
13:     add  $TL_i$  to  $TT|_{row_i}$ 
14:   end for
15:   Run Carpenter( $minsup; TT|_X$ )
16:   Output( $\langle frequent\ closed\ itemsets \rangle$ )
17:   Output( $\langle Transposed\ tables \rangle$ )
18: end procedure

```

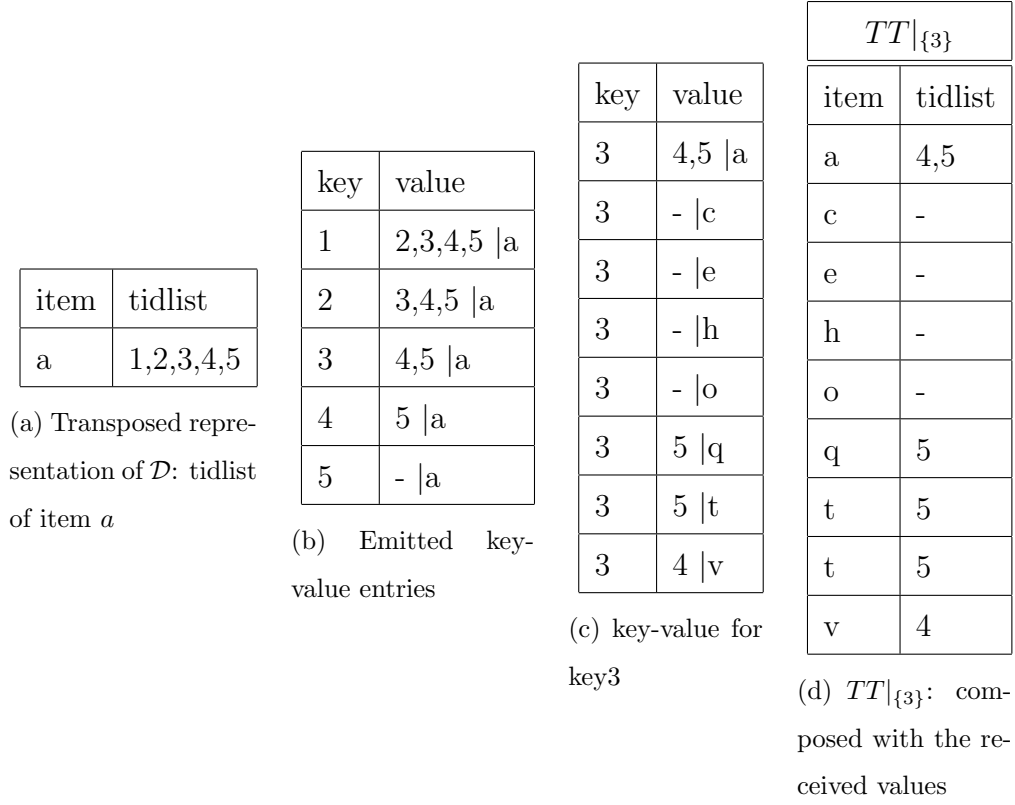


Figure 6: Job 1 applied to the Running example dataset: local Carpenter algorithm is run from the Transposed Table 6d. **troppo larga su due colonne?**

Job 2 Pseudo code

```

1: procedure MAPPER(Frequent Closed itemset / Transposed table)
2:   if Input  $I$  is a table then
3:      $itemset \leftarrow ExtractItemset(I)$ 
4:      $tidlist \leftarrow ExtractTidlist(I)$ 
5:      $Output(<itemset; tidlist + table\ I\ rows>)$ 
6:   else (i.e. input  $I$  is a frequent closed Itemset)
7:      $itemset \leftarrow ExtractItemset(I)$ 
8:      $tidlist \leftarrow ExtractTidlist(I)$ 
9:      $support \leftarrow ExtractSupport(I)$ 
10:     $Output(<itemset; tidlist + support>)$ 
11:   end if
12: end procedure

```

Job 3 Pseudo code

```

1: procedure MAPPER(Frequent Closed itemset / Transposed
   table)
2:   if Input  $I$  is a frequent closed itemset then
3:     save  $I$  to local memory
4:     Output( $\langle itemset; tidlist + support \rangle$ ) Else (i.e. input  $I$ 
   is a Transposed Table)
5:      $itemset \leftarrow ExtractItemset(I)$ 
6:     Run Carpenter( $minsup; TT|_X$ )
7:     Output( $\langle itemset; tidlist + table\ I\ rows \rangle$ )
8:     Output( $\langle itemset; tidlist + support \rangle$ )
9:   end if
10: end procedure
11: procedure REDUCER( $key = itemset; value = itemsets \& tables$ 
    $T[ ]$ )
12:    $oldest \leftarrow null$ 
13:   for each itemset or table  $T$  of  $T[ ]$  do
14:      $tidlist \leftarrow ExtractTidlist(T)$ 
15:     if  $tidlist$  previous of  $oldest$  in a DFS then
16:        $oldest \leftarrow T$ 
17:     end if
18:   end for
19:   Output( $\langle itemset + oldest \rangle$ )
20: end procedure

```

5. Expansion Threshold tuning

Before the comparison with the state of the art distributed algorithms in Section 6, we performed a set of experiments to measure the performance impact of the maximum expansion threshold, evaluating the quality of a set of proposed strategies. This set of experiments, together with the ones in Section 6, are performed on two real life datasets. The first real dataset is the **Kent Ridge Breast Cancer** [?], which contains gene expression data. It is characterized by 97 rows that represent patient samples, and 24,482 attributes related to genes. The attributes are numeric (integers, floating point). Data have been discretized with an equal depth partitioning using 20 buckets (similarly to [3]). The second real dataset is the **PEMS-SF** dataset [?], which describes the occupancy rate of different car lanes of San Francisco bay area freeways. Each transaction represents the daily traffic rate of 963 lanes, sampled every 10 minutes. It is characterized of 440 rows (in some experiments we have utilized only the first 100) and 138672 attributes ($6 \times 24 \times 963$), and it has been discretized in equi-width bins of 0.002. Because of their distribution and their discretization process, the Breast Cancer dataset is more sparse (low correlation among the dataset transactions) than the PEMS-SF dataset. The other two datasets were synthetically generated and tuned to simulate use cases characterized by extremely high-dimensional data, i.e., with massive numbers of features. Both datasets consists of 30 transactions. Dataset #1 has 1,000,000 different items and an average transaction length of 500,000 items, while Dataset #2 is 10 times larger, with 10,000,000 different items and an average transaction length of 5,000,000 items (see Table 1). The discretized version of the

real dataset and the synthetic dataset generator are publicly available at <http://dbdmg.polito.it/PaMPa-HD/>. **TO DO: aggiungere versione finale discretized pemsf**

Table 1: Datasets

Dataset	Number of transactions	Number of different items	Average number of items per transaction
Kent Ridge Breast CancerDataset	97	489,640	24,492
PEMS-SF Dataset	440 (100 rows version)	5,454,414 (3,946,646)	138,672
Synthetic Dataset #1	30	1,000,000	500,000
Synthetic Dataset #2	30	10,000,000	5,000,000

PaMPa-HD is implemented in Java 1.7.0_60 using the Hadoop MR API. Experiments were performed on a cluster of 5 nodes running Cloudera Distribution of Apache Hadoop (CDH5.3.1). Each cluster node is a 2.67 GHz six-core Intel(R) Xeon(R) X5650 machine with 32 Gbyte of main memory running Ubuntu 12.04 server with the 3.5.0-23-generic kernel.

5.1. Impact of the maximum expansion threshold

As already mentioned in Section 4, the maximum expansion threshold (max_exp) parameter indicates the maximum number of nodes to be explored before a preemptive stop of each distributed sub-process is forced. This pa-

parameter strongly affects the enumeration tree exploration, forcing each parallel task to stop before completing the visit of its sub-tree and write partial results on HDFS. This approach allows the synchronization job to globally apply pruning rule 3 and reduce the search space. Low values of *max_exp* threshold decrease the risks of memory issue, because the global problem is split into simpler and less memory-demanding sub-problems, and facilitate the application of pruning rule 3, hence a smaller subspace is searched. However, higher values allow a more efficient execution, by limiting the start and stop of distributed tasks (similarly to the context switch penalty), and the synchronization overheads.

In order to assess the impact of the expansion threshold parameter, we performed two set of experiments. In the first one we have performed the mining on the PEMS-SF (100 transactions) dataset with a Minsup 50, by varying *max_exp* from 10 to 100,000,000. The minsup value was empirically selected in order to let the mining problem being deep enough to show up different performances. In Figure 8 are shown the results in terms of execution time and number of iterations (i.e., the number of jobs).

It is clear how the *max_exp* parameter can influence the performance, with wall-clock times that can be doubled with different configurations. The best performance in terms of execution time is achieved with a maximum expansion threshold equal to 10,000 nodes. With lower values, the execution times are slightly longer, while there is an evident performance degradation with higher *max_exp* values. This result highlights the importance of the synchronization phase. Increasing the *max_exp* parameter makes the number of iterations decreasing, but more useless tree branches are explored, because

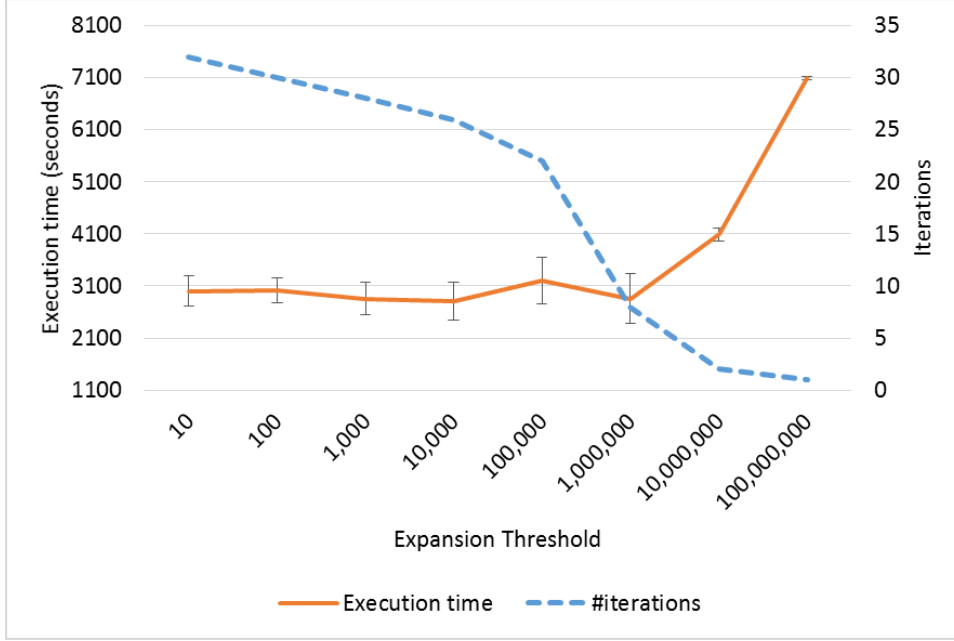


Figure 7: Execution time and number of iterations for different max_exp values on PEMS-SF dataset with $minsup=50$.

pruning rule 3 is globally applied less frequently. Lower values of max_exp , instead, raising the number of iterations, introduce a slight performance degradation caused by iterations overheads.

The same experiment is repeated with the Breast Cancer dataset and a $minsup$ value of 6. As shown in Figure 7, this time the best performances are achieved with max_exp equal to 100,000. In this case, indeed, cannot be noted big differences with the higher values. Instead, the performances strongly decreases with lower values of max_exp parameter.

+ menzione a load balancing.

The max_exp choice has a non-negligible impact on the performances of the algorithm. As demonstrated by the curves in Figures 8 and 7, it is

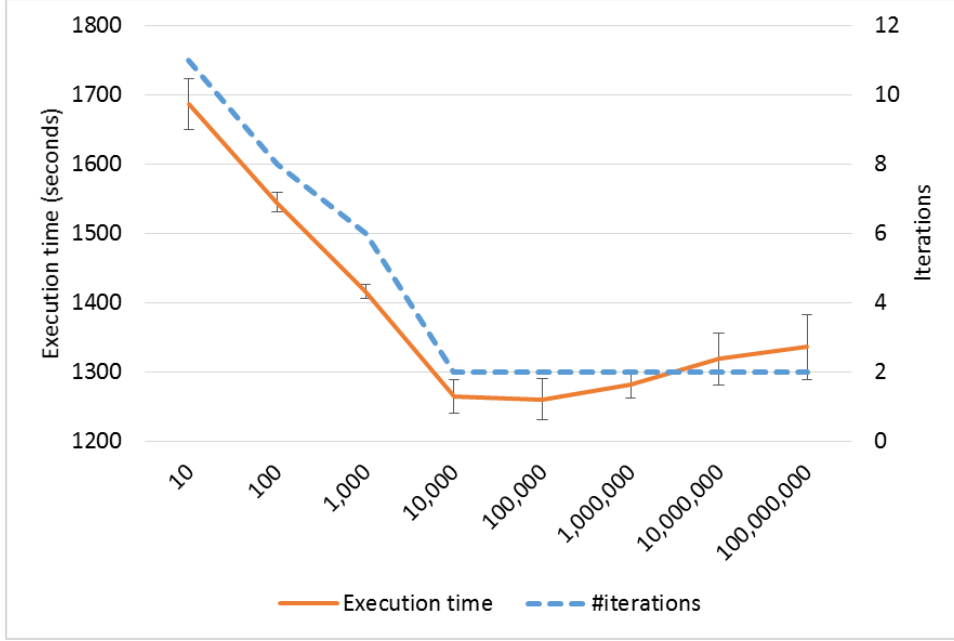


Figure 8: Execution time and number of iterations for different max_exp values on Breast Cancer dataset with $minsup=6$.

very dependent on the use case and distribution of the data. In order to simplify the parameter configuration for the user and, above all, increase the performance of the algorithm, we have implemented some tuning strategies related to max_exp , presented in the next subsection.

5.2. Proposed strategies

This section introduces some heuristic strategies related to the max_exp parameter. The aim of this experiment is to identify an heuristic technique which is able to deliver good performances without the need by the user to tune up the the max_exp parameter. Before the introduction of the techniques, let us motivate the reasons behind their design. Because of the

enumeration tree architecture, the first tables of the tree are the most populated. Each node, in fact, is generated from its parent node as a projection of the parent transposed table on a tid. In addition, the first nodes are, in the average, the ones generating more sub-branches. By construction, their transposed table tidlists are, by definition, longer than the ones of their children nodes. This increases the probability that the table could be projected on a tid. For these reasons, the tables of the initial mining phase are the most heavy to be processed. On the other hand, the number of nodes to process by each local Carpenter iteration tends to increase with the number of iterations. Still, this factor is mitigated by (i) the decreasing size of the tables and (ii) the eventual end of some branches expansion (i.e. when there are not more tids in the node transposed table). These reasons, motivated us to introduce some strategies that assumes a maximum expansion threshold that is increased with the number of iterations. These strategy start with very low values in very initial iterations (i.e. when the nodes are more heavy to be processed) and increase *max_exp* during the mining phases.

The strategy #1 is straightforward: the *max_exp* is increased with a factor of X at each iteration. For instance, if the *max_exp* is set to 10, and X is set to 100 at the second iteration it is raised to 1000 and so on.

Alternatively, we wanted to achieve a balanced growth of the *max_exp* parameter. We want to balance the load of the iterations, but, on the other hand, we want to avoid an overgrowth of the parameter and, therefore, of the execution times of the last iterations. As already said, in fact, increasing the *max_exp* parameter decreases the impact of the synchronization job pruning. In order to monitor the growth, we introduced a set of techniques

based on the execution times and a strategy which monitors the impact of the synchronization jobs. In strategies #2 and #3, we took into exam the execution time of the iterations. Specifically, strategy #2 consists in increasing, at each iteration, the *max_exp* parameter with a factor of $X^{T_{old}/T_{new}}$, given T_{new} and T_{old} the execution time of the previous two jobs. The motivation is to balance the growth of the parameter in order to achieve a stable execution times among the iterations. Strategy #3, instead, is inspired by the congestion control of TCP/IP (a data transmission protocol used by many Internet applications []). Precisely, the *max_exp* is handled like the congestion window size (i.e. the number of packets that are sent without congestion issues). This strategy, called “Slow Start”, assumes two types of growing of the window size: an exponential one and a linear one. In the first phase, the window size is increased exponentially until it reaches a threshold (“ssthresh”, which is calculated empirically from RTT and other values). From that moment, the growth of the window becomes linear, until a data loss occurs. In our case, we are not interested we just inherit the two growth factors. Therefore, our “slow start” strategy consists in increasing the *max_exp* of a factor of X until the last iteration reaches an execution time greater than a given threshold. After that, the growth is more stable, increasing the parameter of a factor of 10 (for this reason X_{10}). We have fixed the threshold to the execution time of the first two jobs (Job 1 and Job 2). These jobs, for the architecture of our algorithm, consists of the very first Carpenter iteration. They are quite different than the others since the first Mapper phase has to build the initial projected transposed tables (first level of the tree) from the input file. We have selected the execution time of the first iteration since it

is consistent with our initial aim, that is to normalize the execution times of the last iterations which are often shorter than the first ones.

The last strategy, #4, is based on the effectiveness of the pruning in terms of closed or tables pruned in the synchronization job. Indeed, the measure of the relative number of tables that are pruned cannot be easily interpreted. An increasing pruning percentage means that there are a lot of tables that are generated uselessly. However, an increasing trend is also normal, since the number of nodes that are processed continues to increase. Given that our intuition is to rise the *max_exp* among the iterations, in strategy #4, we increase the *max_exp* parameter with a factor $X^{Pr_{old}/Pr_{new}}$, given Pr_{new} and Pr_{old} the relative number of pruned tables in the previous two jobs.

In our experiment, we have fixed the initial *max_exp* value to 10. This very low value is motivated by the nature of our strategies, which consist of a balanced increasing of the parameter. We applied the strategies (resumed in Table 2) to the same experiments of Figures 7 and 8, in order to compare execution times the ones obtained with the optimum choice of *max_exp*. In order to assess the impact of the X parameter, we have used values from 10 to 10,000 (except for strategy3 for which we have used values from 100 to 10,000).

The result of the application of the techniques to the PEMS-SF dataset are shown in Figures 9, 10, 11, 12, which represent the relative execution time gain with respect to the best execution time obtained with the fixed *max_exp* of 10,000. In Figure 13 we have grouped the best configuration for each strategy in order to easily compare them. It is clear how the almost all the strategies improve the performance of the algorithm. Only Strategy3 showed to

Table 2: Strategies

Strategy #1(X)	Increasing at each iteration with a factor of X
Strategy #2(X)	Increasing at each iteration with a factor of $X^{T_{old}/T_{new}}$
Strategy #3	Slow start, with a fast increase factor of X
Strategy #4(X)	Increasing at each iteration with a factor of $X^{Pr_{old}/Pr_{new}}$

be slower. Strategy1(1,000), Strategy2(10,000) and Strategy4(1,000) achieve a very similar speedup.

We have repeated the experiment with Breast Cancer dataset and a min-sup value of 6 (Figures 14, 15, 16, 17). We have raised X from 10 to 10,000. Only with Strategy2, as shown in Figure 15, we have raised it to 100,000 , because the experiments could suggest a decreasing execution times trend, but it was not the case. As before, the results are grouped in Figure 18. In this case, all the strategies achieve a positive speedup with respect to the best execution time obtained with a fixed *max_exp* parameter. In addition, the same configurations for each strategy demonstrated to be most performant in both the experiments. The results obtained with the adoption of these strategies are very similar, and the differences are negligible.

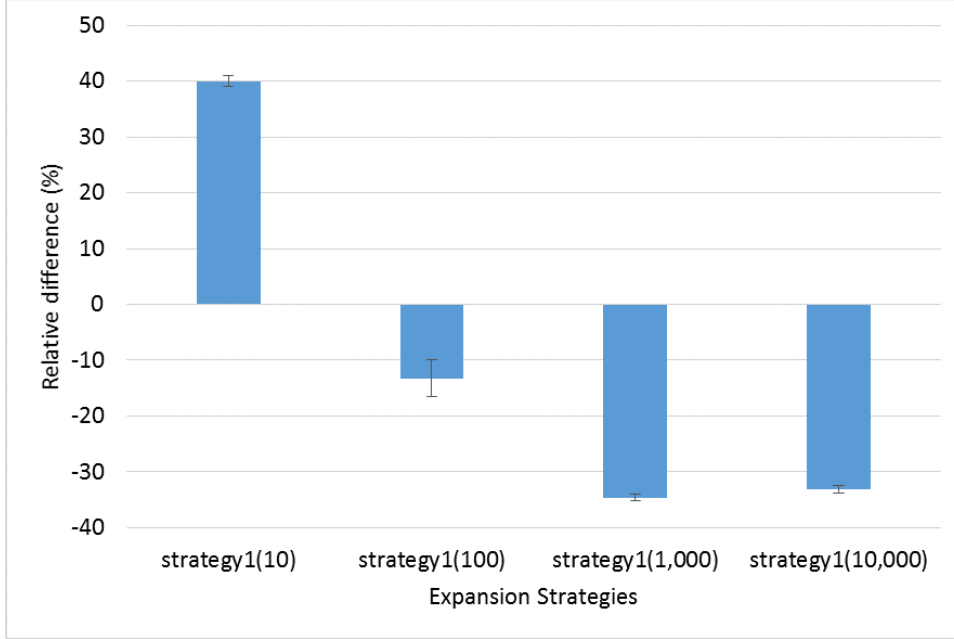


Figure 9: Relative gains on Pems-SF dataset with $minsup=50$, Strategy1 and different X values.

6. Experiments

We performed a set of experiment to evaluate the performance of the proposed algorithm. We measured (i) the efficiency of the proposed algorithm, comparing it with the state of the art distributed approaches (Section 6.1); (ii) performance impact with respect to the number of parallel tasks (see Section 6.2); and, finally (iv) the communication costs and load balancing behavior of the algorithm, very important in such a distributed context.

6.1. Running time

After the identification of a good tradeof strategy in the previous section, we have used it to analyze the efficiency of PaMPa-HD by comparison

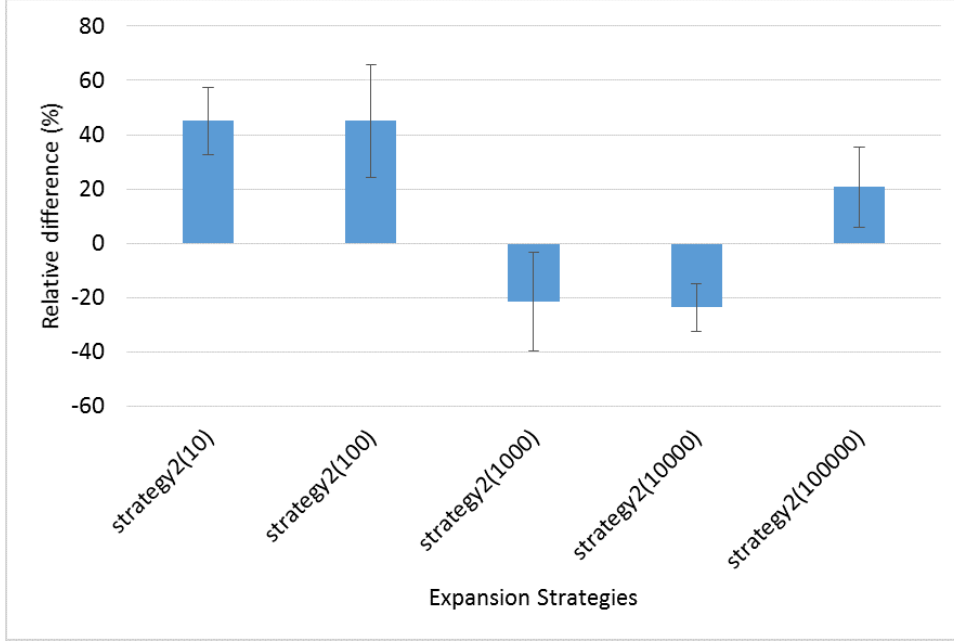


Figure 10: Relative gains on Pems-SF dataset with $minsup=50$, Strategy2 and different X values.

with two state-of-the-art algorithms: a main-memory centralized version of Carpenter, developed in C++ by Sam Borgelt [?], and a MapReduce-based implementation of FP-growth [?] available in Mahout 0.9 [?], (for details see Section 7). The comparison with Carpenter aims to show the ability of PaMPa-HD to address datasets that are not manageable by means of a centralized approach. The comparison with the Hadoop implementation of PFP is a reference benchmark since it represents the most scalable state-of-the-art approach for closed itemset mining, even if PFP was not specifically developed to deal with high-dimensional datasets as PaMPa-HD.

The first set of experiments has been performed with the Kent Ridge Breast Cancer dataset [?]. As reported in Figure ?? (Minsup axis is reversed

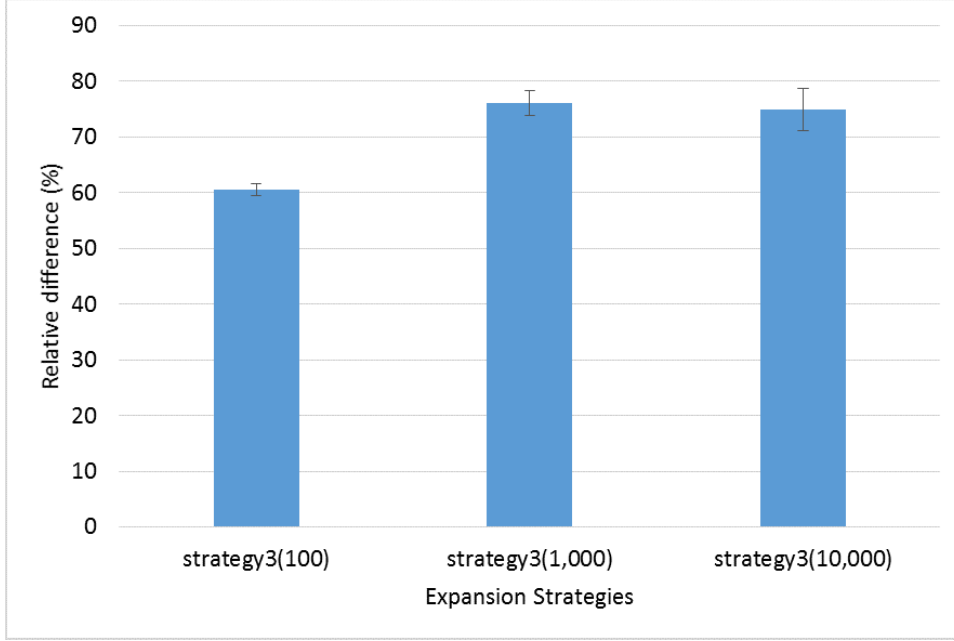


Figure 11: Relative gains on Pems-SF dataset with $minsup=50$, Strategy3 and different X values.

to improve readability), the centralized version of Carpenter is faster than our approach with absolute minimum support thresholds (Minsup) of 6 or higher. However, when the search space to explore becomes larger (i.e., with low support minimum support thresholds) the centralized version is slower or runs out of memory (after days of processing), while PaMPa-HD is able to extract all the frequent closed itemsets. When the problem becomes more complex, and the size of the enumeration three increases, the advantages of the parallelization outperform the overheads of the distributed environment (e.g., communication and synchronization costs). Such overheads can be roughly estimated in the fixed 100 seconds time required by PaMPa-HD to be executed for high Minsup values (8 and 9 in Figure ??). The Hadoop

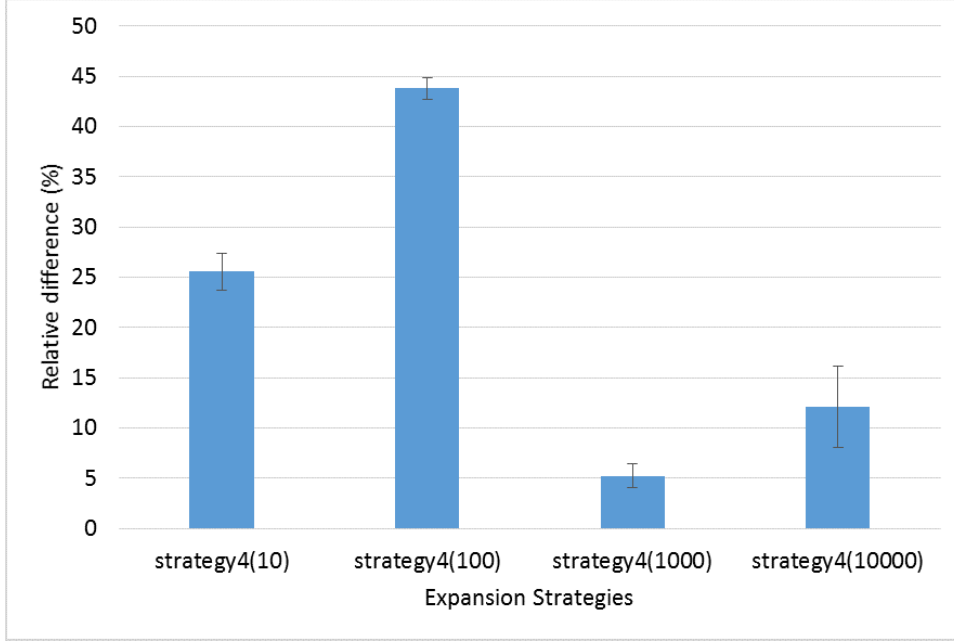


Figure 12: Relative gains on Pems-SF dataset with $minsup=50$, Strategy4 and different X values.

PFP implementation reveals to be unsuitable for the extraction of frequent closed itmsets from the Kent Ridge Breast Cancer dataset: its execution time was more than one day of computation, even for the highest support threshold, so orders of magnitude higher than the others. For this reason, the execution time of PFP is not reported in Figure ???. This result highlights the need for specific algorithms particularly tailored to address high-dimensional data, such as PaMPa-HD, because traditional best-in-class approaches, such as PFP, cannot cope with the curse of dimensionality.

A second set of experiments using the synthetic Dataset #1 has been performed. In Dataset #1, even if the number of different items is only doubled with respect to the breast cancer dataset, the average number of items per

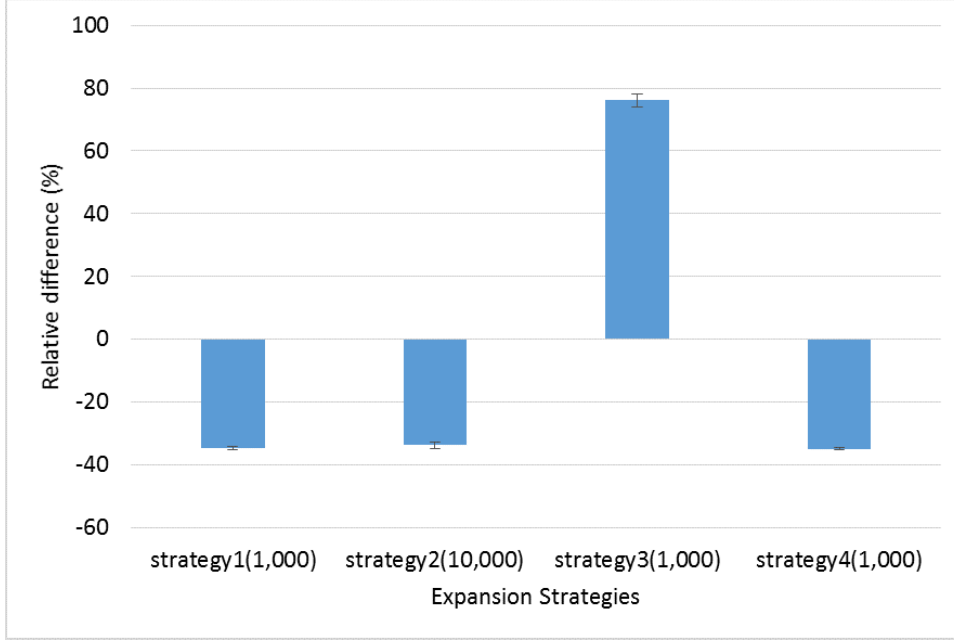


Figure 13: Relative gains of the best configuration for each strategy, on Pems-SF dataset with $minsup=50$.

transaction is much higher (one order of magnitude). This feature makes each conditional transposed table, associated with each node of the enumeration tree, much bigger and heavy to process. As reported in Figure ??, for very high Minsup values (20+) Carpenter scores the lowest execution times, with a Minsup value of 19 the performance of PaMPa-HD and Carpenter are almost the same, but with Minsup values of 18 or lower, Carpenter runs out of memory and the task is unfeasible, whereas PaMPa-HD successfully completes the extraction. Furthermore, the time required for PaMPa-HD to complete the extraction in the latter case and its trend for decreasing Minsup values (19-18-17) is equal or less than the time required by PFP for respectively higher Minsup values (22-21-20). Not only PaMPa-HD is faster

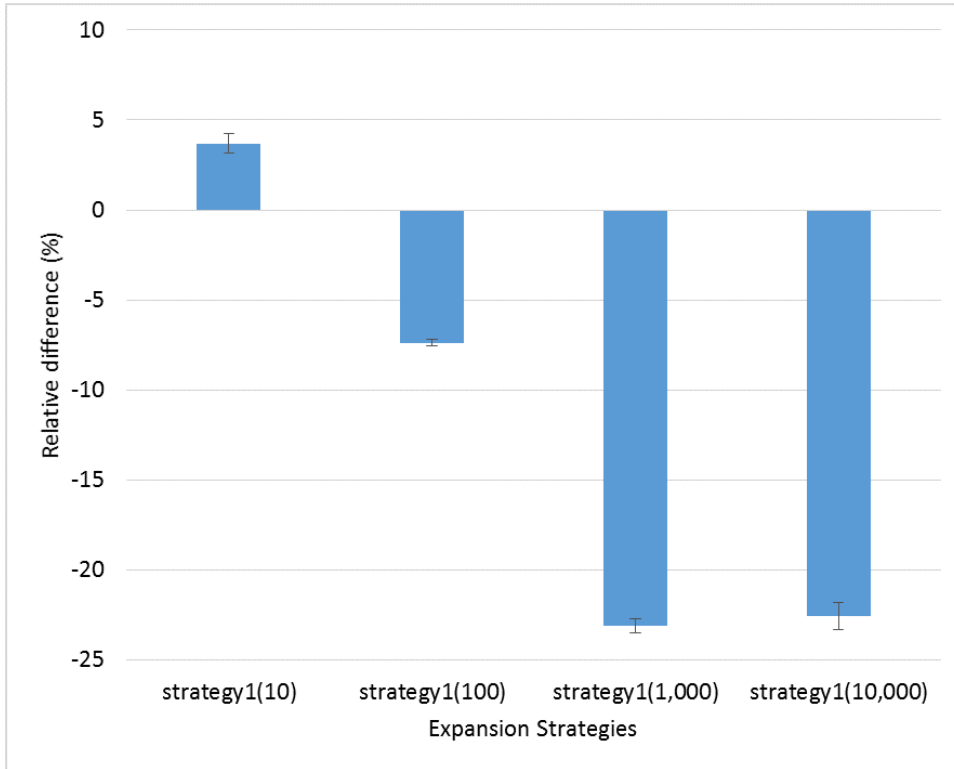


Figure 14: Relative gains on Breast dataset with $minsup=6$, Strategy1 and different X values.

than PFP for high Minsup values (20+), but can it also scale further, since PFP stops at Minsup 20. The results prove that PaMPa-HD is effectively able to address problems that are unfeasible for both centralized state-of-the-art specialized techniques (Carpenter) and scalable state-of-the-art generic itemset mining algorithms (PFP).

The last set of performance experiments is executed on synthetic Dataset #2. Dataset #2 is extremely high-dimensional, achieving an average number of items per transaction of 5,000,000. The results, reported in Figure ??, confirm that (i) for high Minsup values (21+), the problem is conveniently solved

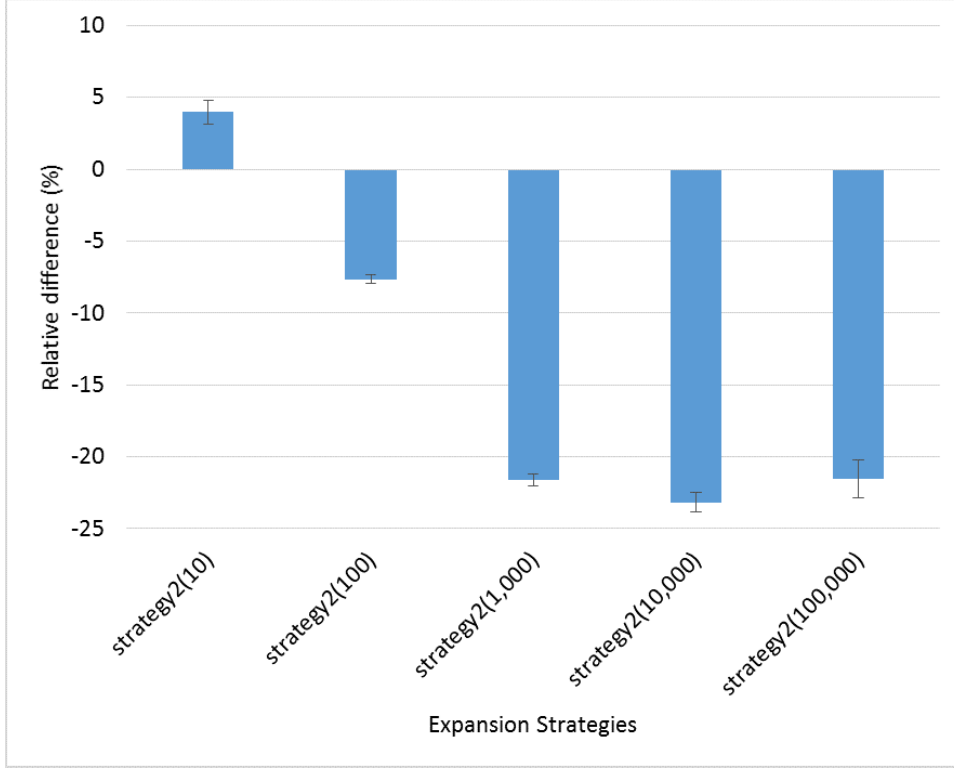


Figure 15: Relative gains on Breast Cancer dataset with $minsup=6$, Strategy2 and different X values.

by current specialized algorithms (Carpenter), (ii) both centralized algorithms and generic distributed approaches cannot cope with high-dimensional datasets with lower Minsup values, being 20 for Carpenter and 25 for PFP the lowest feasible Minsup values, and (iii) the proposed PaMPa-HD algorithm provides both better execution times than PFP by almost an order of magnitude and higher scalability for more complex problems (lower Minsup values). Achieving lower Minsup values allows to discover deeply hidden knowledge, while high Minsup values usually lead to obvious information.

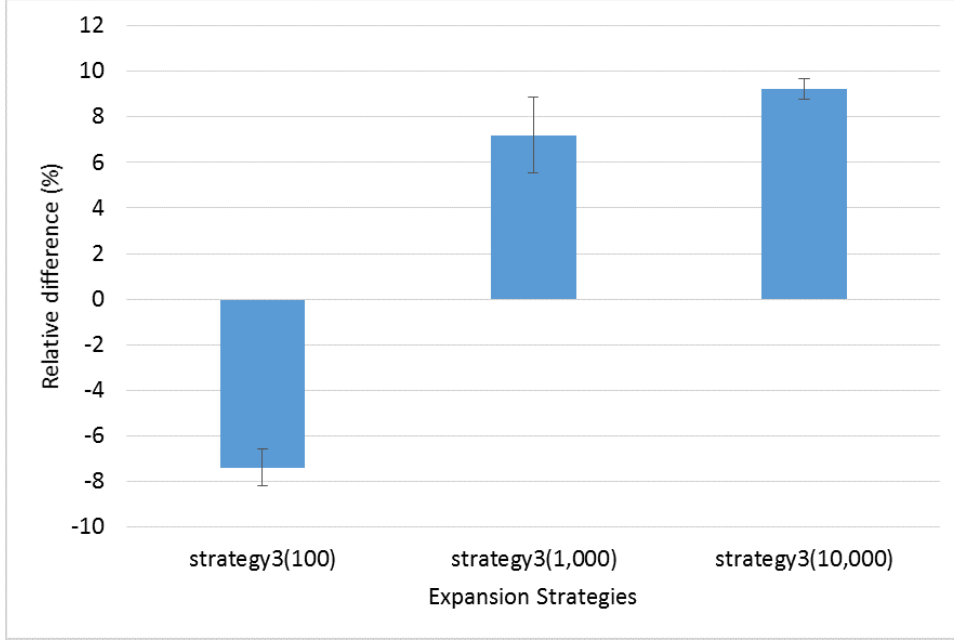


Figure 16: Relative gains on Breast Cancer dataset with $minsup=6$, Strategy3 and different X values.

6.2. Scalability

This section addresses the scalability of PaMPa-HD with respect to the number of reducers, since the most heavy operations (i.e., the execution of the “local” Carpenter) are executed by the reducers. The number of reducers varies from 1 to 18, which is the maximum number of tasks that can be run simultaneously in the commodity cluster at our disposal. The Breast cancer dataset and a minimum absolute support threshold equal to 6 have been used. As shown in figure 6.2, the increase of the number of reducers has a positive impact on the execution time when the number of reducers is less than 10. The marginal benefits decrease as more reducers are added to the computation. Therefore, with more than 10 reducers, the benefits of the

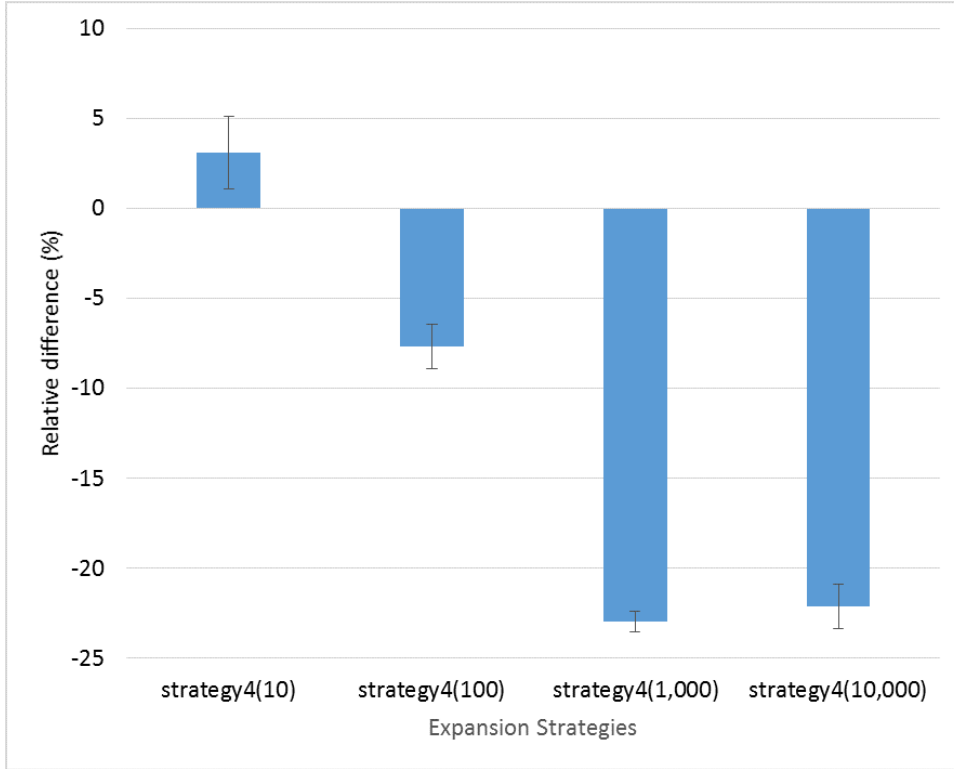


Figure 17: Relative gains on Breast Cancer dataset with $minsup=6$, Strategy4 and different X values.

load distribution are compensated by the decreased effectiveness of the local pruning (i.e., the more the load is distributed, the less effective the local pruning is).

7. Related work

Questo tutto vecchio, rifrasare? Frequent itemset mining represents a very popular data mining technique used for exploratory analysis. Its popularity is witnessed by the high number of approaches and implementations. The most popular techniques to extract frequent itemsets from a transac-

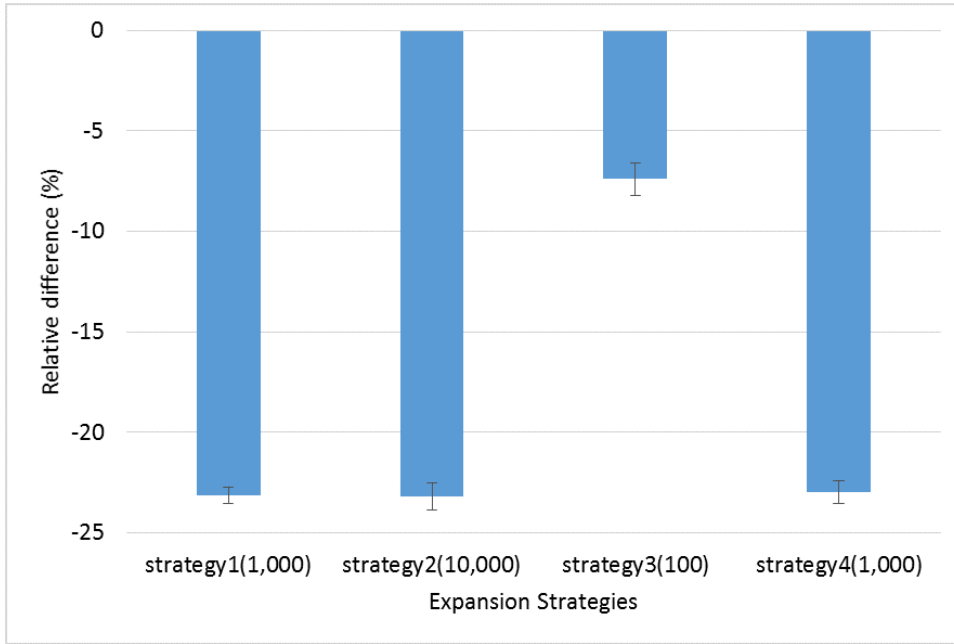


Figure 18: Relative gains of the best configuration for each strategy, on Pems Cancer dataset with $minsup=6$.

tional datasets are Apriori and Fp-growth. Apriori [?] is a bottom up approach: itemsets are extended one item at a time and their frequency is tested against the dataset. FP-growth [?], instead, is based on an FP-tree transposition of the transactional dataset and a recursive divide-and-conquer approach. These techniques explore the search space enumerating the items. For this reason, they work very well for datasets with a small (average) number of items per row, but their running time increases exponentially with higher (average) row lengths [? ?].

In recent years, the explosion of the so called Big Data phenomenon has pushed the implementation of these techniques in distributed environments such as Apache Hadoop [1], based on the MapReduce paradigm [?], and

Apache Spark [2]. Parallel FP-growth [?] is the most popular distributed closed frequent itemset mining algorithm. The main idea is to process more sub-FP-trees in parallel. A dataset conversion is required to make all the FP-trees independent. BigFIM and DistEclat [?] are two recent methods to extract frequent itemsets. DistEclat represents a distributed implementation of the Eclat algorithm [?] an approach based on equivalence classes (groups of itemsets sharing the same prefixes), smartly merged to obtain all the candidates. BigFIM is a hybrid approach exploiting both Apriori and Eclat paradigms. Carpenter [3], which inspired this work, has been specifically designed to extract frequent itemsets from high-dimensional datasets, i.e., characterized by a very large number of attributes (in the order of tens of thousands or more). The basic idea is to investigate the row set space instead of the itemset space. A detailed introduction to the algorithm is presented in section 3. **dovremmo aggiungere il vecchio pampa, qualche aiutino?**. This work is based on the distributed implementation PaMPa-HD [?]. The original algorithm assumes a slightly different architecture, assuming an independent synchronization job at each iteration. As already described in Section 4.1, this implementation includes the synchronization phase in the Job 3. Therefore, the number of MapReduce jobs (with their related overhead) are strongly reduced.

8. Applications

Questo tutto vecchio, rifrasare? Lo eliminiamo? Since PaMPa-HD is able to process extremely high-dimensional datasets we believe it is suitable for many application (scientific) domains. The first example is bioin-

formatics: researchers in this environment often cope with data structures defined by a large number of attributes, which matches gene expressions, and a relatively small number of transactions, which typically represent medical patients or tissue samples. Furthermore, smart cities and computer vision environments are two important application domains which can benefit from our distributed algorithm, thanks to their heterogeneous nature. Another field of application is the networking domain. Some examples of interesting high-dimensional dataset are URL reputation, advertisements, social networks and search engines. One of the most interesting applications, which we plan to investigate in the future, is related to internet traffic measurements. Currently, the market offers an interesting variety of internet packet sniffers like [?], [?]. Datasets, in which the transactions represent flows and the item are flows attributes, are already a very promising application domain for data mining techniques [?], [?], [?].

9. Conclusion

This work introduced PaMPa-HD, a novel frequent closed itemset mining algorithm able to efficiently parallelize the itemset extraction from extremely high-dimensional datasets. Experimental results shos its scalability and its performance in coping with datasets characterized by up to 10 millions different items and, above all, an average number of items per transaction up to 5 millions, on a small commodity cluster of 5 nodes. PaMPa-HD outperforms state-of-the-art algorithms, by showing a better scalability than both Carpenter and PFP, and being more efficient than PFP in the distributed itemset extraction. We plan to apply this approach in the network data anal-

ysis domain and to develop an Apache Spark implementation. We also plan to improve the expansion threshold selection: an auto-tuning mechanism based on the specific data distribution might further unleash the algorithm potential.

Acknowledgement

The research leading to these results has received funding from the European Union under the FP7 Grant Agreement n. 619633 (Project “ONTIC”).

- [1] D. Borthakur, The hadoop distributed file system: Architecture and design, Hadoop Project 11 (2007) 21.
- [2] M. Zaharia, M. Chowdhury, T. Das, A. Dave, J. Ma, M. McCauley, M. J. Franklin, S. Shenker, I. Stoica, Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing, in: NSDI’12, 2012, pp. 2–2.
- [3] F. Pan, G. Cong, A. K. H. Tung, J. Yang, M. J. Zaki, Carpenter: Finding closed patterns in long biological datasets, in: Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’03, ACM, New York, NY, USA, 2003, pp. 637–642. doi:10.1145/956750.956832.
URL <http://doi.acm.org/10.1145/956750.956832>
- [4] M. Cuturi, UCI machine learning repository. PEMS-SF data set (2011).
URL <https://archive.ics.uci.edu/ml/datasets/PEMS-SF>

- [5] J. Leskovec, A. Krevl, SNAP Datasets: Stanford large network dataset collection, <http://snap.stanford.edu/data> (Jun. 2014).
- [6] Pang-Ning T. and Steinbach M. and Kumar V., Introduction to Data Mining, Addison-Wesley, 2006.