



# Efficient mining of the most significant patterns with permutation testing

Leonardo Pellegrina<sup>1</sup> · Fabio Vandin<sup>1</sup>

Received: 24 July 2019 / Accepted: 8 May 2020 / Published online: 9 June 2020

© The Author(s), under exclusive licence to Springer Science+Business Media LLC, part of Springer Nature 2020

## Abstract

The extraction of patterns displaying significant association with a class label is a key data mining task with wide application in many domains. We introduce and study a variant of the problem that requires to mine the top- $k$  statistically significant patterns, thus providing tight control on the number of patterns reported in output. We develop TOPKWY, the first algorithm to mine the top- $k$  significant patterns while rigorously controlling the family-wise error rate of the output, and provide theoretical evidence of its effectiveness. TOPKWY crucially relies on a novel strategy to explore statistically significant patterns and on several key implementation choices, which may be of independent interest. Our extensive experimental evaluation shows that TOPKWY enables the extraction of the most significant patterns from large datasets which could not be analyzed by the state-of-the-art. In addition, TOPKWY improves over the state-of-the-art even for the extraction of *all* significant patterns.

**Keywords** Statistical pattern mining · Hypothesis testing · Top- $k$  patterns

## 1 Introduction

*Frequent patterns mining* is one of the fundamental primitives in data mining, with applications in a large number of domains, ranging from market basket analysis to

---

Responsible editor: M. J. Zaki.

---

A preliminary version of this work appeared in the proceedings of ACM KDD'18 as (Pellegrina and Vandin 2018).

---

✉ Fabio Vandin  
fabio.vandin@unipd.it  
Leonardo Pellegrina  
pellegr@dei.unipd.it

<sup>1</sup> Department of Information Engineering, Università di Padova, Via G. Gradenigo 6/B, 35131 Padua, IT, Italy

biology and medicine (Han et al. 2007). In its original definition (Agrawal et al. 1993) it requires to identify patterns that appear in a fraction at least  $\sigma$  of all the transactions of a transactional dataset. *Significant pattern mining* (Dong and Bailey 2012; Hämäläinen and Webb 2019; Pellegrina et al. 2019a) is an extension of the problem in which each transaction is assigned a binary class label and the goal is to identify patterns having *significant association* with one of the class labels. Significance is commonly assessed using a statistical test [e.g., Fisher exact test (Fisher 1922)], that provides a  $p$  value quantifying the probability that the association observed in real data arises due to chance alone.

Significant pattern mining is crucial in many applications, providing additional information w.r.t. mining patterns that are frequent in the entire dataset: in market basket analysis it serves to identify itemsets that are purchased more frequently by one group of customers than by another one (e.g., married people vs. singles); in social networks, it finds features characterizing users interested in one specific topic; in biology, it identifies sets of genetic variants appearing more frequently in cancer versus normal tissues or in one cancer type versus another one. In all applications, identifying highly reliable associations is of the utmost importance.

One of the critical issues in significant pattern mining is the *multiple hypothesis problem*, due to the huge number of patterns appearing in large datasets. When testing only one pattern, if its  $p$  value is below a fixed threshold  $\alpha$ , one can flag the pattern as significant with the guarantee that the probability of *false discovery* (e.g., flagging the pattern as significant when it is not) is bounded by  $\alpha$ . However, for a fixed significance threshold  $\alpha$ , when  $d$  patterns are tested we expect  $\alpha d$  of them to have  $p$  value below  $\alpha$  even when they are not associated with the class labels. A standard method to correct for multiple hypothesis testing, called *Bonferroni method* (Bonferroni 1936), is to adjust the significance threshold by dividing  $\alpha$  by the number  $d$  of tested patterns. This guarantees that the probability of making one or more false discoveries, called *family-wise error rate* (FWER), is bounded by  $\alpha$ . However, since the number  $d$  of patterns can be huge, this approach results in limited *statistical power*, with very few patterns having  $p$  value passing such small threshold (Webb 2006, 2007, 2008). A naïve solution for the problem is to limit the number of elements in the patterns to be tested, hindering the ability to identify large significant patterns.

A breakthrough in significant pattern mining is the work by Terada et al. (2013a), that proposes LAMP, the first method to identify significant patterns without limiting their size. LAMP is based on the work by Tarone (1990), which shows that patterns that cannot reach statistical significance, called *untestable*, do not need to be taken into account while correcting for multiple hypothesis testing. Subsequent work by Minato et al. (2014) has improved the search strategy employed by LAMP to identify testable patterns. Even so, such methods suffer from limited power, due to the use of Bonferroni correction over the large number of testable patterns.

Recently, methods based on the more powerful Westfall–Young permutation procedure (Westfall and Young 1993) have been proposed, first by Terada et al. (2013b) with FASTWY and then by Llinares-López et al. (2015) with Westfall–Young light (WYlight for short). These methods mine several permuted datasets to identify a threshold such that all patterns with  $p$  value below the threshold can be flagged as statistically significant while controlling the FWER. Such methods achieve a higher

power than methods based on Bonferroni correction, including LAMP, and the state-of-the-art, WYlight, has proved to be more efficient than FASTWY also in terms of runtime and memory.

However, the extraction of significant patterns from large datasets is still challenging, with three crucial issues that are not addressed by currently available methods. First, in several cases the dependency among patterns leads to a huge number of statistically significant patterns even after multiple hypothesis correction. A common approach [used, e.g., by Terada et al. (2013b) and Llinares-López et al. (2015)] to partially alleviate this problem is to consider only closed patterns (Han et al. 2007), discarding patterns with redundant information content in terms of appearance in the dataset and of association with the class label. Even with this restriction the number of significant patterns can be extremely large and when this happens one would like to focus on the most significant ones, without resorting to filtering strategies after the expensive extraction of all significant patterns has been performed. A second and related issue is that current methods work by first identifying the exact corrected threshold for statistical significance, and only subsequently mining the real dataset: when the number of significant patterns is huge, one would like to focus on the most significant ones without the burden of computing the exact significance threshold. Third, all methods may need to process several untestable patterns to identify the correct threshold for significance, resulting in an extremely large running time in particular for datasets with many low frequency patterns. These issues make current methods impractical in many cases, as shown by our experimental evaluation.

## 1.1 Our contribution

In this work we focus on the problem of mining the most statistically significant patterns while rigorously controlling the FWER of the returned set of patterns. In particular, in analogy with frequent pattern mining approaches, we focus on extracting the top- $k$  statistically significant patterns. This problem is more challenging than the extraction of top- $k$  frequent patterns, given that statistical significance does not enjoy the anti-monotonicity property w.r.t. to pattern frequency. In this regards, our contributions are:

- we formally define the problem of mining *Top- $k$  Statistically Significant Patterns*. Our definition allows to properly control the size of the output set while providing guarantees on the FWER of the output.
- we design a novel algorithm, called TOPKWY, for the problem above, which provides guarantees on the FWER by using the Westfall–Young permutation test procedure. TOPKWY adapts to the distribution of significant patterns: it reports *all* significant patterns when their number is small, while it outputs only the most significant patterns when the number of significant patterns is huge. TOPKWY is based on an exploration strategy similar to the one used by TOPKMINER (Pietracaprina and Vandin 2007), an efficient algorithm to identify the top- $k$  frequent patterns. We prove that the use of such strategy guarantees that, in contrast to previous approaches, TOPKWY will never explore *untestable* patterns.

- we introduce several bounds to prune *untestable* patterns that improve over the bound introduced by LAMP and used in WYlight as well. We show that such bounds can be effectively used within the exploration strategy employed by TOPKWY and that it provides a significant speed-up for real datasets.
- we present variants of TOPKWY to mine the top- $k$  significant patterns with control of the Generalized Family-Wise Error Rate ( $g$ -FWER) or the False Discovery Proportion (FDP), which trade an increase in the size of the output with a (potentially) higher, but still controlled, number of false discoveries. These variants lead to an increase in the *statistical power* in situations where the number of results reported controlling the FWER is low.
- we conduct an extensive experimental evaluation of the use of TOPKWY to extract significant itemsets and subgraphs, showing that TOPKWY allows the extraction of statistically significant patterns for large datasets while having reasonable memory requirements. Surprisingly, for many datasets TOPKWY improves over the state-of-the-art even when it is used to find *all* statistically significant patterns.

## 1.2 Additional related work

Since the introduction of the frequent pattern mining problem (Agrawal et al. 1993), a number of methods have been developed to efficiently extract all frequent patterns [see (Han et al. (2007)) for several references]. Given that the number of such patterns can be extremely large and that identifying an appropriate frequency threshold to limit the number of frequent patterns is challenging, methods to identify restricted classes of patterns, e.g., closed patterns (Pasquier et al. 1999) or maximal patterns (Bayardo 1998), have been designed. Methods that directly limit the number of patterns by reporting the  $k$  most frequent closed patterns have been designed (Han et al. 2002; Pietracaprina and Vandin 2007) as well.

Many methods have been developed for *subgroups discovery* [see, e.g., the reviews by Herrera et al. (2011) and by Atzmueller (2015)], that is the task of mining patterns associated with class labels using a *quality score* to quantify the association between a pattern and class labels. Most of the available methods focus on such scores, without assessing the statistical significance of the association or correcting for multiple hypothesis testing. Contrarily to this general trend, a recent work (Li et al. 2014) uses odds ratio as a measure to quantify statistical significance and to reduce redundancies of subgroups, but still does not correct for multiple hypothesis testing. In (Duivesteijn and Knobbe 2011) multiple hypothesis testing is accounted for using a swap randomization technique (Gionis et al. 2007). The methods we describe in this work can be applied directly to the task of discovering the top- $k$  statistically significant subgroups with proper control of the FWER when the *language* defining the subgroups is composed by conjunctive formulas, and could possibly be adapted and extended to more general languages.

In addition to the contributions for significant pattern mining mentioned above (Terada et al. 2013a; Minato et al. 2014; Terada et al. 2013b, 2015; Llinares-López et al. 2015) (described more in depth in Sect. 2), recent work has extended the extraction of statistically sound patterns (Hämäläinen and Webb 2019) in directions that are orthogo-

**Table 1**  $2 \times 2$  contingency table

Variables	$G(S, t) = 1$	$G(S, t) = 0$	Row totals
$c = 1$	$a_S$	$n_1 - a_S$	$n_1$
$c = 0$	$x_S - a_S$	$n - n_1 + a_S - x_S$	$n - n_1$
Col totals	$x_S$	$n - x_S$	$n$

nal to our contributions, for example searching for statistical dependency rules between itemsets and items (Hämäläinen 2012), or using layered critical values (Webb 2008) for correcting for multiple hypothesis testing. Komiyama et al. (2017) have developed efficient techniques for multiple testing correction in the mining of statistical emerging patterns. Terada et al. (2016) and Papaxanthos et al. (2016) have introduced methods to find significant patterns in the presence of covariates.

## 2 Background and problem definition

### 2.1 Significant pattern mining

Let the dataset  $\mathcal{D} = \{t_1, t_2, t_3, \dots, t_n\}$  be a set of  $n$  transactions, where each transaction is an element from a domain  $\mathcal{F}$ . Each transaction  $t_i$  is associated to a binary label  $c_i \in \{0, 1\}$ . We denote by  $n_1$  the number of transactions with label 1, and, without loss of generality, we assume that  $n_1$  is the *minority class*, i.e.  $n_1 \leq n/2$ . We define a *pattern*  $S$  as an element of  $\mathcal{F}$ , potentially with some constraints, and for each transaction  $t$  we define the binary variable  $G(S, t)$  such that  $G(S, t) = 1$  if  $S$  is *contained* in transaction  $t$  and  $G(S, t) = 0$  otherwise. For example: in the case of *itemsets*,  $\mathcal{F}$  is the set of (non-empty) subsets of a universe of binary features  $\mathcal{I}$ , and  $S$  is an element of  $\mathcal{F}$ ; for *subgraphs*,  $\mathcal{F}$  is the set of vertex-labelled graphs, while  $S$  is an element of  $\mathcal{F}$  with the constraint that  $S$  is connected. Given a pattern  $S$ , we define its support  $x_S$  as the number of transactions containing  $S$ , that is  $x_S = \sum_{i=1}^n G(S, t_i)$ . We denote by  $a_S$  the number of class 1 transactions containing  $S$ . In this work we assume that patterns enjoy the *anti-monotonicity property*, such that for any super pattern  $S'$  of  $S$  (i.e.,  $S \subset S'$ ) the support  $x_{S'}$  of  $S'$  is  $x_{S'} \leq x_S$ . This holds for itemsets, subgraphs, and many other kind of patterns.

The objective of *significant pattern mining* is to find patterns with *significant statistical association* to one of the two classes. In order to quantify the statistical association, a rigorous *statistical test* is performed. Such test assesses the association between the *observations*  $G(S, t_1), \dots, G(S, t_n)$  of variable  $G(S, t)$  (describing the presence of  $S$  in the transactions of dataset  $\mathcal{D}$ ) and the observed class labels  $c_i$  of the transactions, representing observations of the variable  $c \in \{0, 1\}$  defining the label of a transaction. Since the variables are binary, for a given pattern  $S$  the  $2 \times 2$  contingency table represented in Table 1 is considered and the popular Fisher exact test (Fisher 1922) is often employed. Such test considers the *marginals*  $(x_S, n_1, n)$  of the contingency table for pattern  $S$  to be fixed; under the *null hypothesis* of independence between variables  $G(S, t)$  and  $c$ , the number  $a_S$  of class 1 transactions containing  $S$  follows a

hypergeometric distribution:

$$\Pr(a_S = a | x_S, n_1, n) = \Pr(a | x_S, n_1, n) = \binom{n_1}{a} \binom{n - n_1}{x_S - a} / \binom{n}{x_S}. \quad (1)$$

Using such distribution, the probability, called  $p$  value, of observing an association that is equally or more extreme than the one observed in the data under the null hypothesis can be computed. Smaller  $p$  values indicate more probable associations. The  $p$  value  $p_S$  for the pattern  $S$  is computed by summing all the probabilities to obtain, under the null hypothesis, contingency tables which are at least as extreme as the one observed in  $\mathcal{D}$ :

$$p_S = p_S(a_S) = \sum_{k: \Pr(k | x_S, n_1, n) \leq \Pr(a_S | x_S, n_1, n)} \Pr(k | x_S, n_1, n). \quad (2)$$

## 2.2 Multiple hypothesis testing

When only one pattern  $S$  is tested, it can be flagged as *significant* when its  $p$  value is smaller than a *significance threshold*  $\alpha$  fixed *a priori*. This guarantees that the probability of a false discovery (i.e., reporting  $S$  as significant when it is not) is bounded by  $\alpha$ . However, if such approach is used when testing  $d$  hypotheses, the expected number of false positives is  $\alpha d$ ; when  $d$  is high, which is typically the case of significant pattern mining, this results in a large number of false positives. Therefore, an appropriate *multiple hypothesis testing correction* of the significance threshold needs to be performed in order to obtain rigorous guarantees in terms of the number of false associations reported in output.

One common approach is to perform a correction in order to bound the *Family-Wise Error Rate* (FWER), which is defined as the probability of reporting at least one false positive. Let  $FP$  be the number of false positives, then:  $\text{FWER} = \Pr(FP > 0)$ . For a given value  $\delta$ , define  $\text{FWER}(\delta)$  as the FWER obtained using  $\delta$  as *corrected* significance threshold, that is by rejecting (i.e., flagging as significant) all null hypotheses (i.e., patterns) with  $p$  value  $\leq \delta$ . Commonly, it is not possible to evaluate  $\text{FWER}(\delta)$  in closed form. One approach to set  $\delta$  is to use the Bonferroni correction, setting  $\delta$  to  $\alpha/d$ . Using the union bound, one can easily show that the resulting  $\text{FWER}(\delta) \leq d\delta = \alpha$ . The problem with this approach is that when  $d$  is high,  $\delta$  is very close to 0, resulting in low *statistical power* with many false negatives.

To increase the statistical power, more sophisticated techniques have been devised. In particular, one can look for an *optimal* corrected significance threshold  $\delta^*$ , which maximizes the statistical power while keeping the FWER bounded by  $\alpha$ :

$$\delta^* = \max\{\delta : \text{FWER}(\delta) \leq \alpha\}.$$

In this regard, a key result from Tarone (1990) is the introduction of *minimum attainable p value*: if we use a corrected significance threshold  $\delta$ , patterns whose  $p$  value cannot be  $\leq \delta$ , called *untestable*, do not need to be considered in the multiple-hypothesis correction. Therefore, if we define  $k(\delta)$  as the number of patterns having minimum

attainable  $p$  value  $\leq \delta$ , then the adjusted Bonferroni correction when threshold  $\delta$  is used can be written as  $\delta = \alpha/k(\delta)$ .

LAMP (Terada et al. 2013a) introduced such concepts into the mining of significant patterns. In particular, it exploits the following observation. For a given pattern  $S$ , its  $p$  value  $p_S$  can be expressed as a function of  $a_S$  only. Since the set of allowable values of  $a_S$  is finite, i.e.  $a_S \in [a_{S,min}, a_{S,max}]$ , there exists a minimum attainable  $p$  value  $\psi(x_S)$ , which depends on  $x_S, n_1$ , and  $n$ , and which corresponds to the most biased case for Eq. 2:

$$\psi(x_S) = \min\{p_S(u) \mid a_{S,min} \leq u \leq a_{S,max}\} .$$

Given a pattern  $S$ , its support  $x_S$ , and a corrected significance threshold  $\delta$ , if  $\psi(x_S) > \delta$  the pattern  $S$  cannot be significant;  $S$  is therefore called *untestable*. To use such observation for significant pattern mining, Terada et al. (2013a) introduced in LAMP a monotonically decreasing lower bound  $\hat{\psi}(x_S)$  on  $\psi(x_S)$ :

$$\hat{\psi}(x_S) = \begin{cases} \psi(x_S) & 0 \leq x_S \leq n_1 \\ 1/\binom{n}{n_1} & n_1 < x_S \leq n . \end{cases}$$

Terada et al. (2013a) showed that identifying a suitable significance threshold  $\delta^*$  translates into finding the maximum support threshold  $\sigma_{max}$  satisfying:

$$\hat{\psi}(\sigma_{max} - 1) > \frac{\alpha}{k(\sigma_{max})}$$

and

$$\hat{\psi}(\sigma_{max}) \leq \frac{\alpha}{k(\sigma_{max} + 1)} .$$

### 2.3 Westfall–Young permutation testing

LAMP significantly increases the statistical power of the over-conservative standard Bonferroni correction. However, it implicitly assumes that hypotheses are independent, that can result in a loss of power when there is dependence between the hypotheses, as in pattern mining. The Westfall–Young (WY) permutation testing method (Westfall and Young 1993) is a multiple hypothesis testing procedure capable of addressing this issue. This method performs *random permutations* of the class labels, creating new datasets for which no pattern  $S$  is truly associated with the permuted class labels. Since every pattern flagged as significant in the permuted datasets is a false positive, the *joint distribution* of the null hypotheses can be directly estimated, resulting in improved statistical power with respect to LAMP.

In detail, the WY method starts by creating  $j_p$  permuted datasets, with  $j_p$  sufficiently large (typically in the order of  $10^3$  or  $10^4$ ). Then for every permuted dataset  $j$  it computes the minimum  $p$  value  $p_{min}^{(j)}$  over all patterns (hypotheses). Then one can estimate  $FWER(\delta)$  as:



$$FWER(\delta) = \frac{1}{j_p} \sum_{j=1}^{j_p} \mathbb{1}[p_{\min}^{(j)} \leq \delta]$$

where  $\mathbb{1}(\cdot)$  is the indicator function (equal to 1 if its argument is true and 0 otherwise). Given a user provided threshold  $\alpha$  for the FWER, the best corrected significance threshold  $\delta^*$  can then be obtained as  $\delta^* = \max_{\delta} \{FWER(\delta) \leq \alpha\}$  with  $FWER(\delta)$  estimated as above.

The WY method does not provide an efficient way of computing the set  $\{p_{\min}^{(j)}\}_{j=1}^{j_p}$  of minimum  $p$  values. Therefore, a naïve implementation requires to exhaustively test *all* the hypothesis on all the  $j_p$  permuted datasets. For significant pattern mining, this means exploring all the patterns appearing in a dataset; since this operation can require exponential time, it is even more challenging to repeat the entire process  $j_p$  times, once for every permuted dataset. Terada et al. (2013b) proposed the first efficient implementation, FASTWY, of the WY procedure for significant pattern mining. The identification of  $\delta^*$  is based on a decremental search scheme, which starts with support  $\sigma = n$  and iteratively decrements  $\sigma$  until an appropriate condition, guaranteeing that all values  $\{p_{\min}^{(j)}\}_{j=1}^{j_p}$  have been computed, is achieved. A more recent method by Terada et al. (2015), HWY, exploits a more efficient mining strategy and parallel computing to accelerate FASTWY.

Llinares-López et al. (2015) proposed WYlight to efficiently compute the optimal value  $\delta^*$  of the corrected significance threshold. The main improvement of WYlight is to avoid the exact computation of all the elements of the set  $\{p_{\min}^{(j)}\}_{j=1}^{j_p}$  and to only produce its exact lower  $\alpha$ -quantile. This result is obtained by maintaining an estimate of the  $\alpha$ -quantile that is only lowered through the mining process. WYlight performs a depth first exploration of the patterns' search tree (Han et al. 2000) in which each pattern has support less or equal than its parent, and performs only one pattern mining instance, testing one pattern at a time and computing its  $p$  value on all the  $j_p$  permuted datasets at the same time. WYlight maintains a threshold  $\sigma$ , initialized at 1, that is raised during the execution of the algorithm pruning patterns whose  $p$  values cannot be in the lower  $\alpha$ -quantile of  $\{p_{\min}^{(j)}\}_{j=1}^{j_p}$ . This is achieved by using the lower bound  $\hat{\psi}(\sigma)$  on the minimum obtainable  $p$  value for patterns of support  $\leq \sigma$ , which allows to effectively prune the search tree. However, during the computation of  $\delta^*$ , some patterns with support  $< \psi^{-1}(\delta^*)$  may be processed, due to the depth first procedure considered by WYlight. After computing  $\delta^*$ , an additional mining of  $\mathcal{D}$  is performed (with minimum support threshold  $\psi^{-1}(\delta^*)$ ) to extract the significant patterns with  $p$  value  $\leq \delta^*$ . As shown in Llinares-López et al. (2015), WYlight significantly improves over FASTWY in particular in terms of memory requirements, allowing the extraction of significant patterns from datasets larger than the ones that can be analyzed by FASTWY.

## 2.4 Problem definition

For a dataset  $\mathcal{D}$ , let  $\delta(\alpha) = \max_{\delta} \{FWER(\delta) \leq \alpha\}$  the threshold obtained through the WY permutation procedure when the bound on the FWER is set to  $\alpha$ . Let  $p^{(k)}$  be the  $p$



value of the  $k$ -th pattern with patterns sorted by (increasing)  $p$  value. Given a dataset  $\mathcal{D}$  and user-provided values  $k$  and  $\alpha$ , our goal is to extract the set  $TSP(\mathcal{D}, k, \alpha)$  of top- $k$  statistically significant patterns with  $FWER \leq \alpha$ , defined as:

$$TSP(\mathcal{D}, k, \alpha) = \left\{ S : p_S \leq \min\{\delta(\alpha), p^{(k)}\} \right\} .$$

Note that when less than  $k$  patterns have  $p$  value below  $\delta(\alpha)$ ,  $TSP(\mathcal{D}, k, \alpha)$  contains all such patterns. In addition, according to our definition more than  $k$  patterns may be in  $TSP(\mathcal{D}, k, \alpha)$ , in case many have the same  $p$  value  $p^{(k)}$ . In particular, for any two patterns  $S, S'$  with  $S' \subset S$  and  $x_{S'} = x_S, a_{S'} = a_S$  we have that  $p_S = p_{S'}$ . For this reason we restrict our interest only to *closed* patterns, i.e. patterns whose supersets have support *strictly lower* than the pattern itself. Since the definition of closed pattern does not depend on the class labels, restricting to closed patterns does not bias any analysis.

The following result establishes the required guarantees on false positives in  $TSP(\mathcal{D}, k, \alpha)$  and it is a direct consequence of the fact that  $TSP(\mathcal{D}, k, \alpha)$  is a subset of all the patterns that would be reported using the WY method.

**Lemma 1** *The set  $TSP(\mathcal{D}, k, \alpha)$  has  $FWER \leq \alpha$ .*

### 3 TOPKWY Algorithm

In this section we present our algorithm TOPKWY for mining the set  $TSP(\mathcal{D}, k, \alpha)$ . We first present its main strategy (Sect. 3.1) that can be applied to any pattern mining problem. We then analyze TOPKWY showing theoretical evidence of the efficiency of its strategy (Sect. 3.2) and introduce improved bounds on the minimum attainable  $p$  value used by TOPKWY (Sect. 4). We also present extensions of TOPKWY (Sect. 5) to control the *generalized* FWER, to control the False Discovery Proportion (FDP), and to employ different exploration strategies on the tree of candidate patterns. Finally, we introduce some crucial implementation details (Sect. 6), focusing on the problem of mining significant itemsets and subgraphs.

#### 3.1 Main strategy

TOPKWY combines two key ideas. First, it maintains an estimate of  $\delta_m = \min\{\delta(\alpha), p^{(k)}\}$  that is updated during the exploration of the patterns and maintains a corresponding minimum support threshold  $\sigma = \psi^{-1}(\delta_m)$  that is raised during the exploration of the patterns. Analogously to the strategy employed by WYlight (Linares-López et al. 2015) the updates of  $\delta_m$  and  $\sigma$  depend on  $\alpha$ , but in addition TOPKWY updates them also depending on the  $p$  values of members of  $TSP(\mathcal{D}, k, \alpha)$ . Second, the search tree of all possible patterns is explored in order of decreasing support, analogously to the strategy used by TOPKMINER (Pietracaprina

and Vandin 2007) for mining top- $k$  frequent patterns, which guarantees that only patterns of support greater or equal to the final value of  $\sigma$  (i.e.,  $\psi^{-1}(\delta_m)$ ) are explored.<sup>1</sup>

TOPKWY is described in Algorithm 1. In line 1, the threshold  $\delta_m$  is initialized to  $\alpha$  (the threshold with no correction for multiple hypothesis) and  $\sigma$  is initialized accordingly to  $\hat{\psi}^{-1}(\delta_m)$ . All the elements of the set of minimum  $p$  values  $\{p_{\min}^{(j)}\}_{j=1}^{j_p}$  observed on the permuted datasets are initialized to 1 (their maximum achievable value) in line 2. The labels of the permuted datasets are generated in line 3. The pattern exploration is organized using a priority queue  $Q$  where each entry represents a pattern  $S$ , with key equal to the support  $x_S$  and value representing all the information needed by the algorithm regarding  $S$  (e.g.,  $a_S$ ) and also with relevant information regarding the parent  $f_S$  of pattern  $S$  in the search tree (see Sect. 4).  $Q$  is initialized in line 4 and stores the frontier of unexplored patterns, keeping them accessible by non-increasing support. TOPKWY stores patterns having  $p$  value  $\leq \delta_m$  in a priority queue  $P$ , keeping them accessible by non-decreasing  $p$  value. This is the set of candidates for  $TSP(\mathcal{D}, k, \alpha)$ , which are collected and produced in output as soon as possible during the exploration. This allows to reduce the memory requirements and to start analyzing the results during the exploration, without the need of waiting for the algorithm's termination. The first patterns in  $Q$  are obtained by the `expand`( $S, Q$ ) operation on line 5 called on the empty pattern  $S = \emptyset$ : this procedure generates all patterns children of the pattern  $S$  in the search tree (and their corresponding projected datasets), and inserts the ones of support  $\geq \sigma$  in the queue  $Q$ . The details of efficient implementations of `expand` are described in Sect. 6. The `while` loop (lines 6–21) implements the main step of the exploration strategy: the most frequent pattern  $S$ , its support  $x_S$ , its support  $a_S$  in the minority class of  $\mathcal{D}$ , and the relevant information for its parent  $f_S$  are extracted from  $Q$  in line 7. If the  $p$  value  $p_S$  of  $S$  is  $\leq \delta_m$ , then  $S$  is inserted in  $P$  in line 10. In line 8,  $\sigma'$  is set to  $x_S$ , which is an upper bound to the support of all elements stored in  $Q$ . This quantity is used to identify patterns surely in the set  $TSP(\mathcal{D}, k, \alpha)$  without waiting for the final corrected significance threshold  $\delta_m$  to be found, done in lines 11 and 12.  $k$  is updated accordingly, reducing it to the number of patterns which still need to be found. In order to compute the corrected significance threshold  $\delta_m$ , the algorithm computes the  $p$  values of pattern  $S$  in the  $j_p$  permuted datasets, updating the values of  $\{p_{\min}^{(j)}\}_{j=1}^{j_p}$  if needed. This operation is done with the `test` procedure. Similarly to WYlight, our algorithm processes all the  $j_p$  permutations for every pattern  $S$  at once, computing only the needed exact lower quantile of the set of minimum  $p$  values of the WY permutations, and not the minimum  $p$  values of every permuted dataset. Differently from WYlight, we use an improved lower bound  $\psi'(x_S, f_S)$  to the minimum attainable  $p$  value of  $S$  to decide (in line 13) whether to test  $S$  on the permuted datasets or not (see Sect. 4). This allows to skip the expensive computation of the  $j_p$  supports  $\{a_S^{(j)}\}_{j=1}^{j_p}$  of  $S$  on the minority class of the permuted datasets for several patterns  $S$ .

The significance threshold  $\delta_m$  is decreased during the exploration in two cases: when the estimated  $FWER(\delta_m)$  for the current threshold  $\delta_m$  increases above  $\alpha$  (line 15), or

<sup>1</sup> This assumes that the search tree for patterns has the property that the children of a node have support not greater than the node itself, which is a usual property of pattern mining algorithms (Han et al. 2007; Uno et al. 2005; Nijssen and Kok 2004) and is required by WYlight as well.

when more than  $k$  patterns with  $p$  value  $\leq \delta_m$  are observed (line 17). The corresponding minimum support threshold  $\sigma$  is then updated accordingly in line 18. The correctness of these steps are proved in Sect. 3.2. After the update of  $\delta_m$  and  $\sigma$ , elements which have become untestable are removed from  $Q$  in line 19, and elements which are not significant are removed from  $P$  in line 20. The current pattern  $S$  is expanded in line 21, and all its children having support  $\geq \sigma$  are inserted into  $Q$ . The exploration ends when  $Q$  gets empty. When this happens, all elements still contained in  $P$  with  $p$  value at most  $\delta_m$  are reported as significant in line 22.

The strategy employed by TOPKWY can be adapted to incrementally update  $k$  for the same  $\alpha$ , providing an interactive mining process. This can be achieved by providing a maximum value  $k^*$  in input to in Algorithm 1 to definitely prune untestable patterns, but freezing the computation after  $k$  patterns with  $p$  value below the current value of  $\hat{\psi}(\sigma)$  have been found. If the user wants to increase  $k$ , the exploration can continue without restarting the entire mining instance.

---

**Algorithm 1:** TOPKWY

---

```

Input: Transaction dataset  $\mathcal{D}$  with class labels  $c$ , number of permutations  $j_p$ , target FWER  $\alpha$ ,
number of results  $k$ 
Output: Set of top- $k$  significant patterns with FWER  $\leq \alpha$ 
1  $\delta_m \leftarrow \alpha; \sigma \leftarrow \hat{\psi}^{-1}(\delta_m);$ 
2  $p_{\min}^{(j)} \leftarrow 1, \forall j \in [1, j_p];$ 
3 generate  $j_p$  permuted class labels;
4  $Q, P \leftarrow$  empty priority queues;
5  $\text{expand}(\emptyset, Q);$ 
6 while  $Q \neq \emptyset$  do
7    $(S, x_S, a_S, f_S) \leftarrow Q.\text{removeMax}();$ 
8    $\sigma' \leftarrow x_S;$ 
9   if  $p_S \leq \delta_m$  then
10     $P.\text{insert}(S, p_S);$ 
11    /*  $\mathcal{O} =$  patterns surely in  $TSP(\mathcal{D}, k, \alpha)$  */
12     $\mathcal{O} \leftarrow \{S' \in P : p_{S'} < \hat{\psi}(\sigma')\};$  produce  $\mathcal{O}$  in output;
13    remove patterns in  $\mathcal{O}$  from  $P; k \leftarrow k - |\mathcal{O}|;$ 
14    if  $\psi'(x_S, f_S) \leq \delta_m$  then
15      $\text{test}(S, \{p_{\min}^{(j)}\}_{j=1}^{j_p});$ 
16     /* update  $\delta_m$  based on estimate of  $\delta^*$  */
17      $\delta_m \leftarrow \min\{\delta_m, \max\{\delta : FWER(\delta) \leq \alpha\}\};$ 
18     /* update  $\delta_m$  based on top- $k$  patterns in  $P$  */
19      $p^{(k)} \leftarrow k\text{-th largest } p \text{ value in } P;$ 
20      $\delta_m \leftarrow \min\{\delta_m, p^{(k)}\};$ 
21     /* update  $\sigma$  */
22      $\sigma \leftarrow \psi^{-1}(\delta_m);$ 
23     /* remove untestable patterns from  $Q$  */
24     remove from  $Q$  all patterns  $S'$  with  $x_{S'} < \sigma;$ 
25     /* remove non-significant patterns from  $P$  */
26     remove from  $P$  all patterns  $S'$  with  $p_{S'} > \delta_m;$ 
27      $\text{expand}(S, Q);$ 
28 produce in output  $\{S' \in P : p_{S'} \leq \delta_m\};$ 

```

---

### 3.2 Analysis

Some important properties of TOPKWAY algorithm can be formally stated. The first regards the correctness of the algorithm.

**Theorem 1** (Correctness of TOPKWAY) *TOPKWAY outputs the set  $TSP(\mathcal{D}, k, \alpha)$  of top- $k$  significant patterns with  $FWER \leq \alpha$ .*

**Proof** The correctness of TOPKWAY follows from two observations: first, the final threshold  $\delta_m$  obtained by the algorithm is correct; second, only patterns with  $p$  value less or equal than the final value of  $\delta_m$  are produced in output. We start by proving the first statement.  $\delta_m$  is initialized to the value  $\alpha$ , that is the uncorrected threshold for significance and is always  $\geq \delta^*$ .  $\delta_m$  is decreased (and the corresponding minimum support threshold  $\sigma$  is increased) during the exploration in two cases. The first case (line 15) is when the estimated  $FWER(\delta_m)$  for the current threshold  $\delta_m$  increases above  $\alpha$ . This means that more than  $\alpha j_p$   $p$  values of  $\{p_{\min}^{(j)}\}_{j=1}^{j_p}$  are below the current significance threshold  $\delta_m = \hat{\psi}(\sigma)$ , which allows for too many false positives, and the FWER is not correctly controlled to the level  $\alpha$ .  $\delta_m$  is then updated to the highest value of  $\delta$  for which  $FWER(\delta) \leq \alpha$ . The second case is when more than  $k$  patterns with  $p$  value  $\leq \delta_m$  are observed (line 17). In this case, let  $\tilde{p}$  be the highest  $p$  value of the  $k$  most significant patterns observed up to this point. Then all patterns of support  $< \hat{\psi}^{-1}(\tilde{p})$  cannot result in a  $p$  value  $< \tilde{p}$  and therefore we need to consider (both in  $\mathcal{D}$  and in the permuted datasets) only patterns of support at least  $\hat{\psi}^{-1}(\tilde{p})$ . That is, the minimum support threshold  $\sigma$  can be safely increased to  $\hat{\psi}^{-1}(\tilde{p})$  with a corresponding significance threshold  $\tilde{p}$ . When  $\delta_m$  is last updated, its value will then be equal to the minimum between  $\delta(\alpha)$  and  $p^{(k)}$ .

We now prove the second statement. This is trivially correct for patterns produced in output by line 22. We then consider patterns produced in output in line 11. Note that the current pattern  $S$  has support  $\sigma'$  and the search strategy employed by TOPKWAY guarantees that all patterns with support  $> \sigma'$  have already been explored. Therefore, from this point on the algorithm will never encounter  $p$  values  $< \hat{\psi}(\sigma')$  and therefore the corrected significance threshold  $\delta_m$  will be  $\geq \hat{\psi}(\sigma')$ . Thus all patterns in  $P$  with  $p$  value  $< \hat{\psi}(\sigma')$  can be safely produced in output (and removed from  $P$ ).  $\square$

The following result provides theoretical guarantees on which patterns will be explored by TOPKWAY, providing analytical evidence of the efficiency of our strategy.

**Theorem 2** (Optimality of TOPKWAY) *TOPKWAY expands only patterns of support  $\geq \hat{\psi}^{-1}(\delta_m)$ .*

**Proof** Similarly to the proof of Theorem 1, when pattern  $S$  of support  $\sigma'$  is extracted from  $\mathcal{Q}$ , we are guaranteed that the algorithm will never encounter  $p$  values  $< \hat{\psi}(\sigma')$  again. Therefore the corrected significance threshold  $\delta_m$  will be  $\geq \hat{\psi}(\sigma')$ , that is  $\sigma' \geq \hat{\psi}^{-1}(\delta_m)$  (i.e.,  $S$  is testable).  $\square$

We now show that in a simplified model for how  $p$  values are obtained, there exists a family of datasets for which the expected difference between the number of patterns

explored by a DFS strategy and the number of patterns explored by TOPKWY is exponential in the size of the dataset. In the simplified model, the  $p$  values obtained by random permutations are uniformly distributed in  $(0, 1]$ . (Note that we do not assume independence among  $p$  values from different itemsets.) Consider now the family of datasets  $\mathcal{D}_n = \{t_1, \dots, t_n\}$  defined on the set of (binary) features  $\mathcal{I} = \{i_1, \dots, i_n\}$  where  $t_j = \mathcal{I} \setminus \{i_j\}$ . Moreover, half of transactions in  $\mathcal{D}_n$  have label 0 while the other half have label 1.

**Theorem 3** Consider a dataset  $\mathcal{D}_n$  from the family described above. Let  $\varepsilon$  be a constant such that  $0 < \varepsilon \leq \frac{1}{2}$  and  $\varepsilon n \in \mathbb{N}$ . Assume that the choice of  $\alpha$  and  $k$  is such  $\delta^* = \hat{\psi}(\varepsilon n) = \binom{\frac{n}{2}}{\varepsilon n} / \binom{n}{\varepsilon n}$ , and that  $j_p$  random permutations are used. Let  $X$  be the difference between the number of patterns explored by a DFS strategy and the number of patterns explored by TOPKWY in the simplified model above. Then  $\mathbb{E}[X] = \Omega(2^{\varepsilon n} / j_p)$ .

**Proof** Note that in such dataset all patterns are closed. Let  $W$  be the set of patterns explored by TOPKWY. The DFS strategy has to explore all the patterns in  $W$  (since they are testable). We now show that it will also explore, in expectation,  $\Omega(2^{\varepsilon n} / j_p)$  additional patterns, that proves the statement.

It is easy to show that since  $\varepsilon n \in \mathbb{N}$  and  $\delta^* = \binom{\frac{n}{2}}{\varepsilon n} / \binom{n}{\varepsilon n}$ , the set of testable patterns  $W$  is given by all patterns of size  $\leq \varepsilon n$ . For each pattern  $S_i$  and each  $j$  with  $1 \leq j \leq j_p$ , let  $X_{ij}$  be the random variable that is 1 if  $S_i$  has  $p$  value  $\leq \delta^*$  in the  $j$ -th permuted dataset, and 0 otherwise. Note that  $S_{ij}$  is a Bernoulli random variable of parameter  $\delta^*$ , for all  $i$  and  $j$ . Let  $Y$  be the number of  $p$  values lower than  $\delta^*$ , assuming that the DFS has explored  $\varepsilon n + m$  patterns. The expectation of  $Y$  is

$$\mathbb{E}[Y] = \mathbb{E} \left[ \sum_{i=1}^{\varepsilon n + m} \sum_{j=1}^{j_p} X_{ij} \right] = \sum_{i=1}^{\varepsilon n + m} \sum_{j=1}^{j_p} \mathbb{E}[X_{ij}] = (\varepsilon n + m) j_p \delta^*. \tag{3}$$

Note that requiring  $\mathbb{E}[Y] \geq 1$  provides a lower bound to the number of  $p$  values needed for the DFS to establish that  $\delta^* = \binom{\frac{n}{2}}{\varepsilon n} / \binom{n}{\varepsilon n}$ , since  $\alpha j_p$   $p$  values below such threshold must be observed.

Let  $v = \varepsilon n$ . The condition  $\mathbb{E}[Y] = (\varepsilon n + m) j_p \delta^* \geq 1$  implies

$$\begin{aligned} (\varepsilon n + m) &\geq \frac{1}{j_p \delta^*} = \frac{1}{j_p} \frac{\binom{n}{v}}{\binom{\frac{n}{2}}{v}} = \frac{1}{j_p} \frac{n! (\frac{n}{2} - v)!}{(\frac{n}{2})! (n - v)!} \\ &= \frac{1}{j_p} \frac{(n - v)! \prod_{j=0}^{v-1} (n - j)}{(n - v)!} \frac{(\frac{n}{2} - v)!}{(\frac{n}{2} - v)! \prod_{j=0}^{v-1} (\frac{n}{2} - j)} \\ &= \frac{1}{j_p} \prod_{j=0}^{v-1} \frac{(n - j)}{(\frac{n}{2} - j)} \geq \frac{1}{j_p} \prod_{j=0}^{v-1} 2 = \frac{2^v}{j_p} = \frac{2^{\varepsilon n}}{j_p} \end{aligned}$$

or, equivalently

$$m \geq \frac{2^{\varepsilon n}}{j_p} - \varepsilon n \in \Omega(2^{\varepsilon n}/j_p).$$

Note that since the set of testable patterns includes all patterns of size  $\leq \varepsilon n$ , all the  $\Omega(2^{\varepsilon n}/j_p)$  patterns explored by the DFS *after* exploring the first  $\varepsilon n$  patterns and before observing the first  $p$  value below  $\delta^*$  have size  $> \varepsilon n$  and, thus, are not in  $W$ , which proves the statement.  $\square$

Compared to a depth first search (DFS) exploration strategy (i.e., the one employed by WYlight), the *best first* exploration strategy followed by TOPKWY has the additional costs required by operations involving data structures  $P$  and  $Q$ .  $P$  can be implemented as a heap of entries  $(p, \ell_p)$ , where the key  $p$  is a  $p$  value and the value  $\ell_p$  is a list of patterns which contains all the patterns with the same  $p$  value  $p$ . With such implementation, operations involving  $P$  (lines 10, 11, 16, and 20) can be performed with  $O(\log k)$  operations, since  $P$  will contain at most  $k$  entries. Analogously,  $Q$  can be implemented as a heap of entries  $(x, \ell_x)$ , where the key  $x$  is a support and the value  $\ell_x$  is a list of patterns with the same value of  $\hat{\psi}(x)$  (that is, having the same support  $x$ ). With such implementation, operations involving  $Q$  (lines 7, 19, and 21) can be performed with  $O(\log n)$  operations ( $n$  is the number of transactions in  $\mathcal{D}$ ). Alternatively,  $P$  and  $Q$  can be implemented so that all operations require time  $O(1)$  by storing references to lists  $\ell_p$  and  $\ell_x$  in arrays of size  $O(n^2)$  (since the  $p$  value of  $S$  is a function of the support  $x_S$  of  $S$  and the number  $a_S$  of transactions with label 1 and containing  $S$ ) and  $O(n)$ , respectively. (Note that this additional space requirement is not impractical, since  $O(nj_p)$  space is needed to store the random permutations.) We also note that computing the improved lower bound  $\psi'(x_S, f_S)$  (line 13), has the same cost as computing the lower bound  $\hat{\psi}(x_S)$  to the  $p$  value used in WYlight (see Sect. 4). Even with the additional costs required by  $P$  and  $Q$ , the best first strategy of TOPKWY leads to significant improvements in running time, as demonstrate by our experimental evaluation (Sect. 7).

#### 4 Improved bounds on minimum attainable $p$ value

In this section we prove novel and efficiently computable lower bounds on the minimum  $p$  value achievable by a pattern  $S$  that are tighter than the ones introduced by LAMP (Terada et al. 2013a) and are of particular interest in the context of WY permutation testing. These bounds are based on information computed when processing a *parent pattern*  $Y$  of  $S$ ; in the case of itemsets,  $Y$  is a parent of  $S$  (or, alternatively,  $S$  is a *child* or a *super pattern* of  $Y$ ) when  $Y \subset S$ . Such bounds can be used to *skip* the expensive processing of the permutations for  $S$  when they ensure that it is not possible to improve the current estimate of the corrected significance threshold. While we present these bounds as a critical component of TOPKWY, they may be of independent interest since can be employed in WYlight or similar algorithms to speed-up WY permutation testing.

Let the pattern  $S$  be a super pattern of  $Y$ , that is  $S \supset Y$ . Then  $x_Y \geq a_Y \geq 0$  and  $x_Y \geq x_S$ . Since the set of transactions (i.e., the *conditional dataset*) containing  $S$  is a subset of the set of transactions containing  $Y$ , we can bound the support  $a_S$  of  $S$  in the class  $c_1$  with the following relations:

$$\max(a_Y - (x_Y - x_S), 0) \leq a_S \leq \min(x_S, a_Y).$$

Considering the  $j_p$  permuted class labels, let  $a_Y^{(j)}$  be the number of transactions containing  $Y$  and in the minority class (i.e.,  $a_Y^{(j)}$  is the value of  $a_Y$  when the class labels are given by the  $j$ -th permutation). An analogous relation holds between  $a_S^{(j)}$  and  $a_Y^{(j)}$ , for all  $j$ :

$$\max(a_Y^{(j)} - (x_Y - x_S), 0) \leq a_S^{(j)} \leq \min(x_S, a_Y^{(j)}).$$

An immediate consequence of these bounds on  $a_S^{(j)}$  are *lower bounds* to  $p_S(a_S^{(j)})$ .

**Lemma 2** Let  $\check{a}_S^{(j)} = \max(a_Y^{(j)} - (x_Y - x_S), 0)$  and  $\hat{a}_S^{(j)} = \min(x_S, a_Y^{(j)})$ . Then, for all  $j \in [1, j_p]$ ,

$$p_S(a_S^{(j)}) \geq \min \left\{ p_S(\check{a}_S^{(j)}), p_S(\hat{a}_S^{(j)}) \right\}.$$

This result suggests that if we have already computed  $a_Y^{(j)}, \forall j \in [1, j_p]$  while processing the permutations of  $Y$ , we could skip the expensive computation of  $a_S^{(j)}$ , and, therefore,  $p_S(a_S^{(j)}), \forall j \in [1, j_p]$ , in situations when the lower bounds to  $p_S(a_S^{(j)})$  are *greater* than the current value of the corrected significance threshold. In the following, we present a bound valid for all  $p_S(a_S^{(j)})$  simultaneously, that is a function of only the minimum and maximum elements of  $a_Y^{(j)}$ , instead of all of them. Let

$$a_{Y_{\min}} = \min \left\{ a_Y^{(j)} : j \in [1, j_p] \right\}$$

and

$$a_{Y_{\max}} = \max \left\{ a_Y^{(j)} : j \in [1, j_p] \right\}.$$

Then,  $\forall j \in [1, j_p]$  we bound  $a_S^{(j)}$  as:

$$a_{S_{\min}} = \max(a_{Y_{\min}} - (x_Y - x_S), 0) \leq a_S^{(j)} \leq \min(x_S, a_{Y_{\max}}) = a_{S_{\max}}.$$

This allows to compute a bound  $\psi'(x_S, x_Y, a_{Y_{\min}}, a_{Y_{\max}})$  to the minimum attainable  $p$  value of  $S$  that is tighter than  $\psi(x_S)$ :

$$\psi'(x_S, x_Y, a_{Y_{\min}}, a_{Y_{\max}}) = \min(p_S(a_{S_{\min}}), p_S(a_{S_{\max}})). \tag{4}$$



The bound in Eq. 4 is evaluated in constant time, assuming  $p_S(a)$  is pre-computed for all valid values of  $a$ , as done in WYlight and in TOPKWY to efficiently implement the `test` function. The following are simple consequences of Lemma 2 and the fact that  $a_{S_{\min}}$  and  $a_{S_{\max}}$  are always equally or more tight than the naive bounds on  $a_S$  assumed by  $\psi(x_S)$ .

**Lemma 3**  $\min \left\{ p_S \left( a_S^{(j)} \right) : j \in [1, j_p] \right\} \geq \psi'(x_S, x_Y, a_{Y_{\min}}, a_{Y_{\max}}) \geq \psi(x_S)$ .

If for the current value of the significance threshold  $\delta_m$  it holds that  $\psi'_S > \delta_m$ , then we can infer, without computing  $\{a_S^{(j)}\}_{j=1}^{j_p}$ , that none of the  $j_p$   $p$  values of  $S$  in the permuted datasets will improve the estimate of the current lower-quantile of the set  $\{p_{\min}^{(j)}\}_{j=1}^{j_p}$  and therefore cannot contribute to the computation of  $\delta(\alpha)$  or  $\delta_m$ . That is, all the computation on the permuted datasets can be skipped for the current pattern  $S$ . For all children of  $S$ , if  $S$  is not tested the bounds  $a_{S_{\min}}$  and  $a_{S_{\max}}$  can be propagated to compute bounds also on their class distribution; if  $S$  is tested, then we propagate the actual minimum and maximum values of  $\{a_S^{(j)}\}_{j=1}^{j_p}$ . In Algorithm 1 we use the bound above with the values propagated by the parent  $f_S$  of  $S$  and use  $\psi'(x_S, f_S)$  to highlight this fact. This optimization is particularly effective when patterns have a high degree of correlation, i.e., when patterns share many transactions.

Note that even if  $S$  does not need to be tested, descendants of  $S$  may need to be tested. However, using the bound  $\psi'(\cdot)$  we can quickly identify cases in which none of the descendants of  $S$  need to be explored and therefore the entire subtree can be pruned. In particular, since all the descendants of  $S$  will have support  $\leq x_S - 1$ , considering  $a_S$  (i.e., the number of transactions containing  $S$  and in the minority class in the dataset  $\mathcal{D}$ ), the algorithm can find  $\min\{\psi'(i, x_S, a_{S_{\min}}, a_{S_{\max}}) : i \in [\sigma, x_S - 1]\}$ , and if such value is  $> \delta_m$  we can prune all the search subtrees rooted in the children of  $S$ . This optimization is part of the `expand` operation in TOPKWY. These novel bounds consider the information of one common ancestor pattern to avoid useless computations for many of its children: in practice, the number of tests to perform across the permuted datasets can be significantly smaller than the number of testable patterns, leading to a significant computational speed-up. The approach above can be extended to bound  $\min \left\{ p_S \left( a_S^{(j)} \right) : j \in [1, j_p] \right\}$  by considering the information computed on the intersection of the conditional datasets of any pair of patterns  $S$  and  $Y$ , even if  $S \not\supseteq Y$ .

Another possible extension of the techniques we derive in this section is to consider not only *one pair*  $(a_{Y_{\min}}, a_{Y_{\max}})$ , but  *$v$  pairs*  $(a_{Y_{\min}}^i, a_{Y_{\max}}^i)$ , for all  $i \in [1, v]$ : for each  $i$ , we define the set  $J_i$  as a subset of  $\{1, \dots, j_p\}$ , with  $\bigcup_{i=1}^v J_i = \{1, \dots, j_p\}$ . For all  $i$ , we bound all values  $a_Y^{(j)}$  for all  $j \in J_i$  with  $a_{Y_{\min}}^i$  and  $a_{Y_{\max}}^i$ . If we define

$$a_{Y_{\min}}^i = \min \left\{ a_Y^{(j)} : j \in J_i \right\}, \quad a_{Y_{\max}}^i = \max \left\{ a_Y^{(j)} : j \in J_i \right\},$$

$$a_{S_{\min}}^i = \max(a_{Y_{\min}}^i - (x_Y - x_S), 0), \quad a_{S_{\max}}^i = \min(x_S, a_{Y_{\max}}^i),$$

then we simply obtain

$$\min \left\{ p_S \left( a_S^{(j)} \right) : j \in J_i \right\} \leq \min \left\{ p_S \left( a_{S_{\min}}^i \right), p_S \left( a_{S_{\max}}^i \right) \right\}.$$

This provides a trade-off between the memory (required to store the  $2v$  bounds on the expanded nodes of the search space) and the time to evaluate the bound (that is linear in the number  $v$  of sets instead of constant), and the time the algorithm saves by skipping computations of the permutations, that is a direct consequence of how tight the bounds on the  $p$  values are; in fact, the permutations that have to be processed are only the ones belonging to the sets  $J_i$  such that  $\min \left\{ p_S \left( a_{S_{\min}}^i \right), p_S \left( a_{S_{\max}}^i \right) \right\}$  is not higher than the current value of the significance threshold.

## 5 Extensions of TOPKWY

### 5.1 Controlling the generalized FWER

While the main focus of TOPKWY is to control the FWER of the output, a simple modification provides an algorithm to control the *generalized* FWER ( $g$ -FWER Lehmann and Romano 2012). The  $g$ -FWER is defined as the probability that at least  $g$  false positives are reported in output. In several applications one may be willing to tolerate a small amount of false discoveries in order to increase the power of detecting significant patterns, provided the number of false discoveries can be controlled. In such cases methods to discover significant patterns while controlling the  $g$ -FWER are preferred to methods controlling FWER.

Let  $FP$  be the number of false positives, then:  $g$ -FWER =  $\Pr(FP \geq g)$ . Let  $g$ -FWER( $\delta$ ) be the  $g$ -FWER obtained using  $\delta$  as *corrected* significance threshold, that is by flagging as significant patterns with  $p$  value  $\leq \delta$ . Note that the Westfall–Young procedure can be used to estimate  $g$ -FWER( $\delta$ ) as

$$g\text{-FWER}(\delta) = \frac{1}{J_p} \sum_{j=1}^{j_p} \mathbb{1}[p_g^{(j)} \leq \delta]$$

where  $p_g^{(j)}$  is the  $g$ -th smallest  $p$  value (over all patterns) in the permuted dataset  $j$ .

Algorithm 1 can be simply modified to obtain the set of top- $k$  significant patterns with  $g$ -FWER  $\leq \alpha$ . To achieve this, it is sufficient to perform the following changes: replace  $\text{test}(S, \{p_{\min}^{(j)}\}_{j=1}^{j_p})$  (line 14) with  $\text{test}(S, \{p_g^{(j)}\}_{j=1}^{j_p})$ , where  $\text{test}(S, \{p_g^{(j)}\}_{j=1}^{j_p})$  computes the  $p$  value of pattern  $S$  in the  $j_p$  permuted datasets and updates the values of  $\{p_g^{(j)}\}_{j=1}^{j_p}$  if needed; replace  $\max\{\delta : \text{FWER}(\delta) \leq \alpha\}$  (line 15) with  $\max\{\delta : g\text{-FWER}(\delta) \leq \alpha\}$ . Let TOPKWY- $g$  be such modified algorithm. We have the following.

**Lemma 4** TOPKWY- $g$  outputs the set  $TSP(\mathcal{D}, k, \alpha)$  of top- $k$  significant patterns with  $g$ -FWER  $\leq \alpha$ .

The proof is analogous to the proof of Theorem 1.

In addition to finding the top- $k$  most significant pattern with bounded  $g$ -FWER, with  $g$  provided in input, TOPKWY- $g$  can be adapted to a different scenario. In this

case, one may want to retrieve the  $k$  most significant patterns, using the random permutations to obtain a rigorous *estimate* of how many of such results are likely to be false positives. More formally, for  $k$  and  $\alpha$  provided by the user, one may be interested in computing the quantity  $g^*$  defined as

$$g^* = \min \left\{ g : g\text{-FWER}(p^k) \leq \alpha \right\},$$

that is the *minimum* value of  $g$  such that the  $g\text{-FWER}$  is controlled when the significance threshold is  $p^{(k)}$ , that is the highest  $p$  value of the set of top- $k$  results we extracted.  $g^*$  provides useful knowledge on the *quality* of the set of top- $k$  results provided to the user. In this situation the user is not required to fix a-priori  $g$  before examining the data. Further simple modifications to TOPKWY- $g$  are sufficient to obtain such variant, that are: remove line 15 and replace line 22 with “produce in output  $\{S' \in P : p_{S'} \leq \delta_m\}$  and  $g^* = \min \{g : g\text{-FWER}(p^k) \leq \alpha\}$ ”. Let TOPKWY\* be such modified algorithm; we obtain the following guarantees.

**Lemma 5** TOPKWY\* outputs the set of patterns  $\{S : p_S \leq p^k\}$  of top- $k$  significant patterns and  $g^* = \min \{g : g\text{-FWER}(p^k) \leq \alpha\}$ .

## 5.2 Bounding the proportion of false discoveries

The False Discovery Proportion (*FDP*) (van der Laan et al. 2004; Lehmann and Romano 2012) of a set of hypotheses  $\mathcal{P}$  is defined as the ratio  $FD/|\mathcal{P}|$ , where  $FD$  is the (unknown) number of false discoveries  $\in \mathcal{P}$ ; note that when  $|\mathcal{P}| = 0$ , the *FDP* is assumed to be 0. Let  $\zeta, \alpha \in (0, 1)$  and  $k, g \in [1, +\infty)$ . Define the set  $\mathcal{P}$  such that all the following hold:

$$\begin{aligned} \max \{p_S : (S, p_S) \in \mathcal{P}\} &\leq \delta^*, \\ \delta^* &= \max \{\delta : g\text{-FWER}(\delta) \leq \alpha\}, \\ g &\leq \zeta |\mathcal{P}|, \quad |\mathcal{P}| \leq k. \end{aligned}$$

It is possible to prove that a set  $\mathcal{P}$  satisfying the above guarantees has size at most  $k$  and  $FDP \leq \zeta$  with probability  $\geq 1 - \alpha$ . Simple modifications to TOPKWY- $g$  lead to an algorithm that outputs  $\mathcal{P}$  with the aforementioned guarantees: remove lines 11 and 12, replace line 15 with “ $\delta_m \leftarrow \max\{\delta : (\lfloor k\zeta \rfloor)\text{-FWER}(\delta) \leq \alpha\}$ ” and line 22 with “produce in output  $\mathcal{P}(\delta^*) = \{S' \in P : p_{S'} \leq \delta^*\}$  where  $\delta^* = \max\{\delta : (\lfloor \zeta |\mathcal{P}(\delta)| \rfloor)\text{-FWER}(\delta) \leq \alpha\}$ ”. Let such algorithm be TOPKWY- $\zeta$ . We obtain the following result.

**Lemma 6** TOPKWY- $\zeta$  outputs the set of patterns  $\mathcal{P}$  of size at most  $k$  with False Discovery Proportion  $\leq \zeta$  with probability  $\geq 1 - \alpha$ .

### 5.3 Alternative exploration strategies

While TOPKWY builds on examining the search tree of all possible patterns in order of decreasing support, i.e. with a *best first strategy* analogous to the one used by TOPKMINER (Pietracaprina and Vandin 2007) for mining top- $k$  frequent patterns, it can also be modified to efficiently obtain the set  $TSP(\mathcal{D}, k, \alpha)$  using different exploration strategies, e.g. a level-wise exploration of the search tree [performed, e.g., by the Apriori algorithm (Agrawal and Srikant 1994) for itemsets] or a depth first search on the tree of all possible patterns (Uno et al. 2005; Nijssen and Kok 2004). This can be achieved by setting  $\sigma'$  to  $\max \{x_{S'} : S' \in Q\}$  (instead that to  $x_S$ ) in line 8 of Alg. 1 and by an appropriate choice of the priority for patterns in the priority queue  $Q$ , that stores the frontier of unexplored patterns: to obtain a level-wise exploration for itemsets, the priority of pattern  $S$  is set to the total number  $|\mathcal{I}|$  of items minus  $|S|$ ; to obtain a depth first search, the priority of patterns  $S$  is set to its level in the search tree.

While for strategies other than the *best first* one the optimality (Theorem 2) is not guaranteed, the possibility to employ other strategies allows to obtain the top- $k$  significant patterns starting from efficient implementations of frequent pattern mining algorithms that build on such strategies for various types of patterns [e.g., subgraphs (Nijssen and Kok 2004)].

## 6 Implementation details

An efficient implementation of *expand* and *test* procedures is critical for the efficiency of TOPKWY. This crucially depends on the representation of  $\mathcal{D}$  and the permuted class labels, and both depend on the type of patterns of interest. In Sects. 6.1 and 6.2 we now describe in more details the implementations for significant itemsets mining and for significant subgraphs mining; in particular, for significant itemsets we discuss the implementation of TOPKWY as described in Sect. 3.1 as well as the variant, described in Sect. 5.3 of TOPKWY based on the DFS strategy, which we denote by TOPKWY-dfs; for significant subgraphs we only consider the implementation of TOPKWY-dfs.

### 6.1 Significant itemset mining

Our implementation of TOPKWY is based upon TOPKMINER (Pietracaprina and Vandin 2007), which mines top- $k$  frequent closed itemsets. As for TOPKMINER, TOPKWY uses a PatriciaTrie (Zandolin and Pietracaprina 2003) to store a compact representation of the dataset  $\mathcal{D}$  in which transactions sharing the same prefix are represented by the same node in the tree. The conditional dataset of (i.e., the set of transactions containing) an itemset  $Y$  is stored as a list  $m_Y$  of nodes of the PatriciaTrie. An additional counter is added to every node, representing how many transactions with prefix represented by the node belong to the class 1. The same is done for the  $j_p$  permutations adding  $j_p$  counters to each node in the trie. Since the PatriciaTrie is built adding one transaction at a time, a technique similar to reservoir sampling is used to generate the  $j_p$  permuted labels of every transaction. Let  $\mathbf{r}$  be a vector of length  $j_p$ ,

with all components  $r^{(j)}$ ,  $j = 1, \dots, j_p$  of  $\mathbf{r}$  initialized to  $n_1$ , the number of transactions to assign to the class  $c_1$  for every permutation of index  $j$ . For every transaction  $t_i$ , with  $i \in [1, n]$ , the  $j$ -th label of  $t_i$  is assigned to  $c_1$  with probability  $r^{(j)}/(n-i+1)$ ,  $c_0$  otherwise. If  $c_1$  is chosen, then  $r^{(j)}$  is decreased of one. This method guarantees that the total number of transactions with label  $c_1$  will be  $n_1$  for every  $j \in [1, j_p]$  and that the labels of the  $j$ -th permuted dataset are obtained by a random permutation of the class labels.

Our implementation of TOPKWY-dfs is based upon LCM 3 (Uno et al. 2005), which mines frequent closed itemsets using a depth first strategy. The third version of LCM combines various techniques and data structures to accelerate the generation and the computation of the frequencies of frequent closed itemsets.

## 6.2 Significant subgraph mining

Our implementation of TOPKWY-dfs relies on GASTON (Nijssen and Kok 2004) to mine significant subgraphs. GASTON first considers simple patterns, such as paths and trees, since efficient techniques for isomorphism checking are available for such acyclic structures. Only after this first phase, denoted as “quickstart”, general subgraphs, containing cycles, are evaluated. The search strategy of GASTON relies on a depth first enumeration of subgraphs. We do not provide a subgraph variant of TOPKWY because no competitive algorithms based on a best first exploration strategy are currently available (Wörlein et al. 2005; Nijssen and Kok 2006).

## 7 Experimental evaluation

We implemented and tested TOPKWY and TOPKWY-dfs for the extraction of significant itemsets and significant subgraphs. Our experimental evaluation has three goals. First, to assess the number of significant patterns found in real datasets. Second, to evaluate the performance of TOPKWY: since no other tool for the extraction of top- $k$  significant patterns exists, we compare TOPKWY and TOPKWY-dfs with the state-of-the-art tool for significant pattern mining, WYlight (Llinares-López et al. 2015). While the techniques introduced in this work can be extended to other multiple hypothesis testing procedures, such as LAMP (Terada et al. 2013a), we do not compare with LAMP or derived strategies (Minato et al. 2014) since (Llinares-López et al. 2015) shows that WY permutation testing results in higher power. Third, to assess the impact of our improved bounds and implementation choices on performances.

In Sect. 7.1 we describe the implementation and computational environment for our experiments. In Sect. 7.2 we describe the datasets we used. In Sect. 7.3 we describe the experiments we have performed and our choice of parameters. Finally, in Sect. 7.4 we report and discuss the results of our experiments.

**Table 2** Itemset Datasets statistics. For each dataset the table reports: the number  $|\mathcal{D}|$  of transactions; the number  $|I|$  of items; the average transaction length  $avg$ ; the fraction  $n_1/n$  of transactions in the minority class; the number  $SP(0.05)$  of significant patterns for  $FWER = 0.05$

Dataset	$ \mathcal{D} $	$ I $	$avg$	$n_1/n$	$SP(0.05)$
svmguide3(L)	1,243	44	21.9	0.23	36,736
chess(U)	3,196	75	37	0.05	$> 10^7$
mushroom(L)	8,124	118	22	0.48	71,945
phishing(L)	11,055	813	43	0.44	$> 10^7$
breast cancer(L)	12,773	1,129	6.7	0.09	6
a9a(L)	32,561	247	13.9	0.24	348,611
pumb-star(U)	49,046	7117	50.5	0.44	$> 10^7$
bms-web1(U)	58,136	60,978	2.51	0.03	704,685
connect(U)	67,557	129	43	0.49	$> 10^8$
bms-web2(U)	77,158	330,285	4.59	0.04	289,012
retail(U)	88,162	16,470	10.3	0.47	3,071
ijcnn1(L)	91,701	44	13	0.10	607,373
T10I4D100K(U)	100,000	870	10.1	0.08	3,819
T40I10D100K(U)	100,000	942	39.6	0.28	5,986,439
codrna(L)	271,617	16	8	0.33	4,088
accidents(U)	340,183	467	33.8	0.49	$> 10^7$
bms-pos(U)	515,597	1,656	6.5	0.40	26,366,131
covtype(L)	581,012	64	11.9	0.49	542,365
susy(U)	5,000,000	190	43	0.48	$> 10^7$

## 7.1 Implementation and environment

We implemented TOPKWY in C++ as an extension of the TOPKMINER algorithm (Pietracaprina and Vandin 2007). For TOPKWY-dfs we modified the C implementation of WYlight (based on LCM (Uno et al. 2005) and GASTON (Nijssen and Kok 2004)) made available by the authors at <https://github.com/flinares/wylight>. All implementations were compiled with gcc 4.8.4. Our experiments have been performed on a 2.30 GHz Intel Xeon CPU machine with 512 GB of RAM, running on Ubuntu 14.04. Our code and scripts to replicate all experiments described in the paper are available at <https://github.com/VandinLab/TopKWY>.

## 7.2 Datasets

For itemsets mining, we performed our experiments using 19 datasets: the 10 largest ones used in Llinares-López et al. (2015) and available at FIMI'04<sup>2</sup> and UCI<sup>3</sup>, all the

<sup>2</sup> <http://fimi.ua.ac.be>.

<sup>3</sup> <https://archive.ics.uci.edu/ml/index.php>.

**Table 3** Subgraph Datasets statistics

Dataset	$ \mathcal{D} $	$V_{avg}$	$E_{avg}$	$n_1/n$	$SP(0.05)$
MUTAG	188	17.93	19.79	19.79	70,184
BZR(30)	405	35.75	38.36	0.21	80,425
COX2	467	41.22	43.45	0.22	$> 10^6$
ENZYMES(10)	600	32.63	62.14	0.17	112,158
DHFR(30)	753	42.43	44.54	0.61	$> 10^6$
DD	1,178	284.32	715.66	0.58	80,256
AIDS(30)	2,000	15.69	16.20	0.2	566,727
NCI1	4,110	29.87	32.30	0.5	$> 10^6$
NCI109	4,127	29.68	32.13	0.5	$> 10^6$
Mutagenicity	4,337	30.32	30.77	0.44	$> 10^6$
Tox_21_AHR(30)	8,169	18.09	18.50	0.12	98,398

For each dataset the table reports: the number  $|\mathcal{D}|$  of graphs; the average number of nodes  $|V_{avg}|$ ; the average number of edges  $|E_{avg}|$ ; the fraction  $n_1/n$  of graphs in the minority class; the number  $SP(0.05)$  of significant patterns for  $FWER = 0.05$ . The number in brackets report the maximum number of vertexes of the explored subgraphs, that has been limited to allow practical running times for the experiments

datasets used in Komiyama et al. (2017), available from the libSVM repository<sup>4</sup>, and 4 additional ones (a9a, bms-web1, accidents, susy) available from libSVM, FIMI'04, and SPMF<sup>5</sup>. The datasets' statistics are in Table 2. For each dataset, we also note if it already contained class labels ( $L$ ) or not ( $U$ ). For unlabeled datasets we simulated a typical analysis requiring to find itemsets correlated with a given item (feature) in a dataset. For every unlabeled dataset we selected the single item whose frequency is closer from below to 0.5, removed the corresponding item from every transaction, and use its appearance to define the target class label. The reported ratio  $n_1/n$  for the minority class of unlabeled datasets refers to the output of this labeling process. For real-valued features we obtained two bins by thresholding at the mean value and using one item for each bin [analogously to Komiyama et al. (2017)].

For subgraphs mining, we considered 11 of the largest datasets with binary target labels available from a repository<sup>6</sup> of benchmark datasets; most of them are also analysed by Llinares-López et al. (2015). The datasets' statistics are in Table 3. In some cases, we bounded the maximum number of vertexes of the subgraphs that are explored by GASTON, in order to obtain practical running times for the experiments. Such limits are reported in the parenthesis after the dataset's name in Table 3.

### 7.3 Parameters and experiments

For TOPKWY and TOPKWY-dfs we considered  $k = 10^i$  for  $i \in [1, 6]$ . For all the datasets we analyzed, we ran TOPKWY and TOPKWY-dfs, for all such values of  $k$ , and

<sup>4</sup> <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>.

<sup>5</sup> <http://www.philippe-fournier-viger.com/spmf>.

<sup>6</sup> <https://ls11-www.cs.tu-dortmund.de/staff/morris/graphkerneldatasets>.



WYlight. We fixed the number of permutations  $j_p = 10^4$ , shown to be a good choice in Llinares-López et al. (2015), and fixed the commonly used value  $\alpha = 0.05$  as FWER threshold. For the comparison between TOPKWY, TOPKWY-dfs, and WYlight, we repeated every experiment 10 times, recording the running time and peak memory provided by the operating system; we report the averages over the 10 runs, standard deviations are negligible and therefore not shown. The measures reported for TOPKWY include the time and space to retrieve statistically significant patterns and write them on file, while for TOPKWY-dfs and WYlight we only report the time and space needed to find the optimal significance threshold, which corresponds to the first step of the method, therefore reporting a lower bound to their runtimes. We stopped the execution of an algorithm if it did not conclude after (at least) one month of computation; for these cases, the indicated time and peak memory are lower bounds. For experiments testing the impact of parameters or implementation choices on TOPKWY we used only one execution.

## 7.4 Results

### 7.4.1 Itemsets mining

Table 2 reports the number of significant patterns for  $\alpha = 0.05$  in the datasets we considered, obtained by running TOPKWY (with  $k = +\infty$ ) or WYlight. For some datasets we stopped the computation after 1 month, so only a lower bound is available. In most cases, the number of significant patterns is extremely large: for 11 out of 19 datasets there are  $> 5 \times 10^5$  significant patterns and in 7 datasets there are  $> 10^7$  significant patterns. Therefore a direct way to limit the number of significant patterns in output, as provided by the top- $k$  significant patterns, is required.

Figure 1 compares the running time of TOPKWY and WYlight. Note that for the 11 datasets in which the number of significant patterns is  $< 10^6$ , TOPKWY with  $k = 10^6$  identifies all the significant patterns and produces the same patterns found with WYlight. For 15 out of 19 datasets, TOPKWY (with  $k = 10^6$ ) is faster than WYlight by a factor at least 2. For 9 datasets TOPKWY is faster than WYlight by at least one order of magnitude, and for 6 datasets WYlight requires  $> 11$  days while TOPKWY identifies up to the  $10^6$  most significant patterns within one day and the  $10^4$  most significant ones in few hours. Even for the datasets where TOPKWY identifies all significant patterns, producing the same patterns as WYlight, TOPKWY is always faster than WYlight, with up to one order of magnitude speed-up in some cases. This shows that TOPKWY is an effective tool to identify all significant patterns whenever possible and enables the analysis of significant patterns when their number is extremely high. For datasets in which the number of significant patterns is  $> 10^6$  we ran TOPKWY with  $k = \infty$  to compare its strategy for finding the corrected significance threshold for *all* significant patterns with the one used by WYlight. The runtime of TOPKWY is always lower than the runtime of WYlight by at least 20%, with a significant speed-up in some case (e.g., for chess, TOPKWY terminates in 2 days, while WYlight needs more than 10 days). These results show that TOPKWY outperform the state-of-the-art even for this task.

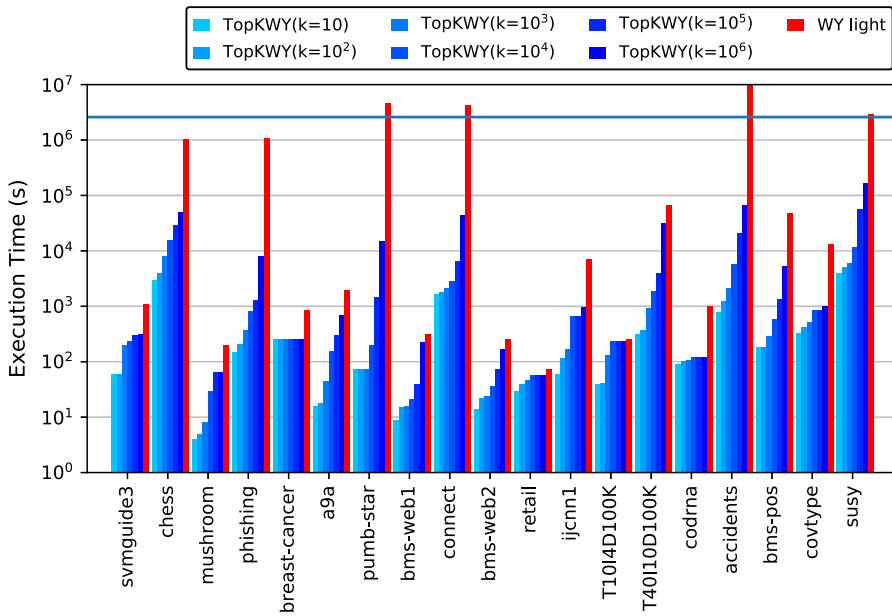


Fig. 1 Running time for TOPKWY (with various values of  $k$ ) and WYlight. The blue horizontal line corresponds to 1 month of computation

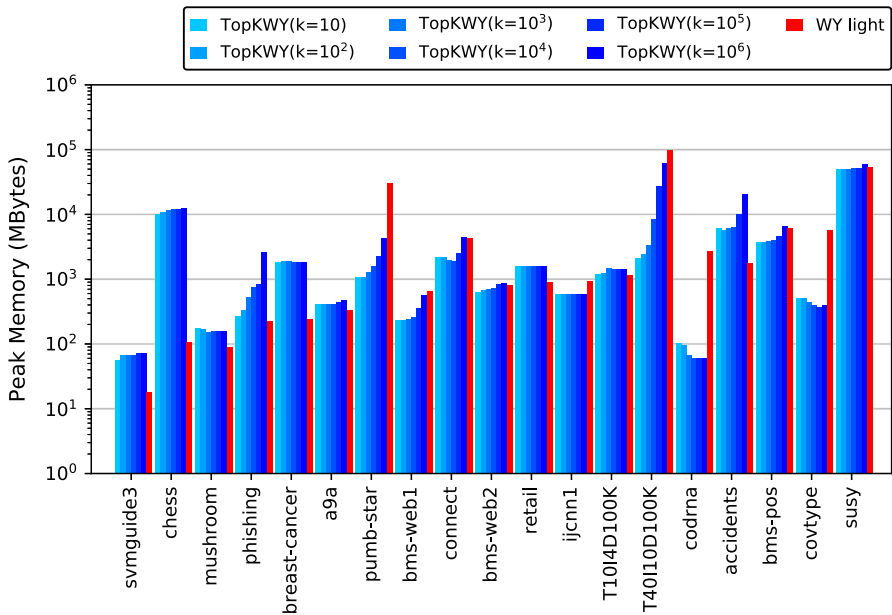


Fig. 2 Peak memory for TOPKWY (with various values of  $k$ ) and WYlight

Figure 2 compares the peak memory required by TOPKWY and WYlight. Given the *best first* strategy employed by TOPKWY, we expected its memory requirement could be higher than WYlight, that follows a depth first strategy. Interestingly, only in three cases TOPKWY required 1 order of magnitude more memory than WYlight and in both such cases the requirements are reasonable ( $\leq 20$  GBs) for current machines. However, in such cases WYlight required  $> 11$  days to complete, while TOPKWY terminated in  $< 1$  day, showing that, by using a reasonably larger amount of memory than WYlight, TOPKWY renders the identification of significant patterns feasible. In all other cases the memory requirement of TOPKWY is either the same or within few GBs of WYlight. For some datasets TOPKWY requires significantly less memory than WYlight: surprisingly this happens for datasets (cod-rna, covtypes) on which TOPKWY reports the same significant patterns as WYlight (i.e., *all* significant patterns). In some cases, memory usage decreases slightly when  $k$  increases, due to our dynamical allocation of the  $p$  values lookup table that may require less space when the minimum support decreases.

We investigated the impact of our implementation choices on the memory requirement of TOPKWY (Fig. 3). We compared the space required to store the permuted labels on all the nodes of the PatriciaTrie used by TOPKWY (see Sect. 6.1) with the space required by storing the permuted labels for each transaction (as done for example by WYlight). Since TOPKWY stores, for each node of the Patricia Trie, a list of  $j_p$  values (i.e., the number of transactions with minority label among the ones sharing the prefix corresponding to the node), one transaction may have more than  $j_p$  values associated to its nodes. In most cases the space required by the two methods is essentially the same, but in three cases the use of the Patricia Trie corresponds to a significant reduction in the memory used. In particular, these three cases are for datasets in which TOPKWY identifies all the significant patterns using less memory than WYlight, providing strong evidence of the importance of our encoding of the permuted class labels.

We compared the exploration strategies used by TOPKWY and by WYlight by recording the number of patterns they test (Fig. 4), restricting to datasets in which WYlight terminates. In all cases, TOPKWY tests a lower number of patterns than WYlight, with differences of almost two orders of magnitude for some datasets. This shows the effectiveness of our exploration strategy and of our novel bounds  $\psi'(\cdot)$  (see Sect. 4) on reducing the number of tests to perform.

We then directly investigated the impact of our novel bounds on the runtime of TOPKWY. We compared the running time of WYlight with the running time of two variants of TOPKWY: one using our improved bound  $\psi'(\cdot)$  and one using the LAMP bound  $\hat{\psi}(\cdot)$  (i.e., the same bound used by WYlight). The results for some representative datasets are in Fig. 5a. The results for the other datasets are similar. We observed that, for all datasets other than chess, the exploration strategy employed by TOPKWY to extract only the top- $k$  significant patterns already provides a substantial (up to more than one order of magnitude) improvement in the running time of TOPKWY with respect to WYlight, even using the same LAMP bounds. When our novel bound  $\psi'(\cdot)$  is used in TOPKWY we observe additional speed-ups, for a total up to more than two orders of magnitude. Therefore, the reduction in the number of patterns that need to

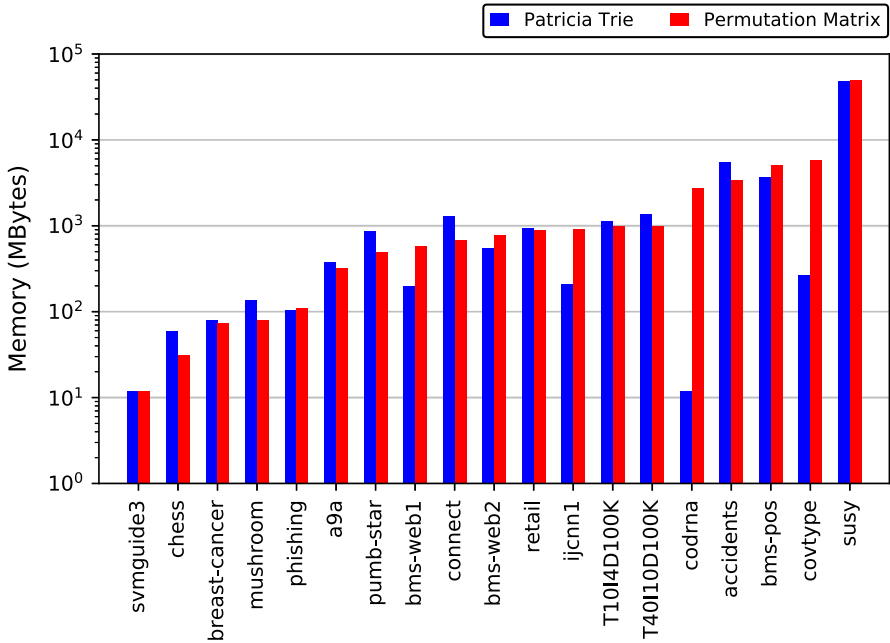


Fig. 3 Memory requirement for permuted class labels using PatriciaTrie and permutation matrix

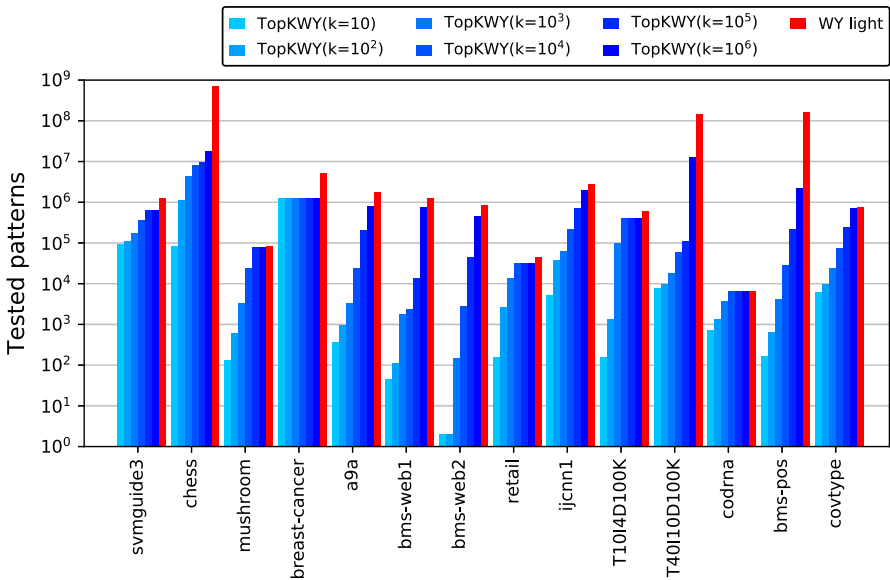
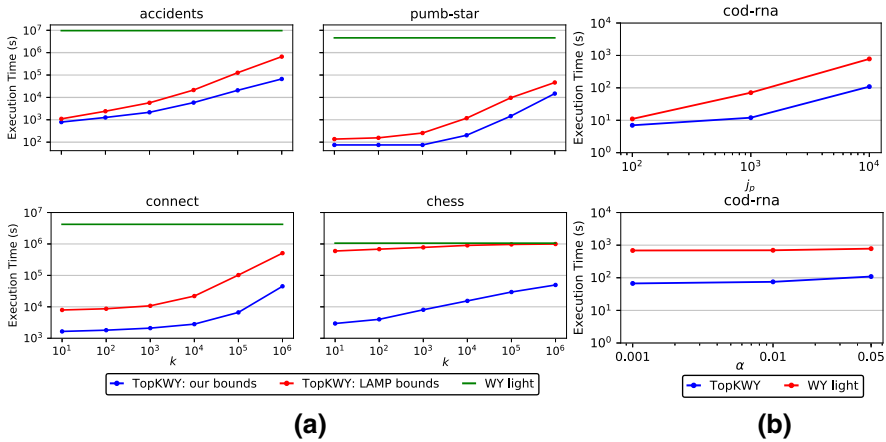


Fig. 4 Comparison between the number of tested patterns on the permuted datasets by TOPKWY and WYlight

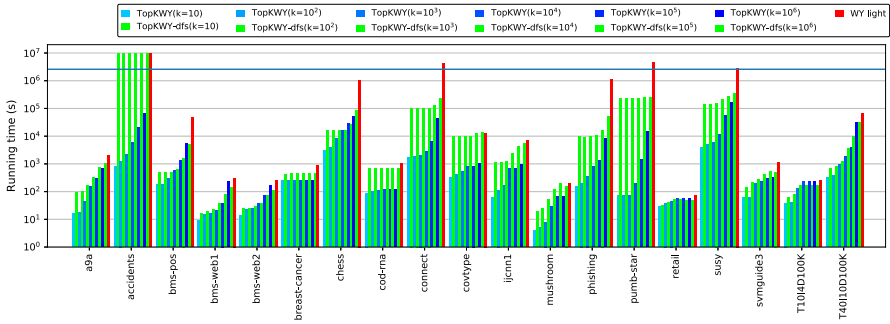


**Fig. 5** **a** Comparison between the running time of WYlight and the running time of TOPKWY using our improved bound  $\psi'(\cdot)$  and the LAMP bound  $\hat{\psi}(\cdot)$ . **b** Running time for different values of  $\alpha$  and  $j_p$

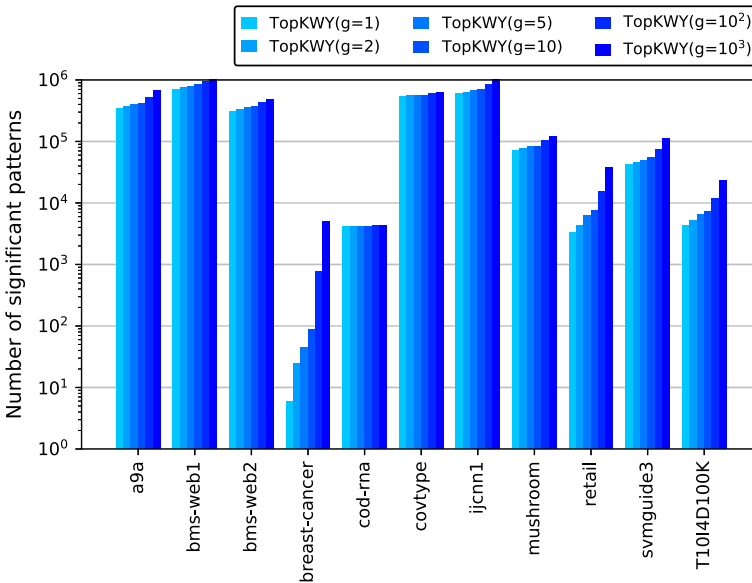
be tested on the permuted datasets, obtained by the exploration strategy of TOPKWY and our improved bound, is a crucial component for the performance of TOPKWY.

Finally, we assessed the impact of  $\alpha$  and  $j_p$  on the running time of TOPKWY and WYlight on two representative datasets, cod-rna and accidents, which are representative for the two scenarios of a small number of significant patterns (cod-rna) and of a large number of significant patterns (accidents). In these experiments we fixed  $k = 10^4$ . Figure 5b reports the results for cod-rna. Results for accidents are not reported since the running time of accidents remained essentially the same for all values of  $\alpha$  and  $j_p$ . This means that for accidents using the bounds introduced in Sect. 4 the computational effort is dominated by the pattern space exploration (and not the evaluation of the permuted datasets): considering only the top- $k$  significant patterns is therefore crucial to analyze such dataset. For cod-rna, we observe that varying  $\alpha$  has some but small impact on the runtime of both methods while there is a linear dependence of the running time of WYlight on  $j_p$  and a similar but less pronounced dependence of TOPKWY. In all cases, TOPKWY is faster than WYlight (for accidents WYlight does not terminate within 1 month) showing the efficiency of TOPKWY for different ranges of the  $\alpha$  and  $j_p$  parameters.

*Comparison between Best First and Depth First strategies.* We investigate the impact of the best first strategy adopted by TOPKWY on the computational performances of the mining tasks. To do so, we compared TOPKWY with TOPKWY-dfs, the variant of TOPKWY, which explores patterns in depth first order (see Sect. 5.3). We ran TOPKWY-dfs on the same set of experiments described in Sect. 7.3. Fig. 6 shows the running times of TOPKWY, TOPKWY-dfs, and WYlight. We can clearly see that, for 9 datasets out of 19, there is a significant difference in the running times of TOPKWY and TOPKWY-dfs: this means that, in particular for smaller values of  $k$ , the exploration strategy is a critical component of TOPKWY. For accidents, one of the most challenging dataset to analyze, both TOPKWY-dfs and WYlight can not



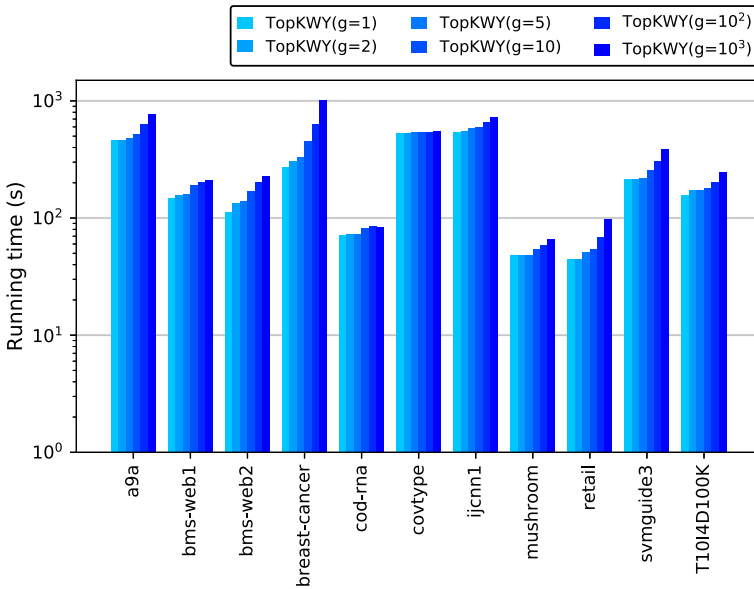
**Fig. 6** Running time for TOPKWY, TOPKWY-dfs (with various values of  $k$ ) and WYlight. The blue horizontal lines corresponds to 1 month of computation



**Fig. 7** Number of results found using TOPKWY when controlling the  $g$ -FWER, for various values of  $g$ ,  $j_p = 10^4$ ,  $k = 10^6$ , and  $\alpha = 0.05$

complete their execution in less than one month, even for  $k = 10$ . Therefore, for such dataset the best first strategy adopted by TOPKWY is crucial.

*Results for  $g$ -FWER.* We investigate the increase in statistical power of TOPKWY when controlling the generalized-Family-Wise Error Rate ( $g$ -FWER) by analyzing datasets described in Sect. 7.3 having less than  $10^6$  results for  $\alpha = 0.05$  when controlling the FWER. As we can see in Fig. 7, for the breast-cancer dataset the number of significant patterns increased by more then two orders of magnitude as the value of  $g$  increases (i.e., when more false positives are allowed). Figure 8 show the running time of TOPKWY when controlling the  $g$ -FWER at different values of  $g$ . As expected, the required time slightly increases but it stays practical for all datasets. We do not compare with



**Fig. 8** Running time for TOPKWY when controlling the  $g$ -FWER, for various values of  $g$ ,  $j_p = 10^4$ ,  $k = 10^6$ , and  $\alpha = 0.05$

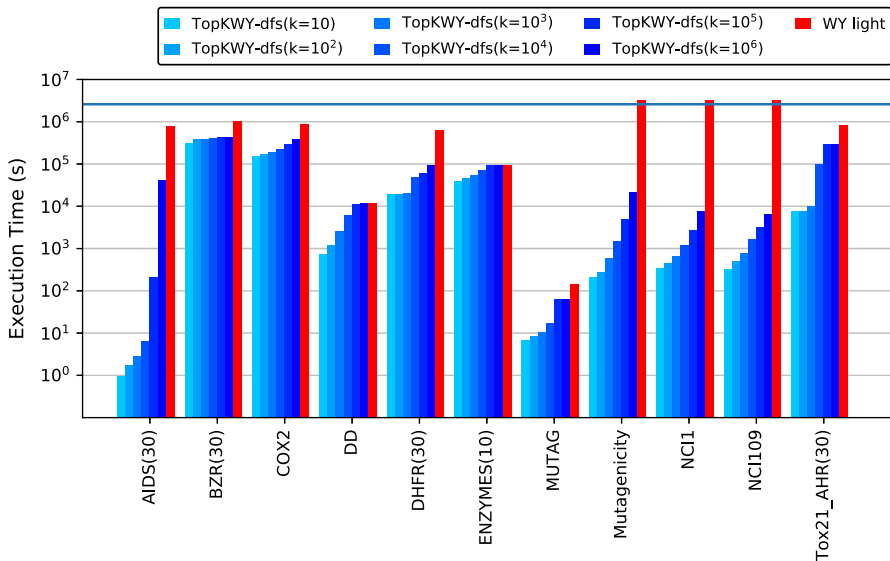
**Table 4 a** Values of  $g^*$  (second row of the Table) computed by TOPKWY\* for different values of  $k$  (first row of the Table) on breast-cancer dataset with  $j_p = 10^4$  and  $\alpha = 0.05$ . **b** Number of results  $|\mathcal{P}|$  (second row of the Table) found by TOPKWY- $\zeta$  on breast-cancer dataset with  $j_p = 10^4$ ,  $\alpha = 0.05$ , and  $k = 10^4$ , for different values of  $\zeta$  (first row of the Table)

<b>a)</b>				
$k$	10	$10^2$	$10^3$	$10^4$
$g^*$	2	11	163	2396
<b>b)</b>				
$\zeta$	0.01	0.05	0.1	0.25
$ \mathcal{P} $	0	24	624	$10^4$

other methods since TOPKWY is the first algorithm to discover significant patterns with a rigorous control on the  $g$ -FWER.

In Table 4.a we show the computed values of  $g^*$  by TOPKWY\* (see Sect. 5.1 for its definition) on the breast-cancer dataset for  $k \in \{10, 10^2, 10^3, 10^4\}$ . We can see that TOPKWY\* is able to provide informative estimates of the quality of the reported set of  $k$  most significant patterns, in terms of the minimum  $g$  such that the  $g$ -FWER( $p^k$ ) is  $\leq \alpha$ , without the need of fixing  $g$  a-priori. In Table 4.b we show the number of results found using TOPKWY- $\zeta$  for  $k = 10^4$  on the breast-cancer dataset, varying  $\zeta \in \{0.01, 0.05, 0.1, 0.25\}$ . From these results we can see that TOPKWY- $\zeta$  is a very flexible tool to discover significant patterns with bounds on both the output size and the





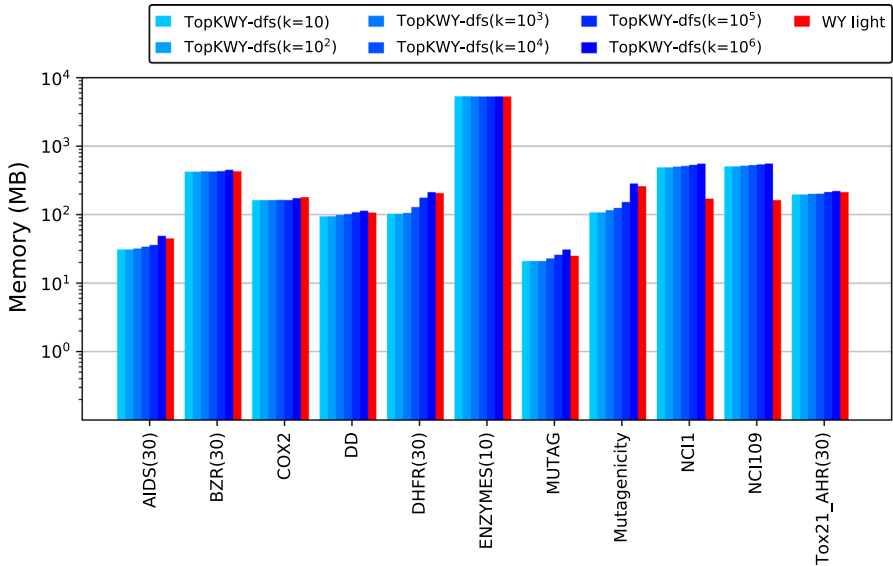
**Fig. 9** Running time for TOPKWAY-dfs (with various values of  $k$ ) and WYlight for significant subgraph mining. The blue horizontal line corresponds to 1 month of computation

maximum ratio of false discoveries, providing improved statistical power in situations where the number of significant patterns when controlling the  $FWER$  is very low. (We do not show the running times for TOPKWAY\* and TOPKWAY- $\zeta$  since those are very similar to the ones reported in Fig. 8.)

#### 7.4.2 Subgraphs mining

We ran TOPKWAY-dfs on the datasets described in Sect. 7.2. Table 3 reports the number of significant patterns for  $\alpha = 0.05$  in the datasets we considered, obtained by running TOPKWAY-dfs (with  $k = +\infty$ ) or WYlight. For some datasets we stopped the computation after 1 month, so only a lower bound is available. In most cases, the number of significant patterns is extremely large: for 6 out of 11 datasets there are  $> 5 \times 10^5$  significant patterns and in 5 datasets there are  $> 10^6$  significant patterns. This shows that a direct way to limit the number of significant patterns in output is required for subgraphs mining as well.

We then compared the running time and memory requirement of TOPKWAY-dfs and of WYlight. Figure 9 compares the running times of TOPKWAY-dfs and WYlight. As for itemsets mining, when the number of significant patterns is lower than  $k$ , TOPKWAY-dfs finds all of them, obtaining the same output as WYlight. We can see that TOPKWAY-dfs is, in all cases, faster than WYlight: for 9 datasets out of 11 and for  $k = 10^6$ , TOPKWAY-dfs improves the running time by a factor at least 2. It is interesting to note that for 5 datasets the number of significant results is  $< 10^6$ , therefore TOPKWAY-dfs is faster even if its output is the same of WYlight. For 3 datasets the running time is reduced by more than two orders of magnitude, and WYlight is not able to terminate in less than 1 month. We can observe that these three datasets contains more



**Fig. 10** Memory usage for TOPKWY-dfs (with various values of  $k$ ) and WYlight for significant subgraph mining

than  $10^6$  significant results; this clearly shows that focusing on the most significant patterns leads to significant computational advantages and enables the analysis of such datasets.

Figure 10 compares the memory usage of TOPKWY-dfs and WYlight. Both algorithms are very memory efficient and, while TOPKWY-dfs usually requires more memory than WYlight, the difference is small: for 7 of the 8 datasets where WYlight terminates, TOPKWY-dfs never requires more than 8% of the memory of WYlight, and 23% in the case of MUTAG, where the difference is of few MBs. (For the three datasets where WYlight does not terminate, we only report a lower bound to its memory usage.)

## 8 Conclusion

In this work we introduce TOPKWY, an efficient algorithm to identify the top- $k$  significant patterns with rigorous guarantees on the FWER and provide theoretical evidence of its effectiveness. Our extensive experimental evaluation shows that TOPKWY enables the identification of significant patterns on large datasets and that it significantly improves over the state-of-the-art.

Our notion of top- $k$  significant patterns and our algorithm TOPKWY could be relevant to other mining problems, for example statistical emerging pattern mining (Komiya et al. 2017), while providing a bound on the FWER. While we focus on bounding the FWER, a different approach would be to bound the false discovery rate (FDR) (Benjamini and Hochberg 1995), that is the expected ratio of false discoveries among all reported patterns. Bounding the FDR is crucial in cases where the number of significant results obtained bounding the FWER is low. In addition, fully

processing extremely large datasets may not be feasible: the combination of the techniques we develop in this work with sampling (e.g., the recently developed techniques by Pellegrina et al. (2019b) to extend TOPKWY for the analysis of random samples) is a promising direction that we will investigate in future work.

**Acknowledgements** This work is supported, in part by the National Science Foundation grant IIS-1247581 ([https://www.nsf.gov/awardsearch/showAward?AWD\\_ID=1247581](https://www.nsf.gov/awardsearch/showAward?AWD_ID=1247581)), by the University of Padova grants SID2017 and STARS: Algorithms for Inferential Data Mining, and by MIUR, the Italian Ministry of Education, University and Research, under PRIN Project n. 20174LF3T8 AHeAD (Efficient Algorithms for HARnessing Networked Data).

## References

- Agrawal R, Imieliński T, Swami A (1993) Mining association rules between sets of items in large databases. *SIGMOD Rec* 22:207–216
- Agrawal R, Srikant R (1994) Fast algorithms for mining association rules in large databases. In: Proceedings of the 20th international conference on very large data bases (VLDB '94), San Francisco, CA, USA. Morgan Kaufmann Publishers Inc, pp 487–499
- Atzmueller M (2015) Subgroup discovery. *Wiley Interdiscip Rev Data Min Knowl Discov* 5(1):35–49
- Bayardo RJ Jr (1998) Efficiently mining long patterns from databases. *ACM Sigmod Rec* 27(2):85–93
- Benjamini Y, Hochberg Y (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J Ro Stat Soc Ser B (Methodol)* 57:289–300
- Bonferroni C (1936) Teoria statistica delle classi e calcolo delle probabilita. *Pubblicazioni del R Istituto Superiore di Scienze Economiche e Commerciali di Firenze* 8:3–62
- Dong G, Bailey J (2012) Contrast data mining: concepts, algorithms, and applications. CRC Press, Boca Raton
- Duivesteyn W, Knobbe A (2011) Exploiting false discoveries—statistical validation of patterns and quality measures in subgroup discovery. In: 2011 IEEE 11th international conference on data mining. IEEE, pp 151–160
- Fisher RA (1922) On the interpretation of  $\chi^2$  from contingency tables, and the calculation of p. *J R Stat Soc* 85(1):87–94
- Gionis A, Mannila H, Mielikäinen T, Tsaparas P (2007) Assessing data mining results via swap randomization. *ACM Trans Knowl Discov Data (TKDD)* 1(3):14
- Hämäläinen W (2012) Kingfisher: an efficient algorithm for searching for both positive and negative dependency rules with statistical significance measures. *Knowl Inf Syst* 32(2):383–414
- Hämäläinen W, Webb GI (2019) A tutorial on statistically sound pattern discovery. *Data Min Knowl Disc* 33(2):325–377
- Han J, Pei J, Yin Y (2000) Mining frequent patterns without candidate generation. In: Chen W, Naughton JF, Bernstein PA (eds) SIGMOD conference. ACM, New York, pp 1–12
- Han J, Wang J, Lu Y, Tzvetkov P (2002) Mining top-k frequent closed patterns without minimum support. In: Proceedings 2002 IEEE international conference on data mining, 2002. ICDM 2003. IEEE, pp 211–218
- Han J, Cheng H, Xin D, Yan X (2007) Frequent pattern mining: current status and future directions. *Data Mining Knowl Discov* 15:55–86
- Herrera F, Carmona CJ, González P, Del Jesus MJ (2011) An overview on subgroup discovery: foundations and applications. *Knowl Inf Syst* 29(3):495–525
- Komiyama J, Ishihata M, Arimura H, Nishibayashi T, Minato S-I (2017) Statistical emerging pattern mining with multiple testing correction. In: Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining. ACM, pp 897–906
- Lehmann EL, Romano JP (2012) Generalizations of the familywise error rate. In: Selected works of EL Lehmann. Springer, pp 719–735
- Li J, Liu J, Toivonen H, Satou K, Sun Y, Sun B (2014) Discovering statistically non-redundant subgroups. *Knowl-Based Syst* 67:315–327

- Llinares-López F, Sugiyama M, Papaxanthos L, Borgwardt K (2015) Fast and memory-efficient significant pattern mining via permutation testing. In: Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining. ACM, pp 725–734
- Minato S-i, Uno T, Tsuda K, Terada A, Sese J (2014) A fast method of statistical assessment for combinatorial hypotheses based on frequent itemset enumeration. In: Joint European conference on machine learning and knowledge discovery in databases. Springer, pp 422–436
- Nijssen S, Kok JN (2004) A quickstart in frequent structure mining can make a difference. In: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining. ACM, pp 647–652
- Nijssen S, Kok JN (2006) Frequent subgraph miners: runtimes don't say everything. In: MLG 2006, p 173
- Papaxanthos L, Llinares-López F, Bodenham D, Borgwardt K (2016) Finding significant combinations of features in the presence of categorical covariates. In: Advances in neural information processing systems, pp 2279–2287
- Pasquier N, Bastide Y, Taouil R, Lakhal L (1999) Discovering frequent closed itemsets for association rules. In: International conference on database theory. Springer, pp 398–416
- Pietracaprina A, Vandin F (2007) Efficient incremental mining of top-K frequent closed itemsets. In: Discovery science, volume 4755 of lecture notes in computer science. Springer, Berlin Heidelberg, pp 275–280
- Pellegrina L, Vandin F (2018) Efficient mining of the most significant patterns with permutation testing. In: Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining. ACM, pp 2070–2079
- Pellegrina L, Riondato M, Vandin F (2019a) Hypothesis testing and statistically-sound pattern mining. In: Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining. ACM, pp 3215–3216
- Pellegrina L, Riondato M, Vandin F (2019b) Spumante: Significant pattern mining with unconditional testing. In: Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining-KDD, vol 19
- Tarone R (1990) A modified bonferroni method for discrete data. *Biometrics* 515–522
- Terada A, Okada-Hatakeyama M, Tsuda K, Sese J (2013a) Statistical significance of combinatorial regulations. *Proc Nat Acad Sci* 110(32):12996–13001
- Terada A, Tsuda K, Sese J (2013b) Fast westfall-young permutation procedure for combinatorial regulation discovery. In: 2013 IEEE international conference on bioinformatics and biomedicine (BIBM). IEEE, pp 153–158
- Terada A, Kim H, Sese J (2015) High-speed westfall-young permutation procedure for genome-wide association studies. In: Proceedings of the 6th ACM conference on bioinformatics, computational biology and health informatics. ACM, pp 17–26
- Terada A, Tsuda K et al (2016) Significant pattern mining with confounding variables. In: Pacific-Asia conference on knowledge discovery and data mining. Springer, pp 277–289
- Uno T, Kiyomi M, Arimura H (2005) Lcm ver. 3: collaboration of array, bitmap and prefix tree for frequent itemset mining. In: Proceedings of the 1st international workshop on open source data mining: frequent pattern mining implementations. ACM, pp 77–86
- van der Laan MJ, Dudoit S, Pollard KS (2004) Augmentation procedures for control of the generalized family-wise error rate and tail probabilities for the proportion of false positives. *Stat Appl Genet Mol Biol* 3(1):1–25
- Webb GI (2006) Discovering significant rules. In: Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining. ACM, pp 434–443
- Webb GI (2007) Discovering significant patterns. *Mach Learn* 68(1):1–33
- Webb GI (2008) Layered critical values: a powerful direct-adjustment approach to discovering significant patterns. *Mach Learn* 71(2–3):307–323
- Westfall PH, Young SS (1993) Resampling-based multiple testing: examples and methods for p-value adjustment. Wiley Series in Probability and Statistics, Hoboken
- Wörlein M, Meinel T, Fischer I, Philippsen M (2005) A quantitative comparison of the subgraph miners mofa, gspan, fsm, and gaston. In: European conference on principles of data mining and knowledge discovery. Springer, pp 392–403
- Zandolin D, Pietracaprina A (2003) Mining frequent itemsets using patricia tries. In: Proceedings of FIMI03, vol 90

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.