

SPIRIT: Sequential Pattern Mining with Regular Expression Constraints

Minos N. Garofalakis

Bell Laboratories
minos@bell-labs.com

Rajeev Rastogi

Bell Laboratories
rastogi@bell-labs.com

Kyuseok Shim

Bell Laboratories
shim@bell-labs.com

Abstract

Discovering sequential patterns is an important problem in data mining with a host of application domains including medicine, telecommunications, and the World Wide Web. Conventional mining systems provide users with only a very restricted mechanism (based on minimum support) for specifying patterns of interest. In this paper, we propose the use of Regular Expressions (REs) as a flexible constraint specification tool that enables user-controlled focus to be incorporated into the pattern mining process. We develop a family of novel algorithms (termed SPIRIT – Sequential Pattern Mining with Regular expression constraints) for mining frequent sequential patterns that also satisfy user-specified RE constraints. The main distinguishing factor among the proposed schemes is the degree to which the RE constraints are enforced to prune the search space of patterns during computation. Our solutions provide valuable insights into the tradeoffs that arise when constraints that do not subscribe to nice properties (like anti-monotonicity) are integrated into the mining process. A quantitative exploration of these tradeoffs is conducted through an extensive experimental study on synthetic and real-life data sets.

1 Introduction

Discovering *sequential patterns* from a large database of sequences is an important problem in the field of knowledge discovery and data mining. Briefly, given a set of data sequences, the problem is to discover subsequences that are *frequent*, in the sense that the percentage of data sequences containing them exceeds a user-specified minimum *support* [3, 11]. Mining frequent sequential patterns has found a host of potential application domains, including retailing (i.e., market-basket data), telecommunications, and, more recently, the World Wide Web (WWW). In market-basket databases, each data sequence corresponds to items bought by an individual customer over time and frequent patterns can be useful for predicting future customer behavior. In telecommunications, frequent sequences of alarms output

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the VLDB copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment.

Proceedings of the 25th VLDB Conference,
Edinburgh, Scotland, 1999.

by network switches capture important relationships between alarm signals that can then be employed for online prediction, analysis, and correction of network faults. Finally, in the context of the WWW, server sites typically generate huge volumes of daily log data capturing the sequences of page accesses for thousands or millions of users¹. Discovering frequent access patterns in WWW logs can help improve system design (e.g., better hyperlinked structure between correlated pages) and lead to better marketing decisions (e.g., strategic advertisement placement).

As a more concrete example, the Yahoo! Internet directory (www.yahoo.com) enables users to locate interesting WWW documents by navigating through large *topic hierarchies* consisting of thousands of different document classes. These hierarchies provide an effective way of dealing with the abundance problem present in today's keyword-based WWW search engines. The idea is to allow users to progressively refine their search by following specific *topic paths* (i.e., sequences of hyperlinks) along a (predefined) hierarchy. Given the wide variety of topics and the inherently fuzzy nature of document classification, there are numerous cases in which distinct topic paths lead to different document collections on very similar topics. For example, starting from Yahoo!'s home page users can locate information on hotels in New York City by following *either* `Travel:Yahoo!Travel:North America:United States:New York:New York City:Lodging:Hotels` or `Travel:Lodging:Yahoo!Lodging:New York:New York Cities:New York City:Hotels and Motels`, where “:” denotes a parent-child link in the topic hierarchy. Mining user access logs to determine the most frequently accessed topic paths is a task of immense marketing value, e.g., for a hotel or restaurant business in New York City trying to select a strategic set of WWW locations for its advertising campaign.

The design of effective algorithms for mining frequent sequential patterns has been the subject of several studies in recent years [3, 4, 7, 8, 11, 12]. Ignoring small differences in the problem definition (e.g., form of input data, definition of a subsequence), a major common thread that runs through the vast majority of earlier work is the *lack of user-controlled focus in the pattern mining process*. Typically,

¹In general, WWW servers only have knowledge of the IP address of the user/proxy requesting a specific web page. However, *referrers* and *cookies* can be used to determine the sequence of accesses for a particular user (without compromising the user's identity).

the interaction of the user with the pattern mining system is limited to specifying a lower bound on the desired support for the extracted patterns. The system then executes an appropriate mining algorithm and returns a very large number of sequential patterns, only some of which may be of actual interest to the user. Despite its conceptual simplicity, this “unfocused” approach to sequential pattern mining suffers from two major drawbacks.

1. *Disproportionate computational cost for selective users.* Given a database of sequences and a fixed value for the minimum support threshold, the computational cost of the pattern mining process is fixed for any potential user. Ignoring user focus can be extremely unfair to a highly selective user that is only interested in patterns of a very specific form.
2. *Overwhelming volume of potentially useless results.* The lack of tools to express user focus during the pattern mining process means that selective users will typically be swamped with a huge number of frequent patterns, most of which are useless for their purposes.

The above discussion clearly demonstrates the need for novel pattern mining solutions that enable the incorporation of user-controlled focus in the mining process. There are two main components that any such solution must provide. First, given the inadequacy of simple support constraints, we need a *flexible constraint specification language* that allows users to express the specific family of sequential patterns that they are interested in. For instance, returning to our earlier “New York City hotels” example, a hotel planning its ad placement may only be interested in paths that (a) begin with `Travel`, (b) end in either `Hotels` or `Hotels and Motels`, and (c) contain at least one of `Lodging`, `Yahoo!Lodging`, `Yahoo!Travel`, `New York`, or `New York City`, since these are the only topics directly related to its line of business. Second, we need novel pattern mining algorithms that can exploit user focus by *pushing user-specified constraints deep inside the mining process*. The abstract goal here is to exploit pattern constraints to prune the computational cost and ensure system performance that is *commensurate* with the level of user focus (i.e., constraint selectivity).

We should note that even though recent work has addressed similar problems in the context of association rule mining [9, 10], the problem of incorporating a rich set of user-specified constraints in sequential pattern mining remains, to the best of our knowledge, unexplored. Furthermore, as we will discover later in the paper, pattern constraints raise a host of new issues specific to sequence mining (e.g., due to the explicit ordering of items) that were not considered in the subset and aggregation constraints for itemsets considered in [9, 10]. For example, our pattern constraints do not satisfy the property of *anti-monotonicity* [9]; that is, the fact that a sequence satisfies a pattern constraint does not imply that all its subsequences satisfy the same constraint. These differences mandate novel solutions that are completely independent of earlier

results on constrained association rule mining² [9, 10].

In this paper, we formulate the problem of mining sequential patterns with *regular expression constraints* and we develop novel, efficient algorithmic solutions for pushing regular expressions inside the pattern mining process. Our choice of regular expressions (REs) as a constraint specification tool is motivated by two important factors. First, REs provide a simple, natural syntax for the succinct specification of families of sequential patterns. Second, REs possess sufficient expressive power for specifying a wide range of interesting, non-trivial pattern constraints. These observations are validated by the extensive use of REs in everyday string processing tasks (e.g., UNIX shell utilities like `grep` or `ls`) as well as in recent proposals on query languages for sequence data (e.g., the Shape Definition Language of Agrawal et al. [1]). Returning once again to our “New York City hotels” example, note that the constraint on topic paths described earlier in this section can be simply expressed as the following RE: `Travel(Lodging|Yahoo!Lodging|Yahoo!Travel|New York|New York City)(Hotels|Hotels and Motels)`, where “|” stands for disjunction. We propose a family of novel algorithms (termed SPIRIT – Sequential Pattern Mining with Regular expression constraints) for mining frequent sequential patterns that also belong to the language defined by the user-specified RE. Our algorithms exploit the equivalence of REs to deterministic finite automata [6] to push RE constraints deep inside the pattern mining computation. The main distinguishing factor among the proposed schemes is the *degree* to which the RE constraint is enforced within the generation and pruning of candidate patterns during the mining process. We observe that, varying the level of user focus (i.e., RE enforcement) during pattern mining gives rise to certain interesting tradeoffs with respect to computational effectiveness. Enforcing the RE constraint at each phase of the mining process certainly minimizes the amount of “state” maintained after each phase, focusing only on patterns that could potentially be in the final answer set. On the other hand, minimizing this maintained state may not always be the best solution since it can, for example, limit our ability to do effective support-based pruning in later phases. Such tradeoffs are obviously related to our previous observation that RE constraints are *not* anti-monotone [9]. We believe that our results provide useful insights into the more general problem of constraint-driven, ad-hoc data mining, showing that there can be a whole spectrum of choices for dealing with constraints, even when they do not subscribe to nice properties like anti-monotonicity or succinctness [9]. An extensive experimental study with synthetic as well as real-life data sets is conducted to explore the tradeoffs involved and their impact on the overall effectiveness of our algorithms. Our results indicate that incorporating RE constraints into the pattern mining computation can some times yield more than an order of magnitude

²Due to space constraints, we omit a detailed discussion of earlier work. The interested reader is referred to the full version of this paper [5].

improvement in performance, thus validating the effectiveness of our approach. Our experimentation with real-life WWW server log data also demonstrates the versatility of REs as a user-level tool for focusing on interesting patterns. The work reported in this paper has been done in the context of the *SERENDIP* data mining project at Bell Laboratories (www.bell-labs.com/projects/serendip).

2 Problem Formulation

2.1 Definitions

The main input to our mining problem is a database of sequences, where each sequence is an ordered list of *elements*. These elements can be either (a) *simple items* from a fixed set of literals (e.g., the identifiers of WWW documents available at a server [4], the amino acid symbols used in protein analysis [12]), or (b) *itemsets*, that is, non-empty sets of items (e.g., books bought by a customer in the same transaction [11]). The list of elements of a data sequence s is denoted by $\langle s_1 s_2 \dots s_n \rangle$, where s_i is the i^{th} element of s . We use $|s|$ to denote the *length* (i.e., number of elements) of sequence s . A sequence of length k is referred to as a *k-sequence*. (We consider the terms “sequence” and “sequential pattern” to be equivalent for the remainder of our discussion.) Table 1 summarizes the notation used throughout the paper with a brief description of its semantics. Additional notation will be introduced when necessary.

| Symbol | Semantics |
|-----------------------------|--|
| s, t, u, \dots | Generic sequences in the input database |
| $\langle s t \rangle$ | Sequence resulting from the concatenation of sequences s and t |
| $ s $ | Length, i.e., number of elements, of sequence s |
| s_i | i^{th} element of sequence s |
| s_i^* | Zero or more occurrences of element s_i (Kleene closure operator) |
| $s_i s_j$ | Select one element out of s_i and s_j (disjunction operator) |
| \mathcal{R} | Regular expression (RE) constraint |
| $\mathcal{A}_{\mathcal{R}}$ | Deterministic finite automaton for RE \mathcal{R} |
| b, c, d, \dots | Generic states in automaton $\mathcal{A}_{\mathcal{R}}$ |
| a | Start state of automaton $\mathcal{A}_{\mathcal{R}}$ |
| $b \xrightarrow{s_i} c$ | Transition from state b to state c in $\mathcal{A}_{\mathcal{R}}$ on element s_i |
| $b \xrightarrow{s} c$ | Transition path from state b to state c in $\mathcal{A}_{\mathcal{R}}$ on the sequence of elements s |
| C_k | Set of candidate k -sequences |
| F_k | Set of frequent k -sequences |

Table 1: Notation.

Consider two data sequences $s = \langle s_1 s_2 \dots s_n \rangle$ and $t = \langle t_1 t_2 \dots t_m \rangle$. We say that s is a *subsequence* of t if s is a “projection” of t , derived by deleting elements and/or items from t . More formally, s is a subsequence of t if there exist integers $j_1 < j_2 < \dots < j_n$ such that $s_1 \subseteq t_{j_1}, s_2 \subseteq t_{j_2}, \dots, s_n \subseteq t_{j_n}$. Note that for sequences of simple items the above condition translates to

$s_1 = t_{j_1}, s_2 = t_{j_2}, \dots, s_n = t_{j_n}$. For example, sequences $\langle 1 3 \rangle$ and $\langle 1 2 4 \rangle$ are subsequences of $\langle 1 2 3 4 \rangle$, while $\langle 3 1 \rangle$ is not. Srikant and Agrawal [11] observe that, when mining market-basket sequential patterns, users often want to place a bound on the *maximum distance* between the occurrence of adjacent pattern elements in a data sequence. For example, if a customer buys bread today and milk after a couple of weeks then the two purchases should probably not be seen as being correlated. Following [11], we define sequence s to be a *subsequence with a maximum distance constraint* of δ , or alternately *δ -distance subsequence*, of t if there exist integers $j_1 < j_2 < \dots < j_n$ such that $s_1 \subseteq t_{j_1}, s_2 \subseteq t_{j_2}, \dots, s_n \subseteq t_{j_n}$ and $j_k - j_{k-1} \leq \delta$ for each $k = 2, 3, \dots, n$. That is, occurrences of adjacent elements of s within t are not separated by more than δ elements. As a special case of the above definition, we say that s is a *contiguous* subsequence of t if s is a 1-distance subsequence of t , i.e., the elements of s can be mapped to a contiguous segment of t .

A sequence s is said to *contain* a sequence p if p is a subsequence of s . We define the *support* of a pattern p as the fraction of sequences in the input database that contain p . Given a set of sequences \mathcal{S} , we say that $s \in \mathcal{S}$ is *maximal* if there are no sequences in $\mathcal{S} - \{s\}$ that contain it.

A RE constraint \mathcal{R} is specified as a RE over the alphabet of sequence elements using the established set of RE operators, such as disjunction ($|$) and Kleene closure ($*$) [6]. Thus, a RE constraint \mathcal{R} specifies a language of strings over the element alphabet or, equivalently, a regular family of sequential patterns that is of interest to the user. A well-known result from complexity theory states that REs have exactly the same expressive power as *deterministic finite automata* [6]. Thus, given any RE \mathcal{R} , we can always build a deterministic finite automaton $\mathcal{A}_{\mathcal{R}}$ such that $\mathcal{A}_{\mathcal{R}}$ accepts exactly the language generated by \mathcal{R} . Informally, a deterministic finite automaton is a finite state machine with (a) a well-defined *start* state (denoted by a) and one or more *accept* states, and (b) deterministic transitions across states on symbols of the input alphabet (in our case, sequence elements). A transition from state b to state c on element s_i is denoted by $b \xrightarrow{s_i} c$. We also use the shorthand $b \xrightarrow{s} c$ to denote the sequence of transitions on the elements of sequence s starting at state b and ending in state c . A sequence s is *accepted* by $\mathcal{A}_{\mathcal{R}}$ if following the sequence of transitions for the elements of s from the start state results in an accept state. Figure 1 depicts the state diagram of a deterministic finite automaton for the RE $1^* (2 2 | 2 3 4 | 4 4)$ (i.e., all sequences of zero or more 1’s followed by 2 2, 2 3 4, or 4 4). Following [6], we use double circles to indicate an accept state and $>$ to emphasize the start state (a) of the automaton. For brevity, we will simply use “automaton” as a synonym for “deterministic finite automaton” in the remainder of the paper.

2.2 Problem Statement

Given an input database of sequences, we define a sequential pattern to be *frequent* if its support in the database ex-

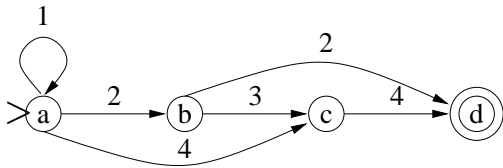


Figure 1: Automaton for the RE $1^* (2\ 2 \mid 2\ 3\ 4 \mid 4\ 4)$.

ceeds a user-specified minimum support threshold. Prior work has focused on efficient techniques for the discovery of frequent patterns, typically ignoring the possibility of allowing and exploiting flexible structural constraints during the mining process. In this paper, we develop novel, efficient algorithms for mining frequent sequential patterns in the presence of user-specified RE constraints. Due to space constraints, the discussion in this paper focuses on the case of *sequences of simple items with no maximum distance constraints*. The necessary extensions to handle itemset sequences and distance constraints for pattern occurrences are described in detail in the full version of this paper [5]. The following definitions establish some useful terminology for our discussion.

Definition 2.1 A sequence s is said to be *legal with respect to state b* of automaton $\mathcal{A}_{\mathcal{R}}$ if every state transition in $\mathcal{A}_{\mathcal{R}}$ is defined when following the sequence of transitions for the elements of s from b .

Definition 2.2 A sequence s is said to be *valid with respect to state b* of automaton $\mathcal{A}_{\mathcal{R}}$ if s is legal with respect to b and the final state of the transition path from b on input s is an accept state of $\mathcal{A}_{\mathcal{R}}$. We say that s is *valid* if s is valid with respect to the start state a of $\mathcal{A}_{\mathcal{R}}$ (or, equivalently, if s is accepted by $\mathcal{A}_{\mathcal{R}}$).

Example 2.1 : Consider the RE constraint $\mathcal{R} = 1^* (2\ 2 \mid 2\ 3\ 4 \mid 4\ 4)$ and the automaton $\mathcal{A}_{\mathcal{R}}$, shown in Figure 1. Sequence $\langle 1\ 2\ 3 \rangle$ is legal with respect to state a and sequence $\langle 3\ 4 \rangle$ is legal with respect to state b , while sequences $\langle 1\ 3\ 4 \rangle$ and $\langle 2\ 4 \rangle$ are not legal with respect to any state of $\mathcal{A}_{\mathcal{R}}$. Similarly, sequence $\langle 3\ 4 \rangle$ is valid with respect to state b (since $b \xrightarrow{3} c \xrightarrow{4} d$ and d is an accept state), however it is not valid, since it is not valid with respect to the start state a of $\mathcal{A}_{\mathcal{R}}$. Examples of valid sequences include $\langle 1\ 1\ 2\ 2 \rangle$ and $\langle 2\ 3\ 4 \rangle$. ■

Having established the necessary notions and terminology, we can now provide an abstract definition of our constrained pattern mining problem as follows.

- **Given:** A database of sequences \mathcal{D} , a user-specified minimum support threshold, and a user-specified RE constraint \mathcal{R} (or, equivalently, an automaton $\mathcal{A}_{\mathcal{R}}$).
- **Find:** All *frequent and valid* sequential patterns in \mathcal{D} .

Thus, our objective is to efficiently mine patterns that are not only frequent but also belong to the language of sequences generated by the RE \mathcal{R} ³. To this end, the next section introduces the SPIRIT family of mining algorithms for

³Our algorithms can readily handle a *set* of RE constraints by collapsing them into a single RE [6].

pushing user-specified RE constraints to varying degrees inside the pattern mining process.

3 Mining Frequent and Valid Sequences

3.1 Overview

Figure 2 depicts the basic algorithmic skeleton of the SPIRIT family, using an input parameter \mathcal{C} to denote a generic user-specified constraint on the mined patterns. The output of a SPIRIT algorithm is the set of frequent sequences in the database \mathcal{D} that satisfy constraint \mathcal{C} . At a high level, our algorithmic framework is similar in structure to the general Apriori strategy of Agrawal and Srikant [2]. Basically, SPIRIT algorithms work in passes, with each pass resulting in the discovery of longer patterns. In the k^{th} pass, a set of candidate (i.e., potentially frequent and valid) k -sequences C_k is generated and pruned using information from earlier passes. A scan over the data is then made, during which the support for each candidate sequence in C_k is counted and F_k is populated with the frequent k -sequences in C_k . There are, however, two crucial differences between the SPIRIT framework and conventional Apriori-type schemes (like GSP [11]) or the Constrained Apriori (CAP) algorithm [9] for mining associations with anti-monotone and/or succinct constraints.

1. *Relaxing \mathcal{C} by inducing a weaker (i.e., less restrictive) constraint \mathcal{C}' (Step 1).* Intuitively, constraint \mathcal{C}' is weaker than \mathcal{C} if every sequence that satisfies \mathcal{C} also satisfies \mathcal{C}' . The “strength” of \mathcal{C}' (i.e., how closely it emulates \mathcal{C}) essentially determines the degree to which the user-specified constraint \mathcal{C} is pushed inside the pattern mining computation. The choice of \mathcal{C}' differentiates among the members of the SPIRIT family and leads to interesting tradeoffs that are discussed in detail later in this section.
2. *Using the relaxed constraint \mathcal{C}' in the candidate generation and candidate pruning phases of each pass.* SPIRIT algorithms maintain the set F of frequent sequences (up to a given length) that satisfy the relaxed constraint \mathcal{C}' . Both F and \mathcal{C}' are used in:
 - (a) the candidate generation phase of pass k (Step 6), to produce an initial set of candidate k -sequences C_k that satisfy \mathcal{C}' by appropriately extending or combining sequences in F ; and,
 - (b) the candidate pruning phase of pass k (Steps 8-9), to delete from C_k all candidate k -sequences containing at least one subsequence that satisfies \mathcal{C}' and does not appear in F .

Thus, a SPIRIT algorithm maintains the following *invariant*: at the end of pass k , F_k is exactly the set of all frequent k -sequences that satisfy the constraint \mathcal{C}' . Note that incorporating \mathcal{C}' in candidate generation and pruning also impacts the terminating condition for the **repeat** loop in Step 15. Finally, since at the end of the loop, F contains frequent patterns satisfying the induced relaxed constraint \mathcal{C}' , an additional filtering step may be required (Step 17).

Procedure SPIRIT(\mathcal{D}, \mathcal{C})**begin**

1. let $\mathcal{C}' :=$ a constraint *weaker* (i.e., less restrictive) than \mathcal{C}
 2. $F := F_1 :=$ frequent items in \mathcal{D} that satisfy \mathcal{C}'
 3. $k := 2$
 4. **repeat** {
 5. // candidate generation
 6. using \mathcal{C}' and F generate $C_k := \{$ potentially frequent k -sequences that satisfy $\mathcal{C}' \}$
 7. // candidate pruning
 8. let $P := \{s \in C_k : s$ has a subsequence t that satisfies \mathcal{C}' and $t \notin F\}$
 9. $C_k := C_k - P$
 10. // candidate counting
 11. scan \mathcal{D} counting support for candidate k -sequences in C_k
 12. $F_k :=$ frequent sequences in C_k
 13. $F := F \cup F_k$
 14. $k := k + 1$
 15. } **until** TerminatingCondition(F, \mathcal{C}') holds
 16. // enforce the original (stronger) constraint \mathcal{C}
 17. output sequences in F that satisfy \mathcal{C}
- end**

Figure 2: SPIRIT constrained pattern mining framework.

Given a set of candidate k -sequences C_k , counting support for the members of C_k (Step 11) can be performed efficiently by employing specialized search structures, like the *hash tree* [11], for organizing the candidates. The implementation details can be found in [11]. The candidate counting step is typically the most expensive step of the pattern mining process and its overhead is directly proportional to the size of C_k [11]. Thus, at an abstract level, the goal of an efficient pattern mining strategy is to employ the minimum support requirement and any additional user-specified constraints to restrict as much as possible the set of candidate k -sequences counted during pass k . The SPIRIT framework strives to achieve this goal by using two different types of pruning within each pass k .

- *Constraint-based pruning* using a relaxation \mathcal{C}' of the user-specified constraint \mathcal{C} ; that is, ensuring that all candidate k -sequences in C_k satisfy \mathcal{C}' . This is accomplished by appropriately employing \mathcal{C}' and F in the candidate generation phase (Step 6).
- *Support-based pruning*; that is, ensuring that all subsequences of a sequence s in C_k that satisfy \mathcal{C}' are present in the current set of discovered frequent sequences F (Steps 8-9). Note that, even though *all* subsequences of s must in fact be frequent, we can only check the minimum support constraint for subsequences that satisfy \mathcal{C}' , since only these are retained in F .

Intuitively, constraint-based pruning tries to restrict C_k by (partially) enforcing the input constraint \mathcal{C} , whereas support-based pruning tries to restrict C_k by checking the minimum support constraint for qualifying subsequences. Note that, given a set of candidates C_k and a relaxation \mathcal{C}' of \mathcal{C} , the amount of support-based pruning is maximized

when \mathcal{C}' is *anti-monotone* [9] (i.e., all subsequences of a sequence satisfying \mathcal{C}' are guaranteed to also satisfy \mathcal{C}'). This is because support information for *all* of the subsequences of a candidate sequence s in C_k can be used to prune it. However, when \mathcal{C}' is *not* anti-monotone, the amounts of constraint-based and support-based pruning achieved vary depending on the specific choice of \mathcal{C}' .

3.1.1 Pushing Non Anti-Monotone Constraints

Consider the general problem of mining all frequent sequences that satisfy a user-specified constraint \mathcal{C} . If \mathcal{C} is anti-monotone, then the most effective way of using \mathcal{C} to prune candidates is to push \mathcal{C} “all the way” inside the mining computation. In the context of the SPIRIT framework, this means using \mathcal{C} *as is* (rather than some relaxation of \mathcal{C}) in the pattern discovery loop. The optimality of this solution for anti-monotone \mathcal{C} stems from two observations. First, using \mathcal{C} clearly maximizes the amount of constraint-based pruning since the strongest possible constraint (i.e., \mathcal{C} itself) is employed. Second, since \mathcal{C} is anti-monotone, all subsequences of a frequent candidate k -sequence that survives constraint-based pruning are guaranteed to be in F (since they also satisfy \mathcal{C}). Thus, using the full strength of an anti-monotone constraint \mathcal{C} maximizes the effectiveness of constraint-based pruning as well as support-based pruning. Note that this is exactly the methodology used in the CAP algorithm [9] for anti-monotone itemset constraints. An additional benefit of using anti-monotone constraints is that they significantly simplify the candidate generation and candidate pruning tasks. More specifically, generating C_k is nothing but an appropriate “self-join” operation over F_{k-1} and determining the pruned set P (Step 8) is simplified by the fact that all subsequences of candidates are guaranteed to satisfy the constraint.

When \mathcal{C} is *not* anti-monotone, however, things are not that clear-cut. A simple solution, suggested by Ng et al. [9] for itemset constraints, is to take an anti-monotone relaxation of \mathcal{C} and use that relaxation for candidate pruning. Nevertheless, this simple approach may not always be feasible. For example, our RE constraints for sequences do not admit any non-trivial anti-monotone relaxations. In such cases, the degree to which the constraint \mathcal{C} is pushed inside the mining process (i.e., the strength of the (non anti-monotone) relaxation \mathcal{C}' used for pruning) impacts the effectiveness of both constraint-based pruning and support-based pruning in different ways. More specifically, while increasing the strength of \mathcal{C}' obviously increases the effectiveness of constraint-based pruning, it can also have a negative effect on support-based pruning. The reason is that, for any given sequence in C_k that survives constraint-based pruning, *the number of its subsequences that satisfy the stronger, non anti-monotone constraint \mathcal{C}' may decrease*. Again, note that only subsequences that satisfy \mathcal{C}' can be used for support-based pruning, since this is the only “state” maintained from previous passes (in F).

Pushing a non anti-monotone constraint \mathcal{C}' in the pattern discovery loop can also increase the computational com-

plexity of the candidate generation and pruning tasks. For candidate generation, the fact that \mathcal{C}' is not anti-monotone means that some (or, all) of a candidate’s subsequences may be absent from F . In some cases, a “brute-force” approach (based on just \mathcal{C}') may be required to generate an initial set of candidates C_k . For candidate pruning, computing the subsequences of a candidate that satisfy \mathcal{C}' may no longer be trivial, implying additional computational overhead. We should note, however, that candidate generation and pruning are inexpensive CPU-bound operations that typically constitute only a small fraction of the overall computational cost. This fact is also clearly demonstrated in our experimental results (Section 4). Thus, the major tradeoff that needs to be considered when choosing a specific \mathcal{C}' from among the spectrum of possible relaxations of \mathcal{C} is the extent to which that choice impacts the effectiveness of constraint-based and support-based pruning. The objective, of course, is to strike a reasonable balance between the two different types of pruning so as to minimize the number of candidates for which support is actually counted in each pass.

3.1.2 The SPIRIT Algorithms

The four SPIRIT algorithms for constrained pattern mining are points spanning the entire spectrum of relaxations for the user-specified RE constraint $\mathcal{C} \equiv \mathcal{R}$. Essentially, the four algorithms represent a natural progression, with each algorithm pushing a stronger relaxation of \mathcal{R} than its predecessor in the pattern mining loop⁴. The first SPIRIT algorithm, termed SPIRIT(N) (“N” for Naive), employs the weakest relaxation of \mathcal{R} – it only prunes candidate sequences containing elements that do not appear in \mathcal{R} . The second algorithm, termed SPIRIT(L) (“L” for Legal), requires every candidate sequence to be *legal* with respect to some state of $\mathcal{A}_{\mathcal{R}}$. The third algorithm, termed SPIRIT(V) (“V” for Valid), goes one step further by filtering out candidate sequences that are not *valid with respect to any state of $\mathcal{A}_{\mathcal{R}}$* . Finally, the SPIRIT(R) algorithm (“R” for Regular) essentially pushes \mathcal{R} “all the way” inside the mining process by counting support only for *valid* candidate sequences, i.e., sequences accepted by $\mathcal{A}_{\mathcal{R}}$. Table 2 summarizes the constraint choices for the four members of the SPIRIT family within the general framework depicted in Figure 2. Note that, of the four SPIRIT algorithms, SPIRIT(N) is the only one employing an anti-monotone (and, trivial) relaxation \mathcal{C}' . Also, note that the progressive increase in the strength of \mathcal{C}' implies a subset relationship between the frequent sequences determined for each pass k ; that is,

$$F_k^{SPIRIT(R)} \subseteq F_k^{SPIRIT(V)} \subseteq F_k^{SPIRIT(L)} \subseteq F_k^{SPIRIT(N)}.$$

The remainder of this section provides a detailed discussion of the candidate generation and candidate pruning

⁴The development of the SPIRIT algorithms is based on the equivalent automaton form $\mathcal{A}_{\mathcal{R}}$ of the user-specified RE constraint \mathcal{R} . Algorithms for constructing $\mathcal{A}_{\mathcal{R}}$ from \mathcal{R} can be found in the theory literature [6].

| Algorithm | Relaxed Constraint \mathcal{C}' ($\mathcal{C} \equiv \mathcal{R}$) |
|-----------|--|
| SPIRIT(N) | all elements appear in \mathcal{R} |
| SPIRIT(L) | legal wrt some state of $\mathcal{A}_{\mathcal{R}}$ |
| SPIRIT(V) | valid wrt some state of $\mathcal{A}_{\mathcal{R}}$ |
| SPIRIT(R) | valid, i.e., $\mathcal{C}' \equiv \mathcal{C} \equiv \mathcal{R}$ |

Table 2: The four SPIRIT algorithms.

phases for each of the SPIRIT algorithms. Appropriate terminating conditions (Step 15) are also presented. The quantitative study of the constraint-based vs. support-based pruning tradeoff for the SPIRIT algorithms is deferred until the presentation of our experimental results (Section 4).

3.2 The SPIRIT(N) Algorithm

SPIRIT(N) is a simple modification of the GSP algorithm [11] for mining sequential patterns. SPIRIT(N) simply requires that all elements of a candidate sequence s in C_k appear in the RE \mathcal{R} . This constraint is clearly anti-monotone, so candidate generation and pruning are performed exactly as in GSP [11].

Candidate Generation. For every pair of $(k - 1)$ -sequences s and t in F_{k-1} , if $s_{j+1} = t_j$ for all $1 \leq j \leq k - 2$, then $\langle s t_{k-1} \rangle$ is added to C_k . This is basically a self-join of F_{k-1} , the join attributes being the last $k - 2$ elements of the first sequence and the first $k - 2$ elements of the second.

Candidate Pruning. A candidate sequence s is pruned from C_k if at least one of its $(k - 1)$ -subsequences does not belong to F_{k-1} .

Terminating Condition. The set of frequent k -sequences, F_k , is empty.

3.3 The SPIRIT(L) Algorithm

SPIRIT(L) uses the automaton $\mathcal{A}_{\mathcal{R}}$ to prune from C_k candidate k -sequences that are not *legal* with respect to any state of $\mathcal{A}_{\mathcal{R}}$. In our description of SPIRIT(L), we use $F_k(b)$ to denote the set of frequent k -sequences that are legal with respect to state b of $\mathcal{A}_{\mathcal{R}}$.

Candidate Generation. For each state b in $\mathcal{A}_{\mathcal{R}}$, we add to C_k candidate k -sequences that are legal with respect to b and have the potential to be frequent.

Lemma 3.1: Consider a k -sequence s that is legal with respect to state b in $\mathcal{A}_{\mathcal{R}}$, where $b \xrightarrow{s_1} c$ is a transition in $\mathcal{A}_{\mathcal{R}}$. For s to be frequent, $\langle s_1 \cdots s_{k-1} \rangle$ must be in $F_{k-1}(b)$ and $\langle s_2 \cdots s_k \rangle$ must be in $F_{k-1}(c)$. ■

Thus, the candidate sequences for state b can be computed as follows. For every sequence s in $F_{k-1}(b)$, if $b \xrightarrow{s_1} c$ is a transition in $\mathcal{A}_{\mathcal{R}}$, then for every sequence t in $F_{k-1}(c)$ such that $s_{j+1} = t_j$ for all $1 \leq j \leq k - 2$, the candidate sequence $\langle s t_{k-1} \rangle$ is added to C_k . This is basically a join of $F_{k-1}(b)$ and $F_{k-1}(c)$, on the condition that the

$(k - 2)$ -length suffix of $s \in F_{k-1}(b)$ matches the $(k - 2)$ -length prefix of $t \in F_{k-1}(c)$ and $b \xrightarrow{s_1} c$ is a transition in $\mathcal{A}_{\mathcal{R}}$.

Candidate Pruning. Given a sequence s in C_k , the candidate generation step ensures that both its prefix and suffix are frequent. We also know that in order for s to be frequent, every subsequence of s must also be frequent. However, since we only count support for sequences that are legal with respect to some state of $\mathcal{A}_{\mathcal{R}}$, we can prune s from C_k only if we find a *legal* subsequence of s that is not frequent (i.e., not in F). The candidate pruning procedure computes the set of maximal subsequences of s with length less than k that are legal with respect to some state of automaton $\mathcal{A}_{\mathcal{R}}$. If any of these maximal subsequences is not contained in F , then s is deleted from C_k .

We now describe an algorithm for computing the maximal legal subsequences of a candidate sequence s . Let $\text{maxSeq}(b, s)$ denote the set of maximal subsequences of s that are legal with respect to state b of $\mathcal{A}_{\mathcal{R}}$. Then, if we let $t = \langle s_2 \cdots s_{|s|} \rangle$, a superset of $\text{maxSeq}(b, s)$ can be computed from $\text{maxSeq}(b, t)$ using the fact that: (a) $\text{maxSeq}(b, s) \subseteq \text{maxSeq}(b, t) \cup \{ \langle s_1 u \rangle : u \in \text{maxSeq}(c, t) \} \cup \{s_1\}$, if $b \xrightarrow{s_1} c$ is a transition in $\mathcal{A}_{\mathcal{R}}$; and, (b) $\text{maxSeq}(b, s) \subseteq \text{maxSeq}(b, t)$, otherwise. The intuition is that for a subsequence $v \in \text{maxSeq}(b, s)$, either v does not involve s_1 , in which case v is a maximal subsequence of t that is legal with respect to b , or $v_1 = s_1$ and $\langle v_2 \cdots v_{|v|} \rangle$ is a maximal subsequence of t with respect to state c . Based on the above observation, we propose a *dynamic programming* algorithm, termed **FINDMAXSUBSEQ**, for computing $\text{maxSeq}(b, s)$ for all states b of $\mathcal{A}_{\mathcal{R}}$ (Figure 3). Intuitively, **FINDMAXSUBSEQ** works by computing the set maxSeq for successively longer suffixes of the input sequence s , beginning with the suffix consisting of only the last element of s .

More specifically, given an input sequence s and two sets of states in $\mathcal{A}_{\mathcal{R}}$ ($Start$ and End), algorithm **FINDMAXSUBSEQ** returns the set of all maximal subsequences t of s such that (a) the length of t is less than $|s|$, and (b) t is legal with respect to a state b in $Start$ and if $b \xrightarrow{t} c$, then $c \in End$. In each iteration of the for loop spanning Steps 3–17, for each state b in $\mathcal{A}_{\mathcal{R}}$, maximal legal subsequences for the suffix $\langle s_l \cdots s_{|s|} \rangle$ are computed and stored in $\text{maxSeq}[b]$. At the start of the l^{th} iteration, $\text{maxSeq}[b]$ contains the maximal subsequences of $\langle s_{l+1} \cdots s_{|s|} \rangle$ that are both legal with respect to state b and result in a state in End . Thus, if a transition from b to c on element s_l is in $\mathcal{A}_{\mathcal{R}}$, then the maximal legal subsequences for b comprise those previously computed for $\langle s_{l+1} \cdots s_{|s|} \rangle$ and certain new sequences involving element s_l . These new sequences containing s_l are computed in the body of the for loop spanning Steps 5–9 and stored in $\text{tmpSeq}[b]$. A point to note is that, since we are only interested in maximal legal subsequences that result in a state in End , we add s_l to $\text{tmpSeq}[b]$ only if $c \in End$ (Step 7).

After the new maximal subsequences involving s_l are

Procedure FINDMAXSUBSEQ($Start, End, s$)

begin

1. **for each** state b in automaton $\mathcal{A}_{\mathcal{R}}$ **do**
2. $\text{maxSeq}[b] := \emptyset$
3. **for** $l := |s|$ **down to** 1 **do** {
4. **for each** state b in automaton $\mathcal{A}_{\mathcal{R}}$ **do** {
5. $\text{tmpSeq}[b] = \emptyset$
6. **if** (there exists a transition $b \xrightarrow{s_l} c$ in $\mathcal{A}_{\mathcal{R}}$) {
7. **if** ($c \in End$) $\text{tmpSeq}[b] := \{s_l\}$
8. $\text{tmpSeq}[b] := \text{tmpSeq}[b] \cup \{ \langle s_l t \rangle : t \in \text{maxSeq}[c] \}$
9. }
10. }
11. **for each** state b in automaton $\mathcal{A}_{\mathcal{R}}$ **do** {
12. $\text{maxSeq}[b] := \text{maxSeq}[b] \cup \text{tmpSeq}[b]$
13. **for each** sequence t in $\text{maxSeq}[b]$ **do**
14. **if** (there exists u in $\text{maxSeq}[b] - \{ \langle s_l \cdots s_{|s|} \rangle \}$ such that t is a subsequence of u)
15. delete t from $\text{maxSeq}[b]$
16. }
17. }
18. **return** $\bigcup_{b \in Start} \text{maxSeq}[b] - \{s\}$ (after deleting non-maximal sequences)

end

Figure 3: Algorithm for finding maximal subsequences.

stored in $\text{tmpSeq}[b]$ for every state b of $\mathcal{A}_{\mathcal{R}}$, they are added to $\text{maxSeq}[b]$, following which, non-maximal subsequences in $\text{maxSeq}[b]$ are deleted (Steps 11–16)⁵. Finally, after maximal legal subsequences for the entire sequence s have been computed for all the states of $\mathcal{A}_{\mathcal{R}}$, only those for states in $Start$ are returned (Step 18).

To recap, the candidate pruning procedure of SPIRIT(L) invokes **FINDMAXSUBSEQ** to determine all the maximal legal subsequences of each candidate s in C_k , and deletes s from C_k if any of these subsequences is not frequent. For SPIRIT(L), algorithm **FINDMAXSUBSEQ** is invoked with $Start$ and End both equal to the set of all states in $\mathcal{A}_{\mathcal{R}}$.

Terminating Condition. The set of frequent k -sequences that are legal with respect to the start state a of $\mathcal{A}_{\mathcal{R}}$ is empty; that is, $F_k(a)$ is empty.

Time Complexity. Consider the candidate pruning overhead for a candidate k -sequence s in C_k . Compared to the candidate pruning step of SPIRIT(N), which has a time complexity of $O(k)$ (to determine the k subsequences of s), the computational overhead of candidate pruning in SPIRIT(L) can be significantly higher. More specifically, the worst-case time complexity of computing the maximal legal subsequences of s using algorithm **FINDMAXSUBSEQ** can be shown to be $O(k^2 * |\mathcal{A}_{\mathcal{R}}| * |\text{maxSeq}(s)|)$, where $|\mathcal{A}_{\mathcal{R}}|$ is the number of states in $\mathcal{A}_{\mathcal{R}}$ and $|\text{maxSeq}(s)|$ is the number of maximal legal subsequences for s . To see this, note that the outermost **for** loop in Step 3 of **FINDMAXSUBSEQ** is executed k times. The time complexity of the first **for** loop in Step 4 is $O(|\mathcal{A}_{\mathcal{R}}| * |\text{maxSeq}(s)|)$,

⁵In Steps 13–15, we have to be careful not to consider $\langle s_l \cdots s_{|s|} \rangle$ to delete other sequences in $\text{maxSeq}[b]$ since we are interested in maximal sequences whose length is less than $|s|$.

while that of the second **for** loop in Step 11 is $O(k * |\mathcal{A}_{\mathcal{R}}| * |\maxSeq(s)|)$, since $\maxSeq[b]$ can be implemented as a trie, for which insertions, deletions, and subsequence checking for k -sequences can all be carried out in $O(k)$ time.

We must point out that the higher time complexity of candidate pruning in SPIRIT(L) is not a major efficiency concern since (a) the overhead of candidate generation and pruning is typically a tiny fraction of the cost of counting supports for candidates in C_k , and (b) in practice, $|\maxSeq(s)|$ can be expected to be small for most sequences. In the worst case, however, for a k -sequence, $|\maxSeq(s)|$ can be $O(2^k)$. This worst case scenario can be avoided by imposing an a-priori limit on the size of $\maxSeq[b]$ in FINDMAXSUBSEQ and using appropriate heuristics for selecting *victims* (to be ejected from $\maxSeq[b]$) when its size exceeds that limit.

Space Overhead. SPIRIT(N) only utilizes F_{k-1} for the candidate generation and pruning phases during the k^{th} pass. In contrast, the candidate pruning step of SPIRIT(L) requires F to be stored in main memory since the maximal legal subsequences of a candidate k -sequence may be of any length less than k . However, this should not pose a serious problem since each F_k computed by SPIRIT(L) contains only frequent and legal k -sequences, which are typically few compared to all frequent k -sequences. In addition, powerful servers with several gigabytes of memory are now fairly commonplace. Thus, in most cases, it should be possible to accommodate all the sequences in F in main memory. In the occasional event that F does not fit in memory, one option would be to only store F_{k-l}, \dots, F_{k-1} for some $l \geq 1$. Of course, this means that maximal subsequences whose length is less than $k - l$ cannot be used to prune candidates from C_k during the candidate pruning step.

3.4 The SPIRIT(V) Algorithm

SPIRIT(V) uses a stronger relaxed constraint \mathcal{C}' than SPIRIT(L) during candidate generation and pruning. More specifically, SPIRIT(V) requires every candidate sequence to be *valid* with respect to some state of $\mathcal{A}_{\mathcal{R}}$ ⁶. In our description of SPIRIT(V), we use $F_k(b)$ to denote the set of frequent k -sequences that are valid with respect to state b of $\mathcal{A}_{\mathcal{R}}$.

Candidate Generation. Since every candidate sequence s in C_k is required to be valid with respect to some state b , it must be the case that the $(k - 1)$ -length suffix of s is both frequent and valid with respect to state c , where $b \xrightarrow{s_1} c$ is a transition in $\mathcal{A}_{\mathcal{R}}$. Thus, given a state b of $\mathcal{A}_{\mathcal{R}}$, the set of potentially frequent and valid k -sequences with respect to b can be generated using the following rule: for every transition $b \xrightarrow{s_i} c$, for every sequence t in $F_{k-1}(c)$, add $\langle s_i t \rangle$ to the set of candidates for state b . The set C_k is

⁶Note that an alternative approach would be to require candidates to be legal with respect to the start state of $\mathcal{A}_{\mathcal{R}}$. This approach is essentially symmetric to SPIRIT(V) and is not explored further in this paper.

simply the union of these candidate sets over all states b of $\mathcal{A}_{\mathcal{R}}$.

Candidate Pruning. The pruning phase of SPIRIT(V) is very similar to that of SPIRIT(L), except that only valid (rather than legal) subsequences of a candidate can be used for pruning. More specifically, given a candidate sequence s in C_k , we compute all maximal subsequences of s that are valid with respect to some state of $\mathcal{A}_{\mathcal{R}}$ and have length less than k . This is done by invoking algorithm FINDMAXSUBSEQ with *Start* equal to the set of all states of $\mathcal{A}_{\mathcal{R}}$ and *End* equal to the set of all *accept* states of $\mathcal{A}_{\mathcal{R}}$. If any of these subsequences is not contained in F , then s is deleted from C_k .

Terminating Condition. The set of frequent k -sequences F_k is empty. Unlike SPIRIT(L), we cannot terminate SPIRIT(V) based on just $F_k(a)$ becoming empty (where a is the start state of $\mathcal{A}_{\mathcal{R}}$). The reason is that, even though there may be no frequent and valid sequences of length k for a , there could still be longer sequences that are frequent and valid with respect to a .

3.5 The SPIRIT(R) Algorithm

SPIRIT(R) essentially pushes the RE constraint \mathcal{R} “all the way” inside the pattern mining computation, by requiring every candidate sequence for which support is counted to be valid (i.e., $\mathcal{C}' \equiv \mathcal{R}$).

Candidate Generation. Since F contains only valid and frequent sequences, there is no efficient mechanism for generating candidate k -sequences other than a “brute force” enumeration using the automaton $\mathcal{A}_{\mathcal{R}}$. The idea is to traverse the states and transitions of $\mathcal{A}_{\mathcal{R}}$ enumerating all paths of length k that begin with the start state and end at an accept state. Obviously, each such path corresponds to a valid k -sequence containing the elements that label the transitions in the path. (The terms “path” and “sequence” are used interchangeably in the following description.)

We employ two optimizations to improve the efficiency of the above exhaustive path enumeration scheme. Our first optimization uses the observation that, if a path of length less than k corresponds to a sequence that is valid but not frequent, then further extending the path is unnecessary since it cannot yield frequent k -sequences. The second optimization involves exploiting *cycles* in $\mathcal{A}_{\mathcal{R}}$ to reduce computation.

Lemma 3.2: Suppose for a path $\langle t u \rangle$ (of length less than k), both t and $\langle t u \rangle$ result in the same state from the start state a . (That is, u corresponds to a cycle in $\mathcal{A}_{\mathcal{R}}$.) Then, if the path $\langle t u v \rangle$ obtained as a result of extending $\langle t u \rangle$ with v is to yield a candidate k -sequence, it must be the case that $\langle t v \rangle$ is both frequent and valid. ■

Consider the generation of candidate k -sequences C_k . Given a path $\langle t u \rangle$ satisfying the assumptions of Lemma 3.2, we only need to extend $\langle t u \rangle$ with sequences v for which $\langle t v \rangle$ belongs to $F_{|\langle t v \rangle|}$ (since the length of $\langle t v \rangle$ is less than k). Due to space constraints, we have omitted the detailed definition of the can-

didate generation algorithm for SPIRIT(R) and examples of its operation. The interested reader is referred to [5].

Candidate Pruning. A candidate sequence s in C_k can be pruned if a valid subsequence of s is not frequent. The maximal valid subsequences of s can be computed by invoking algorithm FINDMAXSUBSEQ with *Start* equal to $\{a\}$ and *End* equal to the set of all accept states of $\mathcal{A}_{\mathcal{R}}$.

Terminating Condition. For some iteration j , sets $F_j, \dots, F_{j+|\mathcal{A}_{\mathcal{R}}|-1}$ are all empty, where $|\mathcal{A}_{\mathcal{R}}|$ is the number of states in automaton $\mathcal{A}_{\mathcal{R}}$. To see this, consider any frequent and valid sequence s whose length is greater than $j + |\mathcal{A}_{\mathcal{R}}| - 1$. Obviously, s contains at least one cycle of length at most $|\mathcal{A}_{\mathcal{R}}|$ and, therefore, s must contain at least one frequent and valid subsequence of length at least j . However, no valid sequence with length greater than or equal to j is frequent (since $F_j, \dots, F_{j+|\mathcal{A}_{\mathcal{R}}|-1}$ are all empty). Thus, s cannot be a frequent and valid sequence.

4 Experimental Results

In this section, we present an empirical study of the four SPIRIT algorithms with synthetic and real-life data sets. The objective of this study is twofold: (1) to establish the effectiveness of allowing and exploiting RE constraints during sequential pattern mining; and, (2) to quantify the constraint-based vs. support-based pruning tradeoff for the SPIRIT family of algorithms (Section 3.1).

In general, RE constraints whose automata contain fewer transitions per state, fewer cycles, and longer paths tend to be more *selective*, since they impose more stringent restrictions on the ordering of items in the mined patterns. Our expectation is that for RE constraints that are more selective, constraint-based pruning will be very effective and the latter SPIRIT algorithms will perform better. On the other hand, less selective REs increase the importance of good support-based pruning, putting algorithms that use the RE constraint too aggressively (like SPIRIT(R)) at a disadvantage. Our experimental results corroborate our expectations. More specifically, our findings can be summarized as follows.

1. The SPIRIT(V) algorithm emerges as the overall winner, providing consistently good performance over the entire range of RE constraints. For certain REs, SPIRIT(V) is more than an order of magnitude faster than the “naive” SPIRIT(N) scheme.
2. For highly selective RE constraints, SPIRIT(R) outperforms the remaining algorithms. However, as the RE constraint becomes less selective, the number of candidates generated by SPIRIT(R) explodes and the algorithm fails to even complete execution for certain cases (it runs out of virtual memory).
3. The overheads of the candidate generation and pruning phases for the SPIRIT(L) and SPIRIT(V) algorithms are negligible. They typically constitute less than 1% of the total execution time, even for complex

REs with automata containing large numbers of transitions, states, and cycles.

Thus, our results validate the thesis of this paper that incorporating RE constraints into the mining process can lead to significant performance benefits. All experiments reported in this section were performed on a Sun Ultra-2/200 workstation with 512 MB of main memory, running Solaris 2.5. The data sets were stored on a local disk.

4.1 Synthetic Data Sets

We used a synthetic data set generator to create a database of sequences containing items. The input parameters to our generator include the number of sequences in the database, the average length of each sequence, the number of distinct items, and a Zipf parameter z that governs the probability of occurrence, $\frac{1}{i^z} / \sum_i \frac{1}{i^z}$, of each item i in the database. The length for each sequence is selected from a Poisson distribution with mean equal to the average sequence length. Note that an item can appear multiple times in a single data sequence.

In addition, since we are interested in a sensitivity analysis of our algorithms with respect to the RE constraint \mathcal{R} , we used an RE generator to produce constraints with a broad range of selectivities. Each RE constraint output by the generator consists of *blocks* and each block in turn contains *terms* with the following structure. A term T_i is a disjunction of items and has the form $(s_1 | s_2 | \dots | s_l)$. Each block B_i is simply a concatenation of terms, $T_1 T_2 \dots T_m$. Finally, the constraint \mathcal{R} is constructed from blocks and has the form $(B_1 | B_2 | \dots | B_n)^*$ – thus, every sequence that satisfies \mathcal{R} is a concatenation of one or more sequences satisfying the block constraints. The generic structure of the automaton $\mathcal{A}_{\mathcal{R}}$ for \mathcal{R} is shown in Figure 4. RE constraints with different selectivities can be generated by varying the number of items per term, the number of terms per block, and the number of blocks in \mathcal{R} . Note that, in terms of the automaton $\mathcal{A}_{\mathcal{R}}$, these parameters correspond to the number of transitions between a pair of states in $\mathcal{A}_{\mathcal{R}}$, the length of each cycle, and the number of cycles contained in $\mathcal{A}_{\mathcal{R}}$, respectively.

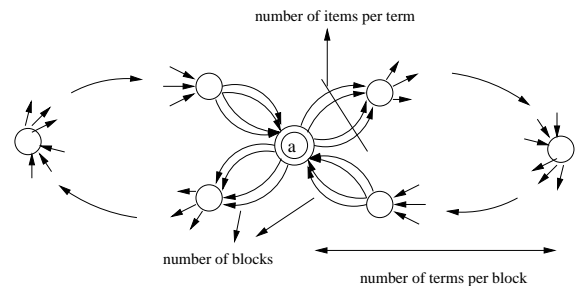


Figure 4: Structure of automaton for RE generation.

The RE generator accepts the maximum number of items per term, the number of terms per block, and the number of blocks as input parameters. In the RE constraint that it outputs, the number of items per term is uniformly distributed between 1 and the maximum specified value.

The items in each term of \mathcal{R} are chosen using the same Zipfian distribution that was used to generate the data set. The RE generator thus enables us to carry out an extensive study of the sensitivity of our algorithms to a wide range of RE constraints with different selectivities.

Table 3 shows the parameters for the data set and the RE constraint, along with their default values and the range of values for which experiments were conducted. The default value of $z = 1.0$ was chosen to model an (approximate) 70-30 rule and to ensure that the item skew was sufficient for some interesting patterns to appear in the data sequences. In each experiment, one parameter was varied with all other parameters fixed at their default values. Once again, due to space constraints, we only present a subset of our experimental results. The full set of results (including scaleup and maximum distance experiments) can be found in [5].

| Parameter | Default | Range |
|----------------------------|---------|---------------------------------|
| No. of Sequences | 10^5 | $5 \cdot 10^4 - 2.5 \cdot 10^5$ |
| Avg. Sequence Length | 10 | |
| No. of Items | 1000 | |
| Zipf Value | 1.0 | |
| Max. No. of Items per Term | 10 | 2 – 30 |
| No. of Terms Per Block | 4 | 2 – 10 |
| No. of Blocks | 4 | 2 – 10 |
| Min. Support | 1.0 | 0.5 – 2.0 |
| Max. Distance | 2 | 0 – 15 |

Table 3: Synthetic data and RE constraint parameters.

4.2 Performance Results with Synthetic Data Sets

Maximum Number of Items Per Term. Figure 5(a) illustrates the execution times of the SPIRIT algorithms as the maximum number of items per term in \mathcal{R} is increased. As expected, as the number of items is increased, the number of transitions per state in $\mathcal{A}_{\mathcal{R}}$ also increases and so do the numbers of legal and valid sequences. Thus, constraint-based pruning becomes less effective and the performance of all SPIRIT algorithms deteriorates as more items are added to each term. As long as the number of items per term does not exceed 15, \mathcal{R} is fairly selective; consequently, constraint-based pruning works well and the SPIRIT algorithms that use \mathcal{R} to prune more candidates perform better. For instance, when the maximum number of items per term is 10, the SPIRIT(N), SPIRIT(L), SPIRIT(V), and SPIRIT(R) algorithms count support for 7105, 1418, 974, and 3822 candidate sequences, respectively. SPIRIT(R) makes only two passes over the data for valid candidate sequences of lengths 4 and 8. The remaining algorithms make 8 passes to count supports for candidates with lengths up to 8, a majority of which have lengths 4 and 5.

However, beyond 15 items per term, the performance of the algorithms that rely more heavily on constraint \mathcal{R} for pruning candidates degenerates rapidly. SPIRIT(R) sustains the hardest hit since it performs very little support-based pruning and its exhaustive enumeration approach for candidate generation results in an enormous number of can-

didates of length 4. In contrast, since SPIRIT(N) only uses \mathcal{R} to prune sequences not involving items in \mathcal{R} , and few new items are added to terms in \mathcal{R} once the number of items per term reaches 15, the execution times for the SPIRIT(N) algorithm hold steady. Beyond 25 items per term, the running times of SPIRIT(L) and SPIRIT(V) also stabilize, since decreases in the amount of constraint-based pruning as \mathcal{R} becomes less selective are counterbalanced by increases in support-based pruning. At 30 items per term, SPIRIT(V) continues to provide a good balance of constraint-based and support-based pruning and, thus, performs the best.

Number of Terms Per Block. The graph in Figure 5(b) plots the running times for the SPIRIT algorithms as the number of terms per block is varied from 2 to 10. Increasing the number of terms per block actually causes each cycle (involving the start state a) to become longer. The initial dip in execution times for SPIRIT(L), SPIRIT(V), and SPIRIT(R) when the number of terms is increased from 2 to 4 is due to the reduction in the number of candidate sequences of lengths 4 and 5. This happens because with short cycles of length 2 in $\mathcal{A}_{\mathcal{R}}$, sequences of length 4 and 5 visit the start state multiple times and the start state has a large number of outgoing transitions. But when $\mathcal{A}_{\mathcal{R}}$ contains cycles of length 4 or more, the start state is visited at most once, thus causing the number of candidate sequences of lengths 4 and 5 to decrease. As cycle lengths grow beyond 4, the number of legal sequences (with respect to a state in $\mathcal{A}_{\mathcal{R}}$) starts to increase due to the increase in the number of states in each cycle. However, the number of valid sequences (with respect to a state in $\mathcal{A}_{\mathcal{R}}$) does not vary much since each of them is still required to terminate at the start state a .

Note that when the number of terms exceeds 6, the number of candidates generated by SPIRIT(R) simply explodes due to the longer cycles. On the other hand, SPIRIT(V) provides consistently good performance throughout the entire range of block sizes.

Number of Blocks. Figure 6(a) depicts the performance of the four algorithms as the number of blocks in \mathcal{R} is increased from 2 to 10. The behavior of the four algorithms has similarities to the “number of items per term” case (Figure 5(a)). The only difference is that, as the number of blocks is increased, the decrease in \mathcal{R} ’s selectivity and the increase in the number of legal and valid sequences in $\mathcal{A}_{\mathcal{R}}$ are not as dramatic. This is because the number of blocks only affects the number of transitions associated with the start state – the number of transitions for other states in $\mathcal{A}_{\mathcal{R}}$ stays the same. Once again, SPIRIT(V) performs well consistently, for the entire range of numbers of blocks. An interesting case is that of SPIRIT(R) whose execution time does degrade beyond SPIRIT(V)’s, as the number of blocks is increased, but it still manages to do better than SPIRIT(L), even when \mathcal{R} contains 10 blocks. This can be attributed predominantly to the effectiveness of the optimization for cycles in $\mathcal{A}_{\mathcal{R}}$ that is applied during SPIRIT(R)’s candidate generation phase. In general, due to

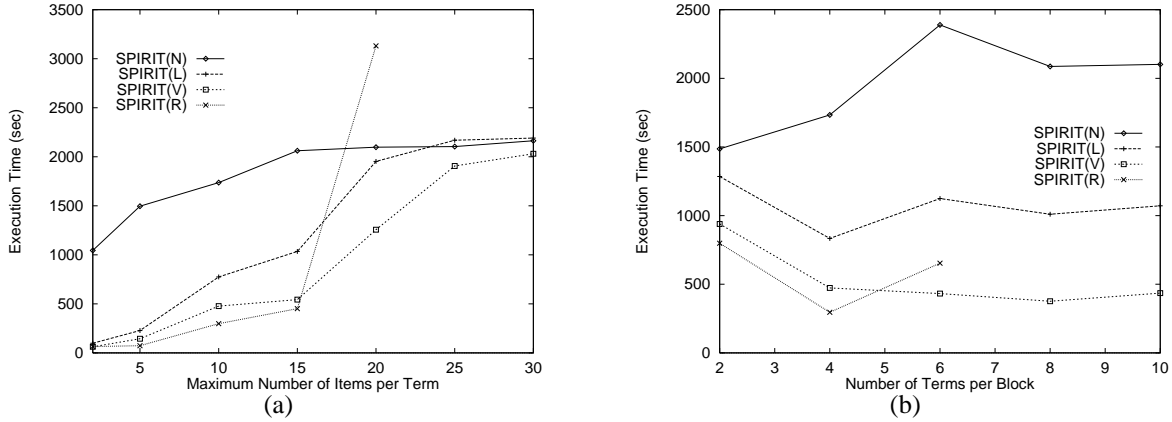


Figure 5: Performance results for (a) number of items and (b) number of terms.

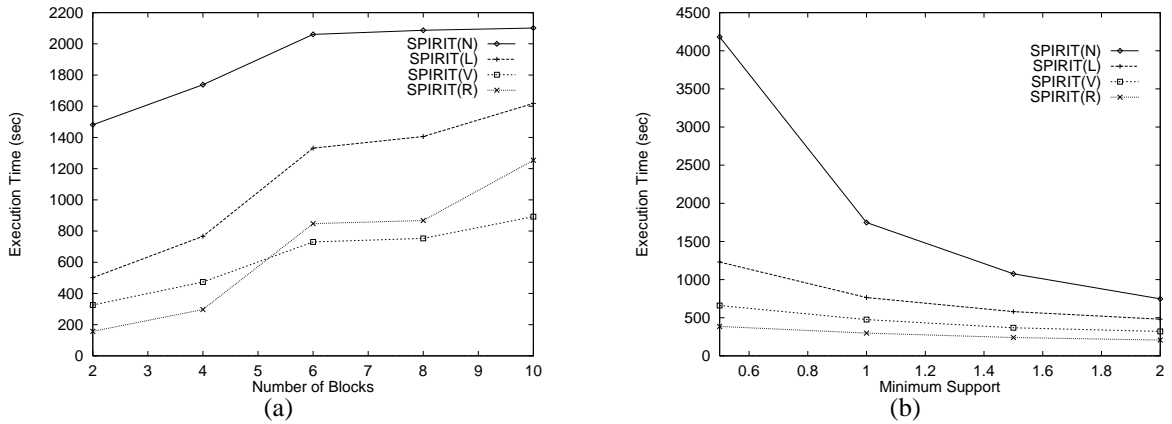


Figure 6: Performance results for (a) number of blocks and (b) minimum support.

our cycle optimization, one can expect the SPIRIT(R) algorithm to perform reasonably well, even when $\mathcal{A}_{\mathcal{R}}$ contains a large number of cycles of moderate length.

Minimum Support. The execution times for the SPIRIT algorithms as the minimum support threshold is increased from 0.5 to 2.0 are depicted in Figure 6(b). As expected, the performance of all algorithms improves as the minimum support threshold is increased. This is because fewer candidates have the potential to be frequent for higher values of minimum support. Furthermore, note that the running times of algorithms that rely more heavily on support-based pruning improve much more rapidly.

4.3 Real-life Data Set

For our real-life data experiments, we used the WWW server access logs from the web site of an academic CS department⁷. The logs contain the sequences of web pages accessed by each user⁸ starting from the department’s web site, for the duration of a week. The department’s

home page contains links to a number of topics, including Academics, Admissions, Events, General information, Research, People, and Resources. There are additional links to the university and college home pages to which the CS department belongs, but we chose not to use these links in our RE constraint. Users navigate through the web pages by clicking on links in each page, and the sequences of pages accessed by a user are captured in the server logs.

We used a RE constraint to focus on user access patterns that start with the department’s home page (located at `/main.html`) and end at the web page containing information on the M.S. degree program (located at `/academics/ms-program.html`). In addition, we restricted ourselves to patterns for which the intermediate pages belong to one of the aforementioned 7 topics (e.g., Academics). Thus, the automaton $\mathcal{A}_{\mathcal{R}}$ contains three states. There is a transition from the first (start) state to the second on `/main.html` and a transition from the second state to the third (accept) state on `/academics/ms-program.html`. The second state has 15 transitions to itself, each labeled with the location of a web page belonging to one of the above 7 topics. We used a minimum support

⁷At the department’s request, we do not disclose its identity.

⁸We use IP addresses to distinguish between users.

| Size | Frequent and Valid Sequences |
|------|--|
| 2 | < /main.html/academics/ms-program.html > |
| 3 | < /main.html/general/contacts.html/academics/ ms-program.html > < /main.html/general/nav.html/academics/ms-program.html > < /main.html/academics/academics.html/academics/ ms-program.html > < /main.html/academics/nav.html/academics/ ms-program.html > < /main.html/admissions/nav.html/academics/ ms-program.html > < /main.html/admissions/admissions.html/academics/ ms-program.html > |
| 4 | < /main.html/general/nav.html/general/contacts.html/ academics/ms-program.html > < /main.html/academics/nav.html/academics/academics.html/ academics/ms-program.html > < /main.html/admissions/nav.html/admissions/ admissions.html/academics/ms-program.html > |

Table 4: Interesting patterns discovered in the WWW logs.

| Algorithm | Exec. Time (sec) | Candidates | Passes |
|-----------|------------------|------------|--------|
| SPIRIT(N) | 1562.8 | 5896 | 13 |
| SPIRIT(L) | 32.77 | 1393 | 10 |
| SPIRIT(V) | 16.0 | 59 | 5 |
| SPIRIT(R) | 17.67 | 52 | 7 |

Table 5: Execution statistics for the SPIRIT algorithms.

threshold of 0.3%. The number of access sequences logged in the one week data set was 12868.

The mined frequent and valid access patterns are listed in increasing order of size in Table 4. Note that there is a number of distinct ways to access the M.S. degree program web page by following different sequences of links (e.g., via admissions, academics). The execution times and the numbers of candidates generated by the four SPIRIT algorithms are presented in Table 5. As expected, since the RE constraint is fairly selective, both SPIRIT(V) and SPIRIT(R) have the smallest running times. SPIRIT(L) is about twice as slow compared to SPIRIT(V) and SPIRIT(R). The execution time for SPIRIT(N) is almost two orders of magnitude worse than SPIRIT(V) and SPIRIT(R), since it generates a significantly larger number of candidate sequences with lengths between 5 and 9 (almost 4000). We believe that our results clearly demonstrate the significant performance gains that can be attained by pushing RE constraints inside a *real-life* pattern mining task.

5 Conclusions

In this paper, we have proposed the use of Regular Expressions (REs) as a flexible constraint specification tool that enables user-controlled focus to be incorporated into the pattern mining process. We have developed a family of novel algorithms (termed SPIRIT) for mining frequent sequential patterns that also satisfy user-specified RE constraints. The main distinguishing factor among the proposed schemes is the degree to which the RE constraints

are enforced to prune the search space of patterns during computation. The SPIRIT algorithms are illustrative of the tradeoffs that arise when constraints that do not subscribe to nice properties (like anti-monotonicity) are integrated into the mining process. To explore these tradeoffs, we have conducted an extensive experimental study on synthetic and real-life data sets. The experimental results clearly validate the effectiveness of our approach, showing that speedups of more than an order of magnitude are possible when RE constraints are pushed deep inside the mining process. Our experimentation with real-life data also illustrates the versatility of REs as a user-level tool for focusing on interesting patterns.

Acknowledgments: We would like to thank Narain Gehani, Hank Korth, and Avi Silberschatz for their encouragement. Without the support of Yesook Shim, it would have been impossible to complete this work.

References

- [1] R. Agrawal, G. Psaila, E. L. Wimmers, and M. Zait. "Querying Shapes of Histories". In *Proc. of the 21st Intl. Conf. on Very Large Data Bases*, September 1995.
- [2] R. Agrawal and R. Srikant. "Fast Algorithms for Mining Association Rules". In *Proc. of the 20th Intl. Conf. on Very Large Data Bases*, September 1994.
- [3] R. Agrawal and R. Srikant. "Mining Sequential Patterns". In *Proc. of the 11th Intl. Conf. on Data Engineering*, March 1995.
- [4] M.-S. Chen, J. S. Park, and P. S. Yu. "Efficient Data Mining for Path Traversal Patterns". *IEEE Trans. on Knowledge and Data Engineering*, 10(2):209–221, March 1998.
- [5] M. N. Garofalakis, R. Rastogi, and K. Shim. "SPIRIT: Sequential Pattern Mining with Regular Expression Constraints". *Bell Labs Tech. Memorandum BL0112370-990223-03TM*, February 1999.
- [6] H. R. Lewis and C. Papadimitriou. "*Elements of the Theory of Computation*". Prentice Hall, Inc., 1981.
- [7] H. Mannila and H. Toivonen. "Discovering Generalized Episodes Using Minimal Occurrences". In *Proc. of the 2nd Intl. Conf. on Knowledge Discovery and Data Mining*, August 1996.
- [8] H. Mannila, H. Toivonen, and A. I. Verkamo. "Discovering Frequent Episodes in Sequences". In *Proc. of the 1st Intl. Conf. on Knowledge Discovery and Data Mining*, August 1995.
- [9] R. T. Ng, L. V.S. Lakshmanan, J. Han, and A. Pang. "Exploratory Mining and Pruning Optimizations of Constrained Association Rules". In *Proc. of the 1998 ACM SIGMOD Intl. Conf. on Management of Data*, June 1998.
- [10] R. Srikant, Q. Vu, and R. Agrawal. "Mining Association Rules with Item Constraints". In *Proc. of the 3rd Intl. Conf. on Knowledge Discovery and Data Mining*, August 1997.
- [11] R. Srikant and R. Agrawal. "Mining Sequential Patterns: Generalizations and Performance Improvements". In *Proc. of the 5th Intl. Conf. on Extending Database Technology (EDBT'96)*, March 1996.
- [12] J. T.-L. Wang, G.-W. Chirn, T. G. Marr, B. Shapiro, D. Shasha, and K. Zhang. "Combinatorial Pattern Discovery for Scientific Data: Some Preliminary Results". In *Proc. of the 1994 ACM SIGMOD Intl. Conf. on Management of Data*, May 1994.