

Efficiently Summarising Event Sequences with Rich Interleaving Patterns

Apratim Bhattacharyya[°]

Jilles Vreeken[°]

Abstract

Discovering the key structure of a database is one of the main goals of data mining. In pattern set mining we do so by discovering a small set of patterns that together describe the data well. The richer the class of patterns we consider, and the more powerful our description language, the better we will be able to summarise the data. In this paper we propose SQUISH, a novel greedy MDL-based method for summarising sequential data using rich patterns that are allowed to interleave. Experiments show SQUISH is orders of magnitude faster than the state of the art, results in better models, as well as discovers meaningful semantics in the form patterns that identify multiple choices of values.

1 Introduction

Discovering the key patterns from a database is one of the main goals of data mining. Modern approaches do not to ask for *all* patterns that satisfy a local interestingness constraint, such as frequency [2, 10], but instead ask for that *set of patterns* that is optimal for the data at hand. There are different ways to define this optimum. The Minimum Description Length (MDL) principle [14, 5] has proven to be particularly successful [22, 16]. Loosely speaking, by MDL we say that the best set of patterns is the set that compresses the data best. How well we can compress, or better, describe the data depends on the description language we use. The richer this language, the more relevant structure we can identify. At the same time, a richer language means a larger search space, and hence requires more efficient search.

In this paper we consider databases of event sequences, and are after that set of sequential patterns that together describe the data best—as we did previously with SQS [20]. Like SQS we describe a database with occurrences of patterns. Whereas SQS requires these occurrences to be disjoint, however, we allow patterns to *interleave*. This leads to more succinct descriptions as well as better pattern recall. Moreover, we use a richer class of patterns. That is, we do not only allow for gaps in occurrences, but also allow patterns to emit *one out of multiple* events at a certain location. For example,

the pattern ‘*paper [proposes | presents] new*’ discovered in the JMLR abstract database matches two common forms of expressing that a paper presents or proposes something new.

With this richer language, we can obtain much better compression rates with much fewer patterns. To discover good models we propose SQUISH, a highly efficient and versatile search algorithm. Its efficiency stems from re-use of information, partitioning the data, and in particular from considering only the currently relevant occurrences of patterns in the data. It is a natural any-time algorithm, and can be ran for any time budget that is opportune.

Extensive experimental evaluation shows that SQUISH performs very well in practice. It is much better at retrieving interleaving patterns than the very recent proposal by Fowkes and Sutton [4], and obtains much better compression rates than SQS [20], while being orders of magnitude faster than both. The choice-patterns it discovers give insight in the data beyond the state of the art, identifying semantically coherent patterns. Moreover, SQUISH is highly extendable, allowing for richer pattern classes to be considered in the future.

2 Preliminaries

Here we introduce basic notation, and give short introductions to the MDL principle.

2.1 Notation We consider databases of *event sequences*. Such a database D is composed of $|D|$ sequences. A sequence $S \in D$ consists of $|S|$ events drawn from an alphabet Ω . The total number of events occurring in the database, denoted by $||D||$, is simply the sum of lengths of all sequences $\sum_{S \in D} |S|$. We write $S[j]$ to refer to the j^{th} event in sequence S . The support of an event $e \in \Omega$ in a sequence S is simply the number of occurrences of e in S , i.e. $\text{supp}(e | S) = |\{j | S[j] = e\}|$. The support of e in a database D is defined as $\text{supp}(e|D) = \sum_{i=1}^{|D|} \text{supp}(e | S_i)$.

We consider two types of sequential patterns. A **serial episode** $X \in \Omega^{|X|}$ is a sequence of $|X|$ events, and we say that a sequence S contains X if there is a subsequence in S equal to X . We allow noise, or *gap events*, within an occurrence of X . We also consider choice episodes, or **choicisodes**. These are serial episodes with positions matching *one out of multiple* events. For example, serial

[°]Max-Planck Institute for Informatics and Saarland University, Saarland Informatics Campus, Saarbrücken, Germany. {abhattach, jilles}@mpi-inf.mpg.de

episode ac matches an a followed by c , whereas choiciscode $[a, b]c$ matches occurrences of a or b followed by c .

2.2 Brief introduction to MDL The Minimum Description Length principle (MDL) [14, 5] is a practical version of Kolmogorov Complexity [9]. Both embrace the slogan Induction by Compression. We use the MDL principle for model selection.

By MDL, the best model is the model that gives the best lossless compression. More specifically, given a set of models \mathcal{M} , the best model $M \in \mathcal{M}$ is the one that minimizes $L(M) + L(D | M)$, in which $L(M)$ is the length in bits of the description of M , and $L(D | M)$ is the length of the data when encoded with model M . Simply put, we are interested in that model that best compresses the data without loss. MDL as describe above is known as two-part MDL, or crude MDL; as opposed to refined MDL. In refined MDL model and data are encoded together [5]. We use two-part MDL because we are specifically interested in the model: the patterns that give the best description. In MDL we are only concerned with code lengths, not actual code words.

Next, we formalise our problem in terms of MDL.

3 MDL for Event Sequences

To use MDL we need to define a model class \mathcal{M} , and how to encode a model $M \in \mathcal{M}$ and data D in bits.

As models we will consider *code tables* [22, 20]. A code table CT is a dictionary between patterns and associated codes. A code table consists of the singleton patterns $e \in \Omega$, as well as a set \mathcal{P} of non-singleton patterns. We write $code_p(X)$ to denote the *pattern code* that identifies a pattern $X \in CT$. Similarly, we write $code_f(X)$ and $code_g(X)$ for the codes resp. identifying a *fill* resp. a *gap* in the occurrence of a pattern X .

We can encode a sequence database D using the patterns in a code table CT , which generates a *cover* C of the database. A cover C uniquely defines a pattern code stream C_p and a meta code stream C_m . The pattern stream is simply the concatenation of the codes corresponding to the patterns in the cover, in the order of their appearance. Likewise, the meta stream C_m is the concatenation of the gap and fill codes corresponding to the cover. In Fig. 1, we illustrate two example covers and corresponding code tables, the first using only singletons and the second cover with interleaving using patterns from a richer code table with choicisodes.

Before formalising our score, it is helpful to know how to decode a database given a code table and the code streams.

3.1 Decoding a database To decode a database, we start by reading a $code_p(X)$ from the pattern stream C_p . If the corresponding pattern X is a singleton, we append it to our reconstruction of the database D . If it is a non-singleton, we append its first event, $X[1]$, D . To allow for interleaving, we

Data D : a, b, d, c, d, c, a ,

Cover 1: using only singletons

C_p a, b, d, c, d, c, a

CT_1 : a a
 b b
 c c
 d d

Cover 2: using interleaving patterns

C_p p, q, a

C_m $?, ?, !, !, ?, !, ?, !$

alignment a, b, d, c, d, c, a
 p, q, a

CT_2 : a a
 b b
 c c fills
 d d gaps
 acd $p, !, ?$
 $b[d, c]c$ $q, !, ?$
 $r, !, ?$

Figure 1: Example of two possible covers. The first uses only singletons, the second includes interleaving and choicisodes.

have to add a new *context* to context list Λ . A context is a tuple (X, i) consisting of a pattern X , and a pointer i to the next event to be read from the pattern. For an example, let us consider Cover 2 in Fig. 1. We read $code_p(p)$ from $code_p$, append $p[1] = a$ to D , and add $(p, 2)$ to the context list.

Next, if the context list is non-empty, we read as many meta codes from C_m as there are contexts in Λ . If we read a fill code $code_f(X)$ corresponding to one of the contexts $(X, i) \in \Lambda$, we append the next event from X , $X[i]$ to the data D , and increment the pointer. If after this step we have finished reading the pattern, we remove its context from the list. If we only read gap codes $code_g(X)$ for every pattern X in the context list, we read again from the pattern stream. We do this until we reach the end of the pattern stream C_p .

Continuing our example, we read $code_g(p)$ from C_m , which corresponds to a gap in the occurrence of pattern p . We read $code_p(q)$ from C_p , write $q[1] = b$ to D , and insert context $(q, 2)$ to Λ . Next, Λ contains two contexts, and we read two meta codes from C_m , viz. $code_g(p)$ and $code_f(q)$. As for context $(q, 2)$ we read a fill code, we write $q[2] = d$ to D , and increment its pointer to 3. Etc.

3.2 Calculating Encoded Lengths Given the above scheme we know which codes to expect when, and can now formalise our score. We build upon and extend the encoding on Tatti & Vreeken [20] for richer covers and patterns.

Encoded Length of the Database We encode the pattern stream C_p using Shannon optimal prefix codes. The length of the pattern code $L(code_p(X))$ for a pattern X depends on how often it is used in the pattern stream. We write $usage(X)$ to denote the number of times $code_p(X)$ occurs in C_p . The length of the optimal pattern code for X then is

$$L(code_p(X)) = -\log\left(\frac{usage(X)}{\sum_{Y \in CT} usage(Y)}\right).$$

The encoded length of the whole pattern stream C_p is then simply $L(C_p) = \sum_{X \in CT} usage(X)L(code_p(X))$.

To avoid arbitrary choices in the model encoding, we use prequential codes [5] to encode the meta stream. Prequential codes are asymptotically optimal without knowing the distribution beforehand. The idea is that we start with an uniform distribution over the events in the stream and update the counts after every received event. This means we have a valid probability distribution at every point in time, and can hence send optimal prefix codes. The total encoded length of the meta stream pattern is

$$L(C_m) = \sum_{X \in CT} \left(- \sum_{i=1}^{fills(X)} \log \left(\frac{\epsilon + i}{2\epsilon + i} \right) - \sum_{i=1}^{gaps(X)} \log \left(\frac{\epsilon + i}{2\epsilon + fills(X) + i} \right) \right).$$

where $\epsilon = 0.5$ is a constant by which we initialize the distribution [5], $fills(X)$ and $gaps(X)$ are the number of times $code_f(X)$ resp. $code_g(X)$ occurs in C_m .

For lossless decoding of database D , the number of sequences $|D|$ and the length of each sequence $S \in D$ should also be encoded. We do this using $L_{\mathbb{N}}$, the MDL optimal code for integers $n \geq 1$ [15].

Combining the above, for the total encoded length of a database, given a code table CT and cover C , we have

$$L(D | CT) = L_{\mathbb{N}}(|D|) + \sum_{S \in D} L_{\mathbb{N}}(|S|) + L(C_p) + L(C_m).$$

Next we discuss how to encode a model.

Encoded Length of the Code Table Note that the simplest valid code table consists of only the singletons Ω . We refer to this code table as ST , or, the standard code table. We use ST to encode the non-singleton patterns \mathcal{P} of a code table CT . The usage of a singleton $e \in ST$ is simply its support in D , and hence the code length $code_p(e) = -\log \left(\frac{supp(e|D)}{|D|} \right)$. To use these codes the recipient needs to know the supports of the singletons. We encode these using a data to model code—an index over a canonically ordered enumeration of all possibilities [21]; here it is the number of possible supports of $|\Omega|$ alphabets over a database length of $|D|$, $\binom{|D|}{|\Omega|}$. The length of the code is now simply the logarithm over the number of possibilities.

Given the standard code table ST , we can now encode the patterns in the code table. We first encode the length $|X|$ of the pattern, and then number of choice spots in the pattern, $||X|| - |X|$. We encode how many choices we have per location using a data to model code. We finally encode

the events $X[i]$ using the standard code table, ST . That is,

$$L(X | ST) = L_{\mathbb{N}}(|X|) + L_{\mathbb{N}}(||X|| - |X| + 1) + \log \left(\frac{||X|| - 1}{|X| - 1} \right) + \sum_{i=1}^{||X||} L(X[i] | ST).$$

Note that if we do not consider choicisodes, we can simplify the above as we only need to transmit the first and last part of this code. That is, the length and the events in the pattern.

Recall that, pattern codes in the pattern stream C_p are optimal prefix codes. The occurrences of the non-singleton patterns need to be transmitted with the model. We do this again using a data to model code. We encode the sum of pattern usages, $usage(\mathcal{P}) = \sum_{X \in CT \setminus \Omega} usage(X)$, by the MDL optimal code for integers. It is equivalent to use a pattern code per choicisode and then identify the choice-events, or to use a separate pattern code for each instantiation of the choicisode. For simplicity we make the latter choice.

The total encoded size of code table CT given a cover C of database D is then given by

$$L(CT | D, C) = L_{\mathbb{N}}(|\Omega|) + \log \left(\frac{||D||}{|\Omega|} \right) + L_{\mathbb{N}}(|\mathcal{P}| + 1) + L_{\mathbb{N}}(usage(\mathcal{P}) + 1) + L(usage(\mathcal{P}), |\mathcal{P}|) + \sum_{X \in CT} L(X, CT).$$

We are interested in the set of patterns and a corresponding cover C which minimizes the total encoded length of the code table and the database, which is,

$$L(CT, D) = L(CT | C) + L(D | CT).$$

We can now formally define our problem as follows.

Minimal Code Table Problem Let Ω be a set of events and let D be a sequence database over Ω , find the minimal set of serial (choice) episodes \mathcal{P} such that for the optimal cover C of D using \mathcal{P} and Ω , the total encoded cost $L(CT, D)$ is minimal, where CT is the code-optimal code table for C .

For a given database D , we would like to find its optimal pattern set in polynomial time. However, there are exponentially many possible pattern sets, and given a pattern set, there are exponentially many possible covers. For neither problem there exists trivial structure such as monotonicity or sub-modularity that would allow for an optimal polynomial time solution.

Hence, we resort to heuristics. In particular, we split our problem into two parts. We first explain our greedy algorithm to find a good cover given a set of patterns. We describe how to find a set of good patterns in Sec. 5.

4 Covering a Database

Given a pattern set \mathcal{P} and database D , we are after a cover C with interleaving and nesting, that minimises $L(CT, D)$.

Each occurrence of a pattern X in database D , possibly with gaps, defines a *window*. We denote by $S[a, b]$ a window in sequence S that extends from the position a to b . Two windows are non-overlapping if they do not have any events in common which belong to their respective patterns. Two interleaving or nesting windows might have common events, which, as we do not allow overlap, leads to gap events for one of the two windows. Two windows are disjoint if they do not have any events in common. For every event in the database D , there can be many windows with which we can choose to cover it. The optimal cover depends upon the pattern, fill, gap codes of the patterns. The choices grow exponentially with sequence length, with no trivial sub-structure.

To find good disjoint covers, Tatti & Vreeken [20] use an EM-style approach. At each step until convergence, given the pattern, gap and fill codes, the authors use the dynamic programming based algorithm ALIGN to find a cover. ALIGN takes a set of possibly overlapping minimal windows and returns a subset of disjoint minimal windows (i.e. a cover) which maximizes the sum of *gain* (a heuristic measure) of each window. Then, the lengths of the codes are reset based upon the found cover. It is unclear if this scheme can be extended to return a cover with interleaved or nested windows efficiently. Moreover if we extend our model with a new pattern, we have to rerun ALIGN from scratch.

We propose an efficient and easily extendible heuristic for good covers with interleaved and nested windows.

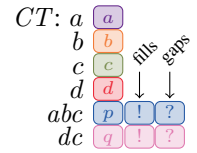
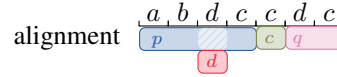
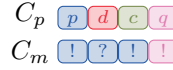
4.1 Window Lengths For a given pattern, as we consider windows with gaps, the length of an window in the database can be arbitrarily long. Tatti & Vreeken therefore consider only *minimal windows*. A window $w = S[i, j]$ is a minimal window of a pattern X if w contains X but no other proper sub-windows of w contain X . If no interleaving or nesting is allowed, it is optimal to consider only minimal windows. Otherwise, it is easy to construct examples where the optimal cover consists of non-minimal windows.

Consider the sequence $abdccdc$ and a code table with the pattern abc , dc and the singletons a , b , c and d . Two possible covers are: $(\underline{ab} \ d \ \underline{c}) \ c \ \underline{dc}$ using only minimal windows and $(\underline{ab}(\underline{dc})\underline{c}) \ \underline{dc}$ where a non-minimal window of abc is used and is nested with a window of dc . It is easy to see that the second cover leads to lower encoded length $L(abdccdc, \{abc, dc, a, b, c, d\})$ (see Fig 2) of about 2.9 bits.

Ideally, we should consider all possible windows. The number of possible windows of a pattern, however, is quadratic in the length of the database. This means that even a search for all windows is computationally inefficient. Therefore, we first search for only the shortest window from each starting position in the database. We consider longer windows when necessary. We do so as follows.

Data D : a, b, d, c, c, d, c

Cover 1: using minimal windows



Cover 2: using non-minimal windows

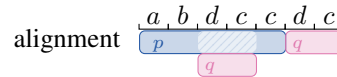
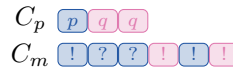


Figure 2: Two possible covers using the same code table. The first cover uses only minimal windows. The second includes nesting and a non-minimal window of abc . The second cover leads to an overall lower encoded length.

4.2 Window Search Given a pattern X , we use the pseudo-code FINDWIN presented as Algorithm 1 to search for its windows in the a sequence or sub-sequence S of database D . It returns us $W_{cand}(X)$ which is a set of candidate windows of the pattern X . It considers only the first window from each starting position in the sequence S . We later choose a subset of these windows (along with those of patterns other than X in CT) to create a cover of the database $L(D | CT)$. To control the ratio of gaps and fills, we maintain a budget variable. This is the number of extra allowed overall gaps. Ideally we would like to have more fills than gaps as it leads to better compression.

To search for windows efficiently in FINDWIN, we use an inverted index: $index^{-1}(x)$ which gives us a list of positions of the event $x \in \Omega$ in the database. We use a priority queue Q to store potential windows sorted by length. Shorter windows means more fills than gaps. We initialize FINDWIN (line 3) by creating potential windows at all the positions where the first alphabet of pattern X occurs in the sequence S_i and pushing these potential windows to Q . Each window w in Q contains the starting position in S_i , its length, and a pointer w_i . This pointer points to a certain event in pattern X which we are search for in S_i . At every step of FINDWIN (line 7) we look at the potential window at the top of the queue Q . We check if the next event in the database equals the character of the pattern X pointed to by w_i and increment the length of the window w . There are now two possibilities i) (line 10) The next database event is the same as the event in X pointed to by w_i . If we have found the full pattern X in the database, we add this window to $W_{cand}(X)$. We can now update our budget if we used more fills than

Algorithm 1: FINDWIN($S_i, X, budget$)

input : sequence S and a pattern X
output : set of windows for X

- 1 $W_{cand} \leftarrow \emptyset$;
- 2 **for** p in $index^{-1}(x)$ **do**
- 3 $w \leftarrow (start, length, w_i = 1)$;
- 4 $push(w, Q)$;
- 5 **while** T is not empty **do**
- 6 $w \leftarrow top(Q)$;
- 7 $(start, length, w_i) \leftarrow w$;
- 8 $length = length + 1$;
- 9 **if** $S_i[start + length]$ equals $X[w_i]$ **then**
- 10 **if** w_i points to the end of X **then**
- 11 $append\ w\ to\ W_{cand}$;
- 12 $b \leftarrow b + 2 \times length(X) - length - 1$;
- 13 **else**
- 14 **if** $b + 2 \times length(X) - length - 1 < 0$ **then**
- 15 $delete\ w\ from\ Q$;

gaps. Using less gaps in one window allows us to use more gaps in another. ii) (line 15) The next database event does not equal the event of X pointed to by w_e . This means that the potential window w has one extra gap. We check if this extra gap is allowed by our budget (line 16). Otherwise, we drop the window.

Now that we can search for windows of patterns, we describe how to choose a subset which generates a good cover C of the database D .

4.3 Candidate Order In the first step of our greedy strategy, we sort the set of patterns in a fixed order, similar to [3]. We call this order the **Candidate Order**. We cover the database using windows of patterns in this order. This order is designed to minimize the code length. This is achieved by putting longer and more frequently occurring patterns higher up in the candidate order. This means we can cover more events while minimizing the code length.

We consider the patterns $X \in CT$ in the order,

1. Decreasing \downarrow in length $|X|$
2. Decreasing \downarrow in support $support(X | D)$
3. Decreasing \downarrow in length of encoding it with the standard table.
4. Increasing \uparrow lexicographically.

4.4 Greedy Cover We now describe our greedy algorithm GREEDYCOVER which we use incrementally build a good cover as pseudo-code in Algorithm 2. We consider patterns in the candidate order. We maintain a set of selected windows W_{sel} . GREEDYCOVER takes this set of selected

Algorithm 2: GREEDYCOVER(W_{sel}, W_{cand})

input : set of selected windows W_{sel} and a set of candidate windows $W_{cand}(X)$ for pattern X .
output : set of selected windows W_{sel} combined with those in $W_{cand}(X)$ not overlapping with W_{sel} .

- 1 $W_{add} \leftarrow \emptyset$;
- 2 $W_{ext} \leftarrow$ create window extends from W_{sel} and $W_{cand}(X)$;
- 3 $last \leftarrow \emptyset$;
- 4 **for** window $w \in W_{ext}$ **do**
- 5 $W_{tmp} \leftarrow$ all windows of X between $last$ and v ;
- 6 **for** v in W_{tmp} , in order of decreasing length **do**
- 7 **if** v does not overlap with W_{tmp} **then**
- 8 $append\ v\ to\ W_{add}$;
- 9 $W_{tmp} \leftarrow W_{tmp} \cup \{ \text{windows of } X \text{ inside } w \text{ in } D \}$;
- 10 $last \leftarrow w$;
- 11 **return** Merge W_{sel} and W_{add} ;

windows W_{sel} and extends it with a subset of candidate windows of pattern X , $W_{cand}(X)$, found with FINDWIN and possibly with (longer, interleaved) windows found on the fly. We assume that both W_{sel} and W_{cand} are sorted.

We refer to a block of windows which are interleaved or nested with each other as an *window extend*. For ease of notation, we refer to windows which are not interleaved or nested also as window extends (containing a single window). We begin GREEDYCOVER by dividing the set W_{sel} into a set of window extends W_{ext} by a linear sweep (if W_{sel} is sorted). For patterns at the top of the candidate order W_{sel} is empty, so we can select all candidate windows $W_{cand}(X)$. For any other pattern, we iterate through the list of window extends W_{ext} (line 4). All the windows of the pattern occurring between any two extends in W_{ext} can be potentially chosen. These are put in W_{tmp} (line 5), a temporary list. It is possible that some windows of X in W_{tmp} overlap. We consider these windows in order of decreasing length (line 6) and discard any window that overlaps with a previously chosen window. We additionally search (on the fly) for interleaved windows occurring within the window extends W_{ext} (line 9).

For example consider the sequence $abcdacbd$ which we want to cover with the patterns ac and bd . Using FINDWIN we get two windows each for the two patterns. If ac is higher up in the candidate order, we first select the two windows of ac ; $abcdacbd$. We now have two window extends in W_{ext} . We search for windows of bd within the first window extend of ac to find one interleaved window: $abcdacbd$ and we select the second window of bd as it is between the two window extends of ac .

Note that, GREEDYCOVER now takes time $O(|W_{sel}| + ||D|| + |W_{cand}| \log(|W_{cand}|))$, in the worst case. Where, $|W_{sel}|$ is the number of windows in W_{sel} and $|W_{cand}|$ is the number of candidate windows. Let, $W_{\max(\mathcal{P})}$

be the maximum number of candidate windows of any pattern in \mathcal{P} . Then GREEDYCOVER takes time $O(|\mathcal{P}|(W_{\max(\mathcal{P})} \log(W_{\max(\mathcal{P})}) + ||D||))$ to construct a cover C of the database D using the patterns \mathcal{P} in the code table CT in the worst case. The maximum number of candidate windows of any pattern $W_{\max(\mathcal{P})}$ is bounded by the size of the database $O(||D||)$. However, GREEDYCOVER makes it computationally more efficient to extend the code table with a new pattern X . We can discard windows of patterns in \mathcal{P} below X in the candidate order from the cover and run GREEDYCOVER for $W_{\text{cand}}(X)$ and the patterns in the code table below X in candidate order. This means that we do not have to recompute the cover from scratch. This is very efficient if the pattern X is near the bottom of the candidate order. As we shall see GREEDYCOVER is very competitive in its execution time compared to SQS [20].

Having presented our greedy approach of covering a database given a set of patterns, we now turn our attention to the task of mining good set patterns.

5 Mining Good Code Tables

Given a pattern set \mathcal{P} we have a greedy algorithm to cover the database D and obtain the encoded length of the model and data $L(CT, D)$. To solve the **Minimal Code Table Problem** we want to find that set \mathcal{P} of patterns which minimizes the total encoded length $L(CT, D)$ of the database. As discussed before, there does not seem to be any trivial substructure in the problem which we can exploit to obtain an optimal set of patterns \mathcal{P} in polynomial time. So, we resort to heuristics. We build upon and extend SQS-SEARCH [20].

5.1 Generating Candidates We build a pattern set \mathcal{P} incrementally. Given a set of patterns \mathcal{P} and a cover C , we aim to find a pattern X and an extension Y , such that $X, Y \in \mathcal{P} \cup \Omega$, whose combination XY would decrease the encoded length $L(\mathcal{P} \cup XY, D)$. We do this until we cannot find any XY that when added to \mathcal{P} reduces the total encode size. Doing is exactly, however, is computationally prohibitive. At every iteration, there would be $O((|\mathcal{P}| + |\Omega|)^2)$ possible candidates. Thus, we again resort to heuristics. We use the heuristic algorithm ESTIMATE from [20] that can find good candidates, with likely decrease in code length if added, in $O(|\mathcal{P}| + |\Omega| + ||D||)$ time. For readability and succinctness, we describe algorithm ESTIMATE in Appendix A.

Candidates are accepted or rejected based on the compression gain. As we can now find richer covers with interleaving and nesting, candidates are potentially more likely to be accepted. However, we want to find a succinct set of patterns which describe the data well. Choicisodes can help in this search for a succinct summary of the data.

5.2 Choicisodes Recall from Sec. 3.2 that we can encode patterns as choicisodes. We have the possibility of combin-

ing a newly discovered non-singleton pattern with a previously discovered non-singleton pattern or choicisode to create or expand a choicisode. Combining non-singleton patterns into a single choicisode may hence lead to savings in the encoded length of the code table $L(CT | C)$ while providing a more succinct representation of the pattern set.

We use a greedy strategy based on MDL for discovering choicisodes. For each newly discovered non-singleton pattern, we consider all previously discovered non-singleton patterns or choicisodes which differ with it at one position. Then, we calculate the increase in code length (of the model) if we encode it as a choicisode with each of these non-singleton patterns or choicisodes. We also consider the increase in code length if we encode it as independently. We choose whichever option with leads to the minimum increase in code length.

Next we present our algorithm SQUISH for mining a succinct and representative pattern set.

5.3 The SQUISH algorithm The present the complete algorithm SQUISH as pseudo-code in Algorithm 3. At each iteration, it considers each pattern $X \in CT$. It creates potential extensions XY , with $Y \in CT$, based on estimated change in the encoded length using ESTIMATE (line 6). SQUISH then considers each of these patterns in the order of the estimated decrease in gain if added to CT (line 7). FINDWIN is used to find the candidate windows of each of these extensions (line 9). GREEDYCOVER is used to cover the data with this candidate pattern XY added to CT . We simultaneously consider the possibility of encoding XY as a choicisode. If XY leads to a decrease in the encoded length of the database D then, we add XY to CT . If XY is to be added, we PRUNE (see Appendix A) the code table to remove redundant patterns. Consider, for example if we decide to add $abcd$, the pattern ab and cd may not be required to construct an effective cover of the database. We also consider the singletons occurring in the gaps of XY , by constructing new extended patterns by using these gap alphabets as intermediate alphabets.

6 Related Work

Discovering sequential patterns is an active research topic. Traditionally there was a focus on mining frequent sequential patterns, with different definitions of how to count occurrences [10, 23, 8]. Mining general patterns, patterns where the order of events are specified by a DAG is surprisingly hard. Even testing whether a sequence contains a pattern is NP-complete [18]. Consequently, research has focused on mining subclasses of episodes, such as, episodes with unique labels [1, 12], strict episodes [19], and injective episodes [1].

Traditional pattern mining typically results in overly many and highly redundant results. Once approach to counter this is mining statistically significant patterns. Com-

Algorithm 3: SQUISH(D)

input : database D
output : pattern set \mathcal{P} with low $L(CT, D)$

- 1 $\mathcal{P} \leftarrow \phi$;
- 2 $C \leftarrow \text{GREEDYCOVER}(\mathcal{P}, D)$;
- 3 **while** *changes* **do**
- 4 $\mathcal{F} \leftarrow \phi$;
- 5 **for** $X \in CT$ **do**
- 6 \lfloor add $\text{ESTIMATE}(X, A, D)$ to \mathcal{F} ;
- 7 **for** $Z \in \mathcal{F}$ ordered by estimated gain **do**
- 8 Sort $\mathcal{P} \cup Z$ in candidate order;
- 9 $W_{cand}(Z) \leftarrow \text{FINDWIN}(D, Z, budget)$;
- 10 $C \leftarrow \text{GREEDYCOVER}(\mathcal{P} \cup Z, D)$;
- 11 **if** $L(D, \mathcal{P} \cup Z) < L(D, \mathcal{P})$ **then**
- 12 \lfloor $\mathcal{P} \leftarrow \text{PRUNE}(\mathcal{P} \cup Z, D)$;

puting the expected frequency of a sequential pattern under a null hypothesis is very complex, however [17, 13].

SQUISH builds upon and extends SQS [20]. Both draw inspiration from the KRIMP [22] and SLIM [16] algorithms. KRIMP pioneered the use of MDL for mining good patterns from transaction databases. Encoding sequential data with serial episodes is much more complicated, and hence SQS uses a much more elaborate encoding scheme. Here, we extend it to discover richer structure in the data. The SLIM algorithm [16] mines KRIMP code table directly from data. SLIM iteratively seeks to improve the current model by considering as candidates joins XY of patterns $X, Y \in CT$. Whereas SLIM considers the full Cartesian product and ranks on the basis of estimated gain, SQS and SQUISH take a batch based approach.

Lam et al. introduced GOKRIMP [7] for mining sets of serial episodes. As opposed to the MDL principle, they use fixed length codes, and do not punish gaps within patterns. This means, their goal is essentially to cover the sequence with as few patterns as possible, which is different from our goal of finding patterns that succinctly summarize the data.

Recently, Fowkes and Sutton proposed the ISM algorithm [4]. ISM is based on a generative probabilistic model of the sequence database, and uses EM to search for that set of patterns that is most likely to generate the database. ISM does not explicitly consider model complexity. Like SQUISH, ISM can handle interleaving and nesting of sequences. We will empirically compare to ISM in the experiments.

7 Experiments

Next we empirically evaluate SQUISH on synthetic and real world data. We compare against SQS [20] and ISM [4]. All algorithms were implemented in C++. We provide the code

Table 1: Database Statistics

Dataset	$ \Omega $	$ D $	$\ D\ $	$L(D, ST)$
Indep	1k	1	10k	103 630
Plant-10	1k	1	10k	103 340
Plant-50	1k	1	10k	102 630
Parallel	25	10k	1M	4 644 290
Sign	267	730	38 689	271 232
Gazelle	497	59k	209 240	1 179 030
Address	5 295	56	62 066	685 593
JMLR	3 846	788	75 646	772 112
Moby	10 277	1	105 719	1 250 149

for research purposes.¹

We evaluate quantitatively on the basis of achieved compression, pattern recall, and execution times. Specifically, we consider the compression gain $\Delta L = L(D, ST) - L(D, CT)$. That is, the gain in compression using discovered patterns versus using the singleton-only code table. Higher scores are better. All experiments were executed single threaded on quad-core Intel Xeon machines with 32GB of memory, running Linux.

Databases We consider four synthetic, and five real databases. We give their base statistics in Table 1.

Indep, *Plant-10*, and *Plant-50* are synthetic data consisting of a single sequence of 10 000 events, over an alphabet of 1000 events. For *Indep*, all events are independent. For *Plant-10*, and *Plant-50* we plant resp. 10 and 50 patterns of 5 events long 10 times each over an otherwise independent sequence, with a 10% probability of having a gap between consecutive events. To evaluate the ability of SQUISH to discover interleaved and nested patterns, we consider the *Parallel* database [4]. Each event in this database is generated by five independent parallel processes chosen at random. Each process i generates the events $\{a_i, b_i, c_i, d_i, e_i\}$ in sequence.

We further consider five real data sets. *Gazelle* is click-stream data from an e-commerce website [6]. The *Sign* database is a list of American sign language utterances [11]. To allow for interpretability we also consider text data. Here the events are the (stemmed) words in the text, with stop words removed. *Addresses* contains speeches of American presidents. *JMLR* contains abstracts from the Journal of Machine Learning research, and *Moby* is the famous novel Moby Dick by Herman Melville.

Synthetic Data As a sanity check we first compare to SQS considering *only* serial episodes and *not* allowing interleaving or nesting. We find that in this setting SQUISH performs

¹<http://eda.mmci.uni-saarland/squish/>

Table 2: Comparing SQS and SQUISH. Results for SQUISH using resp. disjoint serial episodes, interleaving serial episodes, and interleaving serial episodes and choicisodes. Given are the number of non-singleton patterns ($|\mathcal{P}|$), time to finish (t), time reach the same score as SQS (SQS- t), and the gain in compression ΔL (higher is better).

Dataset	SQUISH														
	SQS			Disjoint				Interleaving				Choicisodes			
	$ \mathcal{P} $	t	ΔL	$ \mathcal{P} $	SQS- t	t	ΔL	$ \mathcal{P} $	SQS- t	t	ΔL	$ \mathcal{P} $	SQS- t	t	ΔL
Sign	127	81s	15.5k	157	3.0s	59s	22.5k	156	4.3s	132s	22.7k	93	4.3s	103s	23.5k
Gazelle	934	26m	14.7k	880	1.5s	76m	160.4k	901	0.6s	96m	161.6k	605	1.4s	159m	165.7k
Addresses	155	5m	5.4k	181	3.9s	4m	6.5k	182	3.9s	7m	6.5k	126	3.9s	12m	7.3k
JMLR	580	8m	29.2k	583	5.4s	67m	37.2k	593	6.5s	87m	37.7k	334	5.6s	420m	40.9k
Moby	231	46m	9.6k	231	3m	23m	10.9k	328	270s	39m	10.9k	224	20.3s	66m	12.5k

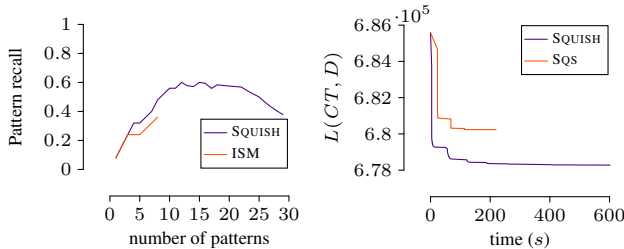


Figure 3: (left) Recall of interleaving patterns (higher is better) for SQUISH and ISM on *Parallel*. (right) Runtime of SQUISH and SQS in seconds vs. encoded length of databases for *Addresses* (lower is better).

on par with SQS in terms of recovering non-interleaving patterns from synthetic data; like SQS it correctly discovers no patterns from *Indep*, it recovers all patterns from *Plant-10*, and recovers 45 patterns exactly from *Plant-50* and fragments of the remaining 5, but does so approximately ten times faster than SQS.

To investigate how well SQUISH retrieves interleaving patterns, we consider the *Parallel* dataset, and compare to ISM. (We also considered SQS but found it did not finish within a day.) To make the comparison fair, we restrict ourselves again to serial episodes, but now do allow for interleaving and nesting. We measure success in terms of *pattern recall*. That is, given a set of patterns \mathcal{P} and a set of target patterns \mathcal{T} , we consider the set \mathcal{T} as the data and cover it with \mathcal{P} (not allowing for gaps). The pattern recall is the ratio of the total number covered events in \mathcal{T} to the maximum of the total number of events in \mathcal{T} or \mathcal{P} .

We give the results in Fig. 3. We find that SQUISH obtains much higher recall scores than ISM. Inspecting the results, we see that SQUISH discovers large fragments of each pattern, whereas ISM retrieves only eight small patterns, most of length 2, and hence does not reconstruct the generating set of patterns well.

Table 3: Sample choicisodes discovered by SQUISH in the *Addresses* and *JMLR* datasets.

Presidential Addresses	
1.	[coordin. execut.] branch govern
2.	fellow [citizen american countrymen]
3.	[discharg perform commenc] duti
4.	god [bless help]
5.	[exercise grant balanc] power
6.	power [grant vest]
7.	[eighteenth fifteenth fourteenth] amendment
8.	[guard war] against
JMLR Abstracts	
1.	[high curse low] dimension.
2.	[empirical structural] risk minimisation
3.	[independent principle] component analysis
4.	paper [proposes presents] new
5.	[Mahalanobis edit Euclidean pairwise] distance
6.	[data train] set
7.	[conditional Markov] random field
8.	[gradient coordinat.] descent

Real data Next we evaluate SQUISH on real data. We compare to SQS in terms of number of patterns, achieved compression, and runtime. We consider three different configurations, 1) disjoint covers of only serial episodes, 2) allowing interleaving and nesting of serial episodes, and 3) allowing interleaving and nesting of serial episodes and choicisodes. We give the results in Table 2.

First of all, the SQS- t columns show that in all setups SQUISH needs only a fraction of the time—up to three orders of magnitude less—to discover a model that is at least good as what SQS returns. To fully converge, SQUISH and SQS take roughly the same amount of time for the disjoint setting, as well as when we do allow interleaving. However, when converged SQUISH discovers models with much better compression rates, i.e. with much higher ΔL , than SQS does.

SQUISH is also significantly faster than ISM, taking only 87 minutes instead of 259 on the *JMLR* database, and on *Gazelle* SQUISH requires only 96 instead of 680 minutes.

SQUISH performs best when we consider our richest description language, allowing both interleaving and choicisodes, discovering much more succinct models that obtain much better scores than if we restrict ourselves. For example, for *Gazelle*, with choicisodes enabled SQUISH needs only 605 instead of 901 patterns to achieve a ΔL of 165.7k instead of 161.6k. Overall, we observe that many choicisodes form semantically coherent groups. We present a number of exemplar choicisode patterns in Table 3. Interesting examples include: data-set and training-set from *JMLR*, god-bless and god-help from Address, cape-horn and cape-cod from *Moby*.

Last, but not least, we report on the convergence of $L(CT, D)$, the encoded length of the database, over time for both SQUISH and SQS in Fig. 3. Both algorithms estimate batches of candidates, and test them one by one tests. We see that the initial candidates are highly effective on increasing compression gain. Candidates generated in the latter iterations lead to only little increase in compression gain. This leads to the possibility of executing SQUISH based upon a time budget, as an any-time algorithm.

8 Conclusion

We considered summarising event sequences. Specifically, we aimed at discovering sets of patterns that capture rich structure in the data. We considered interleaved, nested, and partial pattern occurrences. We proposed the algorithm FINDWIN to efficiently search for pattern occurrences and the greedy algorithm GREEDYCOVER for efficiently covering the data. Experiments show that SQUISH works well in practice, outperforming the state of the art by a wide margin in terms of scores and speed, while discovering pattern sets that are both more succinct and easier to interpret.

As future work we are considering parallel episodes, patterns where certain events are un-ordered e.g. $a\{b, c\}d$ [10]. Discovering such structure presents a significant computational challenges and requires novel scores and algorithms.

Acknowledgements

Apratim Bhattacharyya and Jilles Vreeken are supported by the Cluster of Excellence “Multimodal Computing and Interaction” within the Excellence Initiative of the German Federal Government.

References

[1] A. Achar, S. Laxman, R. Viswanathan, and P. Sastry. Discovering injective episodes with general partial orders. *Data Min. Knowl. Disc.*, 25(1):67–108, 2012.

[2] R. Agrawal and R. Srikant. Fast algorithms for mining association rules. In *VLDB*, pages 487–499, 1994.

[3] R. Bertens, J. Vreeken, and A. Siebes. Keeping it short and simple: Summarising complex event sequences with multivariate patterns. In *KDD*, pages 735–744, 2016.

[4] J. Fowkes and C. Sutton. A subsequence interleaving model for sequential pattern mining. In *KDD*, 2016.

[5] P. Grünwald. *The Minimum Description Length Principle*. MIT Press, 2007.

[6] R. Kohavi, C. Brodley, B. Frasca, L. Mason, and Z. Zheng. KDD-Cup 2000 organizers’ report: Peeling the onion. *SIGKDD Explor.*, 2(2):86–98, 2000. <http://www.ecn.purdue.edu/KDDCUP>.

[7] H. T. Lam, F. Mörchen, D. Fradkin, and T. Calders. Mining compressing sequential patterns. In *SDM*, 2012.

[8] S. Laxman, P. Sastry, and K. Unnikrishnan. A fast algorithm for finding frequent episodes in event streams. In *KDD*, pages 410–419. ACM, 2007.

[9] M. Li and P. Vitányi. *An Introduction to Kolmogorov Complexity and its Applications*. Springer, 1993.

[10] H. Mannila, H. Toivonen, and A. I. Verkamo. Discovery of frequent episodes in event sequences. *Data Min. Knowl. Disc.*, 1(3):259–289, 1997.

[11] P. Papapetrou, G. Kollios, S. Sclaroff, and D. Gunopulos. Discovering frequent arrangements of temporal intervals. In *ICDM*, pages 354–361. IEEE, 2005.

[12] J. Pei, H. Wang, J. Liu, K. Wang, J. Wang, and P. S. Yu. Discovering frequent closed partial orders from strings. *IEEE TKDE*, 18(11):1467–1481, 2006.

[13] F. Petitjean, T. Li, N. Tatti, and G. I. Webb. Skopus: Mining top-k sequential patterns under leverage. *Data Min. Knowl. Disc.*, 30(5):1086–1111, 2016.

[14] J. Rissanen. Modeling by shortest data description. *Automatica*, 14(1):465–471, 1978.

[15] J. Rissanen. A universal prior for integers and estimation by minimum description length. *Annals Stat.*, 11(2):416–431, 1983.

[16] K. Smets and J. Vreeken. SLIM: Directly mining descriptive patterns. In *SDM*, pages 236–247. SIAM, 2012.

[17] N. Tatti. Ranking episodes using a partition model. *Data Min. Knowl. Disc.*, 29(5):1312–1342, 2015.

[18] N. Tatti and B. Cule. Mining closed episodes with simultaneous events. In *KDD*, pages 1172–1180, 2011.

[19] N. Tatti and B. Cule. Mining closed strict episodes. *Data Min. Knowl. Disc.*, 25(1):34–66, 2012.

[20] N. Tatti and J. Vreeken. The long and the short of it: Summarizing event sequences with serial episodes. In *KDD*, pages 462–470. ACM, 2012.

[21] N. Vereshchagin and P. Vitányi. Kolmogorov’s structure functions and model selection. *IEEE TIT*, 50(12):3265–3290, 2004.

[22] J. Vreeken, M. van Leeuwen, and A. Siebes. KRIMP: Mining itemsets that compress. *Data Min. Knowl. Disc.*, 23(1):169–214, 2011.

[23] J. Wang and J. Han. Bide: Efficient mining of frequent closed sequences. *ICDE*, 0:79, 2004.

A Appendix

A.1 Estimating Candidates Here we describe our heuristic strategy for finding new candidates of the form XY as in Sec. 5.1. First, we need two crucial observations.

Constant Time Difference Estimation Given a database D and an cover C . Let P and Q be two patterns. Let $V = \{v_1, \dots, v_N\}$ and $W = \{w_1, \dots, w_N\}$ be two set of windows for P and Q , respectively. Both V and W occur in C . Each of these windows v_i and w_i occur in the same sequence. Given the start positions and end positions of the pattern in sequence k_i , we can write them as $v_i = (a_i, b_i, P, k_i)$ and $w_i = (c_i, d_i, Q, k_i)$. Let U be the set of windows produced by combining them, $U = (a_1, d_1, R, k_1), \dots, (a_N, d_N, R, k_N)$. Let the windows in U be disjoint and the windows in U be disjoint with the windows in $C \setminus (V \cup W)$. Then the difference $L(D, C \cup U \setminus (V \cup W)) - L(D, C)$ depends only N , $gaps(V)$, $gaps(W)$, and $gaps(U)$ and can be computed in constant time from these values.

Shorter Windows in Optimal Cover Given a database D and an cover C . Let $v = (i, j, X, k) \in C$. Assume that there exists a window $S[a, b]$ containing X such that $w = (a, b, X, l)$ does not overlap with any window in C and $b - a < j - i$. Then C is not an optimal cover.

We refer the reader to [20] for detailed proofs.

We present our heuristic procedure ESTIMATE as pseudo-code in Algorithm 4. In this algorithm, given pattern X and a cover C , for a possible extension Y , we enumerate the windows of XY from the shortest to the longest. These windows are constructed by combining two windows in the cover C . We maintain the sets V_Y , W_Y and U_Y (line 1), containing windows of X , windows of Y (to be combined together), and new windows of XY (resulting from the combination) respectively. We do this for every possible extension Y in the code table. At each step we compute the difference in code length of using these windows instead. We maintain d_Y to store this difference. By the observation **Constant Time Difference Estimation**, this can be done in constant time. We prefer patterns XY which are frequently occurring, with more fills than other meta stream characters. Thus, we want to find shorter windows of XY first. Such a set of windows U could potentially lead to a estimated decrease in code length. Therefore, to ensure that we find shorter windows first and efficiency, we search for all windows (all possible Y) simultaneously using a priority queue T and look only at windows in the cover C . For each window of X in the cover C , we look at windows after it to construct windows XY (Y is the pattern of the window following the window of X). We initialize the priority queue T with these windows (line 4-9), sorted based on length. At each step of the candidate generation algorithm, we retrieve

Algorithm 4: ESTIMATE(X, C, D). Heuristic for finding pattern XY with low $L(D, CT \cup XY)$

```

input : database  $D$ , cover  $C$ , and pattern  $X$ 
output : pattern  $XY$  with low  $L(D, CT \cup XY)$ 
1 foreach  $Y \in CT$  do
2    $V_Y \leftarrow \emptyset; W_Y \leftarrow \emptyset; U_Y \leftarrow \emptyset; d_Y \leftarrow 0;$ 
3    $T \leftarrow \emptyset;$ 
4   foreach window  $v$  of  $X$  in cover  $C$  do
5      $(a, b, X, k) \leftarrow v;$ 
6      $d \leftarrow$  end index of window following  $v$  in  $C;$ 
7      $t \leftarrow (v, d, 0); l(t) \leftarrow d - a;$ 
8     add  $t$  into  $T;$ 
9   while  $T$  is not empty do
10     $t \leftarrow \arg \min_{u \in T} l(u);$ 
11     $(v, d, s) \leftarrow t; a \leftarrow$  first index of  $v;$ 
12     $w = (c, d, Y, k) \leftarrow$  active  $w$  of  $Y$  ending at  $d;$ 
13    if  $Y = X$  and (event at  $a$  or  $d$  is marked) then
14      delete  $t$  from  $T;$ 
15      continue;
16    if  $S_k[a, d]$  is a minimal window of  $XY$  then
17      add  $v$  into  $V_Y;$ 
18      add  $w$  into  $W_Y;$ 
19      add  $(a, d, XY, k)$  into  $U_Y;$ 
20       $d_Y \leftarrow \min(div(V, W, U; A) + s, d_Y);$ 
21      if  $|Y| > 1$  then
22         $s \leftarrow s + gain(w);$ 
23      if  $Y = X$  then
24        mark the events at  $a$  and  $d;$ 
25        delete  $t$  from  $T;$ 
26        continue;
27    if  $w$  is the last window in the sequence then
28      delete  $t$  from  $T;$ 
29    else
30       $d \leftarrow$  end index of the active  $w'$  following  $w;$ 
31      update  $t$  to  $(v, d, s)$  and  $l(t)$  to  $d - a;$ 

```

once such window of XY from the priority queue T (line 11) add it to our list U_Y of windows of XY and estimate the change in code length (line 22). As we do not allow overlaps, we need to ensure that windows in U_Y are not overlapping. If a window of XY overlaps with any other window in C , we cannot use both of these windows at the same time. We take this into account by subtracting the $gain(w)$ of this window w overlapping with the window of XY (line 23) [20]. The $gain(w)$ of a window w if a upper bound on the bits gained by encoding the events in the database with this window vs. encoding them as singletons. We define the gain as in [20] for a window w of the pattern Y ($S_k[i, j]$),

$$\begin{aligned}
gain(w) = & -L(code_p(X)) - (j - i - |X|)L(code_g(X)) \\
& - (|X| - 1)L(code_f(X)) + \sum_{x \in X} L(code_p(x)) .
\end{aligned}$$

Overlapping could also happen if $Y = X$. So we simply check if the adjacent scans have already used these two instances of X for creating a window for pattern XX (line 25). We now extend our search by looking at the window following the currently considered window of Y in the cover C (line 34). As we allow interleaving and nesting in our covers, we also look at possible windows Y occurring inside or interleaved with windows of other patterns. That is, we look at singletons inside gaps of windows. For each window X in the cover C , we look at all windows following it, until we reach the window of X or the end of the cover.

A.2 Pruning the Code Table Here, we present the algorithm we use to prune the code table CT , used at line 12 of SQUISH as pseudo-code in Algorithm 5.

Algorithm 5: PRUNE(\mathcal{P}, D)

input : pattern set \mathcal{P} , database D
output pruned pattern set \mathcal{P} ;
:
1 **foreach** $X \in \mathcal{P}$ **do**
2 $CT \leftarrow$ code table corresponding to
 GREEDYCOVER(D, \mathcal{P});
3 $CT' \leftarrow$ code table obtained from CT by deleting
 X ;
4 $g \leftarrow \sum_{w=(i,j,X,k) \in C} gain(w)$;
5 **if** $g < L(CT) - L(CT')$ **then**
6 **if** $L(D, \mathcal{P} \setminus X) < L(D, \mathcal{P})$ **then**
7 $\mathcal{P} \leftarrow \mathcal{P} \setminus X$;
