

# SECRET: Mining Rule Sets from Event Sequences

Aleena Siji,<sup>1</sup> Joscha Cüppers,<sup>2</sup> Osman Mian,<sup>3</sup> Jilles Vreeken<sup>2</sup>

<sup>1</sup> Helmholtz AI, Munich

<sup>2</sup> CISPA Helmholtz Center for Information Security

<sup>3</sup> Institute for AI in Medicine IKIM

aleena.siji@helmholtz-munich.de, joscha.cueppers@cispa.de, osman.mian@uk-essen.de, vreeken@cispa.de

## Abstract

Summarizing event sequences is a key aspect of data mining. Most existing methods neglect conditional dependencies and focus on discovering sequential patterns only. In this paper, we study the problem of discovering both conditional and unconditional dependencies from event sequences. We do so by discovering rules of the form  $X \rightarrow Y$  where  $X$  and  $Y$  are sequential patterns. Such rules are simple to understand and provide a clear description of the relation between the antecedent and the consequent. To discover a succinct and non-redundant set of rules we formalize the problem in terms of the Minimum Description Length principle. As the search space is enormous and does not exhibit helpful structure, we propose the SECRET method to discover high-quality rule sets in practice. Through extensive empirical evaluation we show that unlike the state of the art, SECRET ably recovers the ground truth on synthetic datasets and finds useful rules from real datasets.

## 1 Introduction

In many applications data naturally takes the form of events happening over time. Examples include industrial production logs, the financial market, device failures in a network, medical treatment plans etc. Existing methods for analyzing event sequences primarily focus on mining unconditional, frequent sequential patterns (Agrawal and Srikant 1995; Mannila and Meek 2000; Tatti 2009). Real world processes are often more complex than this, as they often include conditional dependencies. For example, the formation of tropical cyclones ( $C$ ) in the Bay of Bengal is often but not always followed by heavy rainfall ( $R$ ) on the coast. Knowing such a relationship is helpful both in predicting events and in understanding the underlying data generating mechanisms.

In this paper, we are interested in discovering rules of the form  $X \rightarrow Y$  from long event sequences, where  $X$  and  $Y$  are sequential patterns. Existing methods for mining such rules either suffer from pattern explosion, i.e. are prone to returning orders of magnitude more results than we can possibly analyze (Fournier-Viger et al. 2021; Chen et al. 2021),

or have very limited expressivity, e.g. require the constituent events to occur in contiguous order (Bourrand et al. 2021).

We aim to discover succinct sets of rules that generalize the data well and explicitly allow for gaps between the occurrences of the constituent events of these rules. To ensure compact and non-redundant results, we formalize the problem using the Minimum Description Length (MDL) principle (Grünwald 2007). Loosely speaking, we are after that set of sequential rules that together compresses the data best.

However, the problem we so arrive at is computationally challenging. For starters, there exist exponentially many rules, exponentially many rule sets, and then again exponentially many ways to describe the data given a set of rules. Moreover, the search space does not exhibit structure we can use to efficiently obtain the optimal result. To mine good rule sets from data we therefore propose the greedy SECRET algorithm. We introduce two variants. SECRET-CANDIDATES constructs a good rule set from a set of candidate patterns by splitting them into high-quality rules. SECRET-MINE, on the other hand, only requires the data and mines a good rule set from scratch. Starting from a model of singleton rules, it iteratively extends them into more refined rules. To avoid testing all possible extensions, we consider only those extensions that occur significantly more often than expected.

Through extensive evaluation, we show that both variants of SECRET work well in practice. On synthetic data we show that they are robust to noise and recover the ground truth well. On real-world data, we show that SECRET returns succinct sets of rules that give clear insight into the data generating process. This is in stark contrast to existing methods which either return many thousands of rules (Fournier-Viger et al. 2021) or are restricted to rules where events occur contiguously (Bourrand et al. 2021).

To summarize, the main contributions are as follows:

- (a) We define a pattern language for fully ordered sequential rules that accommodates for gaps.
- (b) We present SECRET-CANDIDATES for constructing a high-quality rule set, given a set of sequential patterns.
- (c) We present SECRET-MINE for mining a high-quality rule set given a database of event sequences.
- (d) We extensively evaluate SECRET on synthetic and real-world datasets, comparing it to the state-of-the-art.<sup>1</sup>

<sup>1</sup>Work done while at CISPA Helmholtz Center for Information Security.

Copyright © 2026, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

<sup>1</sup>We make code and data available at [eda.rg.cispa.io/prj/secret/](https://eda.rg.cispa.io/prj/secret/)

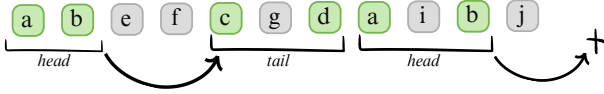


Figure 1: Toy example of rule  $ab \rightarrow cd$  in a sequence. Head  $ab$  triggers the rule twice. Tail  $cd$  follows only once. Hence,  $\text{supp}(ab \rightarrow cd) = 1$  and  $\text{conf}(ab \rightarrow cd) = 0.5$ .

## 2 Preliminaries

In this section we introduce basic notation and give a short introduction to the MDL principle.

### Notation

As data we consider a database  $D$  of  $|D|$  event sequences. A sequence  $S \in D$  consists of  $|S|$  events drawn from a finite alphabet  $\Omega$  of discrete events  $e \in \Omega$ . We denote the total number of events in the data as  $\|D\|$ . We write  $S_t$  for the  $t^{\text{th}}$  sequence in  $D$ , we omit the subscript whenever clear from context. We write  $S[i]$  to refer to the  $i^{\text{th}}$  event in sequence  $S$ , and  $S[i, j]$  for the subsequence from the  $i^{\text{th}}$  up to and including the  $j^{\text{th}}$  event. We denote an empty sequence by  $\epsilon$ .

A serial episode  $X$  is a sequence of  $|X|$  events drawn from  $\Omega$ . A sequential rule  $r$  captures the conditional dependence between serial episodes  $X$  and  $Y$ . Intuitively, it expresses that whenever we see  $X$  in the data it is more likely that  $Y$  will follow. We refer to  $X$  as the *head* of  $r$ , denoted  $\text{head}(r)$ , and to  $Y$  as the *tail* of  $r$ , denoted  $\text{tail}(r)$ . If  $X$  is an empty pattern,  $X = \epsilon$ , we call  $X \rightarrow Y$  an *empty head rule*. Empty head rules where  $|Y| = 1$  are called *singleton* rules.

A subsequence  $S[i, j]$  is a window of pattern  $X$  iff  $X$  is a subsequence of  $S[i, j]$ , and we say  $S[i, j]$  *matches*  $X$  and that  $X$  occurs in  $S[i, j]$ . A pattern window  $S[i, j]$  is *minimal* for  $X$  iff no proper sub-window of  $S[i, j]$  matches  $X$ . A rule window of  $r$  is a tuple of two pattern windows  $S[i, j]$  and  $S[k, l]$  when  $S[i, j]$  matches  $\text{head}(r)$ ,  $j < k$ , and  $S[k, l]$  matches  $\text{tail}(r)$ , we denote it by  $S[i, j; k, l]$ .

We say a window  $S[i, j]$  *triggers* rule  $r$  when it is a minimal window of  $\text{head}(r)$ . A rule window  $S[i, j; k, l]$  *supports* a rule  $r$  if  $S[i, j]$  triggers  $\text{head}(r)$  and  $S[k, l]$  matches  $\text{tail}(r)$ . We call the number of events that occur in a rule window between the rule head and the rule tail,  $k - j - 1$ , the *delay* of the rule instance. We give an example in Fig. 1. We denote the number of windows over all sequences  $S \in D$  that trigger a rule  $r$  as the trigger count  $\text{trigs}(r)$ . We define the *support* of a rule  $r$  as the number of rule windows  $S[i, j; k, l]$  in  $D$  where  $S[k, l]$  is a minimal window of  $\text{tail}(r)$  and follows the head with minimum delay. Finally, we define the confidence of a rule  $r$  as its support relative to its trigger count, formally  $\text{conf}(r) = \text{supp}(r) / \text{trigs}(r)$ .

### Minimum Description Length Principle

The Minimum Description Length (MDL) (Grünwald 2007) is a computable and statistically well-founded approximation of Kolmogorov complexity (Li and Vitányi 1997). For a given model class  $\mathcal{M}$  it identifies the best model  $M \in \mathcal{M}$  as the one that minimizes the number of bits for describing both model and data without loss. Formally, minimize

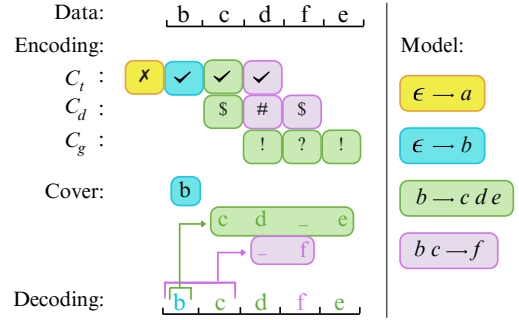


Figure 2: Toy example of encoding sequence  $S$  using rule set  $R$ .  $C_t$  encodes if a triggered rule *hits* or *misses*.  $C_d$  encodes the delay between the trigger and the rule tail.  $C_g$  encodes the gaps in the rule tail. Together, they form a *cover*  $C$ .

$L(M) + L(D | M)$  where  $L(M)$  is the length in bits of model  $M$  and  $L(D | M)$  is the length in bits of data  $D$  given  $M$ . This is known as two-part, or crude MDL—in contrast to one-part, or refined MDL, which is not computable for arbitrary model classes. We use two-part MDL because we are interested in the model, i.e. the set of rules that describes the data best. In MDL we are never concerned with materialized codes but only code lengths. To use MDL, we have to define a model class and code length functions, we do so next.

## 3 MDL for Sequential Rules

We now formally define the problem. We consider sets of sequential rules as our model class  $\mathcal{R}$ . By MDL, we are interested in that set  $R \in \mathcal{R}$  that best describes data  $D$ .

### Decoding an Event Sequence

Before we formally define how we *encode* models and data, we give the main intuition by *decoding* an already encoded sequence. We give an example in Fig. 2.

To decode a symbol, we consider the rules from  $R$  that are currently triggered. For those, we read codes from the trigger stream  $C_t$ . Initially, the context is empty and hence only empty-head rules trigger. The first trigger code is a *miss* for singleton rule  $\epsilon \rightarrow a$ , meaning that the rule tail does not follow. The second trigger code is a *hit* for rule  $\epsilon \rightarrow b$ . As empty-head rules do not incur delays, we can immediately write  $b$  as the first symbol of the sequence.

This triggers rule  $b \rightarrow cde$ . We hence read a code from the trigger stream  $C_t$ , and find that it is a hit. As this rule does not have an empty head, there may be a delay between the head and its tail and we read a code from the delay stream  $C_d$ . It is a start code, so we write the first symbol of the tail ( $c$ ), this creates a minimal window of  $bc$ .

This triggers rule  $bc \rightarrow f$ . We read from  $C_t$  to find that it is a hit, then from  $C_d$  to find that its tail is delayed. To determine if we write the next symbol from tail  $cde$ , we read from the gap stream  $C_g$ . It has a fill code, meaning no gap, so we write  $d$ . This time, no new rule triggers. Tail  $cde$  is not yet completely decoded and  $f$  is delayed. For each delayed tail we read a code from  $C_d$ , and for each incomplete tail we

read a code from  $C_g$ . Here, we read a start code for tail  $f$  and a gap code for tail  $cde$ , we hence write  $f$ . Again, no new rule is triggered. Now only tail  $cde$  is not yet fully decoded. We read a fill code from  $C_g$  and write  $e$  as the last symbol.

To summarize, sequences are encoded from left to right, and rules automatically trigger whenever we observe a minimal window of the head. For each trigger, we encode whether the tail follows using a *hit* or *miss* code. When a rule hits, we encode whether its tail follows immediately or later, using a *start* resp. *delay* code. We encode whether *gaps* occur in the rule tail using *fill* and *gap* codes. Empty-head rules never incur a delay. To avoid unnecessary triggers, we only encode those empty-head rules if no other rule encodes the current symbol (e.g. all active tails say ‘gap’).

### Computing the Description Lengths

Now that we have the intuition, we can formally describe how to encode a model, respectively the data given a model.

**Encoding a Model** A model  $R \in \mathcal{R}$  is a set of rules. To reward structure such as chains where the tail of one rule is the head of another (e.g.  $r_1 = \epsilon \rightarrow AB$ , and  $r_2 = AB \rightarrow CD$ ), we first encode the set  $P$  of all non-empty, non-singleton heads and all non-singleton tails. Formally,

$$P = \{\text{head}(r) \mid \forall r \in R\} \cup \{\text{tail}(r) \mid \forall r \in R\} \setminus (\Omega \cup \epsilon).$$
 The encoded length  $L(P)$ , is defined as

$$L(P) = L_{\mathbb{N}}(|P| + 1) + \sum_{p \in P} L_{\mathbb{N}}(|p|) + |p| \log_2(|\Omega|),$$

where we first encode the number of these patterns using  $L_{\mathbb{N}}$ , the MDL-optimal encoding for integers (Rissanen 1983). It is defined for  $z \geq 1$  as  $L_{\mathbb{N}}(z) = \log^* z + \log c_0$  where  $\log^* z$  is the expansion  $\log z + \log \log z + \dots$  including only the positive terms. To ensure this satisfies the Kraft inequality, we set  $c_0 = 2.865064$  (Rissanen 1983). Since  $P$  can be empty and  $L_{\mathbb{N}}$  is only defined for numbers  $\geq 1$  we offset it by one. Next, we encode each pattern  $p \in P$  where we use  $L_{\mathbb{N}}$  to encode its length and then choose each subsequent symbol  $e \in p$  out of alphabet  $\Omega$ .

Now that we have the set of all heads and tails, we have

$$L(R \mid P) = L_{\mathbb{N}}(|R| + 1) + |R| \cdot (\log_2(|P| + |\Omega| + 1) + \log_2(|P| + |\Omega|)),$$

as the encoded length in bits for a set of rules. We first encode the number of rules. As  $R$  can be empty, we again offset by one. Next, for each rule  $r \in R$ , we choose its head from  $P \cup \Omega$ , and then its tail from  $P \cup \Omega$ .

Putting this together, the number of bits to describe a rule set  $R \in \mathcal{R}$  without loss is

$$L(R) = L(P) + L(R \mid P).$$

**Encoding Data given a Model** As we saw in the example, reconstructing the data requires three code streams  $C_t$ ,  $C_d$ , and  $C_g$ . For an arbitrary database we additionally need to know how many sequences it includes, and their lengths. Formally, the description length of data  $D$  given model  $R$  is

$$L(D \mid R) = L_{\mathbb{N}}(|D|) + \sum_{S \in D} L_{\mathbb{N}}(|S|) + L(C_t) + L(C_d) + L(C_g).$$

To encode the code streams  $C_t, C_d, C_g$  we use prequential codes (Grünwald 2007). These codes are asymptotically optimal without requiring us to make arbitrary choices in how to encode the code distributions. Formally, we have

$$L(C_j) = - \sum_{i=1}^{|C_j|} \log_2 \frac{\text{usg}_i(C_j[i] \mid C_j) + c}{i + \text{unique}(C_j) \cdot c},$$

where  $\text{usg}_i(C_j[i] \mid C_j)$  denotes the number of times  $C_j[i]$  has been used in  $C_j$  up to the  $i^{\text{th}}$  position,  $\text{unique}(C_j)$  denotes the number of unique symbols in  $C_j$ , and  $c$  is a small constant. As common in prequential coding, we set  $c$  to 0.5.

### The Problem, Formally

We can now formalize the problem we aim to solve.

**The Sequential Rule Set Mining Problem** *Given a sequence database  $D$  over alphabet  $\Omega$ , find the smallest rule set  $R \in \mathcal{R}$  and cover  $C$  such that the total encoded size*

$$L(R) + L(D \mid R)$$

*is minimal.*

The resulting search space is enormous. To begin with, there exist super-exponentially many covers of  $D$  given  $R$ . The optimal cover depends on the code lengths, which in turn depend on code usages. Even if the optimal cover is given, the problem of finding the optimal rule set is super-exponential: there exist exponentially many patterns  $p$  in the size of alphabet  $\Omega$ , exponentially many rules  $r$  in the number of patterns, and exponentially many sets of rules. None of these sub-problems exhibit substructure, e.g. monotonicity or submodularity, that we can exploit to find the optimal solution. Hence, we resort to heuristics.

## 4 The Secret Algorithm

In this section we introduce SEQRET, for discovering high-quality **sequential rule-sets** from data. We break the problem down into two parts: optimizing the description of the data given a rule set, and mining good rule sets. For the latter we propose two variants, SEQRET-CANDIDATES for doing so given a set of candidate patterns, and SEQRET-MINE for mining rule sets directly from data.

### Selecting a Good Cover

A lossless description of  $D$  using rules  $R$  correspond to a set of rule windows such that each event  $e$  in  $D$  is *covered* by exactly one window. We are after that cover which minimizes  $L(D \mid R)$ . Finding the optimal cover is infeasible, hence we settle for a good cover and find one greedily.

The idea is to define an order over the rule windows and greedily select the next best window until the data is completely covered. To minimize the encoding length, we prefer covering as many events as possible with a single rule. Therefore, we prefer rules with long tails, high confidence, and high support. Among windows of otherwise equally good rules, we prefer those with lower delays and fewer gaps. We consider the starting position of the rule tail as

a final tie breaker. Combining this, we define WINDOW ORDER as descending on  $|tail(r)|$ ,  $conf(r)$ , and  $supp(r)$ , then ascending on  $l - j - |tail(r)|$  and  $k$ , where  $r$  is a rule and  $S[i, j; k, l]$  is a rule window.

To find a good cover, we start with an empty set and greedily add rule windows to it in WINDOW ORDER. To avoid searching for all possible rule windows in the beginning, we only consider the best window per rule trigger – defined as the one with the fewest gaps in its rule tail window, followed by lowest delay – and look for the next best only if we do not select the former due to conflicts, i.e. when its constituent events are already covered by a previously selected window. To avoid evaluating hopeless windows, we limit ourselves to those within user-set  $max\ delay$  ratio and  $max\ gap$  ratio. We give the pseudocode of COVER algorithm with worst-case time complexity in the appendix.<sup>2</sup>

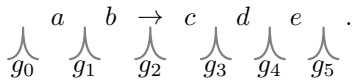
### Selecting Good Rule Sets

Next, we consider the problem of discovering good rule sets. We first propose an approach for doing so given a set of sequential patterns as input. The intuition is that, when the ground truth includes a sequential rule  $X \rightarrow Y$ , a good sequential pattern miner will return  $XY$ . Then we can reconstruct the ground truth by considering splits of candidate patterns  $XY$  into candidate rules  $X \rightarrow Y$ .

We propose SECRET-CANDIDATES, for which we give the pseudocode and worst-case time complexity in the appendix. We first initialize the rule set with all singleton rules to ensure that the data can be encoded losslessly. We then iterate over the candidate patterns in descending order of their contribution to compression (Tatti and Vreeken 2012). We split each pattern into candidate rules – for example, pattern  $abc$  generates candidate rules  $e \rightarrow abc$ ,  $a \rightarrow bc$ , and  $ab \rightarrow c$  – and choose the one whose addition to the rule set minimizes the total description length, with  $L(D|R)$  determined by the COVER procedure. We keep it in our model if it improves the score and discards otherwise. We iterate this process for all candidate patterns.

### Generating Good Rules

Next, we move our attention to mining good rule sets directly from data. The first step is to generate good candidate rules. Suppose we know a rule  $r$ . We consider extending it with events  $e \in \Omega$  that occur significantly more often than expected by chance, either within or directly adjacent to the rule windows of  $r$ . A rule has  $|r| + 1$  gap positions, i.e. before its first event, between its constituent events, and after its last event. For example, rule  $ab \rightarrow cde$  has 6 gap positions,



For each rule  $r$ , we test for every gap position  $g_i$ , if  $e \in \Omega$  is more frequent than expected in its rule windows. Our null hypothesis is

$$H_0 : \sum_{w \in B} \mathbf{1}(e \in g(i, w)) \leq \sum_{w \in B} \Pr(e \in g(i, w))$$

<sup>2</sup>Appendix is available at [eda.rg.cispa.io/prj/secret/](http://eda.rg.cispa.io/prj/secret/)

where  $B$  is the set of best rule windows of  $r$  as used in COVER procedure,  $w \in B$  is one such window of rule  $r$ , and  $g(i, w)$  a function that returns gap  $i$  from window  $w$ . The probability of  $e$  occurring in gap  $g_i$  with length  $|g(i, w_p)|$  in a rule window  $w_p$ , is given by

$$\Pr(e \in g(i, w_p)) = 1 - \left(1 - \frac{supp(\epsilon \rightarrow e)}{|D|}\right)^{|g(i, w_p)|}$$

To test for statistical significance, we model the expected neighborhood as Poisson binomial distribution (Wang 1993), where trials are the rule windows, and success probability per trial is decided by the length of a specific gap position. We correct for multiple hypothesis testing with Benjamini-Hochberg procedure (Benjamini and Hochberg 1995). As computing the exact CDF for high number of trials is expensive (Le Cam 1960; Volkova 1996), we use the fast normal approximation with continuity correction (Hong 2013) for cases where the number of trials, i.e.  $supp(r)$ , is greater than or equals 50. Below that, we use FFT-based exact CDF computation.

If event  $e$  is measured to be significantly more frequent in  $g_i$  than expected, we generate a new rule by inserting  $e$  at the position of  $g_i$  in the rule. We show the pseudo-code for the CANDRULES procedure in the appendix.

### Mining Good Rule Sets

Finally, we describe SECRET-MINE for mining good rule sets directly from data. We provide the pseudocode as Algorithm 1. We initialize rule set  $R$  with all the singleton rules (line 1). Next, we consider adding candidates based on the rules already in the model (line 3). As we want to generate the most promising candidate rules first, we start with rules with high support and high confidence. We define a greedy EXTEND ORDER as 1)  $\uparrow supp(r)$ , 2)  $\uparrow conf(r)$ , 3)  $\uparrow |tail(r)|$  and 4)  $\uparrow |head(r)|$ , where  $\uparrow$  indicates that higher values are preferred. For each, we generate a set of candidate rules as described previously, we test them for addition in the order of their p-values (line 4).

We add those rules into the model whose inclusion results in significant reduction in total encoded size (line 5). We use the no-hypercompression inequality (Bloem and de Rooij 2020; Grünwald 2007) to test for significance at level  $\alpha$ , writing  $\ll_{\alpha}$  for “significantly less”.<sup>3</sup> When adding a candidate rule  $r'$  to the model does not improve compression, we test if replacing the rule  $r$  we generated it from leads to a better compression (line 7). To ensure we can always describe the data without loss, we never remove singleton rules.

After adding a new rule, SECRET-MINE performs a pruning step to remove existing rules that may have become redundant or obsolete (line 10). The PRUNE method iterates over the non-singleton rules in the model and removes those whose exclusion reduces the total encoded size. We do so in PRUNE ORDER where we consider rules in order of lowest usage, highest encoded size, and lowest tail length.

We repeat generating candidate rules, adding them, and pruning redundant rules until convergence. Convergence is

<sup>3</sup>We set  $\alpha = 0.05$ , corresponding to a minimum gain of 5 bits.

guaranteed as our score is lower bounded by 0. We discuss the worst-case time complexity in the appendix.

---

**Algorithm 1: SEQRET-MINE**


---

**Input:** Sequence database  $D$ , Significance level  $\alpha$   
**Output:** Rule set  $R$

```

1  $R \leftarrow \{\epsilon \rightarrow e \mid \forall e \in \Omega\}$ ;
2 do
3   for  $r \in R$  in EXTEND ORDER do
4     for  $r' \in \text{CANDRULES}(D, r)$  in order of
        $p$ -value do
5       if  $L(D, R \cup \{r'\}) \ll_{\alpha} L(D, R)$  :
6          $R \leftarrow R \cup \{r'\}$ ;
7       else if  $r \notin \{\epsilon \rightarrow e \mid \forall e \in \Omega\}$  and
          $L(D, R \cup \{r'\} \setminus \{r\}) \ll_{\alpha} L(D, R)$  :
8          $R \leftarrow (R \setminus \{r\}) \cup \{r'\}$ ;
9       if  $R$  updated :
10         $R \leftarrow \text{PRUNE}(D, R)$ ;
11        continue with next  $r$ 
12 while  $R$  updated;
13 return  $R$ 

```

---

## 5 Related Work

Gaining insight from sequential data is a widely studied topic, ranging from temporal logic learning (Roy et al. 2023) to process mining (Van Der Aalst 2016). While temporal logic learning is concerned with strict, formal logic rules, and process mining with describing the entire, overarching process, our work concentrates specifically on discovering conditional dependencies from event sequences.

Classical methods for sequential pattern mining focus on extracting all frequent patterns and suffer from pattern explosion leading to excessively many, largely redundant, and often spurious results (Fournier-Viger et al. 2017). Mining closed frequent patterns alleviates this, but is sensitive to noise (Yan, Han, and Afshar 2003; Wang and Han 2004).

Pattern set mining avoids the pattern explosion by instead scoring *sets* of patterns. The Minimum Description Length principle has been shown to be a robust criterion for identifying good pattern sets in practice (Galbrun 2022; Cüppers et al. 2024; Tatti and Vreeken 2012). Different pattern languages, scores, and methods have been proposed, such as those allowing gaps (Tatti and Vreeken 2012) and interleaved patterns (Bhattacharyya and Vreeken 2017; Cüppers, Krieger, and Vreeken 2024).

Classical methods for rule mining in event sequences operate similar to frequent pattern mining, but in addition to the frequency requirement also impose a minimum confidence threshold. Various approaches have been proposed to address different data modalities, such as rules over itemsets ordered by time (Fournier-Viger et al. 2014, 2012), or over events in sequences (Dalmas, Fournier-Viger, and Norre 2017; Zaki 2001; Cule and Goethals 2010). The former generally count the number of sequences containing a rule as its support, whereas the latter use sliding or minimal windows to count rule occurrences within a sequence.

Rules can be further categorized into partially ordered rules (Fournier-Viger et al. 2021; Chen et al. 2021) and fully sequential rules (Zaki 2001). These approaches all consider the quality of individual patterns and hence suffer from pattern explosion. To address this, Fournier-Viger and Tseng propose TNS (Fournier-Viger and Tseng 2013), a method that reports top- $k$  non redundant rules, but its notion of redundancy is unable to detect semantically redundant rules.

Most closely related to SEQRET are existing methods which use MDL to select or mine rules. Existing rule set miners for event sequences either filter down an existing set of rules (Chen et al. 2022), do not allow for gaps (Bourrand et al. 2021) or is a supervised method requiring a target (Cüppers, Kalofolias, and Vreeken 2022). Thus, no existing method directly addresses the problem we consider.

## 6 Experiments

In this section we empirically evaluate SEQRET-CANDIDATES and SEQRET-MINE.

We compare to POERMA (Fournier-Viger et al. 2021) and POERMH (Chen et al. 2021) as representative frequent rule miners, to TNS (Fournier-Viger and Tseng 2013) as a top  $k$  non-redundant rule set miner, to COSSU (Bourrand et al. 2021) as an MDL-based rule set miner, and to SQS (Tatti and Vreeken 2012) and SQUISH (Bhattacharyya and Vreeken 2017) as MDL-based sequential pattern miners.

As candidate patterns for SEQRET-CANDIDATES we use the output of SQS (Tatti and Vreeken 2012). For TNS we set  $k$  to the number of rules in the ground truth, and for POERMA and POERMH, the minimum support and minimum confidence values according to the ground truth. For real datasets where the ground truth is unknown, we set  $k$  as the number of rules returned by SEQRET-MINE. Further, we use a minimum support threshold of 10 where feasible, and 20 otherwise. We allow all methods a maximum runtime of 24 hours, except for COSSU, which generally took longer and is allowed upto 48 hours. We provide details in the appendix.

### Experiments on Synthetic Data

We first consider data with known ground truth. For a given alphabet size, number of rules, rule sizes, and confidence, we first generate a rule set  $R$  by selecting events from alphabet  $\Omega$  uniformly at random, with replacement, to form the rule.

**Data Generation** Given a randomly generated rule set  $R$ , we next generate an event sequence  $S$ . We first generate background noise by sampling uniformly at random from  $\Omega$ . Next, we plant patterns, i.e. empty-head rules in  $R$  by sampling them uniformly and writing to  $S$  at random positions. We make sure to not overwrite existing rules. Finally, we go over the generated sequence and wherever a non-empty-head rule is triggered, we sample as per the desired rule confidence whether the trigger is a hit or a miss. If it is a hit, we sample the delay and insert the corresponding rule tail. We provide further details in the appendix. Unless stated otherwise, we generate sequences of length 10 000 over alphabets of size 500, rule sets of size 20 with rule confidence 0.75. We generate 20 datasets per configuration in each experiment.

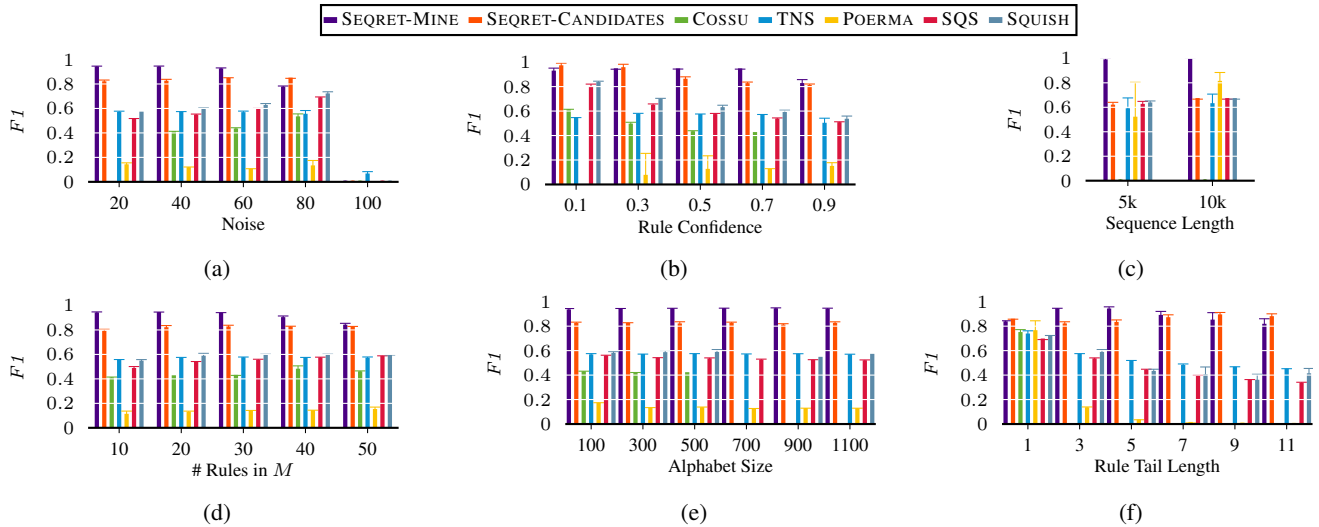


Figure 3: [Higher is better]  $F1$  scores for synthetic data. We observe that SEQRET is robust against (a) high noise, (b) low rule confidence, (d) varying number of true rules, (e) large alphabets, and (f) varying rule tail lengths. In (c) we evaluate rule recovery where heads and tails are only as frequent as by chance, only SEQRET-MINE still picks up the ground truth reliably.

**Evaluation Metric** As evaluation metric, we consider the  $F1$  score. To reward partial discovery, we base the recall and precision on the similarity between discovered rules and ground truth rules, computed using the Levenshtein edit distance without substitution, i.e the longest common subsequence distance (Navarro 2001). The similarity between two rules is a weighted average of the similarity between their respective heads, tails and complete rules. We then compute recall and precision using the similarity scores following Cüppers and Vreeken (2020). We provide full details of the evaluation metric in the appendix.

**Destructive Noise** We first evaluate robustness of SEQRET against destructive noise. To this end, we add noise to generated data by flipping individual events  $e \in S$  with probability between 0.2 and 1. We show the results in Figure 3a. We observe that SEQRET is robust against noise and recovers the ground truth well even at 80% noise. At 100% noise, there is no structure in the data and all methods except TNS correctly discovers no rules. In all experiments, POERMA performs better than or comparable to POERMH, we hence omit POERMH from results. Further, as SQUISH regularly crashes, we report averages over finished runs only.

**Rule Confidence** Next, we evaluate recovery under different rule confidence levels. We vary the rule confidence from 0.1 to 0.9 and show the results in Figure 3b. We observe that SEQRET is robust against low confidence rules.

**Random Rule Triggers** We evaluate whether SEQRET recovers conditional dependencies even when rule heads and tails are infrequent. To test this, we generate data where no rule heads are planted, instead they occur only by chance. To ensure that the rule heads indeed do occur, we limit its size to 1. We then insert rule tails wherever the corresponding rules have triggered. We also limit the size of the tails to 1 to ensure they do not stand out as patterns by themselves.

We show the results in Figure 3c. SEQRET-MINE consistently recovers the rules while SEQRET-CANDIDATES performs worse because SQS does not find good seed patterns.

**Scaling Parameters** Next, we evaluate how SEQRET performs for ground truth models, alphabets, resp. rule tails of different sizes. First, we consider data with 10 to 50 ground truth rules for which we show the results in Fig. 3d. Next, we vary the alphabet from 100 to 1000 unique events and report the results in Fig. 3e. Finally, we vary the length of the rule tails from 1 to 11 and show the results in Figure 3f. We observe for all three scenarios that SEQRET recovers the ground truth well and visibly outperforms the competition.

### Experiments on Real Datasets

Next, we evaluate SEQRET on real world data.

**Datasets** We use nine datasets from different domains. *JMLR* and *Presidential* are two text datasets (Tatti and Vreeken 2012). *POS* contains parts-of-speech tags for a book (Pokou, Fournier-Viger, and Moghrabi 2016). *Ordóñez* and *Lifelog* contain daily activities of a person logged over several days (Ordóñez, De Toledo, and Sanchis 2013). *Rolling Mill* contains process logs from a steel manufacturing plant (Wiegand, Klakow, and Vreeken 2021). *Sepsis* contains treatment logs for sepsis cases from a hospital (Mannhardt 2016). *Ecommerce* contains purchase data of users of an online store spanning across 7 months, obtained from Kaggle. Finally, *Lichess* contains sequences of moves from online chess games. Further details are provided in the appendix. In Table 1 we provide statistics on datasets, as well as the results for different methods.

**General Observations** Overall, we observe that frequency-based methods like POERMA and POERMH discover a high number of rules, making interpretation difficult to impossible. TNS produces largely redundant

Dataset	D	Ω	SEQRET-CANDS			SEQRET-MINE			SQS		POERMA	POERMH	TNS	COSSU	
			P	R	%L	P	R	%L	P	%L	R	R	R	R	%L
JMLR	14501	1920	62	9	1.9	2	250	3.7	116	1.6	127	1282442	205	-	-
Presidential	62010	3973	30	4	0.5	2	167	1.1	58	0.4	57	90552	131	-	-
POS	45531	36	65	6	18.4	36	28	18.0	160	12.6	-	-	71	5	-0.5
Ordonez	739	10	2	0	11.5	1	4	16.7	2	11.5	95923	113337	6	2	-2.4
Lifelog	40520	78	36	6	8.9	16	83	9.8	59	6.5	2001521	1932296	95	5	-1.7
Rollingmill	18416	446	158	50	52.2	10	335	56.7	247	50.5	-	-	247	46	20.6
Sepsis	13114	11	19	5	35.7	8	28	32.2	42	23.5	849299	-	20	9	-1.1
Ecommerce	30875	127	77	10	13.7	92	164	28.1	95	13.6	43513	231824	219	6	-0.1
Lichess	20012	2273	81	18	2.4	11	337	4.6	113	2.1	4326	2068	348	-	-

Table 1: Results on real-world data. We report the number of discovered patterns (non-empty-head rules)  $P$  and rules  $R$ . For SEQRET, SQS and COSSU, we report the percentage of bits saved against the SEQRET null model as  $\%L$ . Failed runs, e.g. because of excessive runtime (COSSU) or out-of-memory errors (POERMA and POERMH), are indicated by ‘-’.

rules. COSSU is limited by its restrictive rule language and discovers very few rules. SQS is a sequential pattern miner that identifies meaningful patterns, but does not capture conditional dependencies that SEQRET successfully models. This difference is evident when comparing the compression achieved by different methods: SQS compresses the data less effectively than SEQRET, likely because SEQRET is more expressive. SEQRET-CANDIDATES improves upon SQS but still compresses worse than SEQRET-MINE.

**Case Studies** Next, we present illustrative examples to highlight how results from SEQRET differ from those of state-of-the-art methods. Consider a phrase from the *JMLR* dataset, ‘*support vector machine*’, that all methods identify in some form. SEQRET-MINE discovers the rules  $\langle \epsilon \rightarrow \text{support, vector} \rangle$  and  $\langle \text{support, vector} \rightarrow \text{machine} \rangle$ , which expresses that *support, vector* is a pattern, and that whenever it occurs it increases the probability of but is *not necessarily* followed by *machine*. In contrast, SEQRET-CANDIDATES and SQS both treat the entire phrase as a single pattern,  $\langle \text{support, vector, machine} \rangle$ , failing to capture the independent existence of *support, vector*. On the other end of the spectrum, POERMA and TNS discover 12 resp. 14 mainly redundant rules involving either *support* or *vector*.

In the *POS* dataset, SEQRET discovers common sentence structures, e.g. the pattern  $\langle \epsilon \rightarrow \text{determiner, cardinal number} \rangle$  capturing phrases such as “the first” or  $\langle \epsilon \rightarrow \text{to, verb-base-form} \rangle$  capturing phrases like “to tell”. SEQRET also captures the rule  $\langle \text{to, verb-base-form} \rightarrow \text{possessive-pronoun} \rangle$ , correctly identifying how possessive pronouns sometimes but not always follow phrases like “to tell”. COSSU, the method closest to our approach, fails to find any of the rules discussed. SQS finds them as several independent patterns disregarding the conditional dependencies.

For the *Lichess* dataset, SEQRET finds the “King’s Pawn Game” opening as  $\langle \epsilon \rightarrow \text{white:e4, black:e5} \rangle$ . It further discovers 12 rules with *white:e4, black:e5* as the head, capturing the different variations that often follow, e.g. the rule  $\langle \text{white:e4, black:e5} \rightarrow \text{white:Nf3} \rangle$ , shown in Figure 4a. Diving deeper into results, SEQRET discovers rules involving “King’s side castling” (shown in Figure 4b) Examples are  $\langle \text{black:O-O} \rightarrow \text{black:Re8} \rangle$  and  $\langle \text{black:O-O} \rightarrow \text{white:Qe2} \rangle$ .

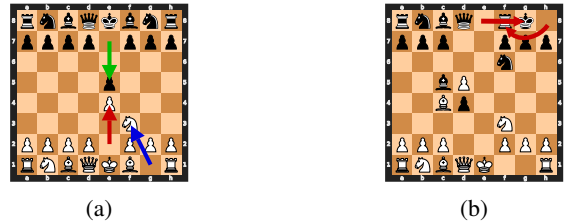


Figure 4: In (a) we show the rule  $\langle \text{white:e4 (red arrow), black:e5 (green arrow)} \rightarrow \text{white:Nf3 (blue arrow)} \rangle$ . In (b) we show black castling,  $\text{black:O-O}$ .

The former captures black moving its rook to e8, a position originally occupied by the King and made available only after castling. The latter captures white moving its queen to e2 following black castling, as a deterrence to black rook. SQS does not find any rule involving castling, and the frequency-based methods find 128 to 1370 mostly redundant rules.

## 7 Conclusion

We considered the problem of mining a succinct and non-redundant set of rules from long event sequences. We formalized the problem in terms of the MDL principle and presented SEQRET-CANDIDATES and SEQRET-MINE algorithms. We evaluated both on synthetic and real-world data. On synthetic data we saw that SEQRET recovers the ground truth well and is robust against noise, low rule confidence, varying alphabet resp. rule set sizes. On real-world data SEQRET provided insights that existing methods could not.

An interesting direction for future work is a differentiable approach to sequential rule mining. Existing work (Gao et al. 2025) focuses on the supervised classification setting, but not unsupervised rules. Additionally, we consider it interesting to study the causal aspects of sequential rules. Our approach lends itself to a causal framework by mapping the rule heads and tails to temporal variables. Moreover, the MDL principle as used in this paper nicely maps to the Algorithmic Markov Condition (Janzing and Schölkopf 2010) criterion to choose a plausible causal model.

## References

- Agrawal, R.; and Srikant, R. 1995. Mining sequential patterns. In *ICDE*, 3–14. Los Alamitos, CA, USA: IEEE Computer Society.
- Benjamini, Y.; and Hochberg, Y. 1995. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1): 289–300.
- Bhattacharyya, A.; and Vreeken, J. 2017. Efficiently Summarising Event Sequences with Rich Interleaving Patterns. In Chawla, N. V.; and Wang, W., eds., *Proceedings of the 2017 SIAM International Conference on Data Mining, Houston, Texas, USA, April 27-29, 2017*, 795–803. SIAM.
- Bloem, P.; and de Rooij, S. 2020. Large-scale network motif analysis using compression. *Data Mining and Knowledge Discovery*, 34: 1421–1453.
- Bourrand, E.; Galarraga, L.; Galbrun, E.; Fromont, E.; and Termier, A. 2021. Discovering Useful Compact Sets of Sequential Rules in a Long Sequence. In *2021 IEEE 33rd International Conference on Tools with Artificial Intelligence (ICTAI)*, 1295–1299. Los Alamitos, CA, USA: IEEE Computer Society.
- Chen, X.; Gan, W.; Wan, S.; and Gu, T. 2022. MDL-based Compressing Sequential Rules. *ArXiv*, abs/2212.10252.
- Chen, Y.; Fournier-Viger, P.; Nouioua, F.; and Wu, Y. 2021. Mining partially-ordered episode rules with the head support. In *Big Data Analytics and Knowledge Discovery: 23rd International Conference, DaWaK 2021, Virtual Event, September 27–30, 2021, Proceedings 23*, 266–271. Springer.
- Cule, B.; and Goethals, B. 2010. Mining association rules in long sequences. In *Advances in Knowledge Discovery and Data Mining: 14th Pacific-Asia Conference, PAKDD 2010, Hyderabad, India, June 21-24, 2010. Proceedings. Part I 14*, 300–309. Springer.
- Cüppers, J.; Kalofolias, J.; and Vreeken, J. 2022. Omen: discovering sequential patterns with reliable prediction delays. *Knowl. Inf. Syst.*, 64(4): 1013–1045.
- Cüppers, J.; Krieger, P.; and Vreeken, J. 2024. Discovering Sequential Patterns with Predictable Inter-event Delays. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 8346–8353. AAAI.
- Cüppers, J.; Schoen, A.; Blanc, G.; and Gimenez, P.-F. 2024. FlowChronicle: Synthetic Network Flow Generation through Pattern Set Mining. *Proceedings of the ACM on Networking*, 2(CoNEXT4): 1–20.
- Cüppers, J.; and Vreeken, J. 2020. Just Wait for it... Mining Sequential Patterns with Reliable Prediction Delays. In *2020 IEEE International Conference on Data Mining*, 82–91. IEEE.
- Dalmas, B.; Fournier-Viger, P.; and Norre, S. 2017. TWINCLE: A constrained sequential rule mining algorithm for event logs. *Procedia computer science*, 112: 205–214.
- Fournier-Viger, P.; Chen, Y.; Nouioua, F.; and Lin, J. C.-W. 2021. Mining partially-ordered episode rules in an event sequence. In *Intelligent Information and Database Systems: 13th Asian Conference, ACIIDS 2021, Phuket, Thailand, April 7–10, 2021, Proceedings 13*, 3–15. Springer.
- Fournier-Viger, P.; Faghihi, U.; Nkambou, R.; and Nguifo, E. M. 2012. CMRules: Mining sequential rules common to several sequences. *Knowledge-Based Systems*, 25(1): 63–76.
- Fournier-Viger, P.; Gueniche, T.; Zida, S.; and Tseng, V. S. 2014. ERMiner: sequential rule mining using equivalence classes. In *Advances in Intelligent Data Analysis XIII: 13th International Symposium, IDA 2014, Leuven, Belgium, October 30–November 1, 2014. Proceedings 13*, 108–119. Springer.
- Fournier-Viger, P.; Lin, J. C.-W.; Kiran, R. U.; Koh, Y. S.; and Thomas, R. 2017. A survey of sequential pattern mining. *Data Science and Pattern Recognition*, 1(1): 54–77.
- Fournier-Viger, P.; and Tseng, V. S. 2013. TNS: mining top-k non-redundant sequential rules. In *Proceedings of the 28th Annual ACM Symposium on Applied Computing*, 164–166.
- Galbrun, E. 2022. The minimum description length principle for pattern mining: A survey. *Data mining and knowledge discovery*, 36(5): 1679–1727.
- Gao, K.; Inoue, K.; Cao, Y.; Wang, H.; and Feng, Y. 2025. Differentiable Rule Induction from Raw Sequence Inputs. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net.
- Grünwald, P. 2007. *The Minimum Description Length Principle*. MIT Press.
- Hong, Y. 2013. On computing the distribution function for the Poisson binomial distribution. *Computational Statistics & Data Analysis*, 59: 41–51.
- Janzing, D.; and Schölkopf, B. 2010. Causal Inference Using the Algorithmic Markov Condition. *IEEE Transactions on Information Theory*, 56(10): 5168–5194.
- Le Cam, L. 1960. An approximation theorem for the Poisson binomial distribution.
- Li, P.; and Vitányi, M. 1997. *An Introduction to Kolmogorov Complexity and Its Applications*. Springer.
- Mannhardt, F. 2016. Sepsis Cases - Event Log.
- Mannila, H.; and Meek, C. 2000. Global Partial Orders From Sequential Data. In *KDD*, 161–168.
- Navarro, G. 2001. A guided tour to approximate string matching. *ACM Comput. Surv.*, 33(1): 31–88.
- Ordóñez, F. J.; De Toledo, P.; and Sanchis, A. 2013. Activity recognition using hybrid generative/discriminative models on home environments using binary sensors. *Sensors*, 13(5): 5460–5477.
- Pokou, Y. J. M.; Fournier-Viger, P.; and Moghrabi, C. 2016. Authorship attribution using small sets of frequent part-of-speech skip-grams. In *The Twenty-Ninth International Flairs Conference*.
- Rissanen, J. 1983. A Universal Prior for Integers and Estimation by Minimum Description Length. *Annals Stat.*, 11(2): 416–431.

- Roy, R.; Gaglione, J.-R.; Baharisangari, N.; Neider, D.; Xu, Z.; and Topcu, U. 2023. Learning interpretable temporal properties from positive examples only. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, 6507–6515.
- Tatti, N. 2009. Significance of Episodes Based on Minimal Windows. In *ICDM*, 513–522.
- Tatti, N.; and Vreeken, J. 2012. The long and the short of it: summarising event sequences with serial episodes. In *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '12*, 462–470. New York, NY, USA: Association for Computing Machinery. ISBN 9781450314626.
- Van Der Aalst, W. 2016. Data science in action. In *Process mining: Data science in action*, 3–23. Springer.
- Volkova, A. Y. 1996. A refinement of the central limit theorem for sums of independent random indicators. *Theory of Probability & Its Applications*, 40(4): 791–794.
- Wang, J.; and Han, J. 2004. BIDE: Efficient Mining of Frequent Closed Sequences. In *ICDE*, 79–90.
- Wang, Y. H. 1993. On the number of successes in independent trials. *Statistica Sinica*, 295–312.
- Wiegand, B.; Klakow, D.; and Vreeken, J. 2021. Mining Easily Understandable Models from Complex Event Data. *SIAM International Conference on Data Mining (SDM)*, SIAM.
- Yan, X.; Han, J.; and Afshar, R. 2003. CloSpan: Mining: Closed Sequential Patterns in Large Datasets. In *SDM*, 166–177. SIAM.
- Zaki, M. J. 2001. SPADE: An efficient algorithm for mining frequent sequences. *Machine learning*, 42: 31–60.