

Discovering Sequential Patterns with Predictable Inter-event Delays

Joscha Cüppers¹, Paul Krieger², Jilles Vreeken¹

¹ CISA Helmholtz Center for Information Security

² Saarland University

joscha.cueppers@cispa.de, s8pakrie@stud.uni-saarland.de, vreeken@cispa.de

Abstract

Summarizing sequential data with serial episodes allows non-trivial insight into the data generating process. Existing methods penalize gaps in pattern occurrences equally, regardless of where in the pattern these occur. This results in a strong bias against patterns with long inter-event delays, and in addition that regularity in terms of delays is not rewarded or discovered—even though both aspects provide key insight.

In this paper we tackle both these problems by explicitly modeling inter-event delay distributions. That is, we are not only interested in discovering the patterns, but also in describing how many times steps typically occur between their individual events. We formalize the problem in terms of the Minimum Description Length principle, by which we say the best set of patterns is the one that compresses the data best. The resulting optimization problem does not lend itself to exact optimization, and hence we propose HOPPER to heuristically mine high quality patterns. Extensive experiments show that HOPPER efficiently recovers the ground truth, discovers meaningful patterns from real-world data, and outperforms existing methods in discovering long-delay patterns.

Introduction

Summarizing event sequences is one of the key problems in data mining. Most existing methods do so in terms of serial episodes and allow for gaps (Tatti and Vreeken 2012) and interleaving (Bhattacharyya and Vreeken 2017) of pattern occurrences. By penalizing every gap equally regardless of where in a pattern it occurs, these methods have a strong bias against long inter-event delays, whereas methods that do not penalize gaps (Fowkes and Sutton 2016) are prone to discover spurious dependencies. What both of these classes lack is a pattern to be able to specify *when* the next symbol is to be expected.

To illustrate, let us consider a toy example of a single event sequence of all national holidays of a given country over the span of multiple years. As is usual, some holidays are ‘fixed’ as they always occur on the same date every year, and others depend on the lunar cycle and hence ‘move’ around. Existing methods have no trouble finding holidays that occur right after each another, e.g. *1st Christmas Day* right before *2nd Christmas Day*, struggle with long

delays, such as *Whit Monday* happening 49 days after *Easter Monday*, and outright fail when the relationship is ‘far’ and ‘loose’ such as *Easter* occurring between 82 to 114 days after *New Year’s*. In this paper, we present a method that can find and describe all these types of dependencies and delays.

To do so, we propose to explicitly model the distributions of inter-event delays in pattern occurrences. That is, as patterns we do not just consider serial episodes, but also discrete distributions that model the number of time-steps between subsequent events of a pattern. This allows us to discover patterns like *New Year* $\xrightarrow{82-114}$ *Easter Monday* $\xrightarrow{49}$ *Whit Monday*, which specify there is a uniformly distributed delay of 82 to 114 days between *New Year’s* and *Easter Monday*, and a fixed delay of 49 days until *Whit Monday*.

We define the problem of mining a succinct and non-redundant set of sequential patterns in terms of the Minimal Description Length Principle (MDL) (Grünwald 2007), by which we are after that model that compresses the data best. Simply put, unlike existing methods we do not plainly prefer patterns with ‘compact’ occurrences but rather those for which the inter-event delays are reliably predictable, no matter if these delays are short or long. This way we can automatically determine which discrete-valued distribution best characterizes the inter-event delays. In practice, we consider Uniform, Gaussian, Geometric, or Poisson distributions, but this set can be trivially extended.

The resulting problem does not lend itself for exact search, which is why we propose the effective HOPPER algorithm to efficiently discover good pattern sets in practice. Starting from just the singletons, HOPPER considers combinations of current patterns as candidates, uses an optimistic estimate to prune out unpromising candidates, explores both short and far dependencies, assigns the best-fitting delay distributions, and greedily chooses the candidate that improves the score most.

Through extensive evaluation, we show that HOPPER works well in practice. On synthetic data we demonstrate that unlike the state-of-the-art, we recover the ground truth well both in terms of patterns and delay distributions even in challenging settings where patterns include delays of hundreds of time steps. On real-world data, we show that HOPPER discovers easily interpretable patterns with meaningful delay distributions. We make all code, synthetic data, and real-world datasets available in the supplementary material.

Preliminaries

In this section, we discuss preliminaries and introduce the notation we use throughout the paper.

Notation

As data D we consider a set of $|D|$ event sequences $S \in D$ each drawn from a finite alphabet Ω of discrete events $e \in \Omega$, i.e. $S \in \Omega^{|S|}$. We write $S[i]$ to refer to the i^{th} event in S , and $|D|$ to denote the total number of events in D .

As patterns we consider serial episodes. A serial episode p is also a sequence drawn over Ω , i.e. $p \in \Omega^{|p|}$. We write $p[i]$ for the i^{th} event in p . We will model the inter-event delays between a subsequent pair of events $p[i]$ and $p[i+1]$ using discrete delay distribution $\pi_{p,i}(\cdot \mid \Theta_{p,i})$. Whenever clear from context we simply write $\pi_{p,i}(\cdot)$.

Finally, a window w^S is an ordered set of indices into S . Two windows a^S and b^S are *in conflict* iff they contain the same index, formally iff $|a^S \cap b^S| > 0$. A window w^S is said to *match* a pattern p if they identify the same events in the same order, i.e. when $\forall_{i \in [1, |p|]} S[w^S[i]] = p[i]$ and $\forall_{i \in [1, |p|-1]} \pi_{p,i}(w^S[i+1] - w^S[i]) > 0$, if p matches we write w_p^S . Whenever S is clear from context, we simply write w_p .

All logarithms are base 2 and we define $0 \log(0) = 0$.

Minimum Description Length

The Minimum Description Length (MDL) principle (Grünwald 2007) is a computable and statistically well-founded model selection principle based on Kolmogorov Complexity (Li and Vitányi 1993). For a given model class \mathcal{M} , it identifies the best model $M \in \mathcal{M}$ as the one minimizing $L(M) + L(D|M)$, where $L(M)$ is the length of model M and $L(D|M)$ the length of the data D given M .

This is known as two-part, or *crude* MDL, in contrast to one-part, or refined MDL (Grünwald 2007), which is not computable for arbitrary model classes. We use two-part MDL because we are particularly interested in the model. In MDL we are never concerned with materialized codes, we only care about code lengths. To use MDL we have to define a model class \mathcal{M} , and code length functions for the model and data given the model. We present these next.

MDL for Patterns with Predictable Delays

In this section we formally define the problem.

Decoding the Database

Before we define how to encode a sequence database using patterns with delay distributions we give the intuition, by explaining how to decode a database from a given cover. A cover C is a description of the data in terms of the patterns p in model M . Formally, a cover is defined as a tuple (C_p, C_d) , where pattern stream C_p describes which pattern (windows) are used in what order, and delay stream C_d consists of the inter-event delays within those windows. Next we explain how to decode a cover C to reconstruct the encoded data.

In Figure 1 we show a toy example. We show a sequence S , a model M , and two covers of S using M .

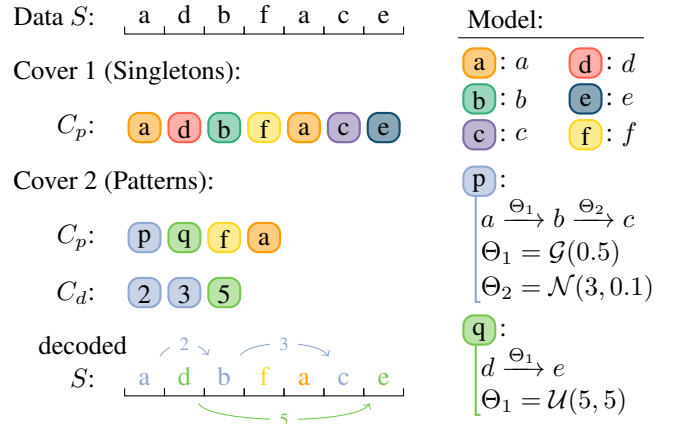


Figure 1: Toy example showing two possible encodings of the same data. Cover 1 uses only singletons, Cover 2 additionally uses two patterns, \boxed{p} and \boxed{q} . A cover consists of the pattern stream C_p encoding the patterns, and the delay stream C_d encoding the inter-event delays. The first gap of pattern \boxed{p} is modeled with a geometric distribution, and the second with a normal distribution. The one gap of \boxed{q} is modeled by a uniform distribution.

We first consider Cover 1. We start by reading the first code from the pattern stream C_p . This is an \boxed{a} which we look up in M and find it encodes event ‘a’. We write this to $S[0]$. We iterate reading and writing until S is decoded.

Next, we consider Cover 2. We again read the first code from C_p , which is now a \boxed{p} . We look up that this stands for pattern p . We write its first symbol, a , to $S[0]$. To know where in S we should write ‘b’ we read a code from the delay stream C_d . We read a 2, which means we write ‘b’ to $S[0+2]$. We continue until we have decoded this instance of pattern p , and then read the next symbol from C_p . This is a \boxed{q} . We start decoding it from the first empty position in S . We iterate until S is fully decoded.

Calculating the Encoding Cost

Now we know how to decode a sequence, we formally define how to compute the encoded sizes of the data and model.

Encoding the data To describe the data without loss, we need in addition to the pattern and delay streams, to know the number and length of sequences in D . We hence have

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

where we encode the numbers using the MDL-optimal encoding for integers $z \geq 1$ (Rissanen 1983). It is defined as $L_{\mathbb{N}}(z) = \log^* z + \log c_0$ where $\log^* z$ is the expansion $\log z + \log \log z + \dots$ where we only include the positive terms. To ensure this is valid encoding, i.e. one that satisfies the Kraft inequality, we set $c_0 = 2.865064$ (Rissanen 1983).

To encode the pattern stream C_p and the delay stream C_d , we use prefix codes, which are codes that are proportional in

length to their probability. For the pattern stream we have,

$$L(C_p) = \sum_{p \in M} -\text{usg}(p) \log \left(\frac{\text{usg}(p)}{\sum_{q \in M} \text{usg}(q)} \right),$$

where usg is the empirical frequency of pattern code $\langle p \rangle$ in the pattern stream C_p . We encode the delay stream C_d similarly, encoding the inter-event delays d_j between events $p[i]$ and $p[i+1]$ of every instance of a pattern p using the corresponding delay distribution $\pi_{p,i}(d_j)$. We hence have

$$L(C_d) = \sum_{p \in M} \sum_{i=1}^{|p|-1} \sum_{j=1}^{\text{usg}(p)} -\log \pi_{p,i}(d_j).$$

Encoding the Model As models we consider sets of patterns M that always include all singletons. We refer to the model that only consists of the singletons as the null model.

For the encoded length of a model we have

$$L(M) = L_{\mathbb{N}}(|\Omega|) + \log \binom{||D|| - 1}{|\Omega| - 1} + L_{\mathbb{N}}(|P| + 1) \\ + L_{\mathbb{N}}(\text{usg}(P)) + \log \binom{\text{usg}(P) - 1}{|P| - 1} + \sum_{p \in P} L(p),$$

where we first encode the size of the alphabet Ω and the supports $\text{supp}(e|D)$ of each singleton event. The latter we do using a so-called data-to-model code — an index over an enumeration of all possible ways to distribute $||D||$ events over alphabet Ω (Vereshchagin and Vítányi 2004). Next, we encode the number $|P|$ of non-singleton patterns $p \in M$ and their combined usage by $L_{\mathbb{N}}()$, and then their individual usages by a data-to-model code. Finally, we encode the non-singleton patterns.

To do so we need to specify how many, and which events a pattern consists of, as well as identify and parameterize its delay distributions. To reward similarities in delay behavior, we allow a distribution to be used for multiple inter-event gaps. As a default, we equip every pattern with one Geometric delay distribution. Formally, the encoded length of a non-singleton pattern $p \in M$ is

$$L(p) = L_{\mathbb{N}}(|p|) + \log(|p| - 1) + \log \binom{|p| - 1}{k} \\ - \sum_{e \in p} \log \left(\frac{\text{supp}(e|D)}{||D||} \right) + \sum_{\Theta \in p} L(\Theta),$$

where we encode the number of events of p , then its number of delay distributions, k , and finally where in the pattern these are used. We encode the events of the pattern using prefix codes based on the supports of events e in D .

To encode a delay distribution $\pi(\cdot | \Theta)$ it suffices to encode Θ . For the non-default delay distributions we first encode its type out of the set $\Psi = \{\text{Geometric}, \text{Poisson}, \text{Uniform}, \text{Normal}\}$ of discrete probability functions under consideration, for which we need $-\log |\Psi|$ bits. We then encode the parameter values $\theta \in \Theta$. We use $L_{\mathbb{N}}(\theta)$ if $\theta \in \mathbb{N}$, and $L_{\mathbb{R}}(\theta)$ if $\theta \in \mathbb{R}$. We have $L_{\mathbb{R}}(\theta) = L_{\mathbb{N}}(d) + L_{\mathbb{N}}(\lceil \theta \cdot 10^d \rceil) + 1$ as the number of bits needed to encode a real number up to user-set precision p (Marx and Vreeken 2019). It does so by shifting θ by d digits, such that $\theta \cdot 10^d \geq 10^p$.

The Problem, Formally

With the above, we can now formally state the problem.

The Predictable Sequential Delay Problem *Given a sequence database D over an alphabet Ω , find the smallest pattern set P and cover C such that the total encoded size, $L(M, D) = L(M) + L(D|M)$ is minimal.*

Considering the complexity of this problem, even when we ignore delay distributions there already exist super-exponential many possible patterns, exponentially many patterns sets over those, as well as, given a pattern set there exist exponentially many covers (Bhattacharyya and Vreeken 2017). Worst of all the search space does not exhibit any structure such as (weak-)monotonicity or submodularity that we can exploit. We hence resort to heuristics.

The HOPPER Algorithm

Now we have formally defined the problem and know how to score a model we need a way to mine good models. We break the problem into two parts, finding a good cover given a model, and finding a good model, and discuss these in turn.

Finding Good Covers

Given a model, we are after that description of the data that minimizes $L(D | M)$. To compute $L(D | M)$, we need a cover C . A cover consists of a set of windows, and hence we first need to find a set of good windows.

Finding Good Windows Mining all possible windows for a pattern p can result in an exponential blow-up. To ensure tractability, we limit ourselves to the 100 windows per starting event with the most likely delays. To avoid wasting time on windows we will never use because they will be too costly, we restrict our search to those for which the delays fall within the 99.7% confidence-interval of the respective probability distribution. For a normal distribution, that corresponds to three standard deviations from the mean. In practice, it is extremely unlikely that we would like to include any of the not considered windows in cover C , hence these restrictions have a negligible to no effect on the results.

Selecting a Good Cover Armed with a set of candidate windows, we next explain how to select a set C of these that together form a good cover. Ideally, we would like to select that cover C that minimizes $L(D | M)$. Finding the optimal cover, however would require testing exponentially many combinations, which would, in turn, result in unfeasible runtime; we hence do it greedily. For a greedy approach we need a way to select the next window for addition. Generally speaking, we prefer long patterns with likely delays. Based on this intuition, we assign each window w_p a score $s(w_p)$. At each step we select the window w_p with the highest $s(w_p)$. If a window conflicts with a previously selected window, we skip it and proceed. We add windows until all events of D are covered. To ensure there always exist a valid cover we always include all singleton windows.

As we prefer long patterns with likely deltas, our window score trades of pattern length ($|p|c$) against the cost of the

individual delays. Formally, we have

$$s(w_p) = |p|c - \sum_{k=1}^{|p|-1} -\log \pi_{p,k}(w_p[k+1] - w_p[k])$$

where c is the average code cost of a singleton event under the null model, that is

$$c = \frac{\sum_{e \in \Omega} -\text{supp}(e|D) \log(\text{supp}(e|D)/||D||)}{||D||}.$$

Mining Good Models

Now that we know how to find a good cover given a set of patterns, we explain how to discover a high-quality pattern set. Since there are super-exponentially many possible solutions, we again take a greedy approach. The general idea is that we use a pattern-growth strategy in which we iteratively combine existing patterns into new longer patterns. Before we explain our method in detail, we explain how we build a pattern candidate given two existing patterns and how to estimate the gain of such a candidate.

Estimating Candidate Gains Computing the total encoded length $L(M \oplus p', D)$ for when we add a new pattern p' to M is costly as this requires covering the data, which in turn requires finding good windows of p' . To avoid doing so for all candidates, we propose to instead use an optimistic estimator to discard those candidates for which we estimate no gain. Specifically, we want to estimate how many bits we will *gain* if we were to add pattern p' to the model.

To do so, we estimate the usage of p' . As we will explain below, every candidate p' is constructed by concatenating two existing patterns $p_1, p_2 \in M$. Assuming that p' will be used maximally, we have an optimistic estimate of its usage as $usg(p') = \min(usg(p_1), usg(p_2))$, or, if $p_1 = p_2$ as $usg(p') = usg(p_1)/2$. We estimate the change in model cost by adding p' by assuming all occurrences of the least frequent parent pattern are now covered by p' . Combined the estimated gain is,

$$\Delta \bar{L}(M \oplus p') = -\hat{L}(p') + L(\underset{p \in \{p_1, p_2\}}{\text{argmin}} usg(p))$$

where $\hat{L}(p')$ is the cost of p' omitting the delay distribution between p_1 and p_2 . We estimate $\Delta L(D | M \oplus p')$ as

$$\Delta \bar{L}(D | M \oplus p') = s \log(s) - s' \log(s') + z \log(z) - x \log(x) + x' \log(x') - y \log(y) + y' \log(y')$$

where s is the sum of all usages, $s = \sum_{p \in M} usg(p)$, and, for readability, we shorten $usg(p')$ to z , $usg(p_1)$ to x , $usg(p_2)$ to y and write x', y', s' for the “updated” usages, that is $x' = x - z$, $y' = y - z$ and $s' = s - z$.

As we do not have any information about the delays between p_1 and p_2 we assume these are encoded for free. Putting the above together gives us an optimistic estimate of the total encoded cost when adding pattern p' to M as

$$\Delta \bar{L}(D, M \oplus p') = \Delta \bar{L}(M \oplus p') + \Delta \bar{L}(D | M \oplus p').$$

Wherever clear from context, we simply write $\Delta \bar{L}(p')$.

Algorithm 1: OPTIMIZEALIGNMENT

Input : pattern candidate p' , alignment A
Output: estimated $gain^*$, optimized alignment A^*

```

1  $gain^* \leftarrow -\infty$ 
2 while  $\Delta \bar{L}_A(p') > gain^*$  do
3    $gain^* \leftarrow \Delta \bar{L}_A(p')$ 
4    $A^* \leftarrow A$ 
5   drop all delays  $d$  with minimal frequency from  $A$ 
6 return  $gain^*, A^*$ 
```

Estimating Candidate Occurrences When we want to evaluate a candidate pattern p' , constructed from patterns p_1 and p_2 , we have to determine its occurrence windows. A simple and crude way to determine candidate windows is by mapping every occurrence of p_1 to the nearest next occurrence of p_2 . We call this procedure ALIGNNEXT. It is particularly good for finding a mapping with the shortest possible delays, but will not do well when delays are relatively long. For this, the ALIGNFAR algorithm by Cüppers, Kalofolias, and Vreeken (2022) provides a better solution. In a nutshell, it efficiently discover that mapping A that minimizes the variance in delays. By a much larger search space it is naturally more susceptible to noise.

As a result, both strategies can give a good starting points, but neither will likely give an alignment that optimizes our MDL score. We propose to greedily optimize these mappings using an optimistic estimate. We first observe that given a mapping, we can trivially compute the delays, on which we can then fit a distribution. We do so for all distributions $\pi \in \Psi$ and choose that $\pi^*, (\cdot | \Theta^*)$ that minimizes the cost of encoding the delays. Second, we observe that a mapping also allows us to better estimate the usage of p' as the number of mapped occurrences of p_1 and p_2 . This gives a gain estimate under alignment A as

$$\Delta \bar{L}_A(p') = -L(p') + \Delta \bar{L}(D | M \oplus p') + \sum_{d \in A} \log \pi(d | \Theta^*).$$

We now use this estimate to identify and remove those mappings with the lowest delay probability (i.e. those with minimal frequency) until $\Delta \bar{L}_A(p')$ no longer increases. We give the pseudocode as Algorithm 1.

Mining Good Pattern Sets Next, we explain how we use the gain estimation and cover strategy to mine good pattern sets P . We give the pseudo-code for our method, HOPPER, as Algorithm 2. The key idea is to use a bottom-up approach and iteratively combine previously found patterns into longer ones.

We iteratively consider the Cartesian product of patterns $p_1, p_2 \in M$ as candidates. We evaluate these in order of potential gain. Events and patterns that occur frequently have the largest potential to compress the data, therefore we consider these combinations first. Specifically, we evaluate combinations of p_1 and p_2 in order of how many events they together currently cover (line 2).

Given a pattern candidate $p' = p_1 \oplus p_2$, we use our optimistic estimator to determine if we expect it to provide any

Algorithm 2: HOPPER

Input : sequence database D , alphabet Ω
Output: model M

```

1  $CT \leftarrow \Omega$ ;  $Cand \leftarrow CT \times CT$ ;
2 forall  $p_1, p_2 \in Cand$  do ordered descending on
    $|p_1|_{usg}(p_1) + |p_2|_{usg}(p_2)$ 
3   if  $\Delta \bar{L}(p_1 \oplus p_2) > 0$  :
4      $gain, p' \leftarrow \text{ALIGNCANDIDATE}(p_1, p_2)$ 
5     if  $gain > 0 \wedge L(D, M) > L(D, M \oplus p')$  :
6        $p' \leftarrow \text{FILLGAPS}(p', |p_1|)$ 
7        $M \leftarrow M \oplus p'$ 
8        $M \leftarrow \text{PRUNE}(M)$ 
9        $Cand \leftarrow Cand \cup \{M \times p', (p_1, p_2)\}$ 
10  $M \leftarrow \text{PRUNEINSIGNIFICANT}(M)$ 
11 return  $M$ 

```

gain in compression. If not, we move on to the next candidate. If we do estimate a gain based on usage of p_1 and p_2 alone, we proceed and optimize the alignment of occurrences of p_1 and p_2 to those of occurrences of p' . We do so using `ALIGNCANDIDATE`, for which we give the pseudocode in the supplementary. In a nutshell, it returns the best optimized result out of `ALIGNNEXT` and `ALIGNFAR`.

If the alignment leads to an estimated gain, we compute our score exactly (l. 5) and if the score improves we are safe to add p' to our model. We do so after we consider augmentations of p' with events that occur between p_1 and p_2 (`FILLGAPS`, line 6) such that we further improve the score. Adding a new pattern to M can make previously added patterns redundant, e.g. when all occurrences of p_1 are now covered by p' . We prune all patterns for which the score improves when we remove them from M (`PRUNE`). Finally, we create new candidates based on the just added pattern, and add (p_1, p_2) back to the candidate set, as we might want to build a different pattern from it in a later iteration.

Before returning the final pattern set, we reconsider all patterns in the model and only keep those that give us a significant gain (Bloem and de Rooij 2020; Grünwald 2007) in compression. We provide further details on the pattern mining procedure in the supplementary.

As we consider the most promising candidates first, the more candidates we evaluate to have no gain, the more unlikely it becomes we will find a candidate that will provide any substantial gain. To avoid evaluating all of those unnecessarily, we propose an early stopping criterion by considering up to $|\Omega|^2/100$, but at least 1 000, unsuccessful candidates in a row. As our score is bounded from below by 0, we know that Hopper will eventually converge.

Related Work

Mining sequential patterns from event sequences has a rich history. Traditional sequential pattern miners focus on finding all frequent patterns (Agrawal and Srikant 1995; Laxman, Sastry, and Unnikrishnan 2007), these suffer from exponentially many patterns, making interpretation hard to

impossible. Closed episodes (Yan, Han, and Afshar 2003; Wang and Han 2004) partially solve this, but are highly sensitive to noise. More recently, research focus shifted to mining patterns with a frequency that is significant with respect to some null hypothesis (Low-Kam et al. 2013; Petitjean et al. 2016; Tonon and Vandin 2019; Jenkins, Walzer-Goldfeld, and Riondato 2022). While this alleviates, it does not solve the pattern explosion.

Pattern set mining solves the pattern explosion by asking for a small and non-redundant set of patterns that generalizes the data well, as instead of asking for *all* patterns that satisfy some individual criterion. There exist different approaches to how to score a pattern set. ISM (Fowkes and Sutton 2016) takes a probabilistic Bayesian approach, unlike us they do not model gaps. SQS (Tatti and Vreeken 2012) is an example of a method that employs the Minimum Description Length principle to identify the best set of serial episodes, which are sequential patterns that allow for gaps. SQUISH (Bhattacharyya and Vreeken 2017) builds upon SQS and additionally allows interleaved and nested patterns. However, SQS and SQUISH, are not capable of finding patterns with long inter-event delays and penalize each individual gap uniformly, regardless where in the pattern it occurs.

Existing methods that enrich patterns with delays can be categorized into two groups, methods that discover frequent patterns that satisfy some user set delay constraints (Yoshida et al. 2000; Giannotti et al. 2006; Dauxais et al. 2017; Cram, Mathern, and Mille 2012), and methods that discovers delay information from the data (Yen and Lee 2013; Nanni and Rigotti 2007). The latter, in contrast to our method, only consider the minimal delay between events, do not work on a single long sequence, and mine all frequent patterns, and hence also suffer from the pattern explosion.

Existing pattern set miners that do model the inter-event delay solve different problems. Galbrun et al. (2018) propose to mine *periodic* patterns, which are patterns that continuously appear throughout the data with near-exact delays. It is therewith well-suited for the holidays example in the introduction, but less so for discovering patterns that only appear more locally. OMEN (Cüppers, Kalofolias, and Vreeken 2022) does discover local patterns and delay distributions, but does so in a supervised setup between a pattern and a target attribute of interest. As such, each of the above methods consider part of the problem we study here, but none address it directly: we aim to discover a small set of sequential patterns where the delays between subsequent events in a pattern are modelled with a probability distribution.

Experiments

In this section we empirically evaluate HOPPER on synthetic and real-world data. We implement HOPPER in Python and provide the source code along with the synthetic data and the real-world data in the supplementary.¹ We compare HOPPER to SKOPUS (Petitjean et al. 2016) as a representative statistically significant sequential pattern miner, SQS (Tatti and Vreeken 2012), SQUISH (Bhattacharyya and Vreeken 2017) and ISM (Fowkes and Sutton 2016) as representatives of the

¹eda.rg.cispa.io/prj/hopper

general class of pattern set miners, and to PPM (Galbrun et al. 2018) as a representative of the periodic pattern miners. For all, we use the implementation by the authors.

HOPPER considers delays up to a user set *max delay*, for all experiments we set it to 200. SKOPUS only works on a set of sequences, when the dataset consists of one sequence, we split the sequence into 100 equally long sequences. We parametrize SKOPUS to report the top 10 patterns of at most length 10, which corresponds to the ground-truth value in our synthetic experiments. PPM only accepts a single sequence as input, to make it applicable on databases of multiple sequences, we concatenate these into one long sequence. We give the full setup description in the supplementary.

Synthetic Data

To evaluate how well HOPPER recovers patterns with known ground, we consider synthetic data. To this end, we generate data as follows. For each synthetic configuration we generate 20 independent datasets. For each dataset we sample uniform at random one sequence of length 10000 over an alphabet of 500 events, we plant 10 unique patterns, uniformly, at random locations while avoiding collisions. The frequency of planted patterns, length and delay distributions between events we vary per experiment.

As evaluation we consider standard F1 score. Where, to reward partial discovery, we weight a reported pattern p_r by the relative edit distance to the planted pattern p_p , that is, $w(p_r, p_p) = \max(1 - \text{lev}(p_r, p_p)/|p_p|, 0)$ where *lev* is the Levenshtein edit distance. Since we do not want to reward redundant discoveries we cap the total reward to one per planted pattern. To illustrate, consider the example where we plant one pattern *abcd*, and discover two patterns, *ab* and *cd*. We value both as 0.5. As we can technically reconstruct the generating pattern we hence have a recall of one, but, as we have to do so using two rather than one pattern, we have a precision of 0.5. This way we reward partial discoveries, which is especially relevant for methods that are designed to pick up events that occur close to one another, but might miss the full pattern if it includes a long delay. We provide additional details on the evaluation in the supplementary.

Sanity Check We start with a sanity check, where we run HOPPER on 20 data sets without structure, generated uniformly at random. It correctly does not report any patterns.

Delay Distributions Next, we test how well HOPPER can recover patterns for varying numbers of delay distributions. We consider the case of no delay distributions up to a pattern including a delay distribution between every subsequent pair of events. We plant 10 unique patterns of length 10 and in total 200 pattern occurrences, that is, on expectation 20 instances per pattern. As delay distribution, we plant Uniform distributions with a delay of between 10 to 20 time steps.

We present the results in the first panel of Fig. 2. We observe that HOPPER performs on par when there are no delay distributions and outperforms the state of the art when we increase their number. We find that SQUISH performs on par with SQS in our experiments and to avoid clutter from here onward postpone its results to the supplementary.

Low Frequency Next, we evaluate performance with low-frequency patterns, we decrease the frequency of the total number of planted patterns. We consider the same setting as above, where we set the number of distributions to four and decrease the total number of planted patterns from 200 to 100, that is, on expectation, from 20 to 10 per pattern. We show the results in the 2nd panel of Fig. 2. We observe that HOPPER outperforms all other methods, ultimately reducing to the performance of SQS in the low-frequency domain.

Long Delays Next, we investigate how robust HOPPER is to long delays, to this end we plant 10 patterns at 200 locations. We plant patterns of length 3, with Normal distributed inter-event delays, with a standard deviation of one, and increase the mean stepwise from 1 to 180. We present the results in the third panel of Fig. 2. We observe HOPPER is very robust against long delays: even with an expected delay of 180 between the individual events it achieves a very high F1 score. In contrast, its competitors do not fare well; SQS and SKOPUS perform well initially but then quickly deteriorate.

High Variance Finally, we evaluate HOPPER under increasing variance of inter-event delays. To this end we plant 400 occurrences of 10 patterns of length 3, with Normally distributed delays with mean 50 and varying the standard deviations. We show the results in the last panel of Fig. 2.

We observe that HOPPER gets near perfect results for lower variance and high F1 score until a standard deviation of 7 at which point 95% percent of the probability mass is distributed over a range of 28 timestamps. In general, we observe that the higher the frequency, the more robust we are against higher variance. We can see that SKOPUS is consistent under increasing variance. This is probably due to the fact that SKOPUS does not care about the distance between events only about the order in which they occur.

Real World Results

Next, we evaluate Hopper on real-world data. We use eight datasets that together span a wide range of use-cases. We consider a dataset of all national *Holidays* in a European country over a century, the playlist a local *Radio* station recorded over a month, the *Lifelog*² of all activities of one person recorded in over seven years, the MIDI data of hundred Bach *Chorales* (Dua and Graff 2017), all commits to the *Samba* project for over ten years (Galbrun et al. 2018), the *Rolling Mill* production log of steel manufacturing plant (Wiegand, Klakow, and Vreeken 2021), the discretized muscle activations of professional ice *Skating* riders (Moerchen and Fradkin 2010), and finally, three text datasets the Gutenberg project, resp. *Romeo* and *Juliet* by Shakespeare, *A Room with a View* by E.M. Forster, and *The Great Gatsby* by F. Scott Fitzgerald. We give the total number of events per dataset in Table 1 and further statistics in the supplementary.

We run HOPPER, SQS, ISM, PPM, and SKOPUS on all datasets. We report the number of patterns ($|P|$), the average expected distance between the first and last event ($\mathbb{E}(w_p[|p|] - w_p[0])$) and for HOPPER, the number of discovered delay distributions ($\#\Theta$). In the interest of space

²<https://quantifiedawesome.com/>

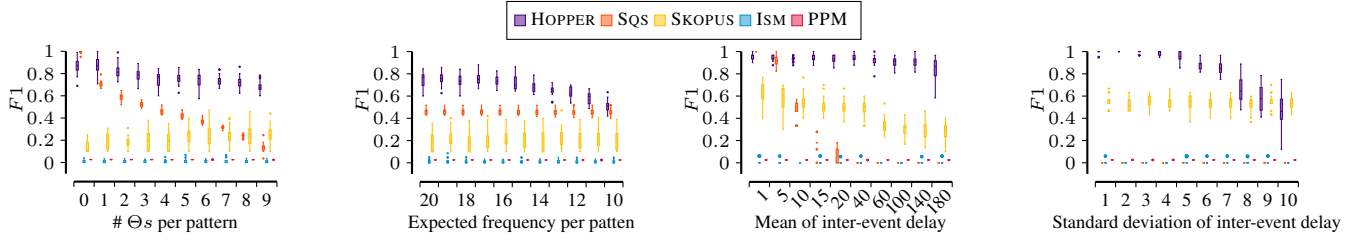


Figure 2: [Higher is better] F1 scores for recovering patterns from synthetic data. From left to right, we evaluate for varying numbers of inter-event distributions, expected frequency of a pattern, mean inter-event delay, resp. different standard deviations for normally distributed delays. We see that HOPPER performs on par with SQS when inter-event delays are few and simply structured, and outperforms the competition with a large margin whenever their structure is more complicated.

Dataset	$ D $	HOPPER			SQS			PPM	
		$ P $	$\#\Theta$	$\mathbb{E}(w)$	$ P $	$\mathbb{E}(w)$		$ P $	$\mathbb{E}(w)$
<i>Holidays</i>	37k	1	7	393	3	19.2		14	51.6
<i>Radio</i>	16k	22	43	48	15	5.8		587	71.9
<i>Lifelog</i>	40k	37	68	129	58	3.9		1.6k	119.1
<i>Samba</i>	29k	40	101	110	221	2.7		1.4k	17.1
<i>Chorales</i>	7k	56	57	4.7	114	2.6		433	2.6
<i>Rolling</i>	54k	237	489	7.4	470	5.0		3.6k	181.9
<i>Skating</i>	26k	86	160	9.1	160	4.0		1.4k	55.5
<i>Romeo</i>	37k	254	284	12.9	254	2.8		2.3k	332.7
<i>Room</i>	87k	565	610	3.1	701	2.5		—	—
<i>Gatsby</i>	64k	439	488	7.3	519	2.6		4.7k	641.9

Table 1: Results on real-world data. For HOPPER, SQS, and PPM we report the number of discovered patterns ($|P|$) and the average expected distance between the first and last symbol of a pattern ($\mathbb{E}(w)$). For HOPPER we additionally give total number of discovered inter-event distributions ($\#\Theta$).

we postpone the results of ISM and SKOPUS, along with the metrics runtime and average events per pattern to the supplementary. HOPPER terminates within seconds to hours, depending on the dataset. We find that while HOPPER and SQS discover similar numbers of patterns, those that HOPPER discovers reveal much longer range dependencies and, in general, include more events. PPM results in an order of magnitude more patterns, most of which are singletons. Next we look at the results for *Holidays* and *Radio* in more detail.

On the *Holidays* dataset, HOPPER finds a single pattern, *May 1st $\xrightarrow{155}$ National Holiday $\xrightarrow{83}$ 1st Christmas Day $\xrightarrow{1}$ 2nd Christmas Day $\xrightarrow{6}$ New Year $\xrightarrow{80-112}$ Good Friday $\xrightarrow{3}$ Easter Monday $\xrightarrow{49}$ Whit Monday*, where all delay distributions are uniform. The pattern precisely describes all fixed and all lunar-calendar dependent holidays within the year. In contrast, the competing methods only find fractions of this pattern, such as *1st Christmas Day*, *2nd Christmas Day*. We show the results for all methods in the supplementary.

The *Radio* dataset includes all the songs played, as well as the ad slots and news segments, for a local radio station over the course of a month. On this data, HOPPER discovers the pattern *Jingle $\xrightarrow{0}$ Ads $\xrightarrow{0}$ News $\xrightarrow{0}$ Jingle $\xrightarrow{U(3,5)}$ Jingle* where the 0-gaps correspond to geometric distributions with

$p = 1$ and the last inter-event delay is a uniform distribution. Other methods find comparable or parts of this patterns, but none give the immediate insight that the first four events follow directly after one another and the last *Jingle* plays between 3 to 5 events after the previous.

More importantly, unlike other methods, HOPPER also picks up patterns such as *Solo Para $\xrightarrow{G(0.02)}$ As It Was $\xrightarrow{N(48,25)}$ I Believe $\xrightarrow{P(24)}$ Anyone for You* that confirm our suspicion that radio stations often play the same sequence of particularly popular songs interspersed with less-well-known songs. No other method finds any comparable patterns. HOPPER discovers much longer patterns than its competitors. Whereas most competitors find patterns of length 2, SQS patterns of at most 4 events, HOPPER discovers patterns of up to 7 events long. Together, this illustrates that HOPPER finds patterns that are not only more detailed in terms of the delay structure, but also in which events they describe.

Conclusion

We consider the problem of summarizing sequential data with a small set of patterns with inter-event delays. We formalized the problem in terms of the Minimum Description Length principle and presented the greedy HOPPER algorithm. On synthetic data we saw that our method recovers the ground truth well and is robust against high delays and variance. On real-world data we observed that HOPPER finds meaningful patterns that go beyond what state of the art methods can capture. While methods that only consider the order of events, can in theory find patterns with long delays, they often do not do this in practice.

We introduce a more powerful pattern language that enables us to discover new structure in data. This comes with the trade-off, of a much larger search space and, in theory, makes us more susceptible to noise, however the experiments have shown that this is not a problem in practice. HOPPER achieves a high F1 score on all experiments in Fig. 2, despite these having 80% or more noise.

Currently, we model the delay between subsequent events in a pattern. In practice, some events may depend on some event earlier in the pattern. We see it as an interesting direction for future work to extend our pattern language to include rule-like dependencies.

References

- Agrawal, R.; and Srikant, R. 1995. Mining sequential patterns. In *ICDE*, 3–14. Los Alamitos, CA, USA: IEEE Computer Society.
- Bhattacharyya, A.; and Vreeken, J. 2017. Squish: Efficiently Summarising Event Sequences with Rich Interleaving Patterns. In *SDM*, 11.
- Bloem, P.; and de Rooij, S. 2020. Large-Scale Network Motif Analysis Using Compression. *DAMI*, 34: 1421–1453.
- Cram, D.; Mathern, B.; and Mille, A. 2012. A Complete Chronicle Discovery Approach: Application to Activity Analysis. *Expert Systems*, 29(4): 321–346.
- Cüppers, J.; Kalofolias, J.; and Vreeken, J. 2022. Omen: Discovering Sequential Patterns with Reliable Prediction Delays. *KAIS*, 64(4): 1013–1045.
- Dauxais, Y.; Guyet, T.; Gross-Amblard, D.; and Happe, A. 2017. Discriminant Chronicles Mining: Application to Care Pathways Analytics. In *AIME*, 234–244. Springer.
- Dua, D.; and Graff, C. 2017. UCI Machine Learning Repository.
- Fowkes, J.; and Sutton, C. 2016. A Subsequence Interleaving Model for Sequential Pattern Mining. In *KDD*.
- Galbrun, E.; Cellier, P.; Tatti, N.; Termier, A.; and Crémilleux, B. 2018. Mining Periodic Patterns with a MDL Criterion. In *ECMLPKDD18*, 535–551. Springer.
- Giannotti, F.; Nanni, M.; Pedreschi, D.; and Pinelli, F. 2006. Mining Sequences with Temporal Annotations. In *Proceedings of the 2006 ACM Symposium on Applied Computing*, 593–597. Dijon France: ACM. ISBN 978-1-59593-108-5.
- Grünwald, P. 2007. *The Minimum Description Length Principle*. MIT Press.
- Jenkins, S.; Walzer-Goldfeld, S.; and Riondato, M. 2022. SPEck: Mining Statistically-Significant Sequential Patterns Efficiently with Exact Sampling. *Data Min Knowl Disc*, 36(4): 1575–1599.
- Laxman, S.; Sastry, P. S.; and Unnikrishnan, K. P. 2007. A Fast Algorithm for Finding Frequent Episodes in Event Streams. In *KDD07*, 410–419. ACM.
- Li, M.; and Vitányi, P. 1993. *An Introduction to Kolmogorov Complexity and its Applications*. Springer.
- Low-Kam, C.; Raissi, C.; Kaytoue, M.; and Pei, J. 2013. Mining Statistically Significant Sequential Patterns. In *ICDM*, 488–497. Dallas, TX, USA: IEEE. ISBN 978-0-7695-5108-1.
- Marx, A.; and Vreeken, J. 2019. Telling Cause from Effect by Local and Global Regression. *KAIS*, 60: 1277–1305.
- Moerchen, F.; and Fradkin, D. 2010. Robust Mining of Time Intervals with Semi-Interval Partial Order Patterns. In *SDM*, 315–326.
- Nanni, M.; and Rigotti, C. 2007. Extracting Trees of Quantitative Serial Episodes. In Džeroski, S.; and Struyf, J., eds., *Knowledge Discovery in Inductive Databases*, volume 4747, 170–188. Berlin, Heidelberg: Springer Berlin Heidelberg. ISBN 978-3-540-75548-7.
- Petitjean, F.; Li, T.; Tatti, N.; and Webb, G. I. 2016. Skopus: Mining Top-k Sequential Patterns under Leverage. *DAMI*, 30(5): 1086–1111.
- Rissanen, J. 1983. A Universal Prior for Integers and Estimation by Minimum Description Length. *Annals Stat.*, 11(2): 416–431.
- Tatti, N.; and Vreeken, J. 2012. The Long and the Short of It: Summarizing Event Sequences with Serial Episodes. In *KDD*, 462–470. ACM.
- Tonon, A.; and Vandin, F. 2019. Permutation Strategies for Mining Significant Sequential Patterns. In *ICDM*, 1330–1335. Beijing, China: IEEE. ISBN 978-1-72814-604-1.
- Vereshchagin, N. K.; and Vitányi, P. M. B. 2004. Kolmogorov’s Structure Functions and Model Selection. *IEEE Transactions on Information Theory*, 50(12): 3265–3290.
- Wang, J.; and Han, J. 2004. BIDE: Efficient Mining of Frequent Closed Sequences. In *ICDE*, 79–90.
- Wiegand, B.; Klakow, D.; and Vreeken, J. 2021. Mining Easily Understandable Models from Complex Event Logs. In *SDM*, 10.
- Yan, X.; Han, J.; and Afshar, R. 2003. CloSpan: Mining: Closed Sequential Patterns in Large Datasets. In *SDM*, 166–177. SIAM.
- Yen, S.-J.; and Lee, Y.-S. 2013. Mining Non-Redundant Time-Gap Sequential Patterns. *Applied Intelligence*, 39(4): 727–738.
- Yoshida, M.; Iizuka, T.; Shiohara, H.; and Ishiguro, M. 2000. Mining Sequential Patterns Including Time Intervals. In Dasarathy, B. V., ed., *AeroSense 2000*, 213–220. Orlando, FL.