

Direct Local Pattern Sampling by Efficient Two-Step Random Procedures

Mario Boley
Fraunhofer IAIS and
University of Bonn
mario.boleym@iais.fhg.de

Claudio Lucchese
I.S.T.I.-C.N.R. Pisa
claudio.lucchese@isti.cnr.it

Daniel Paurat
University of Bonn
daniel.paurat@uni-bonn.de

Thomas Gärtner
Fraunhofer IAIS and
University of Bonn
thomas.gaertner@iais.fhg.de

ABSTRACT

We present several exact and highly scalable local pattern sampling algorithms. They can be used as an alternative to exhaustive local pattern discovery methods (e.g. frequent set mining or optimistic-estimator-based subgroup discovery) and can substantially improve efficiency as well as controllability of pattern discovery processes. While previous sampling approaches mainly rely on the Markov chain Monte Carlo method, our procedures are direct, i.e., non-process-simulating, sampling algorithms. The advantages of these direct methods are an almost optimal time complexity per pattern as well as an exactly controlled distribution of the produced patterns. Namely, the proposed algorithms can sample (item-)sets according to frequency, area, squared frequency, and a class discriminativity measure. We present experimental results demonstrating the usefulness of our procedures for pattern-based model construction as well as their good scalability.

1. INTRODUCTION

This paper presents simple yet effective procedures for local pattern discovery [19] that attack the task from a different algorithmic angle than the standard search approach—namely, by directly generating individual patterns as outcome of a random experiment. Local patterns such as association rules [1] or emerging patterns [12] are used in different application contexts from exploratory data analysis where they constitute units of discovered knowledge to predictive model construction where patterns act as binary features [9, 10, 13]. All applications have in common that usually only a few patterns can be effectively utilized—either due to the limited attention of a data analyst or because too many features can reduce the comprehensibility and performance of a global model. Standard local pattern discovery algorithms, however, are based on exhaustive search within

huge pattern spaces (e.g., frequent set miners [18, 23], or optimistic-estimator-based subgroup and association discovery [16, 21]). Consequently, they tend to either produce a vast amount of output patterns or at least enumerate them implicitly.

This motivates the invention of algorithms that only sample a representative set of patterns without explicitly searching in the pattern space. There are such algorithms in the literature [2, 6, 8] but they provide either no control over the distribution of their output or only asymptotic control by simulating a stochastic process on the pattern space using the Markov chain Monte Carlo method (MCMC). In addition to only offering approximate sampling, MCMC methods have a scalability problem: the number of required process simulation steps is often large and, even more critical, individual simulation steps typically involve support counting and, hence, can be too expensive for large input datasets. Therefore, we present novel pattern generation methods that sample patterns exactly and directly, i.e., without simulating time-consuming stochastic processes. More precisely, given a dataset \mathcal{D} and a number of desired patterns k , the procedures

- produce exactly k patterns each of which is generated following exactly a distribution proportional to either frequency, squared frequency, area (i.e., frequency times size), or discriminativity (i.e., frequency in positive data portion times negative frequency in negative data portion);
- use time $O(\|\mathcal{D}\| + kn)$ respectively $O(\|\mathcal{D}\|^2 + kn)$ in case of squared frequency and discriminativity where n denotes the number of items and $\|\mathcal{D}\|$ the size of the dataset, i.e., the sum of all data record sizes¹.

That is, after a linear respective quadratic preprocessing phase each pattern is produced in a time linear in the number of items. This time complexity appears to be almost optimal, because only reading the data once requires $O(\|\mathcal{D}\|)$ and just printing k patterns without any further computation requires time $O(kn)$.

¹This assumes that $\exp(n) > |\mathcal{D}|$; the actual complexities are $O(\|\mathcal{D}\| + k(n + \ln |\mathcal{D}|))$ and $O(\|\mathcal{D}\|^2 + k(n + \ln^2 |\mathcal{D}|))$.

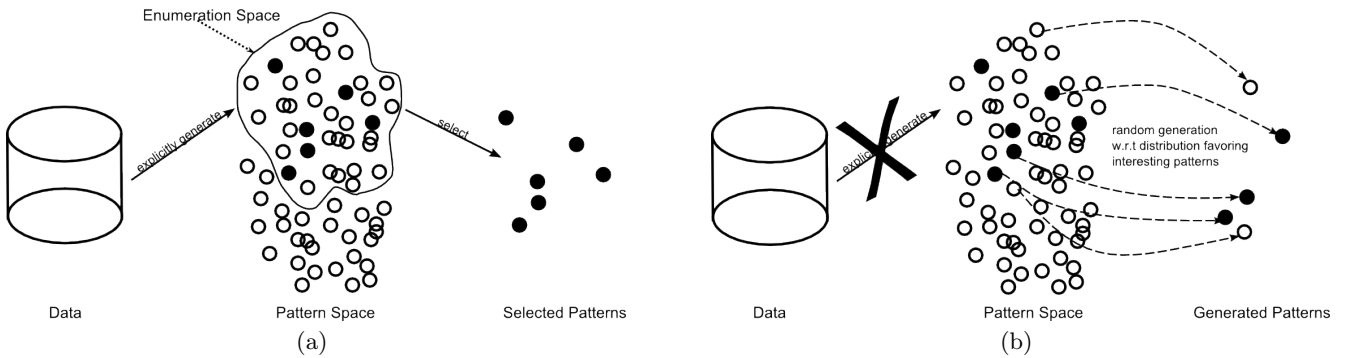


Figure 1: (a) *Exhaustive search*: involves complete generation of an enumeration space guaranteed to contain all interesting patterns; however, size of that space usually has no reasonable bound with respect to input size and is hard to predict. (b) *Controlled pattern sampling*: no explicit construction of potentially huge part of pattern space; instead random generation of small designated number of patterns; no guarantee of finding patterns satisfying hard interestingness threshold, but control over computation time and output size.

After giving some more background on the idea of controlled repeated pattern sampling and reviewing other pattern sampling algorithms, the remainder of this paper is structured as follows. We define formal and notational background (Sec. 2) followed by a detailed description of the sampling procedures (Sec 3). Then we report experimental results showing that sampled patterns are equally useful for pattern-based classification as frequent sets and that pattern sampling can easily outperform exhaustive listing on large datasets (Sec. 4). Finally, we give a summarizing discussion of all results (Sec. 5).

1.1 Pattern Sampling

The data mining literature contains several local pattern discovery algorithms that can efficiently produce large output families. Here, efficiency is defined in an output-sensitive way, i.e., amortized polynomial time per pattern, which is a useful notion assuming that the produced pattern collections are the final output. When viewed from a global (application-driven) perspective though, the enumerated patterns are usually only an intermediate result, from which a final—often much smaller—pattern collection is selected. Hence, enumeration is only the first half of a surrounding local pattern discovery process. This two phase approach, which we refer to as “exhaustive search” is illustrated in Figure 1(a): during the enumeration step a part of the implicitly defined pattern space is physically constructed—we refer to that part as “enumeration space”—and then, during the selection step, the most valuable patterns from this enumeration space are collected with respect to some interestingness measure.

An example for this paradigm is listing frequent sets of an input dataset, but subsequently using only those sets that provide rules with a large lift (or interest) value. A further example is optimistic-estimator-based pruning for subgroup or association discovery. There the enumeration space is the family of all sets having a large enough optimistic estimate of their interestingness and the truly interesting patterns are selected for the result family. Note that, in this example, enumeration and selection are algorithmically interweaved, i.e., sets are already selected throughout the enumeration

phase. Many more examples emerge from the LeGo approach to data mining [20] where patterns are selected according to their utility for constructing global models.

For these approaches, the enumeration step can constitute a severe bottleneck. Even if enumeration is performed by an amortized polynomial time algorithm, its computation time is essentially unpredictable: the size of the enumeration space cannot be directly controlled and its explicit construction takes time at least proportional to that size. On the other hand, if one enforces a maximum computation time by aborting the execution at a certain point, one ends up with an uncontrolled subset of the enumeration space, which depends on the internal search order of the enumeration algorithm.

In contrast, suppose we can access the pattern space \mathcal{L} by an efficient sampling procedure simulating a distribution $\pi: \mathcal{L} \rightarrow [0, 1]$ that is defined with respect to some interestingness measure q , e.g., $\pi(\cdot) = q(\cdot)/Z$ where Z is a normalizing constant. Then it is possible to efficiently generate a pattern collection that consists exactly of as many patterns as are truly required and that is representative for the distribution π and hence for the underlying notion of interestingness q . Figure 1(b) illustrates this alternative approach, which we want to refer to as “controlled repeated pattern sampling”. A potentially positive side-effect of this paradigm is that instead of the usual hard constraints it utilizes parameter-free soft constraints [5]. Hence, the user is freed of the often troublesome task of finding appropriate hard threshold parameters such as a minimum frequency threshold.

1.2 Related Work

In contrast to sampling from the input database (see, e.g., [22, 25]), it is a relatively new development in local pattern discovery to sample from the pattern space. In the context of maximal frequent subgraph mining, Chaoji et al. [8] describes a random process that stops after a number of steps that is bounded by the maximum number of edges present in an input graph and produces a maximal frequent subgraph. A similar process is already applied in Gunopoulos et al. [17] within a Las Vegas variant of the Dualize and Advance

algorithm. More precisely, it is used for the internal randomization of an algorithm with an otherwise deterministic output (all maximal frequent and minimal infrequent sets of a given input database). When applied for the final pattern discovery, however, this random process has the weakness that it provides no control over the generation probabilities of individual patterns.

Several papers propose to overcome this weakness by applying the Markov chain Monte Carlo method. Boley and Grosskreutz [7] proposes frequent set sampling to approximate the effect of specific minimum frequency thresholds. The proposed algorithm simulates a simple Glauber dynamic on the frequent set lattice: starting in the empty set, in each subsequent time step a single item is either removed or added to the current set. A similar Markov chain Monte Carlo (MCMC) method is used in Zaki and Al Hasan [2] for generating a representative set of graph patterns. These MCMC methods provide limited control of the generation probabilities, namely about the infinite limit of the distribution of the current state. The worst-case convergence can, however, be exponentially slow in the size of the input database. For sampling from the family of frequent patterns, this problem appears to be inherent: almost uniform frequent pattern sampling can be used for approximate frequent pattern counting, which one can show to be intractable under reasonable complexity assumptions (see [7]). Similar conclusions can be drawn for enumeration spaces defined by linearly scaled versions of the frequency measure such as the standard optimistic estimator for the binomial test quality function in subgroup discovery [24].

In order to avoid this implication of hard-constraint-based pattern discovery (e.g., using a hard frequency threshold), Boley et al. [6] combines pattern space sampling with soft-constraint-based pattern discovery [5]—resulting in the pattern sampling paradigm described in Section 1.1 above. Still, the underlying method is again MCMC-based, and, despite using a more sophisticated chain defined on the closed set lattice of the input database, it shares the practical weaknesses of this technique. This given paper now retains the idea of controlled pattern sampling without hard constraints, but proposes novel pattern generation methods that are exact and direct, i.e., they do not involve MCMC process simulation. Consequently, the resulting pattern discovery processes are efficient not only theoretically but also on a wide range of real-world benchmark datasets.

2. PRELIMINARIES

Before going into technical details, we fix some basic notions and notation. For a finite set X we denote by $\mathcal{P}(X)$ its power set and by $u(X)$ the uniform probability distribution on X . Moreover, for positive weights $w: X \rightarrow \mathbb{R}^+$ let $w(X)$ denote the distribution on X arising from normalizing w , i.e., the distribution described by $x \mapsto w(x) / \sum_{x' \in X} w(x')$ —assuming that there is an $x \in X$ with $w(x) > 0$.

A binary **dataset** \mathcal{D} over some finite **ground set** E is a bag (multiset) of sets, called **data records**, D_1, \dots, D_m each of which being a subset of $E = \{e_1, \dots, e_n\}$. As **size** of \mathcal{D} , denoted $\|\mathcal{D}\|$ we define the sum of all its data record sizes $\sum_{D \in \mathcal{D}} |D|$. Inspired by the application of market basket analysis the elements of E are often referred to

as “items”. More generally, one can think of E as a set of binary features describing the data records. In particular, a categorical data table can easily be represented as a binary dataset by choosing the ground set as consisting of all attribute/value equality expressions that can be formed from the table. More precisely, a **categorical data table** T consisting of m data row vectors d_1, \dots, d_m with $d_i = (d_i(1), \dots, d_i(n))$ can be represented by the dataset $\mathcal{D}_T = \{D_1, \dots, D_m\}$ with $D_i = \{(j, v) : d_i(j) = v\}$ over ground set

$$E_T = \{(j, d_i(j)) : 1 \leq i \leq m, 1 \leq j \leq n\} .$$

For a given dataset \mathcal{D} over E , the **pattern space** (or *pattern language*) $\mathcal{L}(\mathcal{D})$ considered in this paper is the power set $\mathcal{P}(E)$ of the features and its elements are interpreted conjunctively. That is, the *local* data portion described by a set $F \subseteq E$, called the **support (set)** of F in \mathcal{D} and denoted $\mathcal{D}[F]$, is defined as the multiset of all data records from \mathcal{D} that contain *all* elements of F , i.e., $\mathcal{D}[F] = \{D \in \mathcal{D} : D \supseteq F\}$.

An **interestingness measure** for a pattern language $\mathcal{L}(\cdot)$ is a function

$$q: \{(\mathcal{D}, x) : \mathcal{D} \text{ a binary dataset, } x \in \mathcal{L}(\mathcal{D})\} \rightarrow \mathbb{R} .$$

However, often there is a fixed dataset that is clear from the context. In such cases—and if we want to simplify the notation—we just write q as an unary function $q(\cdot) = q(\mathcal{D}, \cdot)$ and omit the first argument. The most basic measures for set patterns are the **support (count)**, i.e., the size of its support set $q_{\text{supp}}(\mathcal{D}, F) = |\mathcal{D}[F]|$ and the **frequency**, i.e., the relative size of its support with respect to the total number of data records $q_{\text{freq}}(\mathcal{D}, F) = |\mathcal{D}[F]| / |\mathcal{D}|$. For a frequency threshold $t \in [0, 1]$ a set is called **t -frequent** (w.r.t. \mathcal{D}) if $q_{\text{freq}}(\mathcal{D}, F) \geq t$. A further measure considered here is the **area function** [15] $q_{\text{area}}(\mathcal{D}, F) = |F| |\mathcal{D}[F]|$. Intuitively, the area of a set corresponds to the number of 1 entries of the submatrix (of the binary matrix representation of \mathcal{D}) consisting of the columns corresponding to F and the rows corresponding to $\mathcal{D}[F]$.

All measures defined so far are unsupervised measures in the sense that they rely on no further information but the dataset itself. In contrast, there are so-called supervised descriptive rule induction techniques that rely on additional information in the form of **class labels** $l(D) \in C = \{c_1, \dots, c_k\}$ associated to each data record $D \in \mathcal{D}$. For $c \in C$ we denote by \mathcal{D}_c the data portion labeled c , i.e., $\mathcal{D}_c = \{D \in \mathcal{D} : l(D) = c\}$. Examples for this setting are emerging pattern mining [12] and contrast set mining [3], where one is interested in patterns having a high support difference between the positive and the negative portion of the data records, or subgroup discovery [24], where one searches for patterns with a high distributional unusualness of these labels on their support set. In important special case are binary labels, i.e., $C = \{\oplus, \ominus\}$. For this case we consider the following **discriminativity measure**

$$q_{\text{disc}}(F) = |\mathcal{D}_{\oplus}[F]| |\mathcal{D}_{\ominus} \setminus \mathcal{D}_{\oplus}[F]| .$$

A further measure for the discriminative power of a pattern is the **Fisher score** q_{fish} , which is defined for datasets with arbitrary labels C . Intuitively, it measures the relation of the

inter-class variance of a feature to its intra-class variances, i.e.,

$$q_{\text{fish}}(F) = \frac{\sum_{c \in C} |\mathcal{D}_c| (q_{\text{freq}}(\mathcal{D}_c, F) - q_{\text{freq}}(\mathcal{D}, F))^2}{\sum_{c \in C} \sum_{D \in \mathcal{D}_c} (\delta(D \supseteq F) - q_{\text{freq}}(\mathcal{D}_c, F))^2}$$

where $\delta(D \supseteq F)$ is 1 if $D \supseteq F$ and 0 otherwise. This paper does not present a sampling algorithm for this measure. However, the Fisher score is used for post-processing generated patterns in the context of constructing global classification models.

3. SAMPLING ALGORITHMS

After the introduction of set patterns and interestingness measures, we can now present our sampling procedures. A naive approach for sampling a pattern according to a distribution π is to generate a list F_1, \dots, F_N of all patterns with $\pi(F) > 0$, draw an $x \in [0, 1]$ uniformly at random, and then return the unique set F_k with $\sum_{i=1}^{k-1} \pi(F_i) \leq x < \sum_{i=1}^k \pi(F_i)$. However, the exhaustive enumeration of any non-constant part of the pattern space is precisely what we want to avoid. That is, we are interested in *non-enumerative* sampling algorithms. Below we give such algorithms for four quality functions: frequency and squared frequency as well as area and discriminativity.

Note that, in contrast to the frequency measures, for the latter to quality functions it is **NP**-hard to find optimal patterns: Finding a set of maximum area for a given input dataset is equivalent to the **NP**-hard problem of computing a biclique with a maximum number of edges from a given bipartite graph (see [15]). The same hardness result holds for the discriminativity measure because optimizing area can be linearly reduced to optimizing discriminativity: by setting \mathcal{D}_{\oplus} to \mathcal{D} and \mathcal{D}_{\ominus} to $\{E \setminus \{e\} : e \in E\}$ we get $q_{\text{disc}}(\mathcal{D}_{\oplus} \cup \mathcal{D}_{\ominus}, F) = q_{\text{area}}(\mathcal{D}, F)$.

3.1 Frequency and Area

Algorithm 1 Frequency-based Sampling

Require: dataset \mathcal{D} over ground set E ,
Returns: random set $R \sim q_{\text{freq}}(\mathcal{P}(E)) = q_{\text{supp}}(\mathcal{P}(E))$

1. **let** weights w be defined by $w(D) = 2^{|D|}$ for all $D \in \mathcal{D}$
 2. **draw** $D \sim w(\mathcal{D})$
 3. **return** $R \sim u(\mathcal{P}(D))$
-

We start with sampling according to frequency and area, both of which can be achieved by very similar linear time algorithms. The key insight for frequency-based sampling, i.e., $\pi = q_{\text{freq}}(\mathcal{P}(E))$, is that random experiments are good in reproducing frequent events. Namely, if we look at a pattern that is supported by a random data record we are likely to observe a pattern that is supported by many data records altogether. This intuition leads to a two-step non-enumerative sampling routine (see Algorithm 1), which is as fast as simple: First select a data record of the input dataset randomly with a probability proportional to the size of its power set. Then return a uniformly sampled subset of that data record. Using the size of the power set in the first step is important, as otherwise the sampling routine would be biased towards sets occurring in small data records. As

noted in Proposition 1 below, the random set resulting from combining both steps follows the desired distribution.

Regarding the computational complexity of the sampling algorithm we can observe that it is indeed efficient: if one has knowledge of the numbers $|D|$ for all data records $D \in \mathcal{D}$ and, moreover, has index access to all data records, a single random set can be produced in time $O(\log |\mathcal{D}| + |E|)$ (the two terms correspond to producing a random number for drawing a data record in step 1 and of drawing one of its subsets in step 2). Both requirements can be achieved via a single initial pass over the dataset. Thus, we have the following proposition.

PROPOSITION 1. *On input dataset \mathcal{D} over E , a family of k realizations of the random set $\mathbf{R} \sim q_{\text{freq}}(\mathcal{P}(E))$ can be generated in time $O(\|\mathcal{D}\| + k(|E| + \ln |\mathcal{D}|))$.*

PROOF. Let Z be the normalizing constant $\sum_{F \subseteq E} |\mathcal{D}[F]|$ and \mathbf{D} denote the random data record that is drawn in step 2 of the algorithm. For the probability distribution of the returned random set it holds that

$$\begin{aligned} \mathbb{P}[\mathbf{R} = R] &= \sum_{D \in \mathcal{D}} \mathbb{P}[\mathbf{R} = R \wedge \mathbf{D} = D] \\ &= \sum_{D \in \mathcal{D}[R]} \frac{1}{2^{|D|}} \frac{2^{|D|}}{Z} \\ &= \frac{|\mathcal{D}[R]|}{Z} = \frac{q_{\text{supp}}(\mathcal{D}, R)}{Z} \end{aligned}$$

with a normalizing $Z = \sum_{D \in \mathcal{D}} 2^{|D|}$ (which is equal to the desired $\sum_{F \subseteq E} |\mathcal{D}[F]|$). \square

Algorithm 2 Area-based Sampling

Require: dataset \mathcal{D} over ground set E ,
Returns: random set $R \sim q_{\text{area}}(\mathcal{P}(E))$

1. **let** weights w be defined for all $D \in \mathcal{D}$ by

$$w(D) = |D| 2^{|D|-1}$$

2. **draw** $D \sim w(\mathcal{D})$
 3. **draw** $k \sim \text{id}(\{1, \dots, |D|\})$ with weights $\text{id}(i) = i$
 4. **return** $R \sim u(\{F \subseteq D : |F| = k\})$
-

Sampling according to area, i.e., $\pi = q_{\text{area}}(\mathcal{P}(E))$, can be achieved via a slight modification of frequency-based sampling: in step two, instead of drawing a subset uniformly from a data record, draw a subset with probability proportional to its size. As side effect, this modification affects the normalization constants and in particular the data record weights of step one. As for a data record D it holds for the sum of all its subset sizes that

$$\sum_{F \subseteq D} |F| = |D| 2^{|D|-1}$$

we have to modify the data record weights accordingly. The resulting pseudo-code is given in Algorithm 2. Again, after all weights have been computed via an initial pass over the data, an arbitrary number of random sets can be produced in time $O(\log |\mathcal{D}| + |E|)$. Hence, with a similar proof as for Proposition 1 we can conclude:

PROPOSITION 2. *On input dataset \mathcal{D} over E , a family of k realizations of the random set $\mathbf{R} \sim q_{\text{area}}(\mathcal{P}(E))$ can be generated in time $O(\|\mathcal{D}\| + k(|E| + \ln |\mathcal{D}|))$.*

It is an important remark that area can be replaced by *weighted* area relatively easy without changing the asymptotic complexity—where weighted area is defined as

$$q_{\text{ware}}(F) = \left(\sum_{e \in F} w(e) \right) \left(\sum_{D \in \mathcal{D}[F]} w(D) \right)$$

for a set of positive weights $w: (E \cup \mathcal{D}) \rightarrow \mathbb{R}_+$. The same holds for weighted frequency. In this paper, however, for the sake of simplicity we only consider the unweighted case.

3.2 Discriminativity and Squared Frequency

Algorithm 3 Discriminativity-based Sampling

Require: binary labeled dataset \mathcal{D} over ground set E
such that there is an $F \subseteq E$ with $q_{\text{disc}}(F) > 0$
Returns: random set $R \sim q_{\text{disc}}(\mathcal{P}(E))$

1. **let** weights w be defined by

$$w(D_{\oplus}, D_{\ominus}) = (2^{|D_{\oplus} \setminus D_{\ominus}|} - 1) 2^{|D_{\oplus} \cap D_{\ominus}|}$$

for all $(D_{\oplus}, D_{\ominus}) \in \mathcal{D}_{\oplus} \times \mathcal{D}_{\ominus}$

2. **draw** $(D_{\oplus}, D_{\ominus}) \sim w(\mathcal{D}_{\oplus} \times \mathcal{D}_{\ominus})$
3. **draw** $F \sim u(\mathcal{P}(D_{\oplus} \setminus D_{\ominus}) \setminus \emptyset)$ and $F' \sim u(\mathcal{P}(D_{\oplus} \setminus D_{\ominus}))$
4. **return** $R = (F \cup F')$

In order to design a sampling procedure for discriminativity, i.e., $\pi = q_{\text{disc}}(\mathcal{P}(E))$, we can lift the principle of frequency-based sampling to a little more complicated random experiment with the following intuition. If we look at a pattern that is supported by a random positive data record and *not* supported by a random negative data record, we are likely to observe a pattern that is altogether supported by many positive data records and only few negative data records, i.e., we are likely to observe a pattern with a relatively high discriminativity score. Again, in order to control the resulting distribution, it is necessary consider a pair of data records $(D_{\oplus}, D_{\ominus})$ with a probability equal to the number of sets $F \subseteq E$ with $F \subseteq D_{\oplus}$ and $F \not\subseteq D_{\ominus}$. This implies that the increased expressivity of discriminativity compared to frequency comes at a price: due to the necessity of weight computation for all pairs of positive and negative data records, we end up with a quadratic preprocessing phase. Algorithm 3 contains all details of the resulting sampling procedure and leads to the following result.

PROPOSITION 3. *Let \mathcal{D} be a binary labeled input dataset over ground set E such that there is a set $F \subseteq E$ with $q_{\text{disc}}(\mathcal{D}, F) > 0$. A family of k realizations of the random set $\mathbf{R} \sim q_{\text{disc}}(\mathcal{P}(E))$ can be generated in time $O(\|\mathcal{D}\|^2 + k(|E| + \ln^2 |\mathcal{D}|))$.*

PROOF. Let \mathbf{R} denote the random set returned in step 4 of the algorithm and $\mathbf{D}_{\oplus}, \mathbf{D}_{\ominus}$ the data records drawn in step 2. Moreover, for $D \in \mathcal{D}_{\oplus}$ and $D' \in \mathcal{D}_{\ominus}$ let $\delta(D, D')$ denote the family of all sets $F \subseteq E$ that are supported by

D but not supported by D' . We can rewrite this definition as

$$\begin{aligned} \delta(D, D') &= \{F \subseteq E : F \subseteq D, F \not\subseteq D'\} \\ &= \{F \cup F' : \emptyset \subset F \subseteq (D \setminus D'), F' \subseteq (D \cap D')\} . \end{aligned}$$

This form shows that the weights $w(\cdot, \cdot)$ assigned in step 1 are equivalent to $|\delta(\cdot, \cdot)|$ and that, moreover, \mathbf{R} is a set drawn uniformly from $\delta(\mathbf{D}_{\oplus}, \mathbf{D}_{\ominus})$. With this we can conclude similar to the previous algorithms that

$$\begin{aligned} \mathbb{P}[\mathbf{R} = F] &= \sum_{D \in \mathcal{D}_{\oplus}} \sum_{D' \in \mathcal{D}_{\ominus}} \mathbb{P}[\mathbf{R} = F, \mathbf{D}_{\oplus} = D, \mathbf{D}_{\ominus} = D'] \\ &= \sum_{D, D' \in \delta^{-1}[F]} \frac{1}{|\delta(D, D')|} \frac{w(D, D')}{Z} \\ &= \frac{1}{Z} |\{(D, D') \in \mathcal{D}_{\oplus} \times \mathcal{D}_{\ominus} : D \supseteq F, D' \not\supseteq F\}| \\ &= \frac{1}{Z} |\mathcal{D}_{\oplus}[F]| (|\mathcal{D}_{\ominus}| - \mathcal{D}_{\ominus}[F]) \end{aligned}$$

with $Z = \sum_{D, D' \in \mathcal{D}_{\oplus} \times \mathcal{D}_{\ominus}} |\delta(D, D')| = \sum_{F \subseteq E} q_{\text{disc}}(F)$ as required. \square

Algorithm 4 Squared-frequency-based Sampling

Require: dataset \mathcal{D} over ground set E ,
Returns: random set $F \sim q_{\text{freq}}^2(\mathcal{P}(E)) = q_{\text{supp}}^2(\mathcal{P}(E))$

1. **let** weights w be defined by

$$w(D_1, D_2) = 2^{|D_1 \cap D_2|}$$

for all $(D_1, D_2) \in \mathcal{D} \times \mathcal{D}$

2. **draw** $(D_1, D_2) \sim w(\mathcal{D} \times \mathcal{D})$
3. **return** $F \sim u(\mathcal{P}(D_1 \cap D_2))$

It is straightforward to see that the approach of drawing two data records can also be used to express other potentially interesting distributions that are given as the product of two support counts. A basic example is squared frequency. In order to achieve this distribution, one can consider a uniformly² drawn subset of two random data records, i.e., a subset of their intersection. The resulting pseudo-code with appropriate pairwise weights is given in Algorithm 4. Closely following the proof of Prop. 3 this algorithm can be used to show another proposition.

PROPOSITION 4. *On input dataset \mathcal{D} over E , a family of k realizations of the random set $\mathbf{R} \sim q_{\text{freq}}^2(\mathcal{P}(E))$ can be generated in time $O(\|\mathcal{D}\|^2 + k(|E| + \ln^2 |\mathcal{D}|))$.*

In principle, one can design sampling algorithms for an arbitrary power c of the frequency measure by drawing a subset from c random data records. However, the resulting time complexity for computing the weights for each c -tuple of data records gets out of hand quickly.

²For sampling according to the squared area function, draw a subset with probabilities proportional to its squared size instead of uniformly.

4. EVALUATION

The sampling procedures presented in the previous section are provably efficient and correct, i.e., their randomized output follows the specified distributions. It remains to evaluate, beside their practical scalability, how useful these distributions are in the context of local pattern discovery. It is inherently difficult to evaluate pattern discovery methods in an exploratory data analysis context. There one aims to find *interesting* patterns, which relies on an often user-subjective notion of interestingness. Hence, here we resort to the other branch of local pattern discovery applications, i.e., pattern-based global model construction, which allows us to simply use accuracy as objective evaluation metric.

dataset	class	nm/ct	items	rows	density
autos	7	15/10	135	205	0.190
balance-scale	3	4/0	20	625	0.250
breast-cancer	2	0/9	51	286	0.195
colic	3	7/15	84	366	0.271
credit-a	2	6/9	71	690	0.223
diabetes	2	8/0	40	768	0.225
glass	7	9/0	45	214	0.222
heart-c	5	6/7	49	303	0.285
heart-h	5	6/7	46	294	0.246
heart-statlog	2	13/0	55	270	0.254
hepatitis	2	6/13	55	155	0.344
hypothyroid	4	7/22	78	3772	0.364
iris	3	4/0	20	150	0.250
lymph	4	3/15	57	148	0.333
prim.-tumor	22	0/17	37	339	0.468
sonar	2	60/0	300	208	0.203
tic-tac-toe	2	0/9	27	958	0.370
vehicle	4	18/0	90	846	0.211
zoo	7	1/16	135	101	0.133

Table 1: Benchmark datasets with basic statistics: number of classes $|C|$, number of numerical and categorical columns (nm/ct), number $|E_T|$ of items in corresponding binary dataset, number of rows, density $|\mathcal{D}_T| |E_T| / |\mathcal{D}_T|$.

In our experiments we use a variety of databases from the UCI machine learning repository [14] listed in Table 1. In order to apply the pattern discovery algorithms, binary datasets are created from these databases by first converting them into categorical datatables using five bucket frequency discretization of all numeric data columns, and then by considering the corresponding binary datasets (using attribute / value pairs of categorical attributes as items; see Section 2). Implementations of the sampling algorithms are available in the software section of <http://www-kd.iai.uni-bonn.de>.

4.1 Predictive Performance

We start with experiments evaluating the sampling algorithms in the context of pattern-based classification. Here one aims to improve classification accuracy by enriching given labeled training data with pattern-based features. We closely follow the framework of Cheng et al. [9]. In a nutshell it consists of three basic steps: extraction of a collection of patterns (which are subsequently considered as features of the data records supporting them), feature selection based on Fisher score and pattern redundancy, and classification, for which we use a linear support vector machine.

In more detail, for an input data table T with corresponding

binary dataset $\mathcal{D}_T = \mathcal{D}$ and class labels $C = \{1, \dots, c\}$ the following **pattern collections** are considered:

- a frequent sets collection \mathcal{F} consisting of the union of the collection of t -frequent sets of each of the datasets $\mathcal{D}_1, \dots, \mathcal{D}_c$ where $t \in \{1, 0.95, \dots, 0.5\}$ is chosen minimal such that $|\mathcal{F}|$ is less than 200,000 (number is chosen to keep computation times including feature selection below five minutes per training dataset)

as well as three collections each consisting of $k = \alpha |\mathcal{D}|$ independent random sets, namely

- a frequency-based collection \mathcal{R}_{frq} where the i -th set is sampled according to $q_{\text{frq}}(\mathcal{D}_{\lceil ic/k \rceil}, \mathcal{P}(E))$,
- a squared-frequency-based collection \mathcal{R}_{sfr} where the i -th set is sampled according to $q_{\text{frq}}^2(\mathcal{D}_{\lceil ic/k \rceil}, \mathcal{P}(E))$,
- and a discriminativity-based collection \mathcal{R}_{dcr} where the i -th set is sampled according $q_{\text{disc}}(\mathcal{D}_{\oplus} \cup \mathcal{D}_{\ominus}, \mathcal{P}(E))$ with $\mathcal{D}_{\oplus} = \mathcal{D}_{\lceil ic/k \rceil}$ and $\mathcal{D}_{\ominus} = \mathcal{D} \setminus \mathcal{D}_{\lceil ic/k \rceil}$.

Area-based sampling is not considered here, because it is not designed towards providing good features for classification. The parameter α is set to 32 independent of the dataset. With this setting, the pattern selection process is sufficiently stable as captured by the average Fisher score of finally selected features (using the selection procedure described below). See Figure 2 for an illustration on three exemplary datasets.

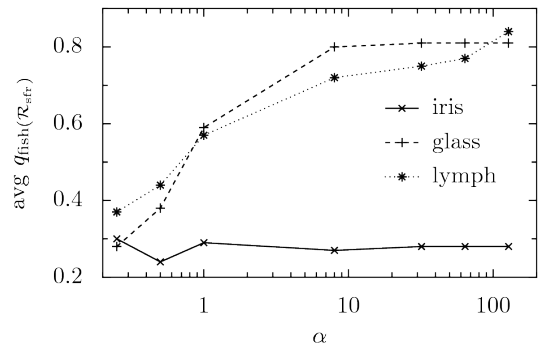


Figure 2: Average Fisher score of the squared-frequency-based random pattern collection resulting from different values of α .

The **feature selection** step for a pattern collection \mathcal{P} is then performed as follows: initialize $\mathcal{P}_0 = \emptyset$ and consider all sets $F_1, \dots, F_l \in \mathcal{P}$ having $q_{\text{frq}}(\mathcal{D}, \cdot) \geq 0.05$ in descending order of their Fisher score $q_{\text{fish}}(\mathcal{D}, \cdot)$. Select F_i if

$$q_{\text{fish}}(F_i) - \max_{F \in \mathcal{P}_i} (r(F_i, F) \min\{q_{\text{fish}}(F_i), q_{\text{fish}}(F)\}) > 0$$

where \mathcal{P}_i is the family of sets already selected when considering F_i and $r(\cdot, \cdot)$ is the redundancy measure given by

$$r(F, F') = \frac{q_{\text{frq}}(\mathcal{D}, F \cup F')}{q_{\text{frq}}(\mathcal{D}, F) + q_{\text{frq}}(\mathcal{D}, F') - q_{\text{frq}}(\mathcal{D}, F \cup F')} .$$

Note that this is a deviation from Cheng et al. where after each selection the remaining patterns are reordered according to the above expression (again, our intention is to avoid computation times of more than five minutes per training dataset). The selection process is stopped after all patterns have been considered or if \mathcal{P}_i covers 80% of the dataset, i.e.,

$$\left| \left(\bigcup_{F \in \mathcal{P}_i} \mathcal{D}[F] \right) / |\mathcal{D}| \right| < 0.2 |\mathcal{D}| .$$

For the final pattern collection \mathcal{P}' the original data table T is then augmented by binary attributes corresponding to the elements of $F_1, \dots, F_k \in \mathcal{P}'$. That is, the augmented table T' has $n + |\mathcal{P}'|$ columns with rows defined by

$$t'_i(j) = \begin{cases} t_i(j), & \text{if } j \leq n \\ 1, & \text{if } j > n \text{ and } D_i \supseteq F_{j-n} \\ 0, & \text{otherwise} \end{cases}$$

where D_i is the data record of \mathcal{D}_T corresponding to the i -th row of T .

dataset	plain	\mathcal{F}	q_{freq}	q_{freq}^2	q_{disc}
autos	79.50	78.25	79.00	79.50	80.00
balance-scale	84.59	84.59	84.59	84.59	84.59
breast-cancer	70.36	74.28	74.11	73.57	73.39
colic	61.39	60.14	61.39	61.39	61.39
credit-a	85.84	86.50	84.96	86.50	86.13
diabetis	74.80	74.14	74.47	74.14	74.40
glass	62.93	68.29	68.78	69.02	66.58
heart-c	79.83	81.87	82.21	81.53	80.00
heart-h	81.48	81.75	81.05	81.75	81.40
heart-statlog	80.76	83.21	82.64	82.64	83.21
hepatitis	81.67	83.00	83.00	83.33	83.00
iris	89.31	88.62	88.62	88.97	88.28
lymph	82.86	82.50	83.93	84.28	83.57
primary-tumor	40.30	46.36	45.61	46.67	45.76
sonar	80.00	78.50	80.00	78.75	80.00
tic-tac-toe	77.05	99.21	99.42	86.26	94.53
vehicle	68.27	68.87	68.69	68.69	69.05
zoo	91.05	91.58	92.10	91.58	91.58

Table 2: Results of SVM classification on plain database and with feature enrichment based on frequent sets and sampled pattern collections (frequency, squared frequency, and discriminativity).

As **classifier** the linear SVM of the LIBSVM software is used wrapped in an optimization layer for its regularization parameter c . That is, the training set is first used to determine the optimal regularization parameter $c \in \{2^i : i = -5, -3, \dots, 14\}$ using 5-fold cross-validation and then a model is trained with the optimal parameter using the complete training set. The complete workflow is validated using two times 5-fold cross-validation for all pattern collections simultaneously.

Table 2 contains the results. A Wilcoxon signed ranks test (see [11]) for our $N = 18$ databases reveals that pattern-based classification with each of the random set collections outperforms the plain SVM significantly at the 2%-level (t-values of 30.5, 22.0, and 21.5 respectively; critical value 33).

Moreover, albeit all random set collections are lying ahead of the frequent set collection on our test databases, it is not significantly outperformed by any of them. We can conclude that pattern-based classification based on all tested sampling algorithms is likely to outperform the plain SVM, and is unlikely to be inferior to standard frequent-pattern-based classification.

4.2 Scalability

Having evaluated the quality of the sampled patterns we now turn to scalability studies. The theoretical potential of the direct sampling procedures is already indicated by the guarantees of Propositions 1-4. In particular for frequency and area-based sampling they suggest applicability on larger to large-scale datasets. Below we investigate to what degree this potential can be transferred into practice.

Regarding a comparison with the MCMC pattern sampling algorithms, the practical advantages are very clear: while for instance the closed set Markov chain simulation [6] takes seven minutes to sample a closed set from a 30K row subset of the US census dataset, frequency-based direct sampling takes only 0.067 seconds for the first sample including preprocessing—afterwards additional samples can be produced in milliseconds.

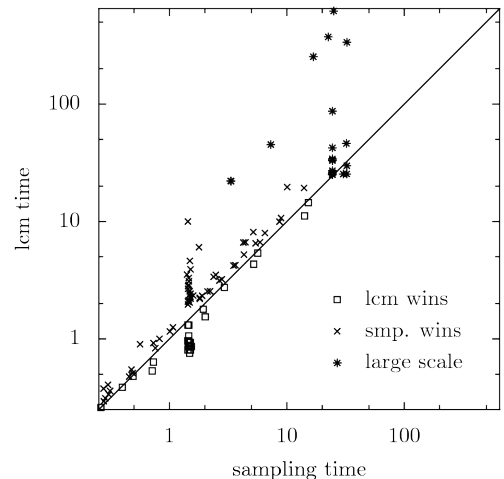


Figure 3: Time of lcm versus time of frequency-based sampling for an identical number of patterns.

For a comparison with exhaustive search algorithms there is a very large number of possible contenders in the literature, each of which in principle requires an individual comparative study. Hence, for this paper we resort to a rather general setting that constitutes a certain worst-case setting: we compare the computation time of the sampling algorithms with that of the linear closed frequent set mining algorithm (lcm) of Uno et al. [23] *per pattern*. This constitutes a worst-case setting because frequent set mining algorithms usually produce the largest output per time unit and lcm is known to be among the fastest of them (winner of the FIMI contest [4]). In addition to the datasets used for the predictive performance study, we also consider several of the larger benchmark datasets of the FIMI workshop, including the 1GB sized “webdocs”, and a 500MB random dataset.

The results are presented as log/log scatter-plot in Figure 3. One can observe that for most configurations lcm and the frequency-based sampling generate their patterns in approximately equal time with the majority of wins going to the sampling. However, focusing on the configurations with large-scale datasets (star symbol) reveals that the sampling algorithm can substantially outperform lcm. For “webdocs” this includes a speed-up factor of 10, for the random dataset even one of 25. We can conclude that frequency-based sampling can substantially outperform (closed) frequent set listing on large datasets and it performs equally well with slight advantage on small-scale data.

hepatitis 0.03	heart-c 0.05	glass.dat 0.03	colic 0.05
bal.-scale 0.1	vehicle 0.15	lymph 0.03	autos 0.03
credit-a 0.2	b-cancer 0.03	tumor 0.05	heart-h 0.05
zoo 0.03	iris 0.03	h-statlog 0.03	hypothyroid 13.69

Table 3: Preprocessing time, i.e., weight computation, for squared-frequency sampling.

While the time performance of frequency-based sampling is also representative for area-based sampling, this is not true for the two sampling procedures with quadratic time weight computation phase. For several datasets the time for this preprocessing step is listed in Table 3. One can observe that for the smaller datasets, weight computation time is only marginal. Moreover, after this phase the performance is essentially equal to frequency-based sampling. However, for large-scale datasets the quadratic complexity is prohibitive; as already indicated on the 4000 row dataset “hypothyroid”. Also, including the preprocessing phase squared-frequency-based sampling can generally not compete with lcm (again measured per pattern). Naturally, the same holds for sampling based on discriminativity.

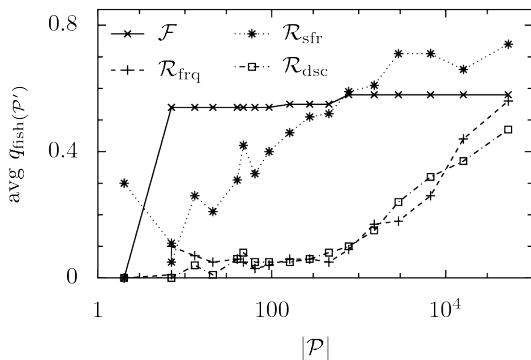


Figure 4: Average Fisher score of selected patterns per collection size for “lymph”.

We conclude this section with a note on pattern quality per time—or equally per collection size because all methods (neglecting the preprocessing phase) use roughly linear time in the size of the pattern collections they produce. Here, we again use the feature selection procedure from Section 4.1 and consider the average Fisher score of selected features

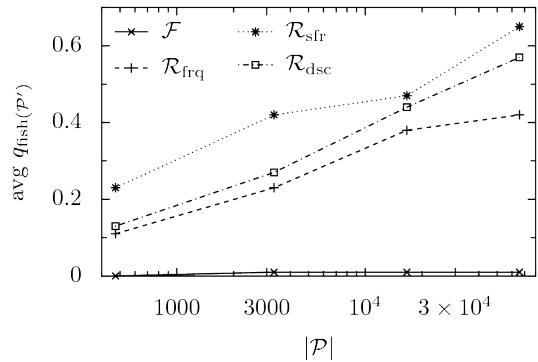


Figure 5: Average Fisher score of selected patterns per collection size for “hypo”.

from frequent set respectively random set families of different sizes. Naturally, the outcome depends on whether the datasets contains high-frequency patterns with discriminative power. If this is the case, frequent sets can quickly provide good features (see, e.g., Figure 4). However, it is not uncommon to observe cases in which there is an extremely large number of high-frequency patterns, non of which possesses any discriminative power. In such cases, the sampling approaches need to generate substantially fewer patterns to provide good features than a frequent set listing algorithm (see, e.g., Figure 5).

5. CONCLUSION

We introduced four simple direct sampling procedures that generate random set patterns distributed according to frequency, squared frequency, area, and a discriminativity measure for binary labels. All procedures come with tight theoretical performance guarantees. Moreover, we described experimental studies demonstrating that the produced patterns are as useful as frequent pattern collections for pattern-based classification, and that direct sampling can compete with and often even outperform the fastest exhaustive mining algorithms when generating an equal number of patterns.

In the context of pattern-based classification there is a large amount of pattern discovery approaches that range from optimistic-estimator-based best-first-search algorithms [21] to methods interweaving model training and pattern discovery [10, 13]. Although such algorithms typically traverse much less patterns per time unit as lcm, their search is more directed towards high quality patterns. This motivates an in-depth comparative study with such methods potentially leading to more sophisticated usage of the sampling algorithms (e.g., applying it within model training just as the cited approaches do with exhaustive mining).

That said, pattern sampling as a paradigm is in no way restricted to pattern-based classification, and should also be evaluated for other in particular unsupervised model construction tasks as well as for exploratory data analysis. This is likely to motivate further variants of pattern sampling procedures. An example is the introduction of column and row weights to the interestingness measure in order to model subjective interest in certain parts of the input data or to decrease the probability of re-discovering redundant patterns.

6. REFERENCES

- [1] R. Agrawal, H. Mannila, R. Srikant, H. Toivonen, and A. I. Verkamo. Fast discovery of association rules. In *Advances in Knowledge Discovery and Data Mining*, pages 307–328. 1996.
- [2] M. Al Hasan and M. J. Zaki. Output space sampling for graph patterns. *PVLDB*, 2(1):730–741, 2009.
- [3] S. D. Bay and M. J. Pazzani. Detecting group differences: Mining contrast sets. *Data Mining and Knowledge Discovery*, 5(3):213–246, 2001.
- [4] R. Bayardo, B. Goethals, and M. J. Zaki, editors. *Proc. of the IEEE ICDM Workshop on Frequent Itemset Mining Implementations, 2004*, volume 126 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2004.
- [5] S. Bistarelli and F. Bonchi. Soft constraint based pattern mining. *Data and Knowledge Engineering*, 62(1):118–137, 2007.
- [6] M. Boley, T. Gärtner, and H. Grosskreutz. Formal concept sampling for counting and threshold-free local pattern mining. In *Proc. of the SIAM Int. Conf. on Data Mining (SDM 2010)*, pages 177–188, 2010.
- [7] M. Boley and H. Grosskreutz. Approximating the number of frequent sets in dense data. *Knowledge and Information Systems*, 21(1):65–89, 2009.
- [8] V. Chaoji, M. A. Hasan, S. Salem, J. Besson, and M. J. Zaki. Origami: A novel and effective approach for mining representative orthogonal graph patterns. *Statistical Analysis and Data Mining*, 1(2):67–84, 2008.
- [9] H. Cheng, X. Yan, J. Han, and C.-W. Hsu. Discriminative frequent pattern analysis for effective classification. In *Proc. of the 23rd Int. Conf. on Data Engineering (ICDE 2007)*, pages 716–725, 2007.
- [10] H. Cheng, X. Yan, J. Han, and P. S. Yu. Direct discriminative pattern mining for effective classification. In *Proc. of the 24th Int. Conf. on Data Engineering (ICDE 2008)*, pages 169–178, 2008.
- [11] J. Demsar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research*, 7:1–30, 2006.
- [12] G. Dong and J. Li. Efficient mining of emerging patterns: Discovering trends and differences. In *Proc. 5th ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining (KDD '99)*, pages 43–52. ACM, 1999.
- [13] W. Fan, K. Zhang, H. Cheng, J. Gao, X. Yan, J. Han, P. S. Yu, and O. Verscheure. Direct mining of discriminative and essential frequent patterns via model-based search tree. In *Proc. of the 14th Int. Conf. on Knowledge Discovery and Data Mining (KDD '08)*, pages 230–238, 2008.
- [14] A. Frank and A. Asuncion. UCI machine learning repository, 2010.
- [15] F. Geerts, B. Goethals, and T. Mielikäinen. Tiling databases. In *Proc. 7th Int. Discovery Science Conf.*, volume 3245 of *LNCS*, pages 278–289. Springer, 2004.
- [16] H. Grosskreutz, S. Rüping, and S. Wrobel. Tight optimistic estimates for fast subgroup discovery. In *Proc. of the European Conf. on Machine Learning and Knowledge Discovery in Databases (ECML/PKDD 2008), Part I*, volume 5211 of *LNCS*, pages 440–456, 2008.
- [17] D. Gunopulos, H. Mannila, and S. Saluja. Discovering all most specific sentences by randomized algorithms. In *Proc. of 6th Int. Conf. of Database Theory (ICDT '97)*, volume 1186 of *LNCS*, pages 215–229, 1997.
- [18] J. Han, J. Pei, Y. Yin, and R. Mao. Mining frequent patterns without candidate generation: A frequent-pattern tree approach. *Data Min. Knowl. Discov.*, 8(1):53–87, 2004.
- [19] D. J. Hand. Pattern detection and discovery. In *Proc. of the ESF Exploratory Workshop on Pattern Detection and Discovery*, volume 2447 of *LNCS*, pages 1–12. Springer, 2002.
- [20] A. J. Knobbe, B. Crémilleux, J. Fürnkranz, and M. Scholz. From local patterns to global models: the lego approach to data mining. In *From Local Patterns to Global Models: Proceedings of the ECML/PKDD 2008 Workshop (LEGO '08)*, 2008.
- [21] S. Morishita and J. Sese. Traversing itemset lattice with statistical metric pruning. In *Proc. of 19th ACM SIGMOD-SIGACT-SIGART Symp. on Principles of Database Systems (PODS)*, pages 226–236, 2000.
- [22] T. Scheffer and S. Wrobel. Finding the most interesting patterns in a database quickly by using sequential sampling. *Journal of Machine Learning Research*, 3:833–862, 2002.
- [23] T. Uno, T. Asai, Y. Uchida, and H. Arimura. An efficient algorithm for enumerating closed patterns in transaction databases. In *Proc. of the 7th Int. Discovery Science Conf. (DS 2004)*, volume 3245 of *LNCS*, pages 16–31. Springer, 2004.
- [24] S. Wrobel. An algorithm for multi-relational discovery of subgroups. In *Proc. 1st Euro. Symp. on Principles of Data Mining and Knowledge Discovery (PKDD '97)*, volume 1263 of *LNCS*, pages 78–87. Springer, 1997.
- [25] M. J. Zaki, S. Parthasarathy, W. Li, and M. Ogihara. Evaluation of sampling for data mining of association rules. In *Proc. of 7th Workshop on Research Issues in Data Engineering (RIDE)*, 1997.