



HAL
open science

Compressing and Querying Skypattern Cubes

Willy Ugarte, Samir Loudni, Patrice Boizumault, Bruno Crémilleux,
Alexandre Termier

► **To cite this version:**

Willy Ugarte, Samir Loudni, Patrice Boizumault, Bruno Crémilleux, Alexandre Termier. Compressing and Querying Skypattern Cubes. IEA/AIE-2019 - 32nd International Conference on Industrial, Engineering & Other Applications of Applied Intelligent Systems, Jul 2019, Graz, Austria. pp.406-421, 10.1007/978-3-030-22999-3_36 . hal-02190788

HAL Id: hal-02190788

<https://hal.science/hal-02190788v1>

Submitted on 22 Jul 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Compressing and Querying Skypattern Cubes

Willy Ugarte¹, Samir Loudni², Patrice Boizumault², Bruno Crémilleux², and
Alexandre Termier³

¹ Peruvian University of Applied Sciences, Peru
willyugarte@gmail.com

² Normandie Univ., UNICAEN, ENSICAEN, CNRS – UMR GREYC, France
{name}. {last-name}@unicaen.fr

³ Univ. Rennes, Inria, CNRS, IRISA, France
alexandre.termier@irisa.fr

Abstract. Skypatterns are important since they enable to take into account user preference through Pareto-dominance. Given a set of measures, a skypattern query finds the patterns that are not dominated by others. In practice, different users may be interested in different measures, and issue queries on any subset of measures (a.k.a subspace). This issue was recently addressed by introducing the concept of skypattern cubes. However, such a structure presents high redundancy and is not well adapted for updating operations like adding or removing measures, due to the high costs of subspace computations in retrieving skypatterns. In this paper, we propose a new structure called Compressed Skypattern Cube (abbreviated CSKYC), which concisely represents a skypattern cube, and gives an efficient algorithm to compute it. We thoroughly explore its properties and provide an efficient query processing algorithm. Experimental results show that our proposal allows to construct and to query a CSKYC very efficiently.

Keywords: skypatterns; pareto-dominance relation; skypattern cubes

1 Introduction

The notion of skyline queries [2] has been quite recently integrated into the pattern discovery paradigm to mine skyline patterns (henceforth called *skypatterns*) [11, 15]. Given a set of measures, skypatterns are based on a Pareto-dominance relation, which means that no measure can be improved without degrading the others. As an example, a user may prefer patterns with a high frequency, large size and a high confidence. Then a pattern x_i dominates another pattern x_j if $\text{freq}(x_j) \geq \text{freq}(x_i)$, $\text{size}(x_j) \geq \text{size}(x_i)$, $\text{conf}(x_j) \geq \text{conf}(x_i)$ where at least one strict inequality holds. The skypattern set contains the patterns that are not dominated by any other pattern. Skypatterns are highly interesting since they do not require thresholds for the measures and the dominance relation gives them global interestingness with a semantics easily understood by the user.

In practice, users do not know the exact role of each measure and it is difficult to select beforehand the most appropriate subset of measures. Users would like to keep all potentially useful measures, look at what happens on skypattern sets when removing

or adding a measure, thus evaluating the impact of measures, and then converge to a convenient skypattern set.

This issue has been first addressed with the notion of a Skypattern Cube [13], which is the lattice of all possible subsets of measures associated with their skypattern sets. More formally, given a set M of n measures, the $2^n - 1$ possible non-empty skypattern subsets should be precomputed to efficiently handle various queries of users. By comparing two neighboring nodes (differentiated by only one measure), users can observe new skypatterns and the ones which disappear, greatly helping to better understand the role of the measures. To sum up, the cube is a structure that enables to discover the most interesting skypattern sets. The skypattern cube has been exploited in various domains such as bioinformatics [10] and mutagenicity [13]. However there are $2^n - 1$ possible non-empty skypattern sets with high redundancy coming from derivations of skypatterns among subspaces of the cube [13].

In this paper, we propose a new structure called the *Compressed Skypattern Cube* (denoted CSKYC). Each subspace stores skypatterns (called *proper skypatterns*) that do not appear in its descendant ones and the compressed skypattern cube contains only non-empty subspaces. Compared to the original skypattern cube [13], the CSKYC has fewer duplicates among subspaces, and does not need to store all of them. Moreover the cube includes unbalanced skypatterns. For instance, let $M = \{\text{freq}, \text{size}\}$ and three patterns x_i , x_j and x_k such that $\text{freq}(x_i) = 10$, $\text{size}(x_i) = 1$, $\text{freq}(x_j) = 2$, $\text{size}(x_j) = 8$, $\text{freq}(x_k) = 4$ and $\text{size}(x_k) = 5$. Clearly, x_i (resp. x_j) is a skypattern for $\{\text{freq}\}$ (resp. $\{\text{size}\}$), thus x_i (resp. x_j) will be instantly a skypattern for M being derived from $\{\text{freq}\}$ (resp. $\{\text{size}\}$). However, x_k is also a skypattern for M , being more balanced over measures than x_i (resp. x_j) which only has an extreme value for $\{\text{freq}\}$ (resp. $\{\text{size}\}$). Proper skypatterns are often well-balanced skypatterns.

Contributions overview. We thoroughly explore interesting properties of the compressed skypattern cube and provide an efficient query processing algorithm. Our contributions can be summarized as follows: (i) we provide the summarization structure CSKYC which concisely represents the whole skypattern cube and preserves its essential information. (ii) We propose a bottom-up approach to construct the CSKYC. (iii) We show how this structure can be used efficiently for query processing. Finally, (iv) we present an extensive set of experiments showing the advantages of our proposals.

Paper organization. The rest of this paper is organized as follows. Section 2 recalls the definitions of the notions used in this paper. Section 3 first introduces the CSKYC, provides algorithms to build it, and shows how the CSKYC can handle various skypattern queries. Section 4 is devoted to related works. Finally, Section 5 shows our experimental results and Section 6 concludes.

2 Preliminaries

Let \mathcal{I} be a set of distinct literals called *items*. A pattern (or itemset) is a non-empty subset of \mathcal{I} . The language of patterns corresponds to $\mathcal{L}_{\mathcal{I}} = 2^{\mathcal{I}} \setminus \emptyset$. A transactional dataset \mathcal{T} is a multiset of patterns in $\mathcal{L}_{\mathcal{I}}$. The traditional example is a supermarket database in which, for each transaction t_i , every item in a transaction is a product bought by the customer i . Tab. 1 summarizes the different notations used throughout the paper.

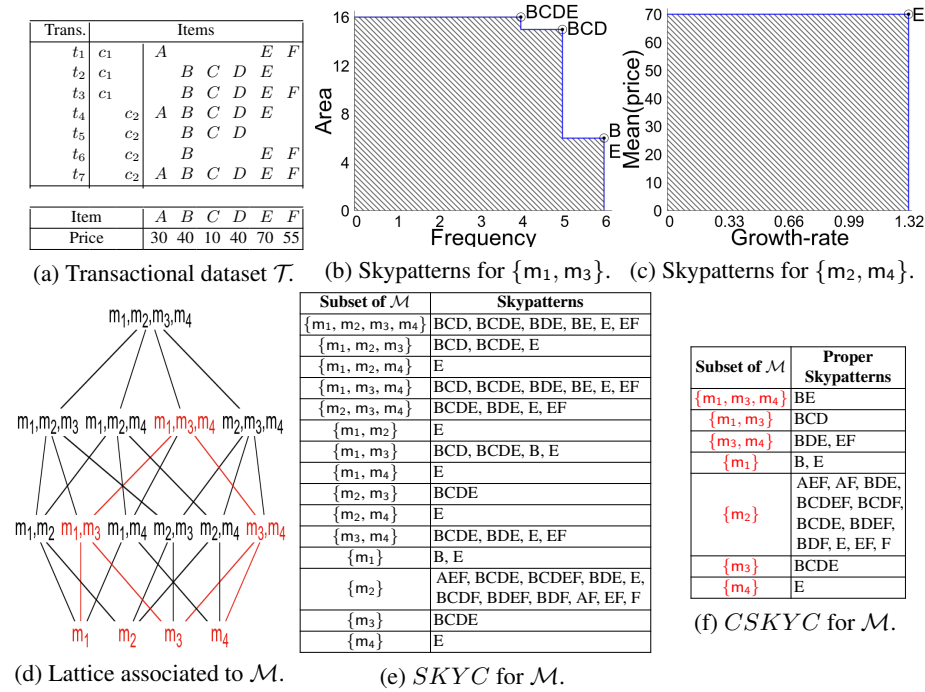


Fig. 1: Running example: $\mathcal{M} = \{ m_1:\text{freq}(x), m_2:\text{gr}_1(x), m_3:\text{area}(x), m_4:\text{mean}(x.\text{price}) \}$.

Example 1. Fig. 1a depicts a transactional dataset \mathcal{T} where items in a transaction t_i are denoted A, \dots, F . It serves as example throughout the paper. An attribute (*price*) is associated to each item. For instance, the *Price* of *A* is \$30. The dataset is partitioned into two classes, class c_1 for clients with loyalty programs and class c_2 for other clients.

Constraint-based pattern mining aims at extracting all patterns $x \in \mathcal{L}_{\mathcal{I}}$ satisfying a query $q(x)$ which is usually called *theory* [7]: $Th(\mathcal{L}_{\mathcal{I}}, q) = \{ x \in \mathcal{L}_{\mathcal{I}} \mid q(x) \text{ is true} \}$. A common example is the minimal frequency constraint ($\text{freq}(x) \geq \theta$) which provides patterns having a number of occurrences exceeding a given minimal threshold θ . Many other measures for patterns can be considered such as:

- $\text{size}(x) = |x|$ is the number of items that x contains.
- $\text{gr}_1(x) = \frac{(|\mathcal{T}| - |\mathcal{T}_1|) \times \text{freq}_1(x)}{|\mathcal{T}_1| \times (\text{freq}(x) - \text{freq}_1(x))}$ where \mathcal{T}_1 is a sub-dataset (i.e a class partition) on \mathcal{T} .
- $\min(x.\text{att}) = \min_{i \in x} \{i.\text{att}\}$ (resp. $\max(x.\text{att}) = \max_{i \in x} \{i.\text{att}\}$) is the lowest (resp. highest) among item values of x for attribute *att*.
- $\text{mean}(x.\text{att}) = (\min(x.\text{att}) + \max(x.\text{att}))/2$.

Example 2. $\text{freq}(BC) = 5$, $\text{mean}(BCD.\text{price}) = 25, \dots$

Skypatterns allow to express a user-preference according to a dominance relation [11].

Table 1: Notations.

Symbol	Definition
\mathcal{T}	Transactional dataset.
\mathcal{I}	Set of items.
$\mathcal{L}_{\mathcal{I}} = 2^{\mathcal{I}} \setminus \emptyset$	Language of patterns.
\mathcal{M}	Set of measures.
$U \subseteq \mathcal{M}$	Subspace U (i.e. subset of measures of \mathcal{M}).
$Skyc(\mathcal{L}_{\mathcal{I}}, U)$	Skypatterns set over $\mathcal{L}_{\mathcal{I}}$ for U.
$P\text{-}Skyc(\mathcal{L}_{\mathcal{I}}, U)$	Proper Skypattern set over $\mathcal{L}_{\mathcal{I}}$ for U.
$\ell\text{-}Skyc(\mathcal{L}_{\mathcal{I}}, U)$	Large Skypattern set over $\mathcal{L}_{\mathcal{I}}$ for U.
$Desc(\mathcal{L}_{\mathcal{I}}, U)$	Union of all proper skypattern sets for all descendant subspaces $V \subset U$.
$SKYC(\mathcal{L}_{\mathcal{I}}, \mathcal{M})$	Skypattern Cube over $\mathcal{L}_{\mathcal{I}}$ for \mathcal{M} .
$CSKYC(\mathcal{L}_{\mathcal{I}}, \mathcal{M})$	Compressed Skypattern Cube over $\mathcal{L}_{\mathcal{I}}$ for \mathcal{M} .
$P \subseteq \mathcal{L}_{\mathcal{I}}$	Set of patterns.
$\mathcal{O}(P)$	Set of data points associated to P.
$Skylines(\mathcal{O}(P), U)$	Set of skyline points on $\mathcal{O}(P)$ for U.

Definition 1 (Pareto-dominance). A pattern x dominates another pattern y w.r.t a measure subset (a.k.a subspace) U , noted by $x \succ_U y$, iff $\forall m_i \in U, m_i(x) \geq m_i(y)$ and $\exists m_j \in U, m_j(x) > m_j(y)$.

Example 3. For $U = \{m_1:\text{freq}(x), m_3:\text{area}(x)\}$, pattern BCD dominates pattern BC since $\text{freq}(BCD) = \text{freq}(BC) = 5$ and $\text{area}(BCD) > \text{area}(BC)$.

The Skypattern Operator [11] extracts the skypattern set w.r.t a subspace U .

Definition 2 (Skypattern Operator). A pattern is a skypattern w.r.t a subspace U iff it is not dominated by any other pattern w.r.t U . The Skypattern Operator returns all the skypatterns w.r.t U : $Skyc(\mathcal{L}_{\mathcal{I}}, U) = \{x \in \mathcal{L}_{\mathcal{I}} \mid \nexists y \in \mathcal{L}_{\mathcal{I}}, y \succ_U x\}$.

Example 4. Consider the dataset of Fig. 1a, users may ask for the skypatterns for every combination of the measures $\{m_1, m_2, m_3, m_4\}$. Figs. 1b and 1c depict skypatterns for $\{m_1, m_3\}$ and $\{m_2, m_4\}$ respectively.

As said above, users may query multiple skypattern sets for different subspaces. Furthermore, for a space \mathcal{M} there are $2^{|\mathcal{M}|-1}$ different skypattern sets. The Skypattern Cube [13] retrieves all skypattern sets for any subspace.

Definition 3 (Skypattern Cube). Given a set of measures \mathcal{M} , the Skypattern Cube of \mathcal{M} is defined as $SKYC(\mathcal{L}_{\mathcal{I}}, \mathcal{M}) = \{(U, Skyc(\mathcal{L}_{\mathcal{I}}, U)) \mid U \subseteq \mathcal{M}, U \neq \emptyset\}$

Example 5. Fig. 1d depicts the lattice associated to \mathcal{M} (power set of \mathcal{M} : $2^{\mathcal{M}} \setminus \emptyset$). Fig. 1e associates to each non-empty subset of \mathcal{M} its skypattern set.

For computing the skypattern cube, [13] has proposed a bottom-up approach using two derivation rules that provide an easy way to automatically infer a large proportion of the skypatterns of a parent node from the skypattern sets of its child nodes without any dominance test (if k measures are associated to a parent node, its child nodes are the nodes defined by the $\binom{k}{k-1}$ subsets of $k-1$ measures).

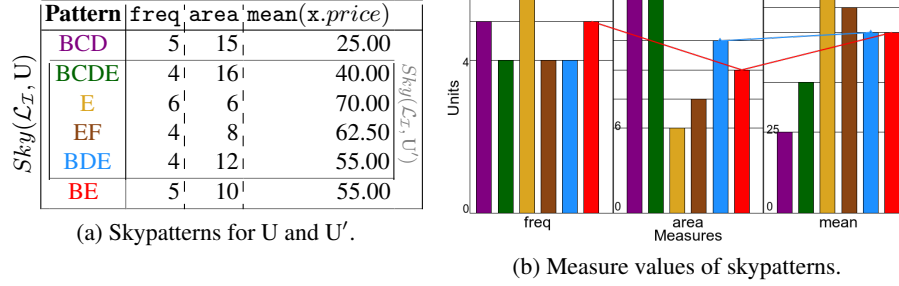


Fig. 2: Example of (Proper) Skypatterns for $U = \{m_1, m_3, m_4\}$ and $U' = \{m_3, m_4\}$.

3 Contributions

This section introduces the CSKYC (Compressed Skypattern Cubes) which concisely represents the entire skypattern cube. The main idea is, for every subspace, to only store its *proper skypatterns*. More precisely, a skypattern x for a subspace U is stored iff $x \in \text{Sky}(\mathcal{L}_{\mathcal{I}}, U)$ and there exists no $V \subset U$ s.t. $x \in \text{Sky}(\mathcal{L}_{\mathcal{I}}, V)$. First, we introduce the CSKYC of a set of measures. Then, we propose a bottom-up approach for building such a CSKYC. Finally, we show how to efficiently query the whole skypattern set for U from the CSKYC.

3.1 The Compressed Skypattern Cube

Definition 4 (Proper Skypattern). *The set of proper skypatterns for a subspace U is the subset of skypatterns on U which are not skypatterns in any subset of U :*

$$\mathbf{P}\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U) = \{ x \in \text{Sky}(\mathcal{L}_{\mathcal{I}}, U) \mid \nexists V \subset U, x \in \text{Sky}(\mathcal{L}_{\mathcal{I}}, V) \}$$

Example 6. Consider $U = \{m_1, m_3, m_4\}$ and $U' = \{m_3, m_4\}$. Tab. 2a shows the measure values for skypatterns for U and U' . Fig. 2b illustrates how proper skypattern BE (in red) (resp. BDE (in light-blue)) is more balanced than other skypatterns for U (resp. U'), having fewer extreme values than the other skypatterns.

Based on this notion, we define the compressed skypattern cube.

Definition 5 (Compressed Skypattern Cube). *Given a set of measures \mathcal{M} , the compressed skypattern cube of \mathcal{M} is defined as*

$$\text{CSKYC}(\mathcal{L}_{\mathcal{I}}, \mathcal{M}) = \{ (U, \mathbf{P}\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U)) \mid U \subseteq \mathcal{M}, U \neq \emptyset, \mathbf{P}\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U) \neq \emptyset \}$$

Example 7. For the dataset shown in Fig. 1a, the $\text{CSKYC}(\mathcal{L}_{\mathcal{I}}, \mathcal{M})$ is depicted in Fig. 1f and its sub-lattice (in red) in Fig. 1d. It contains only 6 non-empty subsets compared to 15 subsets in $\text{SKYC}(\mathcal{L}_{\mathcal{I}}, \mathcal{M})$. Clearly, the CSKYC is much more compact.

3.2 Computing the CSKYC

A first and naive way to get the CSKYC consists in first computing the skypattern cube, and then deriving the CSKYC by removing duplicates from their subspaces. Such an approach is inefficient as the number of subspaces to process is exponential. In this section, we provide a bottom-up algorithm (CSKYC-BUC) for building the CSKYC. Given a set of measures \mathcal{M} of size d , the subspaces are organized into d levels, such that the subspaces of size i are in level i . We only keep *non-empty* subspaces (i.e. those containing proper skypatterns). All descendant skypatterns of a subspace are collected to form a large skypattern set (ℓ -Sky) which are then used as filters, and if no new skypattern is found, the subspace is discarded from the CSKYC.

Let us first give some preliminary definitions in order to compute the CSKYC.

Definition 6 (Indistinct/Incomparable Skypatterns). Let x, y be two skypatterns w.r.t a subspace U : (i) x, y are indistinct, noted $x =_U y$, iff $\forall m_i \in U, m_i(x) = m_i(y)$; (ii) x, y are incomparable, noted $x \prec \succ_U y$, iff $x \not\prec_U y, y \not\prec_U x$ and $x \neq_U y$.

Incomparable skypatterns and indistinct ones for U constitute partitions of $Sky(\mathcal{L}_{\mathcal{I}}, U)$.

Definition 7 (Indistinct Subspace (IS)). A subspace U is an **Indistinct Subspace (IS)** iff all patterns in $Sky(\mathcal{L}_{\mathcal{I}}, U)$ are indistinct from each other.

Example 8. Let $U = \{m_1\}$, and $V = \{m_1, m_3\}$. B and E are indistinct w.r.t. U , while BCDE and BCD are incomparable w.r.t. V .

Lemma 1 states that skypatterns that are common to two different subspaces remain skypatterns in their union.

Lemma 1. $Sky(\mathcal{L}_{\mathcal{I}}, U) \cap Sky(\mathcal{L}_{\mathcal{I}}, V) \subseteq Sky(\mathcal{L}_{\mathcal{I}}, U \cup V)$

Proof (By contradiction). Assume that, for two subspaces U, V s.t. $W = U \cup V$, $\exists x \in \underbrace{Sky(\mathcal{L}_{\mathcal{I}}, U) \cap Sky(\mathcal{L}_{\mathcal{I}}, V)}_{(1)}$, but $x \notin \underbrace{Sky(\mathcal{L}_{\mathcal{I}}, W)}_{(2)}$. From (1): $x \in \underbrace{Sky(\mathcal{L}_{\mathcal{I}}, U)}_{(3)}$ and $x \in \underbrace{Sky(\mathcal{L}_{\mathcal{I}}, V)}_{(4)}$. From (2): $\exists y \in Sky(\mathcal{L}_{\mathcal{I}}, W), y \succ_W x \Rightarrow \underbrace{\forall m_i \in W, m_i(y) \geq m_i(x)}_{(5)}$.

From (3): $y \not\prec_U x \Rightarrow \forall m_i \in U, m_i(y) \leq m_i(x)$. From (5): $x =_U y$
 From (4): $y \not\prec_V x \Rightarrow \forall m_i \in V, m_i(y) \leq m_i(x)$. From (5): $x =_V y$ } $x =_W y$.

Thus, $x \in Sky(\mathcal{L}_{\mathcal{I}}, W)$ leading to a contradiction.

Based on lemma 1, the following theorem enables us to characterize empty subspaces in the CSKYC, i.e. those without proper skypatterns.

Theorem 1 (Empty subspaces in CSKYC). Given two subspaces U and V that are IS, if $Sky(\mathcal{L}_{\mathcal{I}}, U) \cap Sky(\mathcal{L}_{\mathcal{I}}, V) \neq \emptyset$, then $U \cup V$ is an IS and $P\text{-}Sky(\mathcal{L}_{\mathcal{I}}, U \cup V) = \emptyset$.

Proof (By contradiction). Let U, V two IS and $W = U \cup V$.

- Assume that W is not an IS and $\exists x, y \in Sky(\mathcal{T}, U) \cap Sky(\mathcal{T}, V)$. From Lemma 1: $x, y \in Sky(\mathcal{T}, W)$. Since W is not an IS, $x \prec \succ_W y$. As $x, y \in Sky(\mathcal{T}, U), x, y \in Sky(\mathcal{T}, V)$

and U and V are IS, thus, $x =_U y$ and $x =_V y$. Thus, $x =_W y$ leading to a contradiction.
- Assume that $\exists x \in \text{Sky}(\mathcal{T}, U) \cap \text{Sky}(\mathcal{T}, V)$ and $\exists y \in \mathbf{P}\text{-Sky}(\mathcal{T}, W)$. From Lemma 1: $x \in \text{Sky}(\mathcal{T}, W)$. Since W is an IS, $x =_W y$. Thus, $x =_U y$ and $x =_V y$. So, $y \in \text{Sky}(\mathcal{T}, U)$ and $y \in \text{Sky}(\mathcal{T}, V)$. Thus, $y \notin \mathbf{P}\text{-Sky}(\mathcal{T}, W)$ leading to a contradiction.

Example 9. In Fig. 1f, $\mathbf{P}\text{-Sky}(\mathcal{L}_{\mathcal{I}}, \{m_1, m_4\}) = \emptyset$ as $\text{Sky}(\mathcal{L}_{\mathcal{I}}, \{m_1\}) \cap \text{Sky}(\mathcal{L}_{\mathcal{I}}, \{m_4\}) = \{E\}$

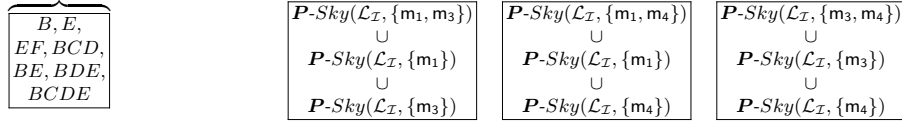
The authors in [13] showed that incomparable skypatterns and some indistinct skypatterns of a child subspace remain also skypatterns in its parent subspace (they are referred to as derivable skypatterns). They also showed that a parent subspace can include non-derivable skypatterns (i.e., those that are not skypatterns in any of its child subspaces). Thus, one can collect the non-empty sets of descendants (which are proper skypatterns) of a subspace to form a large skypattern set and to use them as filters to detect a priori that no proper skypattern exist (see Corollary 1).

Definition 8 (Large Skypattern Set). *The Large Skypattern Set for a subspace U is the union of proper skypattern set of U with all proper skypattern sets of its descendant subspaces $V \subset U$: $\ell\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U) = \mathbf{P}\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U) \cup \text{Desc}(\mathcal{L}_{\mathcal{I}}, U)$,*

where $\text{Desc}(\mathcal{L}_{\mathcal{I}}, U) = \bigcup_{V \subset U} \mathbf{P}\text{-Sky}(\mathcal{L}_{\mathcal{I}}, V) = \bigcup_{V \subset U \wedge |V|=|U|-1} \ell\text{-Sky}(\mathcal{L}_{\mathcal{I}}, V)$

Example 10. For $U = \{m_1:\text{freq}(x), m_3:\text{area}(x), m_4:\text{mean}(x.\text{price})\}$:

$\ell\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U) = \mathbf{P}\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U) \cup \ell\text{-Sky}(\mathcal{L}_{\mathcal{I}}, \{m_1, m_3\}) \cup \ell\text{-Sky}(\mathcal{L}_{\mathcal{I}}, \{m_1, m_4\}) \cup \ell\text{-Sky}(\mathcal{L}_{\mathcal{I}}, \{m_3, m_4\})$



Based on Definition 4, the proper skypatterns of any parent subspace can be computed thanks to the following corollary.

Corollary 1. $\mathbf{P}\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U) = \text{Sky}(\mathcal{L}_{\mathcal{I}}, U) \setminus \text{Desc}(\mathcal{L}_{\mathcal{I}}, U)$.

To compute $\mathbf{P}\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U)$, we first retrieve its descendants (which are proper skypatterns), then we seek for skypatterns that are not in $\text{Desc}(\mathcal{L}_{\mathcal{I}}, U)$. Algorithm 1 gives the pseudo-code of our bottom-up approach. It starts by computing $\mathbf{P}\text{-Sky}(\mathcal{L}_{\mathcal{I}}, m_i)$ for every $m_i \in \mathcal{M}$ (level 1) and then follows a level-wise strategy: from the lower level, each level of the lattice is constructed by applying Theorem 1 and, if needed, computing non-derivable skypatterns (cf. line 21). Two data structures, IS and $\ell\text{-Sky}$ are also maintained during the construction process, storing for each subspace its large pattern set and its status. They allow an incremental computation of $\ell\text{-Sky}$.

3.3 Querying $\text{Sky}(\mathcal{L}_{\mathcal{I}}, U)$ from $\text{CSKYC}(\mathcal{L}_{\mathcal{I}}, \mathcal{M})$

When a skypattern set for a given subspace U is queried, the CSKYC may not have a record for U ; even if it does, the skypattern set that is stored for U is not complete. We propose a straightforward approach to query the complete skypattern set for U from $\text{CSKYC}(\mathcal{L}_{\mathcal{I}}, \mathcal{M})$.

Algorithm 1: CSKYC-BUC: Bottom-up approach for computing CSKYC.

Input: \mathcal{T} : a dataset, \mathcal{M} : a set of measures.
Output: The Compressed Skypattern Cube w.r.t \mathcal{M} .

```

1 CSKYC  $\leftarrow \emptyset$ ; IS[ ]  $\leftarrow \emptyset$ ;  $\ell$ -Sky[ ]  $\leftarrow \emptyset$ ;           // Initialization step
2 foreach  $m_i \in \mathcal{M}$  do
3    $P$ -Sky  $\leftarrow \text{CP+SKY}(\mathcal{L}_{\mathcal{I}}, \{m_i\})$ ;           // Compute  $P$ -Sky( $\mathcal{L}_{\mathcal{I}}, m_i$ )
4    $\ell$ -Sky[ $\{m_i\}$ ]  $\leftarrow P$ -Sky;
5   CSKYC  $\leftarrow \text{CSKYC} \cup \{(\{m_i\}, P\text{-Sky})\}$ ;
6   IS[ $\{m_i\}$ ]  $\leftarrow \text{true}$ ;
7 for  $i \leftarrow 2$  to  $|\mathcal{M}|$  do
8   foreach  $U \subseteq \mathcal{M}$  s.t.  $|U| = i$  do
9      $P$ -Sky  $\leftarrow \text{ComputeProperSky}(U)$ ;           // Compute  $P$ -Sky( $\mathcal{L}_{\mathcal{I}}, U$ )
10    if  $P$ -Sky  $\neq \emptyset$  then
11      CSKYC  $\leftarrow \text{CSKYC} \cup \{(U, P\text{-Sky})\}$ ;
12 return CSKYC;
13 Function ComputeProperSky( $U$ ):
14   children  $\leftarrow \{V \subset U \mid |V| = |U| - 1\}$ ;           // Children of U
15   childrenIS  $\leftarrow \{W \in \text{children} \mid \text{IS}[W] = \text{true}\}$ ;
16   IS[U]  $\leftarrow \text{false}$ ;
17   if  $\exists V, W \in \text{childrenIS}$  s.t. Sky[V]  $\cap$  Sky[W]  $\neq \emptyset$  then // Apply theorem 1
18      $P$ -Sky  $\leftarrow \emptyset$ ;
19     IS[U]  $\leftarrow \text{true}$ ;
20   Desc  $\leftarrow \bigcup_{V \in \text{children}} \ell$ -Sky[V]; // Generate the filter skypatterns
21   if  $\neg \text{IS}[U]$  then
22      $P$ -Sky  $\leftarrow \text{CP+SKY}(\mathcal{L}_{\mathcal{I}} \setminus \text{Desc}, U)$ ;           // Apply corollary 1
23    $\ell$ -Sky[U]  $\leftarrow P$ -Sky  $\cup$  Desc;           // Update  $\ell$ -Sky for U
24   return  $P$ -Sky

```

Our approach is based on the fact that $\text{Sky}(\mathcal{L}_{\mathcal{I}}, U) \subseteq \ell\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U)$ and proceeds in two steps (see Algorithm 2): first, approximating $\text{Sky}(\mathcal{L}_{\mathcal{I}}, U)$ by $\ell\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U)$, and then, performing domination tests to filter dominated patterns.

(i) *Approximating $\text{Sky}(\mathcal{L}_{\mathcal{I}}, U)$.*

Based on Definition 8, we have that: $\forall U \subseteq \mathcal{M}$, $\text{Sky}(\mathcal{L}_{\mathcal{I}}, U) \subseteq \ell\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U)$. The proof is straightforward: $\forall x \in \text{Sky}(\mathcal{L}_{\mathcal{I}}, U)$, either $x \in P\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U)$, or $\exists V \subset U$ s.t. $x \in P\text{-Sky}(\mathcal{L}_{\mathcal{I}}, V)$.

(ii) *Filtering dominated skypatterns.*

To remove dominated skypatterns, we convert the problem into skyline mining operation in $|U|$ dimensions to process it more efficiently. Let f be a mapping function from a set of patterns $P \subseteq \mathcal{L}_{\mathcal{I}}$ to \mathbb{R}^n that associates, to each pattern $x_i \in P$, a data point $f(x_i) \in \mathbb{R}^n$ with coordinates $(m_1(x_i) = v_{i,1}, \dots, m_n(x_i) = v_{i,n})$. Let us note by $\mathcal{O}(P) = \{f(x) \mid x \in P\}$ the set of data points associated to P (see Tab. 2) and $\text{Skyline}(\mathcal{O}(P), U)$ be the set of skyline points of $\mathcal{O}(P)$ w.r.t. U . Thus, $\forall U \subseteq \mathcal{M}$,

Table 2: The multidimensional view for a set of patterns $P \subseteq \mathcal{L}_{\mathcal{I}}$ w.r.t. a subspace U ($|U| = n$).

Pattern	m_1	m_2	\dots	m_n
\mathbf{x}_1	$v_{1,1}$	$v_{1,2}$	\dots	$v_{1,n}$
\mathbf{x}_2	$v_{2,1}$	$v_{2,2}$	\dots	$v_{2,n}$
\vdots	\vdots	\vdots	\vdots	\vdots
\mathbf{x}_{p-1}	$v_{p-1,1}$	$v_{p-1,2}$	\dots	$v_{p-1,n}$
\mathbf{x}_p	$v_{p,1}$	$v_{p,2}$	\dots	$v_{p,n}$

Algorithm 2: Querying $Sky(\mathcal{L}_{\mathcal{I}}, U)$ from $CSKYC(\mathcal{L}_{\mathcal{I}}, \mathcal{M})$

Input: U : a subspace and $CSKYC$: the compressed skypattern cube w.r.t. \mathcal{M} .
Output: $Sky(\mathcal{L}_{\mathcal{I}}, U)$

- 1 $\ell\text{-Sky} \leftarrow \bigcup_{V \subseteq U} P\text{-Sky}(\mathcal{L}_{\mathcal{I}}, V)$; // Approximating $Sky(\mathcal{L}_{\mathcal{I}}, U)$
- 2 $Sky \leftarrow \text{BNL}(\mathcal{O}(\ell\text{-Sky}), U)$; // Filtering dominated skypatterns
- 3 **return** $Sky, \ell\text{-Sky}$

$Sky(P, U) = Skyline(\mathcal{O}(P), U)$. So, applying the skyline operator on $\mathcal{O}(P)$ provides the skypattern set.

Theorem 2. $Sky(\mathcal{L}_{\mathcal{I}}, U) = Skyline(\mathcal{O}(\ell\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U)), U)$.

Proof. Given a subspace U , we prove the two implications:

(\Rightarrow) Assume that $\exists x \in Sky(\mathcal{L}_{\mathcal{I}}, U)$: $\forall y \in \mathcal{L}_{\mathcal{I}}, y \not\prec_U x$. So, $\forall y \in \ell\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U), y \not\prec_U x$ therefore $x \in Sky(\ell\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U), U)$. Thus, $f(x) \in Skyline(\mathcal{O}(\ell\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U)), U)$.

(\Leftarrow) Assume that $\exists f(x) \in Skyline(\mathcal{O}(\ell\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U)), U)$. So, $\forall y_1 \in \ell\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U), y_1 \not\prec_U x$. From Def. 8: $\forall y_2 \in \mathcal{L}_{\mathcal{I}} \setminus \ell\text{-Sky}(\mathcal{L}_{\mathcal{I}}, U), \exists y_3 \in Sky(\mathcal{L}_{\mathcal{I}}, U), y_3 \succ_U y_2$. Thus, $y_2 \not\prec_U x$. Therefore, $x \in Sky(\mathcal{L}_{\mathcal{I}}, U)$.

The second step is performed using a skyline algorithm based on the BNL approach [2].

4 Related Work

Skylines vs Skypatterns. The notion of skyline queries [2, 3] has been recently integrated into pattern discovery to mine skypatterns [11]. Even if these notions seem similar, they correspond to very different extraction tasks. Skyline queries focus on the extraction of dominant tuples of a (point) database (T). The points (objects) are known in advance and then dominance test are applied. The skypattern mining task requires to mine patterns from a dataset (T) that must be Pareto-dominant for a given set of measures. Therefore, the latter problem is much harder since the search space for skypatterns is much larger than the search space for skylines: $O(\mathcal{L}_{\mathcal{I}} = 2^{|\mathcal{I}|})$ instead of $O(|T|)$. Two methods have been designed for mining skypatterns: AETHERIS [11] is a two-step method that benefits from theoretical relationships between condensed representations

and skypatterns, while $CP+SKY$ [15] mines skypatterns using dynamic CSPs. Finally, [12] provides a point-to-point comparison between these two approaches.

Skyline Cube vs Skypattern Cube. To offer the best possible response time for a subspace skyline query, skyline cubes (a.k.a SkyCube) were introduced independently by [17, 9]. They proposed several strategies to share skyline computation in different subspaces. Pei et al. proposed in [8] *Stellar*, which computes seed skylines groups in the full space, then extend it to build the final set of skyline groups and thus avoid the computation of skylines in all the subspaces. Similarly to the notions of skyline/skypattern, the skyline cube differs from the skypattern cube. A SkyCube tackles a point database looking for skyline point sets for a given set of dimensions. The skypattern cube computation has to deal with all the skypattern sets for a given set of measures. As seen in the previous paragraph, even if these notions are close, computing the skypattern cubes is much harder due to the huge search space. Two methods have been proposed to compute the skypattern cube. The first method, called as $CP+SKY+CUBE$ [13], is based on a bottom-up approach and derivation rules exploiting the relation between the nodes in the lattice. The second method [14] proposes an approximation of the skypattern cube and then applies skyline cube mining in $|\mathcal{M}|$ dimensions on that approximation.

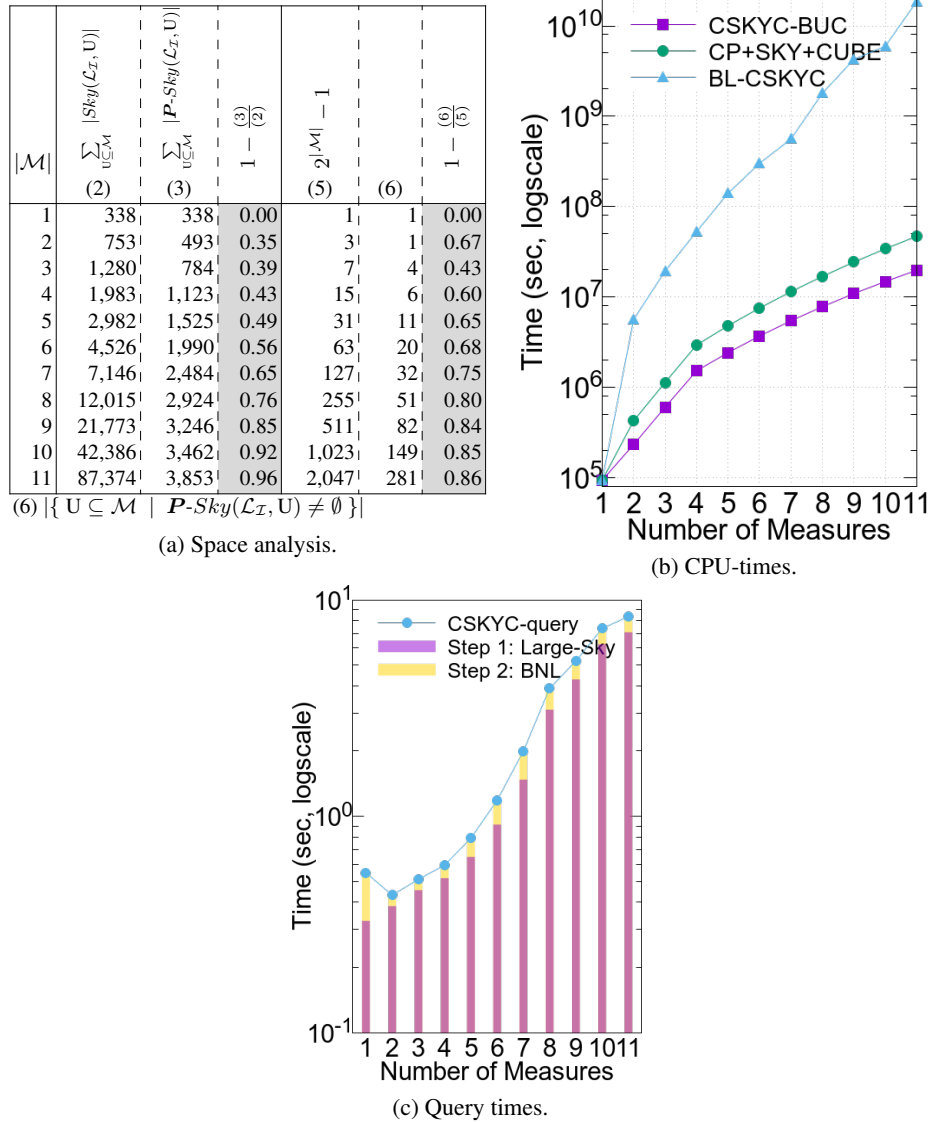
The Compressed Skyline Cube. Probably the closest previous work to our proposal is the so called *compressed skycube* (CSC) [16]. Its compression technique consists in storing for every subspace its partial skyline. It also supports concurrent subspace skyline queries in frequent updated databases. Our CSKYC can be seen as a reshaping of the CSC. However, the compressed skypattern cube computation is much harder due to the huge search space. Indeed, as shown previously, we need to extract patterns from a transactional dataset (\mathcal{T}) in order to determinate *proper* skypatterns for a given subspace. Other skycube summarization techniques have also been introduced. For instance, [1] proposed Hashcube, a structure based on bit-strings for storing the whole skycube. The work described in [6] proposed the negative skycube that returns subspaces where objects are not skylines.

5 Experimental Evaluation

This section evaluates constructing and querying the CSKYC on a real-life dataset and benchmarks. We compare the performances of the CSKYC with those of the original SKYC in terms of running-time and space storage, followed by query performance using CSKYC. The implementation of the different algorithms were carried out in C++. All experiments were conducted on a PC running Linux with a core i3 processor at 2.13 GHz.

5.1 Compressed Skypattern Cubes for Mutagenicity Dataset

We performed experiments on a real-life dataset of large size extracted from mutagenicity data [5] (a major problem in risk assessment of chemicals). This dataset has $|\mathcal{T}|=6,512$ transactions encoding chemicals and $|\mathcal{I}|=1,073$ items encoding frequent closed subgraphs previously extracted from \mathcal{T} with a 2% relative frequency threshold. Chemists use up to $|\mathcal{M}|=11$ measures, five of them are typically used in contrast mining


 Fig. 3: Results on Mutagenicity Dataset with $|\mathcal{M}| = 11$.

(e.g. growth rate) and allow to express different kinds of background knowledge. The other six measures are related to topological and chemical properties of the chemicals.

Space analysis. Fig. 3a shows the storage comparison of CSKYC to skypattern cubes of different dimensionality. Column 1 corresponds to the number of measures. Columns 2 and 3 report the total number of skypatterns for SKYC and CSKYC respectively. Column 4 gives their ratio. Columns 5 and 6 report the total number of subspaces for SKYC

and CSKYC respectively. Column 7 gives their ratio. For each $|\mathcal{M}| = k$, reported values in columns (2), (3), (5) and (6) represent the averages over all $\binom{11}{k}$ possible skypattern cubes. Overall, CSKYC achieves the best compression of the skypattern sets. The effect of duplicate elimination is greatly amplified for $|\mathcal{M}| \geq 6$. CSKYC achieves up to $20.6\times$ compression (in number of skypatterns) and permits using $4 - 7\times$ fewer subspaces. For $|\mathcal{M}| = 11$, the total number of proper skypatterns is 3,853, while for SKYC the total number of skypatterns is 87,374. This lead to a substantial gain greater than 95%.

CPU-time analysis. We compare our approach (CSKYC-BUC) with two methods: (i) a base-line method (BL-CSKYC) for computing CSKYC, and (ii) CP+SKY+CUBE proposed in [13] for computing SKYC. BL-CSKYC follows a bottom-up strategy: from the lower level, for each level and each subspace of the lattice, we compute its skypatterns, collect the skypatterns of its descendant subspaces, and then we remove all the duplicates. Fig. 3b shows the performance of the three methods according to the number of measures $|\mathcal{M}|$. The scale is logarithmic. For CSKYC-BUC (resp. CP+SKY+CUBE) and for $|\mathcal{M}|=k$, the reported CPU-time is the average of CPU-times over all $\binom{11}{k}$ possible CSKYC (resp. SKYC). As we can see, CSKYC-BUC clearly outperforms BL-CSKYC by several orders of magnitude. This is particularly obvious for higher values of $|\mathcal{M}|$ due to the reduced number of skypatterns involved in the construction. For $(2 \leq |\mathcal{M}| \leq 5)$, the average speed-up is 37.3. For $|\mathcal{M}| = 8$, there is an order of magnitude (speed-up value 213.15). For $|\mathcal{M}| = 11$, the speed-up value reaches 949. Finally, CSKYC-BUC is an average 2x faster than CP+SKY+CUBE for building the CSKYC.

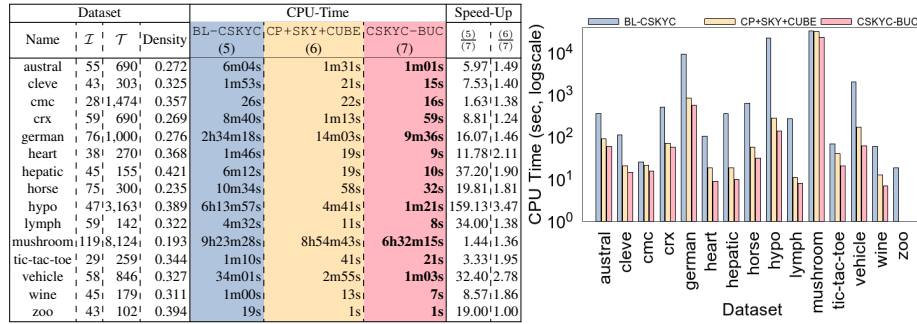
Querying CSKYC. Evaluating the query performance of Algorithm 2, for $|\mathcal{M}|=k$, is performed by dividing the total time to sequentially query every subspace by $2^k - 1$. Each query is extracted from the CSKYC. The reported CPU-time in Fig. 3c are the averages of CPU-times over all $\binom{11}{k}$ possible CSKYCs. By comparing the CPU-times for the two steps of Algorithm 2, overall, the BNL step is negligible as compared to the first step. The scale is logarithmic. Experimental results show that query processing of the CSKYC is fast (less than 10 seconds for $|\mathcal{M}|=11$).

5.2 Compressed Skypattern Cubes for UCI Datasets

Experiments were carried out on 15 datasets from UCI benchmarks [4]. We considered 5 measures $M = \{\text{freq}, \text{max}, \text{area}, \text{mean}, \text{gr}_1\}$. In order to use measures using numeric values, like mean, we generated random values associated to attributes, each value being within the range $[0..1]$. Fig. 4 summarizes the results we obtained.

CPU-time analysis. Fig. 4a compares the performance of the three methods (with a graphical view). Cols. 1-4 give the characteristics of each dataset (name, number of items (\mathcal{I}), number of transactions (\mathcal{T}) and density). CSKYC-BUC clearly dominates the base-line method. On half of the datasets, there is an order of magnitude (speed-up value at least 11.78) (Col. 8). CSKYC-BUC is an average 2 times faster than CP+SKY+CUBE.

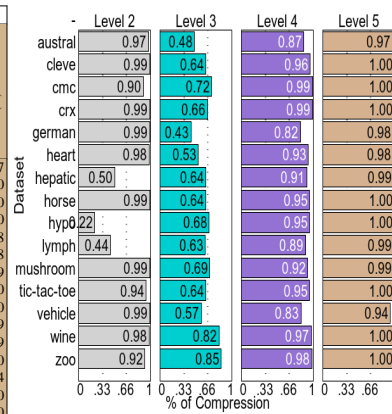
Space analysis. Fig. 4b compares, for each dataset, the number of proper skypatterns vs. the total number of skypatterns at each level of the cube. For each level i ($2 \leq i \leq 5$), the corresponding summarization ratio is also depicted. Fig. 4c shows the graphical



(a) Time Analysis.

Dataset	$ \mathcal{M} $	Space Analysis				
		1	2	3	4	5
		$\sum_{ U =1} S_{\mathcal{H}}(C_x, U)$	$\sum_{ U =2} S_{\mathcal{H}}(C_x, U)$	$\sum_{ U =3} S_{\mathcal{H}}(C_x, U)$	$\sum_{ U =4} S_{\mathcal{H}}(C_x, U)$	$\sum_{ U =5} S_{\mathcal{H}}(C_x, U)$
austral	101,630	1,314 ¹	45 ⁰ 0.97	403 ¹ 211 ⁰ 0.48	648 ¹ 83 ⁰ 0.87	317 ¹ 9 ⁰ 0.97
cleveland	43,504	6,601 ¹	34 ⁰ 0.99	238 ¹ 85 ⁰ 0.64	300 ¹ 12 ⁰ 0.96	111 ¹ 0 ⁰ 1.00
cmc	12,208	559 ¹	57 ⁰ 0.90	269 ¹ 76 ⁰ 0.72	299 ¹ 1 ⁰ 0.99	104 ¹ 0 ⁰ 1.00
crx	107,742	5,216 ¹	58 ⁰ 0.99	307 ¹ 105 ⁰ 0.66	437 ¹ 1 ⁰ 0.99	215 ¹ 0 ⁰ 1.00
german	1,173,279	35,083 ¹	60 ⁰ 0.99	620 ¹ 352 ⁰ 0.43	1,198 ¹ 210 ⁰ 0.82	610 ¹ 13 ⁰ 0.98
heart	34,400	1,218 ¹	29 ⁰ 0.98	283 ¹ 133 ⁰ 0.53	468 ¹ 35 ⁰ 0.93	211 ¹ 5 ⁰ 0.98
hepatic	9,598	189 ¹	95 ⁰ 0.50	537 ¹ 194 ⁰ 0.64	762 ¹ 71 ⁰ 0.91	343 ¹ 4 ⁰ 0.99
horse	129,716	10,352 ¹	22 ⁰ 0.99	177 ¹ 64 ⁰ 0.64	220 ¹ 10 ⁰ 0.95	84 ¹ 0 ⁰ 1.00
hypo	671,961	789 ¹	613 ⁰ 0.22	2,134 ¹ 686 ⁰ 0.68	2,812 ¹ 129 ⁰ 0.95	1,172 ¹ 0 ⁰ 1.00
lymph	45,644	98 ¹	55 ⁰ 0.44	428 ¹ 157 ⁰ 0.63	524 ¹ 59 ⁰ 0.89	268 ¹ 4 ⁰ 0.99
mushroom	650,965	96,164 ¹	45 ⁰ 0.99	216 ¹ 66 ⁰ 0.69	307 ¹ 24 ⁰ 0.92	131 ¹ 1 ⁰ 0.99
tic-tac-toe	16,157	1,301 ¹	77 ⁰ 0.94	259 ¹ 93 ⁰ 0.64	310 ¹ 16 ⁰ 0.95	109 ¹ 0 ⁰ 1.00
vehicle	141,974	15,901 ¹	69 ⁰ 0.99	480 ¹ 208 ⁰ 0.57	827 ¹ 144 ⁰ 0.83	421 ¹ 24 ⁰ 0.94
wine	15,112	1,346 ¹	25 ⁰ 0.98	152 ¹ 27 ⁰ 0.82	170 ¹ 5 ⁰ 0.97	57 ¹ 0 ⁰ 1.00
zoo	4,871	599 ¹	50 ⁰ 0.92	197 ¹ 29 ⁰ 0.85	211 ¹ 5 ⁰ 0.98	74 ¹ 0 ⁰ 1.00

(b) Space Analysis.



(c) Summarization Ratios.

 Fig. 4: Results on UCI datasets with $|\mathcal{M}| = 5$.

view of these ratios. For level 2, on most of the datasets, CSKYC achieves a very high summarization ratios (up to 99%). For level 3, these ratios mostly decrease since both levels (2 and 3) share most subspaces ($\binom{5}{2} + \binom{5}{3} = 20$ in total against 31 for the cube). Finally, for level 4 (resp. 5), CSKYC uses less storage than SKYC by at least 82% (resp. 94%) in size for all datasets we considered. Within these levels, there are almost no *proper* skypatterns since they have few subspaces ($\binom{5}{4} = 5$ for level 4 and $\binom{5}{5} = 1$) for level 5 and each one of these subspaces has a lot of descendant subspaces (14 for level 4 and 30 for level 5).

6 Conclusion

We have presented the compressed skypattern cube which concisely represents the skypattern cube and preserves its essential information. Compared to the original skypattern cube, the compressed skypattern cube has much less duplicates among subspaces. We have provided an efficient algorithm to compute it and to query the skypattern set for any subspace. Our experimental study shows that CSKYC is particularly efficient

in terms of build time and space usage compared to the original skypattern cube. Another interesting property is its ability to efficiently provide the skypattern set of any subspace. As future work, we plan to investigate the incremental maintenance of the CSKYC by allowing to add and/or remove any measure.

References

1. Bøgh, K.S., Chester, S., Sidlauskas, D., Assent, I.: Hashcube: A data structure for space- and query-efficient skycube compression. In: CIKM. pp. 1767–1770 (2014)
2. Börzsönyi, S., Kossmann, D., Stocker, K.: The skyline operator. In: ICDE. pp. 421–430 (2001)
3. Chomicki, J., Godfrey, P., Gryz, J., Liang, D.: Skyline with presorting. In: ICDE. pp. 717–719 (2003)
4. Dua, D., Graff, C.: UCI machine learning repository (2017), <http://archive.ics.uci.edu/ml>
5. Hansen, K., Mika, S., Schroeter, T., Sutter, A., ter Laak, A., Steger-Hartmann, T., Heinrich, N., Müller, K.: Benchmark data set for in silico prediction of ames mutagenicity. *JCIM* 49(9), 2077–2081 (2009)
6. Hanusse, N., Wanko, P.K., Maabout, S.: Computing and summarizing the negative skycube. In: CIKM. pp. 1733–1742 (2016)
7. Mannila, H., Toivonen, H.: Levelwise search and borders of theories in knowledge discovery. *Data Min. Knowl. Discov.* 1(3), 241–258 (1997)
8. Pei, J., Fu, A.W., Lin, X., Wang, H.: Computing compressed multidimensional skyline cubes efficiently. In: ICDE. pp. 96–105 (2007)
9. Pei, J., Jin, W., Ester, M., Tao, Y.: Catching the best views of skyline: A semantic approach based on decisive subspaces. In: VLDB. pp. 253–264 (2005)
10. Pham, H., Lavenier, D., Termier, A.: Identifying genetic variant combinations using skypatterns. In: DEXA Workshops. pp. 44–48. IEEE Computer Society (2016)
11. Soulet, A., Raïssi, C., Plantevit, M., Crémilleux, B.: Mining dominant patterns in the sky. In: ICDM. pp. 655–664 (2011)
12. Ugarte, W., Boizumault, P., Crémilleux, B., Lepailleur, A., Loudni, S., Plantevit, M., Raïssi, C., Soulet, A.: Skypattern mining: From pattern condensed representations to dynamic constraint satisfaction problems. *Artif. Intell.* 244, 48–69 (2017)
13. Ugarte, W., Boizumault, P., Loudni, S., Crémilleux, B.: Computing skypattern cubes. In: ECAI. pp. 903–908 (2014)
14. Ugarte, W., Boizumault, P., Loudni, S., Crémilleux, B.: Computing skypattern cubes using relaxation. In: ICTAI. pp. 859–866 (2014)
15. Ugarte, W., Boizumault, P., Loudni, S., Crémilleux, B., Lepailleur, A.: Mining (soft-) skypatterns using dynamic CSP. In: CPAIOR. pp. 71–87 (2014)
16. Xia, T., Zhang, D.: Refreshing the sky: the compressed skycube with efficient support for frequent updates. In: SIGMOD Conference. pp. 491–502 (2006)
17. Yuan, Y., Lin, X., Liu, Q., Wang, W., Yu, J.X., Zhang, Q.: Efficient computation of the skyline cube. In: VLDB. pp. 241–252 (2005)