

Summarizing high-dimensional feature vectors



Ting Xie

Department of Computer Science and Engineering
University at Buffalo

A dissertation submitted for the degree of

Doctor of Philosophy

Sep 2019

Acknowledgements

todo

Contents

| | | |
|----------|--|----------|
| 1 | Introduction | 1 |
| 2 | Problem Overview | 2 |
| 3 | Application: Query Log Compression for Workload Analytics | 3 |
| 3.1 | Introduction | 5 |
| 3.2 | Problem Definition | 7 |
| 3.2.1 | Preliminaries and Notation | 8 |
| 3.2.2 | Coding Queries | 8 |
| 3.2.3 | Log Compression | 10 |
| 3.2.3.1 | Information Content of Logs | 10 |
| 3.2.3.2 | Communicating Information Content | 12 |
| 3.3 | Information Loss | 12 |
| 3.3.1 | Lossless Summaries | 12 |
| 3.3.2 | Lossy Summaries | 13 |
| 3.3.3 | Idealized Information Loss Measures | 14 |
| 3.4 | Practical Loss Measure | 15 |
| 3.4.1 | Reproduction Error | 15 |
| 3.4.2 | Practical vs Idealized Information Loss | 16 |
| 3.5 | Pattern Mixture Encodings | 17 |
| 3.5.1 | Example: Naive Mixture Encodings | 18 |
| 3.5.2 | Generalized Encoding Fidelity | 19 |
| 3.6 | Pattern Mixture Compression | 20 |

| | | |
|----------|---|-----------|
| 3.6.1 | Constructing Naive Mixture Encodings | 20 |
| 3.6.1.1 | Clustering | 22 |
| 3.6.2 | Approximating Log Statistics | 23 |
| 3.6.3 | Pattern Synthesis & Frequency Estimate | 24 |
| 3.6.4 | Naive Encoding Refinement | 25 |
| 3.7 | Experiments | 26 |
| 3.7.1 | Validating Reproduction Error | 29 |
| 3.7.2 | Feature-Correlation Refinement | 30 |
| 3.7.2.1 | Pattern-based vs Naive Mixture Encodings | 31 |
| 3.7.2.2 | Refining Naive Mixture Encodings | 31 |
| 3.8 | Alternative Applications | 32 |
| 3.8.1 | Experiments | 32 |
| 3.8.1.1 | Error Measures | 33 |
| 3.8.1.2 | Classical Laserlight and MTV | 34 |
| 3.8.1.3 | Generalizing Laserlight and MTV | 35 |
| 3.8.1.4 | Comparison with Naive Mixture Encoding | 35 |
| 3.9 | Related Work | 36 |
| 3.9.1 | Workload Analysis | 36 |
| 3.9.2 | Workload Compression Schemes | 37 |
| 3.10 | Conclusions | 38 |
| 3.11 | Future Work | 39 |
| 3.12 | Acknowledgements | 40 |
| 4 | Application: Interactive Semi-Structured Schema Design | 41 |
| 4.1 | Introduction | 43 |
| 4.1.1 | Extracting Relational Entities | 44 |
| 4.1.2 | Human-Scale S ³ D | 44 |
| 4.1.3 | Overview | 45 |
| 4.2 | Summarization | 46 |
| 4.2.1 | Data Model | 47 |
| 4.2.2 | Paths as Attributes | 48 |

| | | |
|---------------------|--|-----------|
| 4.2.3 | Schema Collections | 48 |
| 4.2.4 | Summarizing Schema Collections | 50 |
| 4.3 | Visualization | 53 |
| 4.3.1 | Schema Segmentation | 53 |
| 4.3.2 | Schema Exploration | 57 |
| 4.3.3 | An Iterative Approach | 58 |
| 4.4 | Related Work | 59 |
| 4.5 | Future Work | 60 |
| 5 | Application: Summarizing Tuple Correlation in Probabilistic Databases | 61 |
| 5.1 | Introduction | 63 |
| 5.2 | Problem Definition | 65 |
| 5.2.1 | Practical Input Representation | 66 |
| 5.2.2 | Summarizing Joint Distribution | 67 |
| 5.2.3 | Pattern-based Summary | 68 |
| 5.3 | Marginal Estimation | 69 |
| 5.4 | Summary Accuracy | 70 |
| 5.4.1 | Query Accuracy | 72 |
| 5.4.1.1 | Accuracy under join | 72 |
| 5.4.1.2 | Accuracy under selection predicate | 73 |
| 5.4.1.3 | Accuracy under project | 74 |
| 5.5 | Constructing Summaries | 75 |
| 5.6 | Related Work | 75 |
| 6 | Findings and Discussions | 77 |
| 7 | Conclusions and Future Work | 78 |
| Bibliography | | 79 |
| A Appendix | | 86 |

List of Figures

| | | |
|-----|---|----|
| 3.1 | Example Encoding Visualizations | 11 |
| 3.2 | Clustering Schemes Comparison | 21 |
| 3.3 | Effectiveness of Naive Mixture Encoding | 24 |
| 3.4 | Validating Reproduction Error | 27 |
| 3.5 | Feature-correlation refinement (US bank) | 32 |
| 3.6 | Naive Mixture v. Laserlight/MTV Mixture | 35 |
| 4.1 | Prototype user interface for SCHEMADRILL. | 46 |
| 4.2 | FP-Tree based schema summaries | 51 |
| 4.3 | RCDC based schema summaries | 53 |
| 4.4 | Segmentation breaks up schema representations into manageable chunks. | 54 |
| 4.5 | Covariance Cloud for the full Yelp dataset. | 55 |
| 4.6 | KNN-PCA cloud with K=6 on the Yelp dataset. | 56 |
| 4.7 | Points containing user selected attributes are highlighted green. | 58 |
| 5.1 | Theoretical v.s. Practical Data-flows. | 66 |

Chapter 1

Introduction

todo

Chapter 2

Problem Overview

todo

Chapter 3

Application: Query Log Compression for Workload Analytics

Abstract

Analyzing database access logs is a key part of performance tuning, intrusion detection, benchmark development, and many other database administration tasks. Unfortunately, it is common for production databases to deal with millions or more queries each day, so these logs must be summarized before they can be used. Designing an appropriate summary encoding requires trading off between conciseness and information content. For example: simple workload sampling may miss rare, but high impact queries. In this paper, we present LOGR, a lossy log compression scheme suitable for use in many automated log analytics tools, as well as for human inspection. We formalize and analyze the space/fidelity trade-off in the context of a broader family of “pattern” and “pattern mixture” log encodings to which LOGR belongs. We show through a series of experiments that LOGR compressed encodings can be created efficiently, come with provable information-theoretic bounds on their accuracy, and outperform state-of-art log summarization strategies.

3.1 Introduction

Automated analysis of database access logs is critical for solving a wide range of problems, from database performance tuning [10], to compliance validation [14], and query recommendation [12]. For example, the Peloton self-tuning database [39] searches for optimal configurations by repeatedly simulating database performance based on statistical properties of historical queries. Unfortunately, query logs for production databases can grow to be large — A recent study of queries at a major US bank for a period of 19 hours found nearly 17 million SQL queries and over 60 million stored procedure executions [30] — and computing these properties from the log itself is slow.

Tracking only a sample of these queries is not sufficient, as rare queries can disproportionately affect database performance, for example, if they benefit from an otherwise unnecessary index. Rather, we need a compressed *summary* of the log on which we can compute aggregate statistical properties. The problems of compression and summarization have been studied extensively (e.g., [47, 48, 21, 33, 24, 8, 42, 29]). However, these schemes either require the use of heavyweight inference to desired statistical measures, or produce unnecessarily large encodings.

In this paper, we adapt ideas from pattern mining and summarization [35, 16] to propose a middle-ground: LOGR, a summarization scheme that facilitates efficient (both in terms of storage and time) approximation of workload statistics. By adjusting a tunable parameter in LOGR, users can choose to obtain a high-fidelity, albeit large summary, or obtain a more compact summary with lower fidelity. Constructing the summary that best balances compactness and fidelity is challenging, as the search space of candidate summaries is combinatorially large [35, 16]. LOGR offers a new approach to summary construction that avoids searching this space, making inexpensive, accurate computation of aggregate workload statistics possible. As a secondary benefit, the resulting summaries are also human-interpretable.

LOGR does not admit closed-form solutions to classical fidelity measures like information loss, so we propose an alternative called *Reproduction Error*. We show

through a combination of analytical and experimental evidence that Reproduction Error is highly correlated with several classical measures of encoding fidelity.

LOGR-compressed data relies on a codebook of structural elements like `SELECT` items, `FROM` tables, or conjunctive `WHERE` clauses [3]. This codebook provides a bi-directional mapping from SQL queries to a bit-vector encoding, reducing the compression problem to one of compactly encoding a collection of feature-vectors. We further simplify the problem by observing that a common theme in use cases like automated performance tuning or query recommendation is the need for predominantly aggregate workload statistics. As these are order-independent, we are able to focus exclusively on compactly representing *bags* of feature-vectors.

LOGR works by identifying groups of co-occurring structural elements that we call patterns. We define a family of *pattern encodings* of access logs, which map patterns to their frequencies in the log. For pattern encodings, we consider two idealized measures of fidelity: (1) Ambiguity, which measures how much room the encoding leaves for interpretation; and (2) Deviation, which measures how reliably the encoding approximates the original log. Neither Ambiguity nor Deviation can be computed efficiently for pattern encodings. Hence we propose a measure called *Reproduction Error* that is efficiently computable and that closely tracks both Ambiguity and Deviation.

In general, the size of the encoding is inversely related with Reproduction Error: The more detailed the encoding, the more faithfully it represents the original log. Thus, log compression may be defined as a search over the space of pattern-based encodings to identify the one that best trades off between these two properties. Unfortunately, searching for such an ideal encoding from the space can be computationally expensive [35, 16]. To overcome this limitation, we reduce the search space by first clustering entries in the log and then encoding each cluster separately, an approach that we call *pattern mixture encoding*. Finally we identify a simple approach for encoding individual clusters that we call *naive mixture encodings*, and show experimentally that it produces results competitive with more general techniques for log compression and summarization.

Concretely, in this paper we make the following contributions: (1) We define two families of compression for query logs: pattern and pattern mixture, (2) We define a computationally efficient measure, Reproduction Error, and demonstrate that it is a close approximation of Ambiguity and Deviation (two commonly used measures), (3) We propose a clustering-based approach to efficiently search for naive mixture encodings, and show how these encodings can be further optimized, and, (4) We experimentally validate LOGR and show that it produces more precise encodings, faster than several state-of-the-art pattern encoding algorithms.

Roadmap. The paper is organized as follows: Section 3.2 formally defines the log compression problem and the summary representation; Section 3.3 then defines information loss of the summaries; Section 3.4 explains the difficulty in computing classical loss measures and provides a practical alternative; Section 3.5 motivates data partitioning and generalizes the practical loss measure to partitioned data; Section 3.6 then introduces the proposed LOGR compression scheme; Section 3.7 empirically validates the practical loss measure and evaluates the effectiveness of LOGR by comparing it with two state-of-the-art summarization methods; Section 3.8 further evaluates LOGR under applications of comparison methods; Section 3.9 discusses related work. Section 3.10 concludes the paper and Section 3.11 discusses future work.

3.2 Problem Definition

In this section, we introduce and formally define the log compression problem. We begin by exploring several applications that need to repeatedly analyze query logs.

Index Selection. Selecting an appropriate set of indexes requires trading off between update costs, access costs, and limitations on available storage space. Existing strategies for selecting a near-optimal set of indexes typically repeatedly simulate database performance under different combinations of indexes, which in turn requires repeatedly estimating the frequency of specific predicates in the workload.

Materialized View Selection. The results of joins or highly selective selection predicates are good candidates for materialization when they appear frequently in the

workload. Like index selection, view selection is a non-convex optimization problem, typically requiring repeated simulation, which in turn requires repeated frequency estimation over the workload.

Online Database Monitoring. In production settings, it is common to monitor databases for atypical usage patterns that could indicate a serious bug or security threat. When query logs are monitored, it is often done retrospectively, some hours after-the-fact [30]. To support real-time monitoring it is necessary to quickly compute the frequency of a particular class of queries in the system’s typical workload.

In each case, the application’s interactions with the log amount to counting queries that have specific features: selection predicates, joins, or similar.

3.2.1 Preliminaries and Notation

Let L be a log, or a finite collection of queries $\mathbf{q} \in L$. We write $f \in \mathbf{q}$ to indicate that \mathbf{q} has some *feature* f , such as a specific predicate or table in its `FROM` clause. We assume (1) that the universe of features in both a log and a query is enumerable and finite, (2) that the features are selected to suit specific applications and (3) optionally that a query is isomorphic to its feature set (motivated in Section 3.2.3.2). We outline one approach to extracting features that satisfies all three assumptions below. We abuse syntax and write \mathbf{q} to denote both the query itself, as well as the set of its features.

Let ω denote some set of features $f \in \omega$. We write these sets using vector notation: $\omega = (x_1, \dots, x_n)$ where n is the number of distinct features in the entire log and x_i indicates the presence (absence) of i th feature with a 1 (resp., 0). For any two patterns ω, ω' , we say that ω' is *contained* or *appears* in ω if $\omega' \subseteq \omega$, or equivalently if $\forall i, x'_i \leq x_i$.

3.2.2 Coding Queries

For this paper, we specifically adopt the feature extraction conventions of a query summarization scheme by Aligon et al. [3]. In this scheme, each feature is one of the

following three query elements: (1) a table or sub-query in the `FROM` clause, (2) a column in the `SELECT` clause, and (3) a conjunctive atom of the `WHERE` clause.

Example 1. Consider the following example query.

```
SELECT _id, sms_type, _time FROM Messages
WHERE status=? AND transport_type=?
```

The query has 6 features: $\langle \text{_id}, \text{SELECT} \rangle$, $\langle \text{sms_type}, \text{SELECT} \rangle$, $\langle \text{_time}, \text{SELECT} \rangle$, $\langle \text{Messages}, \text{FROM} \rangle$, $\langle \text{status}=\text{?}, \text{WHERE} \rangle$, and $\langle \text{transport_type}=\text{?}, \text{WHERE} \rangle$

Although this scheme is simple and limited to conjunctive queries, it fulfills all three assumptions we make on feature extraction schemes. The features of a query (and consequently a log) are enumerable and finite, and the feature set of the query is isomorphic to the original query. Furthermore, even if a query is not itself conjunctive, it may be rewritable into a conjunctive equivalent.

Although we do not explore more advanced feature encoding schemes in detail here, we direct the interested reader to work on query summarization [34, 7, 30]. For example, a scheme by Makiyama et. al. [34] also captures aggregation-related features like group-by columns, while an approach by Kul et. al. [30] encodes partial tree-structures in the query.

3.2.3 Log Compression

As a lossy form of compression, LOGR only approximates the information content of a query log. We next develop a simplified form of LOGR that we call pattern-based encoding, and develop a framework for reasoning about the fidelity of a LOGR-compressed log. As a basis for this framework, we first formulate the information content of a query log to allow us to adapt classical measures of information content.

3.2.3.1 Information Content of Logs

We define the information content of the log as a distribution $p(Q | L)$ of queries Q drawn uniformly from the log.

Example 2. Consider the following query log, which consists of four conjunctive queries.

1. `SELECT _id FROM Messages WHERE status = ?`
2. `SELECT _time FROM Messages WHERE status = ? AND sms_type = ?`
3. `SELECT _id FROM Messages WHERE status = ?`
4. `SELECT sms_type, _time FROM Messages WHERE sms_type = ?`

Drawing uniformly from the log, each entry will appear with probability $\frac{1}{4} = 0.25$. The query $q_1 (= q_3)$ occurs twice, so the probability of drawing it is double that of the others (i.e., $p(q_1 | L) = p(q_3 | L) = \frac{2}{4} = 0.5$)

Treating a query as a vector of its component features, we can define a query $\mathbf{q} = (x_1, \dots, x_n)$ to be an observation of the multivariate distribution over variables $Q = (X_1, \dots, X_n)$ corresponding to features. The event $X_i = 1$ occurs if feature i appears in a uniformly drawn query.

Example 3. Continuing, the universe of features for this query log is (1) $\langle \text{_id}, \text{SELECT} \rangle$, (2) $\langle \text{_time}, \text{SELECT} \rangle$, (3) $\langle \text{sms_type}, \text{SELECT} \rangle$, (4) $\langle \text{status} = ?, \text{WHERE} \rangle$, (5) $\langle \text{sms_type} = ?, \text{WHERE} \rangle$, and (6) $\langle \text{Messages}, \text{FROM} \rangle$. Accordingly, the queries can

```

SELECT sms_type, external_ids, _time,
_id
FROM messages
WHERE (sms_type=?)

```

(a) *Correlation-ignorant*: Features are highlighted independently

| |
|---|
| SELECT sms_type FROM messages WHERE sms_type=? |
| SELECT sms_type FROM messages WHERE status=? |

(b) *Correlation-aware*: Pattern groups are highlighted together.

Figure 3.1: **Example Encoding Visualizations**

be encoded as feature vectors, with fields recording each feature’s presence: $\mathbf{q}_1 = \langle 1, 0, 0, 1, 0, 1 \rangle$, $\mathbf{q}_2 = \langle 0, 1, 0, 1, 1, 1 \rangle$, $\mathbf{q}_3 = \langle 1, 0, 0, 1, 0, 1 \rangle$, $\mathbf{q}_4 = \langle 0, 1, 1, 0, 1, 1 \rangle$

Patterns. Our target applications require us to count the number of times features (co-)occur in a query. For example, materialized view selection requires counting tables used together in queries. Motivated by this observation, we begin by defining a broad class of *pattern-based encodings* that directly encode co-occurrence probabilities. A *pattern* is an arbitrary set of features $\omega = (x_1, \dots, x_n)$ that may co-occur together. Each pattern captures a piece of information from the distribution $p(Q | L)$. In particular, we are interested in the probability of uniformly drawing a query \mathbf{q} from the log that *contains* the pattern \mathbf{b} (i.e., $\mathbf{q} \supseteq \omega$):

$$p(Q \supseteq \omega | L) = \sum_{\mathbf{q} \in L \wedge \mathbf{q} \supseteq \omega} p(\mathbf{q} | L)$$

When it is clear from context, we abuse notation and write $p(\cdot)$ instead of $p(\cdot | L)$. Recall that $p(Q)$ can be represented as a joint distribution of variables (X_1, \dots, X_n) and probability $p(Q \supseteq \omega)$ is equivalent to $p(X_1 \geq x_1, \dots, X_n \geq x_n)$.

Pattern-Based Encodings. Denote by $\mathcal{E}_{max} : \{0, 1\}^n \rightarrow [0, 1]$, the mapping from each pattern (ω) to its frequency in the log: $\mathcal{E}_{max} = \{ (\omega \rightarrow p(\omega)) \mid \omega \in \{0, 1\}^n \}$

A *pattern-based encoding* \mathcal{E} is any such partial mapping $\mathcal{E} \subseteq \mathcal{E}_{max}$. We denote the frequency of pattern ω in encoding \mathcal{E} by $\mathcal{E}[\omega]$ ($= p(Q \supseteq \omega)$). When it is clear from context, we abuse syntax and also use \mathcal{E} to denote the set of patterns it maps (i.e., $domain(\mathcal{E})$). Hence, $|\mathcal{E}|$ is the number of mapped patterns, which we call the

encoding’s *Verbosity*. A *pattern-based encoder* is any algorithm $\text{encode}(L, \epsilon) \mapsto \mathcal{E}$ whose input is a log L and whose output is a set of patterns \mathcal{E} , with Verbosity thresholded at some integer ϵ . Many pattern mining algorithms [35, 16] can be used for this purpose.

3.2.3.2 Communicating Information Content

A side-benefit of pattern-based encodings is that, under the assumption of isomorphism in Section 3.2.1, patterns can be translated to their query representations and used for human inspection of the log. Figure 3.1 shows two examples. The approach illustrated in Figure 3.1a uses shading to show each feature’s frequency in the log, and communicates frequently occurring predicates or columns. This approach might, for example, help a human to manually select indexes. A second approach illustrated in Figure 3.1b conveys correlations, showing the frequency of entire patterns. The accompanying technical report [45] explores visualizations of pattern-based summaries in greater depth.

3.3 Information Loss

Our goal is to encode the distribution $p(Q)$ as a set of patterns: obtaining a less verbose encoding (i.e., with fewer patterns), while also ensuring that the encoding captures $p(Q)$ with minimal information loss. In this section, we define information loss for pattern-based encodings.

3.3.1 Lossless Summaries

To establish a baseline for measuring information loss, we begin with the extreme cases. At one extreme, an empty encoding ($|\mathcal{E}| = 0$) conveys no information. At the other extreme, we have the encoding \mathcal{E}_{max} which is the full mapping from all patterns. Having this encoding is a sufficient condition to exactly reconstruct the original distribution $p(Q)$.

Proposition 1. *For any query $\mathbf{q} = (x_1, \dots, x_n) \in \{0, 1\}^n$, the probability of drawing exactly \mathbf{q} at random from the log (i.e., $p(X_1 = x_1, \dots, X_n = x_n)$) is computable, given \mathcal{E}_{max} .*

3.3.2 Lossy Summaries

Although \mathcal{E}_{max} is lossless, its Verbosity is exponential in the number of features (n). Hence, we will focus on lossy encodings that can be less verbose. A lossy encoding $\mathcal{E} \subset \mathcal{E}_{max}$ may not precisely identify the distribution $p(Q)$, but can still be used to approximate it. We characterize the information content of a lossy encoding \mathcal{E} by defining a *space* (denoted by $\Omega_{\mathcal{E}}$) of distributions $\rho \in \Omega_{\mathcal{E}}$ allowed by an encoding \mathcal{E} . This space is defined by constraints as follows: First, we have the general properties of probability distributions:

$$\forall \mathbf{q} \in \{0, 1\}^n : \rho(\mathbf{q}) \geq 0 \quad \sum_{\mathbf{q}} \rho(\mathbf{q}) = 1$$

Each pattern \mathbf{b} in the encoding \mathcal{E} constrains relevant probabilities in distribution ρ to sum to the target frequency:

$$\forall \mathbf{b} \in domain(\mathcal{E}) : \mathcal{E}[\mathbf{b}] = \sum_{\mathbf{q} \supseteq \mathbf{b}} \rho(\mathbf{q})$$

Note that the dual constraints $1 - \mathcal{E}[\mathbf{b}] = \sum_{\mathbf{q} \not\supseteq \mathbf{b}} \rho(\mathbf{q})$ are redundant under constraint $\sum_{\mathbf{q}} \rho(\mathbf{q}) = 1$.

The resulting space $\Omega_{\mathcal{E}}$ is the set of all query logs, or equivalently the set of all possible distributions of queries, that obey these constraints. From the outside observer's perspective, the distribution $\rho \in \Omega_{\mathcal{E}}$ that the encoding conveys is ambiguous: We model this ambiguity using a random variable $\mathcal{P}_{\mathcal{E}}$ with support $\Omega_{\mathcal{E}}$. The true distribution $p(Q)$ derived from the query log must appear in $\Omega_{\mathcal{E}}$, denoted as $p(Q) \equiv \rho^* \in \Omega_{\mathcal{E}}$ (i.e., $p(\mathcal{P}_{\mathcal{E}} = \rho^*) > 0$). Of the remaining distributions ρ admitted by $\Omega_{\mathcal{E}}$, it is possible that some are more likely than others. For example, a query containing a column (e.g., `status`) is only valid if it also references a table that contains the column (e.g., `Messages`). This prior knowledge may be modeled as a prior on the distribution of $\mathcal{P}_{\mathcal{E}}$ or equivalently by an additional constraint. However,

for the purposes of this paper, we take the uninformed prior by assuming that $\mathcal{P}_{\mathcal{E}}$ is uniformly distributed over $\Omega_{\mathcal{E}}$:

$$p(\mathcal{P}_{\mathcal{E}} = \rho) = \begin{cases} \frac{1}{|\Omega_{\mathcal{E}}|} & \text{if } \rho \in \Omega_{\mathcal{E}} \\ 0 & \text{otherwise} \end{cases}$$

Naive Encodings. One specific family of lossy encodings that treat each feature as being independent (e.g., as in Figure 3.1a) is of particular interest to us. We call this family *naive encodings*, and return to it throughout the rest of the paper. A naive encoding $\ddot{\mathcal{E}}$ is composed of all patterns that have exactly one feature with non-zero frequency.

$$\text{domain}(\ddot{\mathcal{E}}) = \{ (0, \dots, 0, x_i, 0, \dots, 0) \mid i \in [1, n], x_i = 1 \}$$

3.3.3 Idealized Information Loss Measures

Based on the space of distributions constrained by the encoding, the information loss of an encoding can be considered from two related, but subtly distinct perspectives: (1) *Ambiguity* measures how much room the encoding leaves for interpretation and (2) *Deviation* measures how reliably the encoding approximates the target distribution $p(Q)$.

Ambiguity. We define the Ambiguity $I(\mathcal{E})$ of an encoding as the entropy of the random variable $\mathcal{P}_{\mathcal{E}}$. The higher the entropy, the less precisely \mathcal{E} identifies a specific distribution.

$$I(\mathcal{E}) = \sum_{\rho} p(\mathcal{P}_{\mathcal{E}} = \rho) \log(p(\mathcal{P}_{\mathcal{E}} = \rho))$$

Deviation. The deviation from any permitted distribution ρ to the true distribution ρ^* can be measured by the Kullback-Leibler (K-L) divergence $\mathcal{D}_{KL}(\rho^* || \rho)$. We define the Deviation $d(\mathcal{E})$ of a encoding as the expectation of the K-L divergence over all permitted $\rho \in \Omega_{\mathcal{E}}$:

$$d(\mathcal{E}) = \mathbb{E}_{\mathcal{P}_{\mathcal{E}}} [\mathcal{D}_{KL}(\rho^* || \mathcal{P}_{\mathcal{E}})] = \sum_{\rho \in \Omega_{\mathcal{E}}} p(\mathcal{P}_{\mathcal{E}} = \rho) \cdot \mathcal{D}_{KL}(\rho^* || \rho)$$

Limitations. There are two limitations to these idealized measures in practice. First, K-L divergence is not defined on any permitted distribution ρ where the true distribution ρ^* is not *absolutely continuous* (denoted $\rho^* \ll \rho$). Second, neither Deviation nor Ambiguity has a closed-form formula.

3.4 Practical Loss Measure

Computing either Ambiguity or Deviation requires enumerating the entire space of permitted distributions. One approach to approximating either measure is repeatedly sampling from, rather than enumerating the space. However, accurate approximations require a large number of samples, rendering this approach similarly inefficient. In this section, we propose a faster approach to assessing the fidelity of a pattern encoding. Specifically, we select a single representative distribution $\bar{\rho}_{\mathcal{E}}$ from the space $\Omega_{\mathcal{E}}$, and use $\bar{\rho}_{\mathcal{E}}$ to approximate both Ambiguity and Deviation.

3.4.1 Reproduction Error

Maximum Entropy Distribution. The representative distribution is chosen by applying the maximum entropy principle [23] commonly used in pattern-based summarization [35, 16]. That is, we select the distribution $\bar{\rho}_{\mathcal{E}}$ with maximum entropy:

$$\bar{\rho}_{\mathcal{E}} = \arg \max_{\rho \in \Omega_{\mathcal{E}}} \mathcal{H}(\rho) \quad \text{where } \mathcal{H}(\rho) = \sum_{\mathbf{q} \in \{0,1\}^n} -\rho(\mathbf{q}) \log \rho(\mathbf{q})$$

The maximum entropy distribution best represents the current state of knowledge. That is, a distribution with lower entropy assumes additional constraints derived from patterns that we do not know, while one with higher entropy violates the constraints from patterns we do know.

Maximizing an objective function belonging to the exponential family (entropy in our case) under a mixture of linear equalities/inequality constraints is a convex optimization problem [9] which guarantees a *unique* solution and can be efficiently solved using the cvx toolkit [18, 38], and/or by *iterative scaling* [35, 16]. For naive

encodings specifically, we can assume independence between each feature X_i . Under this assumption, $\bar{\rho}_{\mathcal{E}}$ has a closed-form representation:

$$\bar{\rho}_{\mathcal{E}}(\mathbf{q}) = \prod_i p(X_i = x_i) \quad \text{where } \mathbf{q} = (x_1, \dots, x_n)$$

We define *Reproduction Error* $e(\mathcal{E})$ as the entropy difference between the representative and true distributions:

$$e(\mathcal{E}) = \mathcal{H}(\bar{\rho}_{\mathcal{E}}) - \mathcal{H}(\rho^*) \quad \text{where } \bar{\rho}_{\mathcal{E}} = \arg \min_{\rho \in \Omega_{\mathcal{E}}} -\mathcal{H}(\rho)$$

3.4.2 Practical vs Idealized Information Loss

In this section we prove that Reproduction Error closely parallels Ambiguity. We define a partial order lattice over encodings and show that for any pair of encodings on which the partial order is defined, a like relationship is implied for both Reproduction Error and Ambiguity. We supplement the proofs given in this section with an empirical analysis relating Reproduction Error to Deviation in Section 3.7.1.

Containment. We define a partial order over encodings \leq_{Ω} based on *containment* of their induced spaces $\Omega_{\mathcal{E}}$:

$$\mathcal{E}_1 \leq_{\Omega} \mathcal{E}_2 \equiv \Omega_{\mathcal{E}_1} \subseteq \Omega_{\mathcal{E}_2}$$

That is, one encoding (i.e., \mathcal{E}_1) precedes another (i.e., \mathcal{E}_2) when all distributions admitted by the former encoding are also admitted by the latter.

Containment Captures Reproduction Error. We first prove that the total order given by Reproduction Error is a superset of the partial order \leq_{Ω} .

Lemma 1. *For any pair of encodings $\mathcal{E}_1, \mathcal{E}_2$ that induce spaces $\Omega_{\mathcal{E}_1}, \Omega_{\mathcal{E}_2}$ and maximum entropy distributions $\bar{\rho}_{\mathcal{E}_1}, \bar{\rho}_{\mathcal{E}_2}$ it holds that $\mathcal{E}_1 \leq_{\Omega} \mathcal{E}_2 \rightarrow e(\mathcal{E}_1) \leq e(\mathcal{E}_2)$.*

Proof. First we have $\Omega_{\mathcal{E}_2} \supseteq \Omega_{\mathcal{E}_1} \rightarrow \bar{\rho}_{\mathcal{E}_1} \in \Omega_{\mathcal{E}_2}$. Since $\bar{\rho}_{\mathcal{E}_2}$ has the maximum entropy among all distributions $\rho \in \Omega_{\mathcal{E}_2}$, we have $\mathcal{H}(\bar{\rho}_{\mathcal{E}_1}) \leq \mathcal{H}(\bar{\rho}_{\mathcal{E}_2}) \equiv e(\mathcal{E}_1) \leq e(\mathcal{E}_2)$. \square

Containment Captures Ambiguity. Next, we show that the partial order based on containment implies a like relationship between Ambiguities of pairs of encodings.

Lemma 2. *Given encodings $\mathcal{E}_1, \mathcal{E}_2$ with uninformed prior on $\mathcal{P}_{\mathcal{E}_1}, \mathcal{P}_{\mathcal{E}_2}$, it holds that $\mathcal{E}_1 \leq_{\Omega} \mathcal{E}_2 \rightarrow I(\mathcal{E}_1) \leq I(\mathcal{E}_2)$.*

Proof. Given an uninformed prior: $I(\mathcal{E}) = \log |\Omega_{\mathcal{E}}|$, we have $\mathcal{E}_1 \leq_{\Omega} \mathcal{E}_2 \rightarrow |\Omega_{\mathcal{E}_1}| \leq |\Omega_{\mathcal{E}_2}| \rightarrow I(\mathcal{E}_1) \leq I(\mathcal{E}_2)$ \square

3.5 Pattern Mixture Encodings

Thus far we have defined the problem of log compression, treating the query log as a multivariate distribution $p(Q)$ where patterns capture positive frequencies of feature (co-)occurrence. However in cases like logs of *mixed* workloads, there are also many cases of anti-correlation between features. For example, consider a log that includes queries drawn from a mixture of two workloads with disjoint feature sets. Pattern-based summaries can not convey such anti-correlations easily. As a result, patterns including features from both workloads never actually co-occur in the log, but a pattern-based summary of the log will suggest otherwise. Such false positives are especially problematic for use-cases of LOGR involving outlier detection (e.g., [32]). Even in other settings, capturing correlations reduces data dimensionality and improves both runtime and effectiveness of state-of-the-art pattern mining algorithms (See Section 3.8.1).

In this section, we propose a generalization of pattern encodings where the log is modeled not as a single probability distribution, but rather as a mixture of several simpler distributions. The resulting encoding is likewise a mixture: Patterns for each component of the mixture are stored independently. Hence, we refer to it as a *pattern mixture encoding*, and it forms the basis of LOGR compression. We first focus on a simplified form of this problem, where we only mix *naive* encodings (we explore more general mixtures in Section 3.6.4). We refer to the resulting scheme as *naive mixture encodings*, and give examples of the encoding in Section 3.5.1. Then we generalize Reproduction Error and Verbosity to pattern mixture encodings in Section 3.5.2. Finally, with generalized encoding evaluation measures, we evaluate several clustering methods for creating naive mixture encodings.

3.5.1 Example: Naive Mixture Encodings

Consider a toy query log with only 3 conjunctive queries.

1. `SELECT id FROM Messages WHERE status = ?`
2. `SELECT id FROM Messages`
3. `SELECT sms_type FROM Messages`

The codebook of this log includes 4 features: $\langle \text{id}, \text{SELECT} \rangle$, $\langle \text{sms_type}, \text{SELECT} \rangle$, $\langle \text{Messages}, \text{FROM} \rangle$, $\langle \text{status} = ?, \text{WHERE} \rangle$. Re-encoding the three queries as vectors, we get:

$$1. \langle 1, 0, 1, 1 \rangle \quad 2. \langle 1, 0, 1, 0 \rangle \quad 3. \langle 0, 1, 1, 0 \rangle$$

A naive encoding of this log can be expressed as:

$$\left\langle \frac{2}{3}, \frac{1}{3}, 1, \frac{1}{3} \right\rangle$$

This encoding captures that all queries in the log pertain to the `Messages` table, but obscures the relationship between the remaining features. For example, this encoding obscures the anti-correlation between `id` and `sms_type`. Similarly, the encoding hides the correlation between `status = ?` and `id`. Such relationships are critical for evaluating the effectiveness of views or indexes.

Example 4. *The maximum entropy distribution for any naive encoding assumes that features are independent. Assuming independence, the probability of query 1 uniformly drawn from the log is estimated as:*

$$p(\text{id}) \cdot p(\neg \text{sms_type}) \cdot p(\text{Messages}) \cdot p(\text{status}=?) = \frac{4}{27} \approx 0.148$$

This is a significant difference from the true probability of this query (i.e., $\frac{1}{3}$). Conversely queries not in the log, such as the following, have non-zero probability in the encoding.

`SELECT sms_type FROM Messages WHERE status = ?`

$$p(\neg id) \cdot p(sms_type) \cdot p(Messages) \cdot p(status=?) = \frac{1}{27} \approx 0.037$$

To achieve a more faithful representation of the original log, we could partition it into two components, with the corresponding encoding parameters:

| Partition 1 (L_1) | Partition 2 (L_2) |
|---------------------------------------|---------------------------------------|
| (1, 0, 1, 1) | (0, 1, 1, 0) |
| ↓ ↓ | ↓ |
| ⟨ 1, 0, 1, $\frac{1}{2}$ ⟩ | ⟨ 0, 1, 1, 0 ⟩ |

Although there are now two encodings, the encodings are not ambiguous. The feature `status = ?` appears in exactly half of the log entries, and is indeed independent of the other features. All other attributes in each encoding appear in all queries in their respective partitions. Furthermore, the maximum entropy distribution induced by each encoding is exactly the distribution of queries in the partitioned log. Hence, the Reproduction Error is zero for both encodings.

3.5.2 Generalized Encoding Fidelity

We next generalize our definitions of Reproduction Error and Verbosity from pattern-based to pattern mixture encodings. Suppose query log L has been partitioned into K clusters with L_i , \mathcal{E}_i , $\bar{\rho}_{\mathcal{E}_i}$ and ρ_i^* (where $i \in [1, K]$) representing the log of queries, encoding, maximum entropy distribution, and true distribution (respectively) for the i th cluster. First, observe that the distribution for the whole log (i.e., ρ^*) is the sum of distributions for each partition (i.e., ρ_i^*) weighted by the proportion (i.e., $\frac{|L_i|}{|L|}$) of queries:

$$\rho^*(\mathbf{q}) = \sum_{i=1, \dots, K} w_i \cdot \rho_i^*(\mathbf{q}) \quad \text{where } w_i = \frac{|L_i|}{|L|}$$

Generalized Reproduction Error. Similarly, the maximum entropy distribution $\bar{\rho}_{\mathcal{E}}$ for the whole log is:

$$\bar{\rho}_{\mathcal{E}}(\mathbf{q}) = \sum_{i=1, \dots, K} w_i \cdot \bar{\rho}_{\mathcal{E}_i}(\mathbf{q})$$

We define the *Generalized Reproduction Error* of a pattern mixture encoding similarly, as the weighted sum of Reproduction Error for each partition:

$$e(\mathcal{E}) = \mathcal{H}(\bar{\rho}_{\mathcal{E}}) - \mathcal{H}(\rho^*) = \sum_i w_i (\mathcal{H}(\bar{\rho}_{\mathcal{E}_i}) - \mathcal{H}(\rho_i^*)) = \sum_i w_i e(\mathcal{E}_i)$$

When it is clear from context, we refer to Generalized Reproduction Error as Error in the rest of this paper. As in the base case, a pattern mixture encoding with low Error indicates a high-fidelity representation of the original log. A process can infer the frequency of any query $p(Q = \mathbf{q} | L)$ drawn from the original distribution, simply by inferring its frequency in each cluster i (i.e., $p(Q = \mathbf{q} | L_i)$) and taking a weighted average over all inferences.

Generalized Verbosity. We generalize verbosity to pattern mixture encodings as the *Total Verbosity* ($\sum_i |S_i|$), or the total size of the encoded representation.

3.6 Pattern Mixture Compression

We are now ready to describe the LOGR compression scheme. Broadly, LOGR attempts to identify a pattern mixture encoding that optimizes for some target trade-off between Total Verbosity and Error. A naive — though impractical — approach to finding such an encoding would be to search the entire space of possible pattern mixture encodings. Instead, LOGR approximates the same outcome by first identifying the naive mixture encoding that is closest to optimal for the desired trade-off. As we show experimentally, this naive mixture encoding is competitive with more complicated, slower techniques for summarizing query logs. We also explore a hypothetical second stage, where LOGR refines the naive mixture encoding to further reduce Error. The outcome of this hypothetical stage has a slightly lower Error and Verbosity, but does not admit efficient computation of database statistics.

3.6.1 Constructing Naive Mixture Encodings

LOGR searches for a naive mixture encoding that best optimizes for a requested tradeoff between Total Verbosity and Error. As a way to make this search efficient,

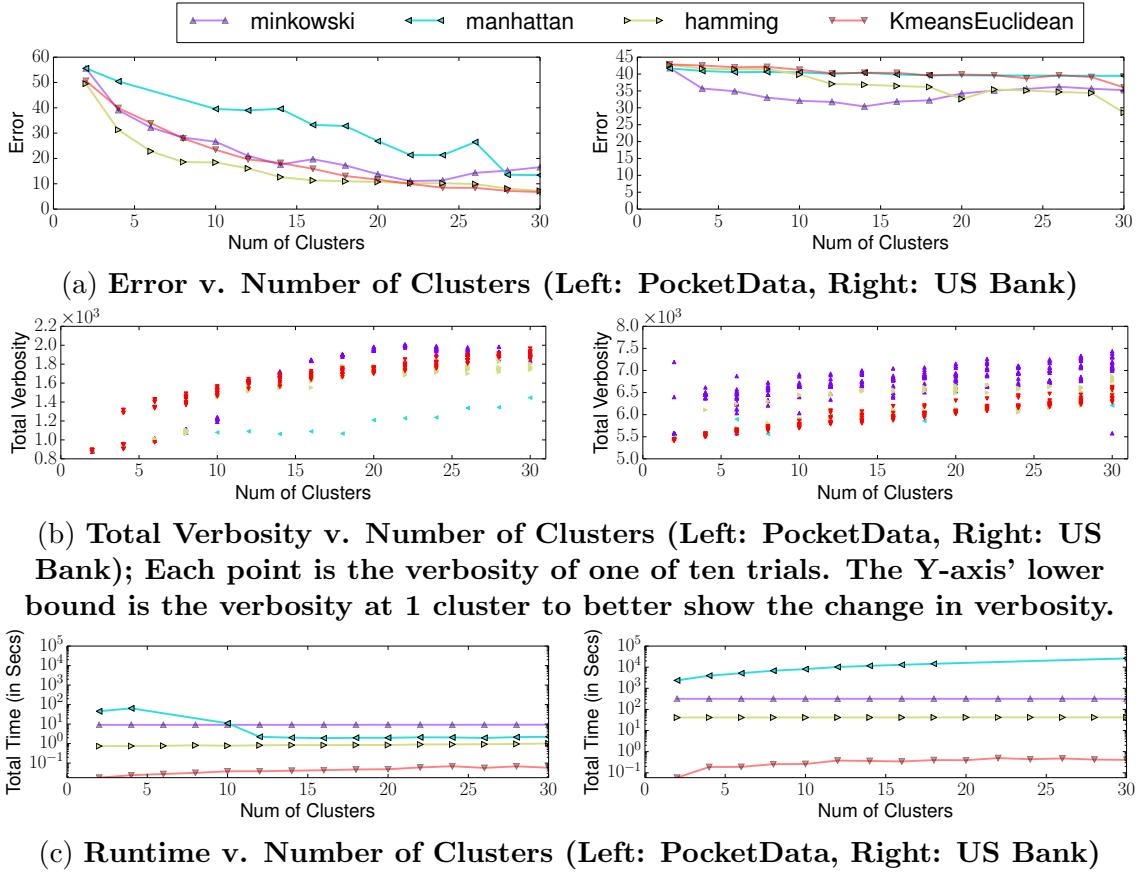


Figure 3.2: Clustering Schemes Comparison

we observe that a log (or log partition) uniquely determines its naive (or naive mixture) encoding. Thus the problem of searching for a naive mixture encoding reduces to searching for the corresponding log partitioning. We further observe that the Error of a naive mixture encoding is proportional to the diversity of the queries in the log being encoded: The more uniform the log (or partition), the lower the Error. Hence, the partitioning problem further reduces to clustering queries in the log by feature overlap. To identify a suitable clustering scheme, we next evaluate four commonly used clustering schemes with respect to their ability to create naive mixture encodings with low Error and Verbosity: (1) KMeans [22] with Euclidean distance (i.e., l_2 -norm) and Spectral Clustering [26] with (2) Manhattan (i.e., l_1 -

norm), (3) Minkowski (i.e., l_p -norm) with $p = 4$, and (4) Hamming distances¹.

Experiment Setup. Spectral and KMeans clustering algorithms are implemented by *sklearn* [40] in Python. We gradually increase K (i.e., the number of clusters) for each clustering scheme to mimic the process of continuously sub-clustering the log, tolerating higher Total Verbosity for lower Error. To reduce randomness in clustering, we run each of them 10 times for each K and averaging the Error of the resulting encodings. We used two datasets: “US Bank” and “PocketData”. We describe both datasets and the data preparation process in detail in Section 3.7. All results for our clustering experiments are shown in Figure 3.2.

3.6.1.1 Clustering

We next show that clustering is an effective way to consistently reduce Error, although no one clustering scheme is ideal for all three of Error, Verbosity, and runtime.

More clusters reduces Error. Figure 3.2a compares the relationship between the number of clusters (x-axis) and Error (y-axis), showing the varying rates of convergence to zero Error for each clustering scheme. We observe that adding more clusters does consistently reduce Error for both data sets, regardless of clustering algorithm or distance measure. We note that the US Bank dataset is significantly more diverse than the PocketData dataset, with respect to the total number of features (See Table 3.1) and that more than 30 clusters may be required for reaching near-zero Error. In general, Hamming distance converges faster than other distance measures on PocketData. Minkowski distance shows faster convergence rate than Hamming within 14 clusters on the US bank dataset.

Adding more clusters increases Verbosity. Figure 3.2b compares the relationship between the number of clusters (x-axis) and Verbosity (y-axis). We observe that Verbosity increases with the number of clusters. This is because when a partition is split, each feature common to both partitions increases the Verbosity by one.

¹We also evaluated Spectral Clustering with Euclidean, Chebyshev and Canberra distances; These did not perform better and we omit them in the interest of conciseness.

Hierarchical Clustering. The clustering schemes produce non-monotonic cluster assignments. That is, Error can occasionally grow as clusters are added (Figure 3.2a). An alternative is to use hierarchical clustering [22], which forces monotonic assignments and offers more dynamic control over the Error/Verbosity tradeoff.

Run Time Comparison. The total runtime (y-axis) in Figure 3.2c includes both distance matrix computation time (if any) and clustering time. Note the log-scale: K-Means is orders of magnitude faster than the others.

Take-Aways. For time-sensitive applications, KMeans algorithm is preferred to Spectral Clustering. With respect to distance measures, minkowski (i.e., l_p -norm) with $p = 4$ provides the best tradeoff between Error and runtime.

Visualizing Naive Mixture Encoding. As with normal pattern summaries, naive mixture summaries are also interpretable. For example a visualization like that of Figure 3.1a can be repeated, once for each cluster. For more details, see our accompanying technical report [45].

3.6.2 Approximating Log Statistics

Recall that our primary goal is estimating statistical properties. In particular, we are interested in counting the occurrences $\Gamma_\omega(L)$ (i.e., $p(Q \supseteq \mathbf{b}) \cdot |L|$) of some pattern ω in the log. Recall that a naive encoding $\ddot{\mathcal{E}}$ includes only single-feature patterns (i.e., patterns exactly encoding $p(X_i \geq x_i)$) and that the closed-form representation for the maximum entropy distribution $\bar{\rho}_{\ddot{\mathcal{E}}}$ arises by independence between features (i.e., $\bar{\rho}_{\ddot{\mathcal{E}}}(Q = \mathbf{q}) = \prod_i p(X_i = x_i)$). Similarly, we use the independence assumption to estimate:

$$est[\Gamma_\omega(L) \mid \ddot{\mathcal{E}}] = \bar{\rho}_{\ddot{\mathcal{E}}}(Q \supseteq \mathbf{b}) \cdot |L| = \prod_i p(X_i \geq x_i) \cdot |L|$$

This process trivially generalizes to naive pattern mixture encodings by mixing distributions. Specifically, given a set of partitions $L_1 \cup \dots \cup L_K = L$, the estimated counts for $\Gamma_\omega(L)$ under each individual partition L_i can be computed based on the

partition's naive encoding $\ddot{\mathcal{E}}_i$, and we then sum up the estimated counts in each partition:

$$est[\Gamma_\omega(L_i) | \ddot{\mathcal{E}}_1, \dots, \ddot{\mathcal{E}}_K] = \sum_{i \in [1, K]} est[\Gamma_\omega(L_i) | \ddot{\mathcal{E}}_i]$$

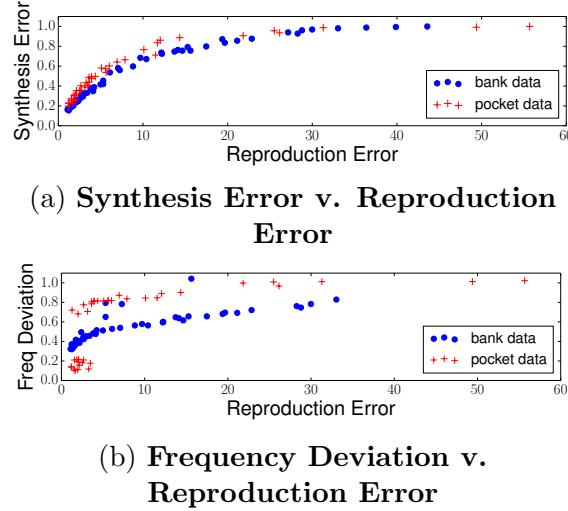


Figure 3.3: Effectiveness of Naive Mixture Encoding

3.6.3 Pattern Synthesis & Frequency Estimate

In this section, we empirically verify the effectiveness of naive mixture encodings in approximating log statistics from two related perspectives. The first perspective focuses on *synthesis error*. It measures whether patterns synthesized by the naive mixture encoding actually appear in the log. From the second perspective, we further investigate the *frequency deviation* of patterns contained in the log. This evaluates whether a naive mixture encoding computes the correct frequency for patterns of interest to client applications. Experimental results are shown in Figure 3.3. Both synthesis error and frequency deviation consistently decrease given more clusters. Furthermore, as we vary the number of clusters, both measures correlate with Reproduction Error.

Synthesis Error is measured by $1 - \frac{m}{n}$ where m out of n randomly synthesized patterns actually appear in the log. Intuitively, when synthesis error grows, it is

more likely that a pattern from the synthesized log will not appear in the original log (i.e., smaller values are better). Figure 3.3a shows synthesis error (y-axis) versus Reproduction Error (x-axis). The figure is generated by synthesizing $n = 10000$ patterns from each cluster of the log. Note that different values of n give similar observations. The overall synthesis error is measured by the average of synthesis errors for all clusters, weighted by the proportion of queries in each cluster.

Frequency Deviation is measured for a pattern by $\frac{|est - t|}{t}$ where t stands for true frequency of a pattern and est is the one estimated by the naive mixture encoding. Since frequency deviation is smaller when evaluated on a pattern contained in the other, as an alternative, we treat each distinct query in the log as a pattern and the frequency deviation on it will be the worst case for all patterns that it contains. Intuitively, this value captures the percentage error of frequency estimates (i.e., smaller values are better). For each cluster, we sum frequency deviations on all of its distinct queries and the final frequency deviation for the whole log is an weighted average (same as synthesis error) over all clusters. Figure 3.3b shows frequency deviation (y-axis) versus Reproduction Error (x-axis).

3.6.4 Naive Encoding Refinement

Naive mixture encodings can already achieve close to near-zero Error (Figure 3.2a), have low Total Verbosity, and admit efficiently computable log statistics $\Gamma_\omega(L)$. Doing so makes estimating statistics more computationally expensive. However, as a thought experiment we consider a hypothetical second pass to enrich naive mixture encodings with non-naive patterns. We start by considering the simpler problem of identifying the *individual* non-naive pattern that maximally reduces the Reproduction Error of a naive encoding.

Feature-Correlation Refinement. Recall that under naive encodings, we have a closed-form estimation $\bar{p}_{\mathcal{E}}(Q \supseteq \mathbf{b})$ of pattern frequencies $p(Q \supseteq \mathbf{b})$. We thus define the *feature-correlation* of pattern \mathbf{b} as the log-difference from its actual frequency to the estimate.

$$fc(\mathbf{b}, \ddot{\mathcal{E}}) = |\log(p(Q \supseteq \mathbf{b})) - \log(\bar{p}_{\mathcal{E}}(Q \supseteq \mathbf{b}))|$$

Intuitively, patterns with higher feature correlations carry more information content of the log that its naive encoding ignores, making them ideal candidates for addition to the naive encoding. For two patterns with the same feature-correlation, the one that occurs more frequently [19] will have greater impact on Reproduction Error. As a result, we compute an overall score for ranking individual patterns:

$$\text{corr_rank}(\mathbf{b}) = p(Q \supseteq \mathbf{b}) \cdot fc(\mathbf{b}, \ddot{\mathcal{E}})$$

We show in Section 3.7.1 that *corr_rank* closely correlates with Reproduction Error. That is, a higher *corr_rank* value indicates that a pattern produces a greater reduction in Reproduction Error if introduced into the naive encoding.

Pattern Diversification. In general, we would like to identify a *set* of patterns. The greedy approach that adds patterns one by one based on their ranking scores *corr_rank* is unreliable, as modifying the naive encoding invalidates the closed-form estimation $\bar{\rho}_{\mathcal{E}}(Q \supseteq \mathbf{b})$ that score *corr_rank* relies on. In other words, we can not sum up *corr_rank* scores of patterns in a set to rank its overall contribution to Reproduction Error reduction, as information content carried by patterns may overlap. To counter such overlap, or equivalently to *diversify* patterns, a search through the space of pattern-sets is needed. This type of diversification is commonly used in pattern mining applications, but can quickly become expensive. As we show experimentally in Section 3.7.2, the benefit of diversification is minimal.

3.7 Experiments

In this section, we design experiments to empirically (1) validate that Reproduction Error correlates with Deviation and (2) evaluate the effectiveness of LOGR compression.

We use two specific datasets in the experiment: (1) SQL query logs of the Google+ Android app extracted from the PocketData public dataset [27] and (2) SQL query logs that capture all query activity on the majority of databases at a major US bank over a period of approximately 19 hours. A summary of these two datasets is given in Table 3.1.

Table 3.1: **Summary of Data sets**

| Statistics | PocketData | US bank |
|---------------------------------|------------|---------|
| # Queries | 629582 | 1244243 |
| # Distinct queries | 605 | 188184 |
| # Distinct queries (w/o const) | 605 | 1712 |
| # Distinct conjunctive queries | 135 | 1494 |
| # Distinct re-writable queries | 605 | 1712 |
| Max query multiplicity | 48651 | 208742 |
| # Distinct features | 863 | 144708 |
| # Distinct features (w/o const) | 863 | 5290 |
| Average features per query | 14.78 | 16.56 |

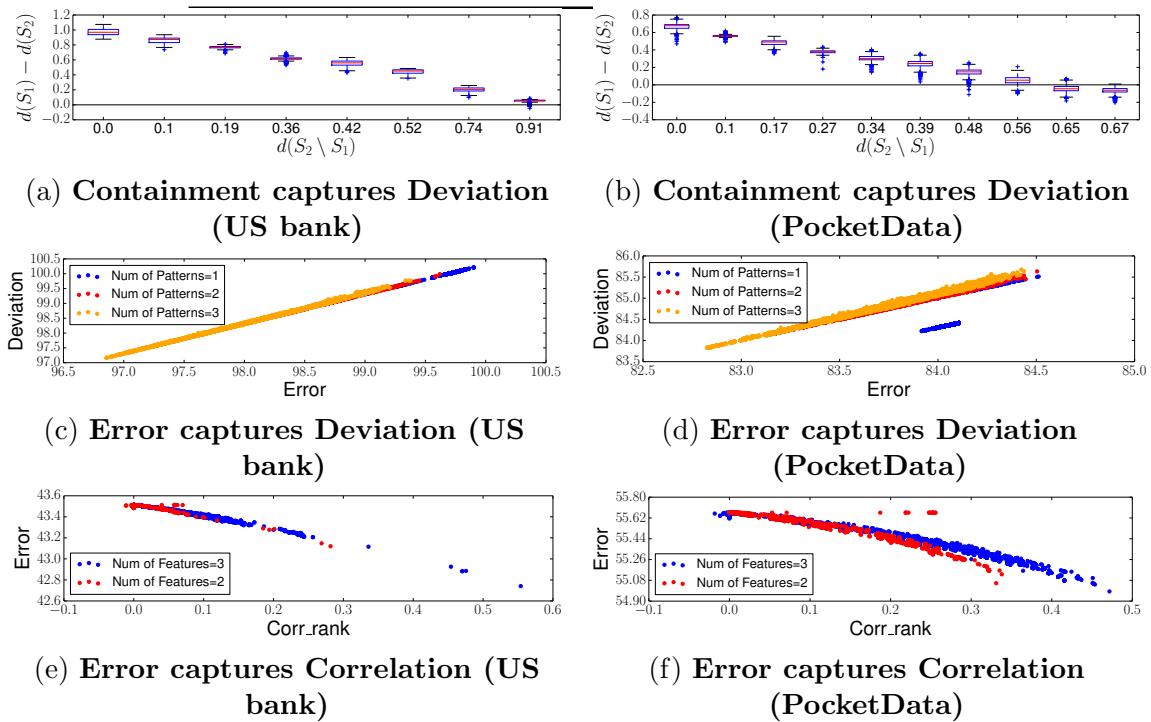


Figure 3.4: **Validating Reproduction Error**

The **PocketData-Google+** query log. The dataset consists of SQL logs that capture all database activities of 11 Android phones. We selected the Google+ ap-

plication for our study since it is one of the few applications where all users created a workload. This dataset is a stable workload of exclusively machine-generated queries.

The US bank query log. This log is an anonymized record of queries processed by multiple relational database servers at a major US bank [30] over a period of 19 hours. Of the nearly 73 million database operations captured, 58 million are not directly queries, but rather invocations of stored procedures. A further 13 million used non-standard SQL features not supported by our SQL parser. Of the remaining of the 2.3 million parsed SQL queries, we base our analysis on the 1.25 million conjunctive SELECT queries. This dataset can be characterized as a diverse workload of both machine- and human-generated queries.

Common Experiment Settings. Experiments were performed on a 2.8 GHz Intel Core i7 CPU with 16 GB 1600 MHz DDR3 memory and a SSD running macOS Sierra.

Constant Removal. A number of queries in the US bank query log differ only in hard-coded constant values. Table 3.1 shows the total number of queries, as well as the number of distinct queries if we ignore constants. By comparison, queries in PocketData all use JDBC parameters. For these experiments, we ignore constant values in queries.

Query Regularization. We apply query rewrite rules (same as [31]) to regularize queries into equivalent conjunctive forms, where possible. Table 3.1 shows that $\frac{135}{605}$ and $\frac{1494}{1712}$ of distinct queries are in conjunctive form for PocketData and US bank respectively. After regularization, all queries in both data sets can be either simplified into conjunctive queries or re-written into a UNION of conjunctive queries compatible with feature scheme of Aligon et al. [3].

Convex Optimization Solving. All convex optimization problems for measuring Reproduction Error and Deviation are solved by the *successive approximation heuristic* implemented by the CVX toolbox [18] with the Sedumi solver.

3.7.1 Validating Reproduction Error

In this section, we validate that Reproduction Error is a practical alternative to Deviation. In addition, we also offer measurements on its correlation with Deviation and score *corr_rank* in Section 3.6.4. As it is impractical to enumerate all possible encodings, we choose a subset of encodings for both datasets. Specifically, we first select all features with frequencies in the range [0.01, 0.99] and use these features to construct patterns. We then enumerate combinations of K (up to 3) patterns as our chosen encodings.

Containment Captures Deviation. Here we empirically verify that containment (Section 3.4.2) captures Deviation (i.e., $\mathcal{E}_1 \leq_{\Omega} \mathcal{E}_2 \rightarrow d(\mathcal{E}_1) \leq d(\mathcal{E}_2)$) to complete the chain of reasoning that Reproduction Error captures Deviation. Figures 3.4a and 3.4b show all pairs of encodings where $\mathcal{E}_2 \supset \mathcal{E}_1$. The y-axis shows the difference in Deviation values (i.e., $d(\mathcal{E}_2) - d(\mathcal{E}_1)$). Deviation $d(\mathcal{E})$ is approximated by drawing 1 million samples from the space $\Omega_{\mathcal{E}}$ induced by the encoding \mathcal{E} . For clarity, we bin pairs of encodings by the degree of overlap between them, measured by the Deviation of the set-difference $d(\mathcal{E}_2 \setminus \mathcal{E}_1)$; Higher $d(\mathcal{E}_2 \setminus \mathcal{E}_1)$ implies less overlap. Y-axis values are grouped into bins and visualized by boxplot where the boxes represent ranges within standard deviation and crosses are outliers. Intuitively, *points above zero* on the y-axis (i.e., $d(\mathcal{E}_2) - d(\mathcal{E}_1) > 0$) are pairs of encodings where the Deviation order agrees with containment order. This is the case for virtually all encoding pairs.

Additive Separability of Deviation. We also observe from Figures 3.4a and 3.4b that agreement between Deviation and containment order is correlated with overlap: More similar encodings are more likely to have agreement. Combined with Proposition 1, this shows first that for similar encodings, Reproduction Error is likely to be a reliable indicator of Deviation. This also suggests that Deviation is additively separable: The information loss (i.e., $d(\mathcal{E}_2) - d(\mathcal{E}_1)$) caused by excluding the encoding $\mathcal{E}_2 \setminus \mathcal{E}_1$ from \mathcal{E}_2 correlates with the quality (i.e., $d(\mathcal{E}_2 \setminus \mathcal{E}_1)$) of the encoding $\mathcal{E}_2 \setminus \mathcal{E}_1$ itself:

$$\mathcal{E}_2 \supset \mathcal{E}_1 \rightarrow d(\mathcal{E}_2) - d(\mathcal{E}_1) < 0 \quad \text{and} \quad d(\mathcal{E}_2 \setminus \mathcal{E}_1) \propto d(\mathcal{E}_2) - d(\mathcal{E}_1)$$

Error correlates with Deviation. As a supplement, Figures 3.4c and 3.4d empirically confirm that that Reproduction Error (x-axis) indeed closely correlates with Deviation (y-axis). Mirroring our findings above, correlation between them is tighter at lower Reproduction Error.

Error and Feature-Correlation. Figure 3.4e and 3.4f show the relationship between Reproduction Error (y-axis) and score $corr_rank$ (x-axis), as discussed in Section 3.6.4. Values of y-axis are Reproduction Error of the naive encodings extended by a non-naive pattern \mathbf{b} containing multiple features (up to 3 for illustrative purposes). One can observe that Reproduction Error of extended naive encodings almost linearly correlates with $corr_rank(\mathbf{b})$. In addition, one can also observe that $corr_rank$ becomes higher when the pattern \mathbf{b} encodes more correlated features.

3.7.2 Feature-Correlation Refinement

In this section, we design experiments serving two purposes: (1) Evaluating the potential reduction of Error from refining naive mixture encodings through state-of-the-art pattern-based summarizers, and (2) Evaluating whether we can replace naive mixture encodings by the encodings created from summarizers that we have plugged-in.

Experiment Setup. To serve both purposes, we construct pattern mixture encodings under three configurations: (1) Naive mixture encodings; (2) Pattern-based encodings and (3) Naive mixture encodings refined into pattern-based encodings. Naive mixture encodings are constructed by K-Means clustering. Pattern-based encodings are generated by two state-of-the-art pattern-based summarizers: (1) *Laserlight* [16] that summarizes multi-dimensional data in order to predict an augmented binary variable and (2) *MTV* [35] that aims at mining maximally informative patterns that summarize binary multi-dimensional data.

The experimental results are shown in Figure 3.5 that contains 3 sub-figures sharing the same x-axis, i.e., the number of clusters. Figure 3.5a compares the Error (y-axis) between naive mixture encodings and pattern mixture encodings that consist of patterns mined from *MTV* or *Laserlight*. Figure 3.5b evaluates the change in Error

(y-axis) through refining naive mixture encodings by adding patterns from *MTV* or *Laserlight*. Figure 3.5c compares the runtime (y-axis) between constructing naive mixture encodings and applying *MTV* or *Laserlight*. We only show the results for US bank query log as results for PocketData give similar observations.

3.7.2.1 Pattern-based vs Naive Mixture Encodings

Figure 3.5a and 3.5c suggest that naive mixture encodings outperform pattern-based encodings in two ways.

Reproduction Error. We observe from Figure 3.5a that the Reproduction Error of naive mixture encodings are orders of magnitude lower than pattern-based encodings generated by *Laserlight* or *MTV* alone.

Computation Efficiency. From Figure 3.5c we observe that the runtime of constructing naive mixture encodings is significantly lower than that of *Laserlight* and *MTV*.

The one way where pattern-based encodings outperform naive mixture encodings is in Total Verbosity. *Laserlight* and *MTV* produce encodings with significantly fewer patterns, as the naive mixture encoding requires at least one pattern for each feature (e.g., 5290 patterns in the US bank query log). Conversely, mining this number of patterns is computationally infeasible (Figure 3.5c).

3.7.2.2 Refining Naive Mixture Encodings

The experiment result is shown in Figure 3.5b. Note that we offset y-axis to show the change in Error. We observe from the figure that reduction of Error contributed by plugging-in pattern-based summarizers is small for both algorithms.

Dimensionality Restriction. For *Laserlight*, the observation is partially due to the fact that we only keep top 100 features (in terms of variability) of the data as its input, since *Laserlight* is implemented in PostgresSQL 9.1 which has a threshold of 100 arguments (one argument for each feature) that can be passed to a function.

Pattern Restriction. For *MTV*, this is due to a runtime error that limits us to 15 or less patterns. We refer the reader to Section 4.5 in [35] that explains the difficulty

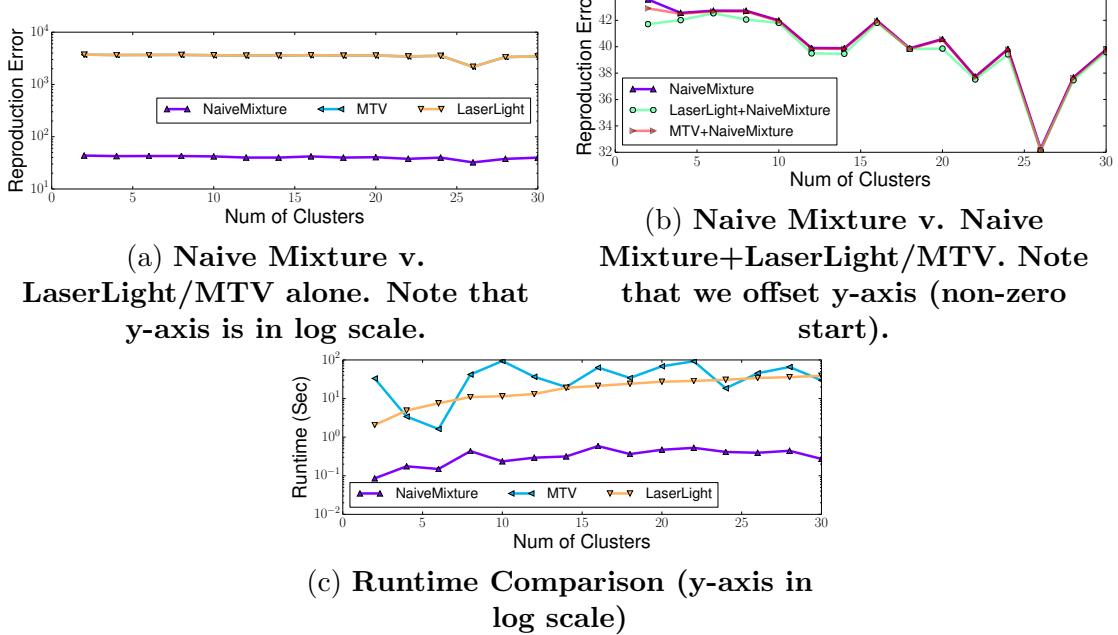


Figure 3.5: Feature-correlation refinement (US bank)

in inferring the maximum entropy distribution constrained by a large number of non-naive patterns.

3.8 Alternative Applications

To fairly evaluate *Laserlight* and *MTV*, we incorporate their own data sets and empirically evaluate them against *naive mixture encoding* under their own applications.

Data Sets. Specifically, we choose *Mushroom* data set used in *MTV* [35] which is obtained from FIMI dataset repository and U.S. Census data on Income or simply *Income* data set, which is downloaded from IPUMS-USA at <https://usa.ipums.org/usa/> and used in *Laserlight* [16]. The basic statistics of the data sets are given in Table 3.2.

3.8.1 Experiments

All experiments involving *Laserlight* and *MTV* will be evaluated under their own Error measures and data sets, unless otherwise stated. The experiments are organized

Table 3.2: Data Sets of Alternative Applications

| Statistics | Income | Mushroom |
|---------------------------------|------------|-----------|
| # Distinct data tuples | 777493 | 8124 |
| # Features per tuple | 9 | 21 |
| Feature Binary-valued? | no | no |
| # Distinct features | 783 | 95 |
| Binary Classification Feature | > 100,000? | Edibility |
| Assumed data tuple multiplicity | 1 | 1 |

as follows: First, we establish baselines by evaluating classical *Laserlight* and *MTV* on their original data; Then we show that classical *Laserlight* and *MTV* can be generalized to partitioned data and that the generalization improves on their Error measures and also runtime; At last, we compare their generalized versions with *naive mixture encoding* to show that *naive mixture encoding* is a reasonable alternative.

3.8.1.1 Error Measures

We first explain how *naive mixture encoding* is evaluated based on Error defined by *Laserlight* and *MTV*.

Evaluating Naive Encoding on Laserlight Error. Algorithm *Laserlight* summarizes data D which consists of feature vectors t augmented by some binary feature v . Denote the valuation of the binary feature v for each feature vector t as $v(t)$. The goal is to mine a summary encoding \mathcal{E} , which is a set of patterns contained in $t \in D$ that offer predictive power on $v(t)$. Denote the estimation (based on \mathcal{E}) of $v(t)$ as $u_{\mathcal{E}}(t) \in [0, 1]$, the *Laserlight* Error is measured by

$$\sum_t (v(t) \log(\frac{v(t)}{u_{\mathcal{E}}(t)}) + (1 - v(t)) \log(\frac{1 - v(t)}{1 - u_{\mathcal{E}}(t)}))$$

Since *naive encoding* $\tilde{\mathcal{E}}$ assumes feature independence, estimation of $v(t)$ is independent of t , namely $u_{\tilde{\mathcal{E}}}(t) = u_{\tilde{\mathcal{E}}} = |\{\tau | v(\tau) = 1, \tau \in D\}| / |D|$. Consequently, the

Laserlight Error of *naive encoding* is

$$-|D|(u_{\mathcal{E}} \log u_{\mathcal{E}} + (1 - u_{\mathcal{E}}) \log(1 - u_{\mathcal{E}}))$$

Evaluating Naive Encoding on MTV Error. Given binary feature vectors D , the *MTV* Error of encoding \mathcal{E} is

$$-|D|H(\bar{\rho}_{\mathcal{E}}) + 1/2|\mathcal{E}| \log |D|$$

where $H(\bar{\rho}_{\mathcal{E}})$ is the entropy of maximum entropy distribution $\bar{\rho}_{\mathcal{E}}$ defined in Section 3.4.1. The second term in *MTV* Error penalizes Verbosity of the encoding \mathcal{E} . Since naive encoding assumes feature independence, we can first compute entropy of the marginal distribution of each individual feature. Entropy $H(\bar{\rho}_{\mathcal{E}})$ is simply the sum of feature entropies.

Evaluating Naive Mixture Encoding. Evaluation of *naive encoding* can be generalized to *naive mixture* by taking a weighted average over resulting clusters (See Section 3.5.2).

3.8.1.2 Classical Laserlight and MTV

Establishing Baselines. To establish baselines, we evaluate *Laserlight* and *MTV* on their own data sets. The take-aways from related experiments are that (1) *naive encoding* is faster and more accurate than classical *Laserlight* and *MTV*; (2) the runtime increases superlinearly with the number of patterns mined from both *Laserlight* and *MTV*. For detailed experiment results, we refer the reader to [45].

Anti-correlation and Dimensionality Reduction. Recall in Section 3.7.2.2 that *Laserlight* is restricted to 100 features. For its own *Income* data set, *Laserlight* can be applied with its full set of 783 features. This is due to the prior knowledge that the 783 features belong to 9 groups. In each group, features are mutually anti-correlated which can be reduced to a single feature. Similarly, *Mushroom* data set can be reduced from 95 to 21 features.

3.8.1.3 Generalizing Laserlight and MTV

We generalize *Laserlight* and *MTV* on partitioned data by applying them on each cluster. We then combine Errors on all clusters by taking a weighted average, as described in Section 3.5.2. Depending on how many patterns are mined from each cluster, *Laserlight* and *MTV* can be generalized into two types: (1) The number of patterns mined from each cluster is scaled to be equal to Verbosity of the *naive encoding*; and (2) The total number of patterns mined from all clusters is fixed to a given number. We name the first type *Laserlight (MTV) Mixture Scaled*, which is comparable to *naive mixture encoding*. We name the second type *Laserlight (MTV) Mixture Fixed*, which is comparable to the classical *LaserLight (MTV)* algorithm.

Take-away. As the data is partitioned into more clusters, both runtime and Error of *Laserlight (MTV) Mixture Fixed* exponentially decrease. This observation can be potentially generalized to other pattern mining algorithms. For experiment details, we refer the reader to [45].

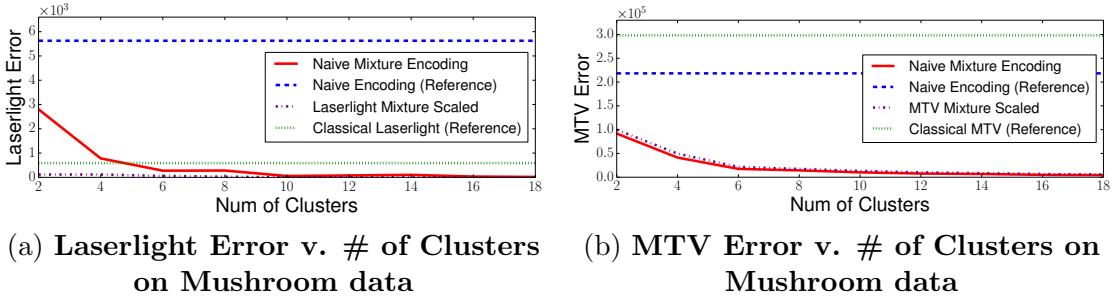


Figure 3.6: Naive Mixture v. Laserlight/MTV Mixture

3.8.1.4 Comparison with Naive Mixture Encoding

At last, we compare *Laserlight (MTV) Mixture Scaled* with *naive mixture encoding*. Note that it is time-consuming for *Laserlight* to mine the same number of patterns as *naive encoding* on *Income* data (See runtime analysis in [45]), we choose *Mushroom* data for *Laserlight Mixture Scaled* instead. The experiment results are given in Figure 3.6. The x-axis for all sub-figures in Figure 3.6 represents the number of clusters and the y-axes stands for *Laserlight* and *MTV* Error respectively. We

incorporate baselines (i.e., *naive encoding*, classical *Laserlight* and *MTV*) as reference lines in Figure 3.6a and 3.6b respectively. We also experienced a limitation of 15 patterns for *MTV*. Hence the comparison between *MTV Mixture Scaled* and *naive mixture encoding* is not strictly on equal footing as *MTV Mixture Scaled* is not able to reach the same Total Verbosity as *naive mixture encoding*. Note that their difference in Verbosity is mitigated by the fact that *MTV* Error measure penalizes encoding Verbosity.

Figure 3.6a shows that both *naive mixture encoding* and *Laserlight Mixture Scaled* have lower Error than their baselines. In addition, *Laserlight Mixture Scaled* has lower Error than *naive mixture encoding* when the number of clusters is less than 4 and they become close after 6 clusters. In other words, *Laserlight* is more accurate on lightly partitioned data. As the data is further partitioned, clusters become easier to summarize, and *naive encoding* becomes more similar to *Laserlight*. Figure 3.6b shows that *naive mixture encoding* marginally outperforms *MTV Mixture Scaled*.

Take-away. *Naive mixture encoding* is faster and has similar (lower) Error than *Laserlight (MTV) Mixture Scaled*.

3.9 Related Work

We aim at compressing query logs for accurately and efficiently computing workload statistics. Before the discussion of compression, we first review usecases and related work for workload analysis.

3.9.1 Workload Analysis

Existing approaches related to workload analysis usually aim at specific tasks like query recommendation [36, 17, 28, 46, 3], performance optimization [7, 11], outlier detection [25] or visual analysis [34].

Query Recommendation. This task aims at tracking historical querying behavior and generating query recommendations. Related approaches [36, 28] flatten a query *abstract syntax tree* as a set of *fragments* [36] or *snippets* [28]. User profiles are then

built by grouping and summarizing queries of specific users in order to make personalized recommendation. Under OLAP systems, profiles are also built for workloads of similar OLAP sessions [3].

Performance Optimization. Index selection [13, 15] and materialized view selection [2, 7, 11] are typical performance optimization tasks. The configuration search space is usually large, but can be reduced with appropriate summaries.

Outlier Detection. Kamra *et al.* [25] aim at detecting anomalous behavior of queries in the log by summarizing query logs into profiles of normal user behavior.

Visual Analysis. Makiyama *et al.* [34] provide a set of visualizations that facilitate further workload analysis on Sloan Digital Sky Survey (SDSS) dataset. QueryScope [20] aims at finding better tuning opportunities by helping human experts to identify patterns shared among queries.

In these approaches, queries are commonly encoded as feature vectors or bit-maps where a bit array is mapped to a list of features with 1 in a position if the corresponding feature appears in the query and 0 otherwise. Workloads under the bit-map encoding must then be compressed before they can be efficiently queried or visualized for analysis.

3.9.2 Workload Compression Schemes

Run-length Encoding. *Run-length encoding (RLE)* is a loss-less compression scheme commonly used in *Inverted Index Compression* [43, 49] and *Column-Oriented Compression* [1]. RLE-based compression algorithms include but not limited to: Byte-aligned Bitmap Code (BBC) used in Oracle systems [6], Word-aligned Hybrid (WAH) [44] and many others [37, 4, 5]. In general, RLE-based methods focus on column-wise compression and requires additional heavyweight inference on frequencies of cross-column (i.e., row-wise) patterns used for workload analysis.

Lempel-Ziv Encoding. Lempel-Ziv [47, 48] is the loss-less compression algorithm used by gzip. It takes variable sized patterns (row-wise in our case) and replaces them with fixed length codes, in contrast to Huffman encoding [21]. Lempel-Ziv encoding

does not require knowledge about pattern frequencies in advance and builds the pattern dictionary dynamically. There are many other similar schemes for compressing files represented as sequential bit-maps, e.g. [41].

Dictionary Encoding. *Dictionary encoding* is a more general form of Lempel-Ziv. It has the advantage that patterns with frequencies stored in the dictionary can be interpreted as workloads statistics useful for analysis. In this paper, we extend dictionary encoding and focus on using a dictionary to infer frequencies of patterns not in it. Mampaey *et al.* proposed *MTV* algorithm [35] that finds the dictionary (of given size) having optimal *Bayesian Information Criterion(BIC)* score. Gebaly *et al.* proposed *Laserlight* algorithm [16] that builds a pattern dictionary for correctly inferring the truth-value of some augmented binary feature.

Generative Models. A generative model is a lossy compressed representation of the original log. Typical generative models are *probabilistic topic models* [8, 42] and *noisy-channel* model [29]. Generative models can infer pattern frequencies but they lack a model-independent measure for efficiently evaluating overall inference accuracy.

Matrix Decomposition. Matrix decomposition methods including Principal Component Analysis (PCA) [24] and Non-negative matrix factorization (NMF) [33] offer lossy data compression. But the resulting matrices after decomposition are not suited for inferring workload statistics.

3.10 Conclusions

In this paper, we introduced the problem of log compression and defined a family of pattern-based log encodings. We precisely characterized the information content of logs and offered three principled and one practical measures of encoding quality: Verbosity, Ambiguity, Deviation and Reproduction Error. To reduce the search space of pattern-based encodings, we introduced the idea of log partitioning, which induces the family of pattern mixture as well as its simplified form: naive mixture encodings. Finally, we experimentally showed that naive mixture encodings

are more informative and can be constructed more efficiently than state-of-the-art pattern-based summarization techniques. We expect that making accurate and efficient inference on pattern frequencies will enable a range of more powerful database tuning and intrusion detection systems.

3.11 Future Work

Multiplicity-aware clustering. As the number of feature vectors can be millions or more, practically we only keep *distinct* feature vectors as input of clustering schemes. We can store feature vector frequencies in a separate column called *multiplicities*. A multiplicity-ignorant clustering scheme assumes a uniform distribution of queries in the log. However, query distributions $p(Q)$ of production database logs are usually skewed. For example, routine queries repeat themselves overwhelmingly in the log but contribute to a minority of distinct queries. We plan to improve naive mixture encodings by exploring *multiplicity-aware* clustering schemes such that distinct feature vectors can be clustered *as if they have been replicated*. The use of mixture models for summarization has potential implications for work on pattern mining; As we show, existing techniques can be substantially improved both in runtime and Error.

Feature Clustering. For the usecase of materialized view selection, computing pattern frequencies may not be enough. We may need to summarize a query log as a limited set of *basis* views such that queries in the log can be represented by a simple join of a subset of basis views. Capturing basis views is not only relevant to data tuning tasks, but also facilitates human inspection of workloads in the log. To achieve the goal, in addition to partitioning queries into separate workload clusters, for each cluster we need to further partition its features into separate clusters where each cluster is equivalent to a *basis view*.

3.12 Acknowledgements

The authors would like to thank the anonymous reviewers for their helpful comments and contributions. This work was supported by NSF Awards SaTC-1409551 and IIS-1750460. The conclusions and opinions in this work are solely those of the authors and do not represent the views of the National Science Foundation.

Chapter 4

Application: Interactive Semi-Structured Schema Design

Abstract

Ad-hoc data models like Json make it easy to evolve schemas and to multiplex different data-types into a single stream. This flexibility makes Json great for generating data, but also makes it much harder to query, ingest into a database, and index. In this paper, we explore the first step of Json data loading: schema design. Specifically, we consider the challenge of designing schemas for existing Json datasets as an *interactive* problem. We present **SCHEMADRILL**, a roll-up/drill-down style interface for exploring collections of Json records. **SCHEMADRILL** helps users to visualize the collection, identify relevant fragments, and map it down into one or more flat, relational schemas. We describe and evaluate two key components of **SCHEMADRILL**: (1) A summary schema representation that significantly reduces the complexity of JSON schemas without a meaningful reduction in information content, and (2) A collection of schema visualizations that help users to qualitatively survey variability amongst different schemas in the collection.

4.1 Introduction

Semi-structured formats like Json allow users to design schemas on-the-fly, as data is generated. For example, adding a new attribute to the output of a system logger does not break backwards compatibility with existing data. This flexibility facilitates the addition of new features and enables low-overhead adaptation of data-generating processes. However, because the data does not have a consistent underlying schema, it can be harder (and slower) to explore than simple tabular data. The logic of each and every query must explicitly account for variations in the schema like missing attributes. Performance also suffers, as there is no one physical data representation that is ideal for all schemas.

To address these problems, a variety of techniques [?, ?, ?, ?, ?] have arisen to generate schemas after-the-fact. The goal of these *semi-structured schema discovery* (S^3D) techniques is to propose a schema for collections of Json records. A common approach to this problem is to bind the Json records to a normalized relational representation, or in other words, to derive a set of flat *views* over the hierarchical Json data.

Existing automated approaches to this problem (e.g., [?, ?]) operate in a single-pass: They propose a schema and consider their job done. Unfortunately these techniques also rely heavily on general heuristics to select from among a set of schema design choices, as a clear declarative specification of a domain would be tantamount to having the schema already. To supplement domain-agnostic heuristics with feedback from domain experts, we propose a new *iterative and interactive* approach to S^3D called **SCHEMADRILL**.

SCHEMADRILL provides a OLAP-style interface (with analogs to roll-up, drill-down, and slice+dice) specialized for exploring collections of Json records. Every state of this interface corresponds to a relational view defined over the Json data. When ready, this view can be exported to a classical RDBMS or similar tool for simplified, more efficient data access. In this paper, we explore several design options for **SCHEMADRILL**, and discuss how each interacts with the challenges of S^3D .

4.1.1 Extracting Relational Entities

The first class of challenges we address involve the nuts and bolts of mapping hierarchical Json schemas to flat relational entities. Fundamentally, this involves a combination of three relational operators: Projection, Selection, and Unnesting.

Projecting Attributes. Json schema discovery can, naively, be thought of as a form of schema normalization [?], where each distinct path in a record is treated as its own attribute. Entities then, are simply groups of attributes projected out of the Json data, and the S³D problem reduces to finding such groups (e.g., by using Functional Dependencies [?]).

Selecting Records. This naive approach fails in situations where the collection of Json records is a mix of different entity types that share properties. As a simple example, Twitter streams mix three entity types: tweets, retweets, and deleted tweets. Although each entity appears in distinct records, they share attributes in common. Hence, entity extraction is not just normalization in the classical sense of partitioning attributes, but rather also a matter of partitioning records by content.

Collapsing Nested Collections. Json specifies two collection types: Arrays and Objects. Typically the former is used for encoding nested collections and the latter for encoding tuples with named attributes. However, this is not a strict requirement. For example, latitude and longitude are often encoded as a 2-element array. Conversely, in some data sets, objects are used as way to index collections by named keys rather than by positions. Hence, simple type analysis can not distinguish between the two cases. This is problematic because treating a collection as a tuple creates an explosion of attributes that make classical normalization techniques incredibly expensive.

4.1.2 Human-Scale S³D

Even in settings where Json data is comparatively well behaved, it is common for it to have dozens, or even hundreds of attributes per record. Similarly, individual Json records can be built from any of the hundreds or thousands (or more) different

permutations of the full set of attributes used across the entire collection. Bringing this information down to human scale requires simultaneously simplifying and summarizing.

Summarization. For the purposes of entity construction, the full set of attributes is often unnecessary. It is often possible to collapse multiple attributes together, or express attributes as equivalent alternatives. As an example, an address might consist of four distinct attributes city, zip code, street, and number when it could conceptually be expressed as just one.

Visualization. In addition to simplifying the underlying problem, it is also useful to give users a coarse “top-down” view of the schema process. Specifically, users need to (1) be able to see patterns of structural similarity between distinct schemas, and (2) understand how much variation exists in the data set as-is.

Iteration. By combining straightforward summarization and data visualization techniques, SCHEMADRILL helps users to quickly identify natural clusters of records and attributes that represent relational entities. SCHEMADRILL facilitates an *iterative* schema design process to allow human experts to better evaluate whether structures in the data indicate conceptual relationships between records or attributes, or are merely data artifacts.

4.1.3 Overview

Figure 4.1 shows the prototype interface of SCHEMADRILL. The pane on the left, discussed in Section 4.2, shows the schema of the currently selected JSON view, highlighting attributes and groups based on relevance. The pane on the right, discussed in Section 4.3, provides a top-down visual sketch of schemas in the currently selected JSON data, and allows users to interactively filter out parts of it. Finally in Section 4.3, we show examples of how SCHEMADRILL facilitates incremental, iterative exploration and mapping of JSON schemas.

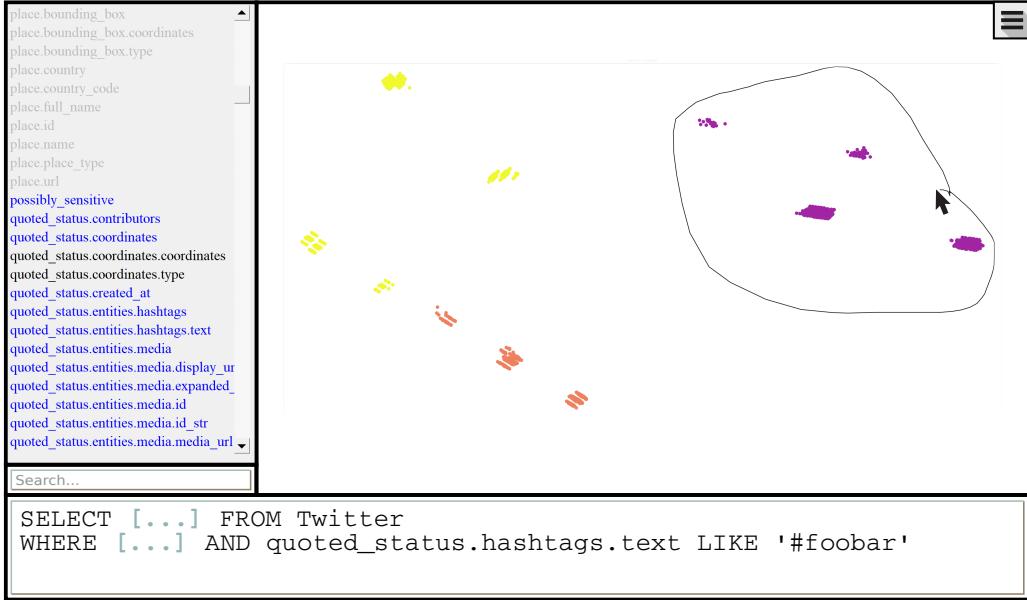


Figure 4.1: Prototype user interface for SCHEMADRILL.

4.2 Summarization

The first component of SCHEMADRILL, the schema pane, shows the relational schema of the extracted view. Initially, this schema consists of one attribute for every path in the Json collection being summarized. Attributes may be deleted or restored, and sets of attributes may be unified.

The core challenge behind implementing this pane is that, depending on data set, the schema could consist of hundreds or thousands of attributes. This can be overwhelming for users who just wants to find account profiles appearing in a Twitter stream. To mitigate this problem, SCHEMADRILL presents the schema in a summarized form.

Specifically, attributes are grouped based on both correlations and anti-correlations between them. Groups of *correlated* attributes, or those that frequently co-occur in Json records are likely to be part of a common structure. Similarly, groups of *anti-correlated* attributes, or those that rarely co-occur in Json records are likely to represent alternatives (e.g., a Street Address vs GPS coordinates). We use correlations

and anti-correlations between attributes to compact the schema for representation. Before describing the summary itself, we first formalize the problem.

4.2.1 Data Model

A Json object is an order-independent mapping from keys to values. A key is a unique identifier of a Json object, typically a string. A value may be atomic or complex. Atomic values in Json may be integers, reals, strings, or booleans. A complex value is either a nested object, or an *array*, an indexed list of values. For simplicity, we model arrays as objects by treating array indexes as keys. A Json record may be any value type, but for simplicity of exposition we will assume that all records are objects.

Example 5. *A fragment of Twitter Tweet encoded as Json*

```
"tweet": {
    "text": "#SIGMOD2018 in Houston this year",
    "user": {
        "name": "Alice Smith", "id": 12345,
        "friends_count": 1023, ...
    },
    "retweeted_status": {
        "tweet": {
            "user": { ... }, "entities": { ... },
            "place": { ... }, "extended_entities": { ... }
        }, ...
    },
    "place": { ... }, "extended_entities": { ... },
    "entities": { "hashtags": [ "SIGMOD2018" ], ... },
    ...
}
```

Json objects are typically viewed as trees with atomic values at the leaves, complex values as inner nodes, and edges labeled with the keys of each child. Our goal is to identify commonalities in the structure of this tree across multiple Json records in a collection. To capture the structure, we define the *schema* S of a Json record as the record with all leaf values replaced by a constant \perp . A *path* P is a sequence of keys

$P = (k_0, \dots, k_N)$. For convenience, we will write paths using dots to separate keys (e.g., `tweet.text`). We say that a path appears in a schema (denoted $P \in S$) if it is possible to start at the root of S and traverse edges in order. If the value reached by this traversal is \perp , we say that P is a terminal path of S (denoted $P \perp S$).

4.2.2 Paths as Attributes

Ultimately, our goal is to create a flat, relational representation suitable for use with an entire collection of Json records. The first step to reaching this goal is to flatten individual Json schemas into collections of attributes. We begin with a naive translation where each attribute corresponds to one terminal path in the schema. We write S^\perp to denote the path set, or relational schema of Json schema S , defined as: $S^\perp = \{ P \mid P \perp S \}$ Since keys are unique, commutative and associative, this representation is interchangeable¹ with the tree representation. Hence, when clear from context we will abuse syntax, using S to denote both a schema and its path set.

Example 6. *The path set of the Json object from Example 5 includes the paths:*

`1. tweet.text 2. tweet.user.friends_count 3. tweet.user.id 4. tweet.entities.hashtags.[0]`

Each terminal appears in the set. Note in particular that single element of the array at `tweet.entities.hashtags` is assigned the key [0].

Path sets make it possible to consider containment relationships between schemas. We say that S_1 is contained in S_2 iff $S_1^\perp \subseteq S_2^\perp$.

4.2.3 Schema Collections

We now return to our main goal, summarizing the schemas of collections of Json records. The starting point for this process is the schemas themselves. Given a collection of Json records, we can extract the set of schemas $\{ S_1, \dots, S_N \}$ of records in the collection, which we call the *source schemas*. One existing technique for summarizing these records, used by Oracle's JSON Data Guides [?, ?], is to simply

¹modulo empty objects or arrays

present the set of all paths that appear anywhere in this collection. We call this the *joint schema* \mathbb{S} :

$$\mathbb{S}^\perp \stackrel{\text{def}}{=} \bigcup_i S_i^\perp$$

Observe that, by definition, each of the source schemas is contained in the joint schema. The joint schema mirrors existing schemes for relational access to JSON data like for example. However, the joint schema can still be very large, with hundreds, thousands, or even tens of thousands of columns². To summarize them we need an even more compact encoding for sets of schemas.

A Schema Algebra. As a basis for compacting schema sets, we define a simple algebra. Recall that we are particularly interested in summarizing cooccurrence and anti-cooccurrence relationships between attributes.

$$\mathbf{A} := \mathbf{P} \mid \emptyset \mid \mathbf{A} \wedge \mathbf{A} \mid \mathbf{A} \vee \mathbf{A}$$

Expressions in the algebra construct sets of schemas from individual attributes. There are two types of leaf expressions in the algebra: A single terminal path P represents a singleton schema ($\{\{P\}\}$), while \emptyset denotes a set containing no schemas $\{\}$. Disjunction acts as union, merging its two input sets:

$$\{S_1, \dots, S_N\} \vee \{S'_1, \dots, S'_M\} \stackrel{\text{def}}{=} \{S_1, \dots, S_N, S'_1, \dots, S'_M\}$$

Disjunction models anti-correlation: The resulting schema set is effectively a collection of schema alternatives. For example, $P_1 \vee P_2$ indicates two alternative schemas: $\{P_1\}$ or $\{P_2\}$. Conjunction combines schema sets by cartesian product:

$$\{S_1, \dots, S_N\} \wedge \{S'_1, \dots, S'_M\} \stackrel{\text{def}}{=} \{S_i \cup S'_j \mid i \in [1, N], j \in [1, M]\}$$

Conjunction models correlations: The resulting schema set mandates that exactly one option from the left-hand-side and one option from the right-hand-side be present. On singleton inputs, the result is also a singleton. For example, $P_1 \wedge P_2$ is a single schema that includes both P_1 and P_2 . For inputs larger than one element,

²One dataset [?] achieved 2.4 thousand paths through nested collections of objects.

the conjunction requires one choice from each input. For example, $(P_1 \vee P_2) \wedge (P_3 \vee P_4)$ is the set of all schemas consisting of one of P_1 or P_2 , and also one of P_3 or P_4 .

Although we omit the proofs for conciseness, both \wedge and \vee are commutative and associative, and \vee distributes over \wedge ³. For conciseness, we use the following syntactic conventions: (1) When clear from context, a schema S denotes its own singleton set, and (2) We write $P_1 P_2$ to denote $P_1 \wedge P_2$.

This schema algebra gives us a range of ways to represent schema sets. At one extreme, the set of source schemas arrives in what is effectively disjunctive normal form (DNF). One schema may be expressed as a conjunction of its elements, and the full set of source schemas can be constructed by disjunction. For example, the source schemas $\{P_1, P_2\}$ and $\{P_2, P_3\}$ may be represented in the algebra as $P_1 P_2 \vee P_2 P_3$.

At the other extreme, the joint schema is a superset of all of the source schemas. It too can be thought of as a schema set, albeit one that loses information about which attributes appear in which schemas. Hence, this joint schema set may be defined as the power set of all attributes in the joint schema.

$$2^S = \bigwedge_{P \in S} (P \vee \emptyset)$$

Observe that at a minimum, each of the source schemas must appear in this schema set ($S_i \in 2^S$). However many other schemas appear in the resulting schema set as well.

4.2.4 Summarizing Schema Collections

These two extreme representations (the raw source schemas and the joint schema set) are bad, but for subtly different reasons. In both cases, the representation is too verbose. In the former case boseness stems from redundancy, with significant overlap in variables between the source schemas. Conversely in the latter case it stems from imprecision, as the schema set encompasses schemas that do not appear in the source schemas. Of the two, the latter is more compact, in particular because each

³To be precise, the structure $\langle \{\{\mathbf{P}\}\}, \vee, \wedge, \emptyset, \{\{\}\} \rangle$ is a semiring.

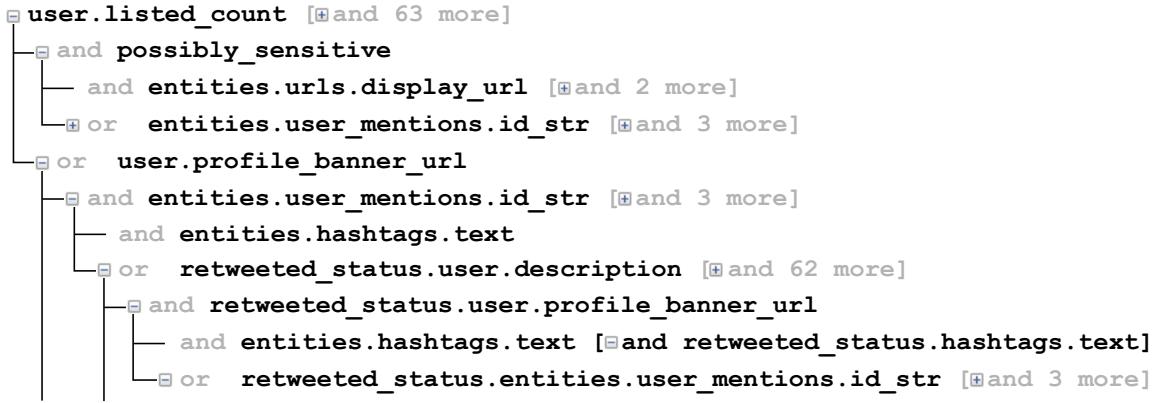


Figure 4.2: FP-Tree based schema summaries

attribute appears no more than once. This is a distinct representational advantage because the joint schema can be displayed simply as a list.

We would like to preserve this only-once property. Our aim then, is to derive an algebraic expression (1) in which each attribute appears exactly once, and (2) that is as tight a fit to the original source schema set as possible. Ultimately, this problem reduces to polynomial factorization and the discovery of read-once formulas [?], a problem that is, in general, worse than P-time [?]. Hence, for this paper, approximations are required. We consider two approaches and allow the user to select the most appropriate one for their needs. The first is based on Frequent Pattern Trees [?] (FPTrees), a data structure commonly used for frequent pattern mining. The second is to limit our search to read-once conjunctions of disjunctions of conjunctions, a form we call RCDC.

FP-Tree Summaries. An FP-Tree is a trie-like data structure that makes it easier to identify common patterns in a query. Every edge in the tree represents the inclusion of one feature, or in our case one attribute. Hence, every node in the tree corresponds to a set of paths (obtained by traversal from the root), and every leaf corresponds to one source schema. We observe that every node in an FP tree corresponds to a disjunction: For a node with 3 children, each subtree represents a different branch. Similarly, every edge corresponds to a conjunction with

a singleton. Although the resulting tree may duplicate some attributes, duplications are minimized [?].

Example 7. *Figure 4.2 illustrates a schema summary based on FP-Trees. Sequences of nodes with a single child are collapsed into single rows of the display (e.g., `user.listed_count` and 63 immediate descendants). A toggle switch allows these entities to be displayed to the user, if desired. Every level of the tree represents a set of alternatives. For example, `possibly_sensitive` never co-occurs with `user.profile_banner_url`.*

RCDC Summaries. Our second visualization is based on correlations and anticorrelations. To construct this visualization, we begin with the joint summary. Recall that the joint summary has the form

$$(P_1 \vee \emptyset) (P_2 \vee \emptyset) (P_3 \vee \emptyset) (P_4 \vee \emptyset) \dots$$

We create a covariance matrix based on the probability of two attributes co-occurring in the schemas of one of our input Json records. Using this covariance matrix, hierarchical clustering [?], and a user-controlled threshold on the covariance, we cluster the attributes by parenthesizing. For example, clustering might group P_1 with P_2 and likewise P_3 with P_4 . We can rewrite this formula as:

$$\approx (P_1 P_2 \vee \emptyset) (P_3 P_4 \vee \emptyset) \dots$$

Observe that this formula omits schemas that the original formula captures (e.g., any schema including P_1 but not P_2). However, because clustering ensures that attributes within a group co-occur frequently, there are comparatively few such schemas.

We next repeat the process with a new covariance matrix built using the frequency of co-occurrence of *groups* (like $P_1 P_2$). As before, we create clusters, but this time we cluster based on extremely negative co-variances. Hence, members of the resulting clusters are unlikely to co-occur. Continuing the example, let us assume that $P_1 P_2$ and $P_3 P_4$ are highly anti-correlated. Approximating and simplifying, we get an expression in RCDC form.

$$\approx (P_1 P_2 \vee P_3 P_4 \vee \emptyset) \dots$$

```

user.listed_count [and 63 more]
entities.user_mentions.id_str [and 3 more] [or 3 more]
possibly_sensitive [ ]
    or user.profile_banner_url
retweeted_status.user.description [and 66 more]

```

Figure 4.3: RCDC based schema summaries

As with the FP-Tree display, we use counts and an example attribute as a summary name for the group, and a toggle button to allow users to expand the group along either the OR or AND axes.

4.3 Visualization

Even within a mostly standardized collection of records like exported Twitter or Yelp data or production system logs, it is possible to find a range of schema usage patterns. Grouping by [anti]-cooccurrence is one step towards helping users understand these usage patterns, but is insufficient for three reasons: 1. Conceptually distinct fragments of the schema may share attributes in common (e.g., delete tweet records share attributes in common with tweet records). 2. Even if they do not co-occur, certain [groups of] attributes may be correlated (e.g., due to mobile phones, tweets with photos are also often geotagged). 3. There is no general way to differentiate Json objects and arrays being used to represent collections from those being used to represent structures (e.g., twitter stores geographical coordinates as a 2-element array). The second part of the SCHEMADRILL interface addresses these issues by presenting top-down visual surveys of the schema. These surveys help users to quickly assess variations in schema usage across the collection, to identify related schema structures, and to “drill down” into finer grained relationships.

4.3.1 Schema Segmentation

Specifically, we want to help the user to focus on particular parts of the joint schema; We want to allow the user to filter out, or segment the schema based on certain



Figure 4.4: Segmentation breaks up schema representations into manageable chunks.

required attributes that we call *subschemas*. We define a subschema s as a set of attributes, where s is contained in one or more source schemas. Further, the s -segment of source schemas S_1, \dots, S_N to be the subset that contain s :

$$\text{segment}(s) \stackrel{\text{def}}{=} \{ S_i \mid i \in [1, N] \wedge s \subseteq S_i \}$$

We are specifically interested in visual representations that can help users to identify subschemas of interest. By then focusing solely on the segments defined by these subschemas can significantly reduce the complexity of the schema design problem, as illustrated in Figure 4.4. Figures 4.4a illustrates the full schema summary as a tree, while Figure 4.4b shows a partial summary identified by the user using the lasso tool we describe shortly.

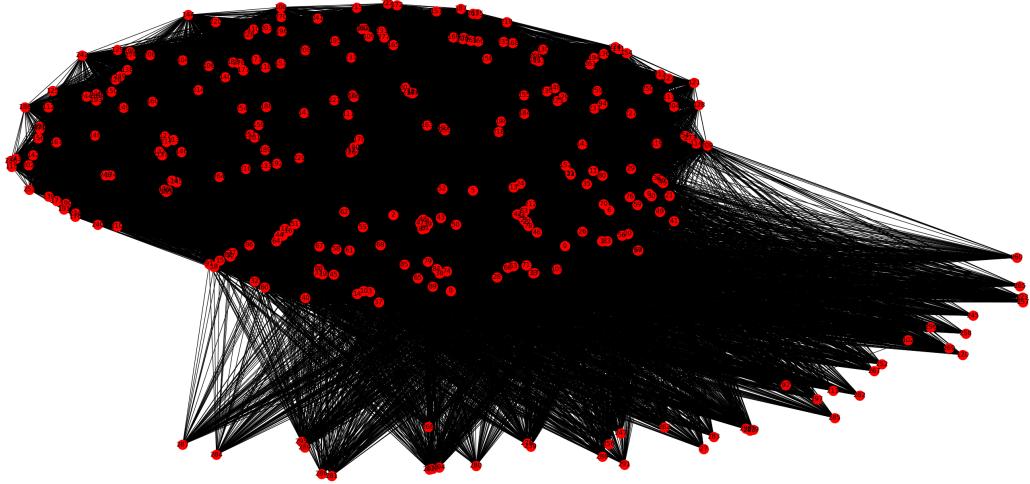


Figure 4.5: Covariance Cloud for the full Yelp dataset.

Covariance Clouds. Our second visual representation, also like schema summaries, uses correlations and anti-correlations to communicate subschemas of interest. To generate a covariance cloud, we create a covariance matrix from the source schemas, using the appearance of each attribute as a variable. Based on a user-controllable threshold, we then construct a graph from the covariance matrix with every attribute as one node, and every covariance exceeding the threshold as an edge. The graph is then displayed to the user as a cloud using standard force-based layout techniques (e.g., those used by GraphViz [?]). Cliques in the graph represent commonly co-occurring subschemas that might form segments of interest. This includes every conjunctive group identified in the schema summary. However, unlike the schema summary, this visual representation more effectively captures subschemas with attributes in common.

KNN-PCA Clouds. While the first visualizaton works on simple schemas, we found that on more complex Json data like Twitter streams [?], or the Yelp open dataset [?] there were too many inter-attribute relationships, and the resulting visu-alizatons were noisy. An approach we settled on is a mixture of Principle Component

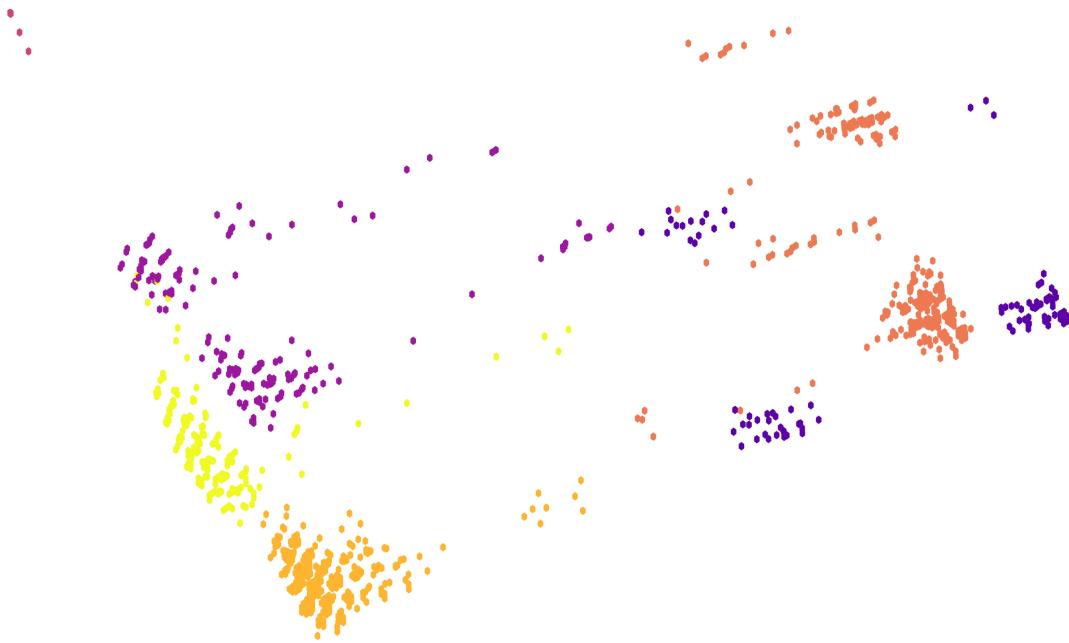


Figure 4.6: KNN-PCA cloud with $K=6$ on the Yelp dataset.

Analysis (PCA) and K Nearest Neighbor clustering (KNN). As before, we treat each source schema as a feature vector with each attribute representing one feature. We then use PCA to plot our source schemas in two-dimensions. The resulting visualization illustrates relationships between source schemas, with greater distances in the visualization representing (approximately) more differences between the schemas. Hence, clustered groupings of schemas represent potentially interesting sub-schemas.

A key limitation with this visualization is that for more complex datasets the somewhat arbitrary choice of 2 dimensions can be too low. Conversely adding more dimensions directly through PCA makes the visualization more complicated and hard to follow. To mitigate these limitations, we use K-Nearest Neighbors (KNN) to colorize the PCA Cloud. In addition to using PCA, we do KNN clustering on the schemas using a user-provided K (number of clusters). Each cluster identified by PCA is assigned a different color. Combined, these two algorithms to provide users an initial insight into the potential correlations that exist in their dataset.

Example 8. Figures 4.1 and 4.6 show an example of the resulting view on Twitter

and Yelp data. Note the much tighter clustering of attributes in the Twitter data: each cluster represents one particular, common type of tweet. Conversely, the Yelp schema includes a nested collection with, for example, hourly checkins at the business. The presence (or absence) of these terminal paths is more variable, and the resulting PCA cloud follows more of a gradient.

These graphs may not initially be intuitive to interpret, but they provide abstract bearings that map directly to phenomenon’s that exist in the data. In lieu of a formal user study, we postulate that the information gained from each algorithm independently will benefit the exploration process that is agnostic to labeled information. This is in contrast to other tools [?], that largely utilize naming cues and attribute nesting to meet user needs.

4.3.2 Schema Exploration

Now that we have shown the user potentially interesting subschemas, our next task is to help them to (1) narrow down on actually interesting subschemas, and (2) use those schemas to drill down to a segment of the schema data. For the first step, it is critical that the user develop a good intuition for what the visual representations encode. One way to accomplish this is to establish a bi-directional mapping between SCHEMADRILL’s two data panes.

To map from visual survey to schema summary, we provide users with a lasso tool. As illustrated in Figure 4.1, users can select regions of the KNN-PCA Cloud and the corresponding schemas within that region. Doing so identifies the maximal subschema contained in all subschemas and regenerates the schema summary pane based only on the segment containing the maximal subschema. The maximal subschema itself is also highlighted in the schema summary pane. On the Covariance cloud, the lasso tool behaves similarly, selecting attributes explicitly rather than a maximal subschema.

The reverse mapping is achieved using highlighting, as illustrated in Figure 4.7. Users can select one or more attributes (or groupings) in the schema summary pane, and the KNN-PCA Cloud (resp., Covariance Cloud) is modified to highlight schemas

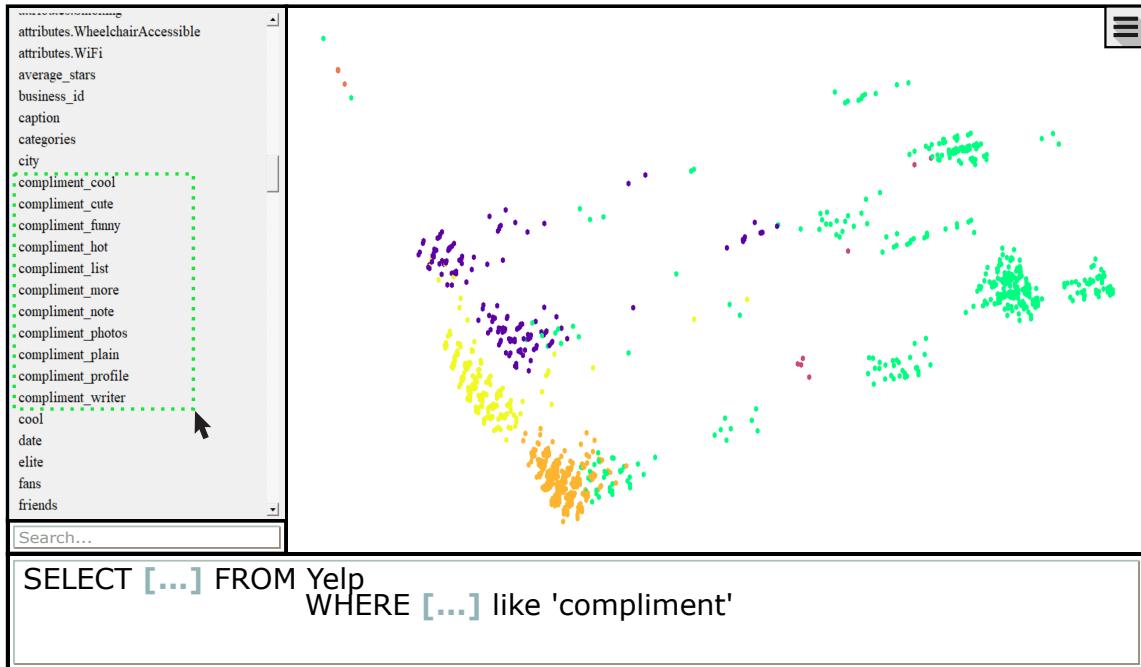


Figure 4.7: Points containing user selected attributes are highlighted green.

in the corresponding segment (resp., to highlight the attributes). In conjunction with their prior knowledge of their tasks and ideal use cases, we use this approach to perform the initial schema segmentation.

In either case, after selecting a set of attributes or schemas, the analyst may choose to drill down into the selected segment, regenerating both views for the now restricted collection of schemas.

4.3.3 An Iterative Approach

At any time in our exploration pipeline analysts may stop where they are, take the knowledge they gained about their dataset, and restart the process from the beginning. Through exploring the Twitter dataset we found retweets and quoted tweets to have a high correlation, with this knowledge we can then start back at the beginning and depending on whether our task allows us to merge these attributes, we can then choose a more appropriate K value for KNN. In addition, attributes

that have no relationship to core attributes, such as a users profile_image_url being present, can easily be pruned to compress our summary output. By incrementally learning about their dataset, an analyst can converge on the views required for their tasks.

4.4 Related Work

Schema Extraction. Schema extraction for Json, as well as for other self-describing data models like XML has seen active interest from a number of sources. An early effort to summarize self-describing *hierarchical* data can be found in the LORE system’s DataGuides [?]. DataGuides view schemas begin with a forest of tree-shaped schemas and progressively merge schemas, deriving a compact encoding of the forest as a DAG. Although initially designed for XML data, similar ideas have been more recently applied for Json data as well [?, ?]. Key challenges in this space involve simply extracting schemas from large, multi-terabyte collections of Json data [?], as well as managing ambiguity in the set of possible factorizations of a schema [?, ?]. The approach taken by Baazizi et. al. [?] in particular adopts a type unification model similar to ours, but lacks the conjunctive operator of our type-system. For non-hierarchical data, interactive tools like Wrangler [?] provide an interactive frameworks for regularizing schemas.

Physical Layout. While schemas play a role in the interpretability of a Json data set, they can also help improve the performance of Json queries. One approach relies on inverted indexes [?] to quickly identify records that make use of sparse paths in the schema. Another approach is to normalize schema elements [?]. Although the resulting schema may not always be interpretable, this approach can result in substantial space savings.

Information Retrieval. From a more general perspective, the schema extraction problem which aims at making large datasets tractable for interactive exploration, is an instance of *categorization* problem that has been repeatedly studied in the literature. More precisely, attributes (metadata) of the datasets can be grouped into

a hierarchy of "facets" (i.e., categories) [?] where the child-level facets are conditioned on the presence of the parent one. In our approach, we adopt the hierarchical data visualization and focus more on the algorithmic essence of the problem: How to (1) balance between the preciseness and conciseness of the visualization and (2) respond to users' data exploration requests in a scalable way.

4.5 Future Work

One challenge that we will need to address in SCHEMADRILL is coping with nested collections. At the moment, the user can manually merge collections of attributes that correspond to disjoint entities. However, we would like to automate this process. One observation is that a typical collection like an array has a schema with the general structure:

$$(P_1 \vee P_1P_2 \vee P_1P_2P_3 \vee \dots) = (P_1 \wedge (\emptyset \vee P_2 \wedge (\emptyset \vee P_3 \wedge (\dots))))$$

The version of this expression on the right hand side is notable as its closure over the semiring $\langle \{\{P\}\}, \vee, \wedge, \emptyset, \{\{\}\} \rangle$ would indicate that the semiring is "quasiregular" or "closed", an algebraic structure best associated with the Kleene star. Hence, we plan to explore the use of the Kleene star to encode nested collections in our algebra. A key challenge in doing so is detecting opportunities for incorporating it into a summary, a more challenging form of the factorization problem.

A further step to increase the capabilities of SCHEMADRILL is to incorporate type information in the summarization. This adds an extra layer of information an analyst can extract from our system, as well as the ability to identify and correct schema errors. As a long term goal we will provide capabilities for linking views, for example by defining functional dependencies. The goal is to create full entity relationship diagrams. In particular, one interesting way to identify potential relationships that exist between entities is by leveraging the overlap between segments.

Chapter 5

Application: Summarizing Tuple Correlation in Probabilistic Databases

Abstract

[[]]

5.1 Introduction

Many tasks in machine learning, information extraction, etc, generate uncertain data. That is, output tuples generated from one probabilistic prediction differ from the other. Though it is a general practice to deterministically output the most probable prediction, it is beneficial to explore other possible predictions, especially when an application needs to compare predictions under different parameter tuning strategies [?]. However, as [?] pointed out, state-of-the-art statistical models of information extraction provide a sound probability distribution over extractions but are challenging to represent and query in relational framework.

One naive strategy would be letting the application to iterate over all observed outputs, together with their probabilities of being observed. This strategy, on one extreme, defines a *probabilistic database* [?] and walks the application through the *entire space* of possible worlds. This is a sufficient condition for the application to discover all information that may be useful: the probability that any set of tuples will be contained in the output. However, this naive strategy puts an unnecessary burden on the application side. That is, there are practically exponential number of possible outputs [?]. Each observed output may contain duplicate tuples and it is possible to deliver the probability distribution over possible outputs in a more compact way.

The other naive strategy would be letting the application view each *individual* tuple together with the probability that the tuple is contained in an output. The probability of any *set* of tuples can then be estimated by multiplying the probability of each content tuple. This strategy, on the other extreme, defines a *tuple-independent probabilistic database* [?] that compactly represents the distribution of outputs by assuming each tuple independently exists in an output. However, for applications where tuples are predominantly correlated in the output, probabilities delivered by tuple-independent probabilistic database may greatly deviate from the truth [?]. Consider an example application where tuple-correlation is dominant by nature.

Example 9. Consider an Entity Resolution (ER) task where the input is a set of tuples where each tuple contains a set of attributes that may help to determine

whether two tuples are different facets of the same entity. The output will be the same set of tuples but augmented with an additional attribute indicating its predicted entity ID.

Since the prediction is uncertain, the output data can be viewed as a space of possible sets of augmented tuples that a predictor may generate. Augmented tuples in the output are correlated by design, that is, given that a tuple \mathbf{t} is assigned entity ID eid , other tuples having attributes similar to \mathbf{t} being assigned the same eid increases and are no longer independent events. This example can be generalized to classification or clustering related tasks where the output is a set of label-augmented tuples where the label assignments are uncertain.

To encode tuple correlation, various kind of auxiliary information (e.g., *c-tables* [?], factors [?, ?, ?], Markov Logic [?]) are introduced for delivering more accurate probabilities. Auxiliary information inevitably increases the model complexity, making the inference of the desired probabilities harder. It also puts a heavier burden on the application side for digesting the representation.

To carefully introduce auxiliary information, we start by considering *c-tables* [?]. A c-tables compactly represents the distribution over outputs by presenting only individual tuples, annotated with a propositional formula. A tuple exists in the output iff its propositional formula evaluates to be true. By walking the application through a c-table as well as the corresponding probability distribution over variables in its formula, we equivalently delivers the distribution of possible outputs. However, linking the probability distribution with its c-table generally requires heavy weight inference. That is, to compute the probability of a set of tuples being co-existent, we need to enumerate all possible evaluations of their variables such that their propositional formula are true. This process is non-trivial especially when we need to repeated compute such probabilities for different sets of tuples. Moreover, in our example application, it is generally difficult to obtain such propositional formula, since tuple correlation is caused by the complex nature of some predictor. As a result, for our target applications, we only expect a tuple in the c-table to be annotated with a single distinct variable X_i with $X_i = 1$ meaning the presence of the i th tuple in

the output. The probability of two tuples $\mathbf{t}_i, \mathbf{t}_j$ being co-existent is directly linked with the marginal probability $p(X_i = 1, X_j = 1)$. To deliver the distribution of outputs, we need to walk an application through the c-table and the joint distribution $\mathbf{P} = p(\dots, X_i, \dots)$. Enumerating all valuations of the variables is not practical. Instead, we group correlated variables together and only enumerate valuation within each group. This essentially creates a *factor* for each group. Mappings in a factor is still exponential with respect to the number of correlated variables. In this paper, we study the problem of encoding the joint distribution \mathbf{P} by only a subset of mappings in factors. We show how desired probabilities of the distribution \mathbf{P} can be computed by the lossy encoding. Due to incomplete mappings, the computed probabilities generally deviates from the truth. However, the encoding allows finer granularity in trading-off between conciseness and probability computation accuracy. That is, by gradually populating the mappings in a factor, an application can choose to obtain an encoding that incurs a high-fidelity in computing probabilities but is verbose, or obtain a more compact encoding that incurs a greater loss of accuracy. To manage the loss-rate, we develop a framework for reasoning about the trade-off between verbosity and fidelity. While the encoding does not admit closed-form solutions to classical information theoretical fidelity measures like information loss, we propose an efficiently computable fidelity measure called Reproduction Error.

Contributions.

Roadmap.

5.2 Problem Definition

Input Definition. As shown in Figure 5.1, the input data consists of a finite set of possible instances $\mathbf{D} = \{D_1, \dots, D_n\}$ where each instance D_i consists of a set of tuples. We use the *possible worlds semantics* [?] and call \mathbf{D} as *possible worlds*. Each world is assigned a probability by the distribution $\mathbf{P} : \mathbf{D} \rightarrow [0, 1]$ such that $\sum_i \mathbf{P}(D_i) = 1$. In other words, our input is defined as a discrete probability space $PDB = (\mathbf{D}, \mathbf{P})$.

Goal. An application theoretically may require scanning through all possible worlds, in order to aggregate the probabilities where the corresponding worlds contain some target *set* of tuples. In the context of Entity Resolution task, an application may require aggregating probabilities of all predictions that contain tuples labeled by a target *eid*. For convenience, we will refer to these aggregated probabilities as *marginal probabilities* or simply *marginals* in the rest of the paper. Our goal is then to help the application efficiently explore and compute marginals in a practical way without iterating through a large number of possible worlds.

System Overview. We begin by giving the reader an overview on how the input flows *theoretically* and *practically* to a client-side application. In Figure 5.1, the red dashed lines encompasses a theoretical data flow that iterates through all possible worlds. In practice, possible worlds in the input will converted into an initial form of representation. Based on the initial form, further compression can be made that provides a materialized view over the input. The materialized view can be efficiently queried and the query result will be simplified before being presented to the application.

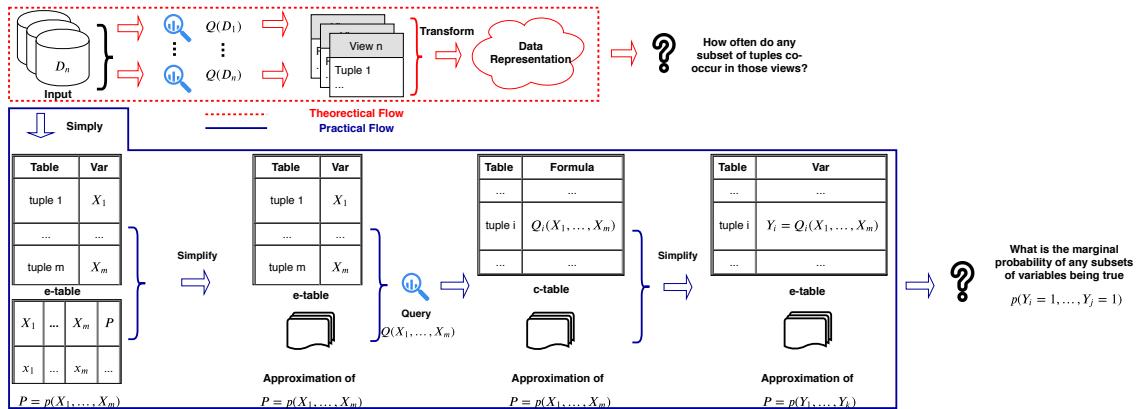


Figure 5.1: Theoretical v.s. Practical Data-flows.

5.2.1 Practical Input Representation

The previous definition of \mathbf{D} does not suggest a practical representation, especially when explicit enumeration of all possible worlds is not feasible when n is very large.

C-tables and pc-tables. We adopt *c-tables* to compactly represent \mathbf{D} . A *c-table* [?] is a set of *distinct* tuples contained in all possible worlds, with each tuple annotated by a propositional formula over random variables. A *valuation* over these variables assigns each propositional formula a truth value, indicating the existence of the corresponding tuple in a world. There is a many-to-one mapping from valuations to a possible world. Extending a c-table by a joint probability distribution over random variables, the resulting *pc-table* can represent any probabilistic database *PDB* [?].

Simplified C-tables. For our example application, there are two limitations by adopting a general c-table representation. Firstly, it is difficult to obtain such propositional formula due to lack of knowledge on how tuples are correlated in *PDB*. That is, due to complex nature of the function that generates *PDB*, it is generally hard to formulate how tuples are correlated. Secondly, the mapping from valuations of variables to possible worlds is many-to-one for a general c-table. To compute desired marginals, it may require an expensive aggregate over all valuations that map to the same world.

To overcome the limitations, we represent our input by a family of simplified c-tables, which we call *existence tables* or simply *e-tables*. Each tuple \mathbf{t}_i in an e-table is annotated with a single random variable X_i whose value is 1, meaning that the i th tuple exists in a world, or 0 otherwise. We refer to these variables as *existence variables*. There is an *one-to-one* mapping from a valuation of existence variables $(x_1, \dots, x_n) \in \{0, 1\}^n$ to a world D . The probability of the world (i.e., $\mathbf{P}(D)$) can be computed directly from the joint probability distribution $p(X_1 = x_1, \dots, X_n = x_n)$. When it is clear from context, we abuse notation and denote the distribution $p(X_1 = x_1, \dots, X_n = x_n)$ also as \mathbf{P} . We call an e-table, together with \mathbf{P} , a *pe-table*.

5.2.2 Summarizing Joint Distribution

It is not practical, both in space and marginal computation efficiency, to present pe-table in a view that iterates through all valuations (i.e., worlds) of existence variables. Hence we need a view that summarizes \mathbf{P} into a more compact form.

Naive Summary. There can be a set of tuples $\{\dots, \mathbf{t}_i, \dots\}$ independent of each other such that,

$$p(\dots, X_i = 1, \dots) = \prod_i p(X_i = 1)$$

In other words, we do not need to store valuations of any combination of those variables in the view, which greatly reduces the number of parameters. We name the set of parameters $p(X_i = 1)$ as a *naive summary* and return to it through out the rest of this paper. The naive summary is often referred to as *tuple-independent* model [?]. When variables are correlated, which is prevalent in our example application, the marginals computed from a naive summary are error-prone.

Loss-less Summary. For tuples whose variables that are correlated, a *factor graph* [?] that is equivalent to the marginal probability distribution over these variables, is a sufficient condition for computing any relevant marginals. However, it requires heavy weight inference to obtain such set of factors and for each factor, the number of parameters is still exponential to the number of variables that the factor contains.

Lossy Summary. Practically an application may not require computing exact marginals for *every* set of tuples and we can *cache* marginals that are frequently computed and offer estimation on others. Specifically, in our example Entity Resolution task, an application may require to explore labelled tuples that are contained in *all* possible predictions, and tuples that are (in-)frequently or otherwise statistically significant. The application may tolerate inexact marginal computation for other sets of tuples. This observation motivates us to study the option of a *lossy* summary containing only a limited number of marginals and offers estimation on marginals not contained.

5.2.3 Pattern-based Summary

Patterns. We name a set of co-existing tuples as a *pattern*. More formally we define a *pattern* as an arbitrary binary vector $\omega = (x_1, \dots, x_n)$ with $x_i = 1$ meaning the i th tuple being present in the pattern. A pattern ω is *contained* in the other $\bar{\omega}$ if $\forall i, x_i \leq \bar{x}_i$, denoted as $\omega \subseteq \bar{\omega}$. Note that a world D is a pattern by definition.

Randomly drawing a world $D \in \mathbf{D}$ from PDB , the probability that it contains pattern ω is defined by:

$$p(\mathbf{D} \supseteq \omega) = \sum_i \mathbf{P}(I_i) \mathbf{1}_\omega(I_i) \text{ where } \mathbf{1}_\omega(I_i) = \begin{cases} 1 & \text{if } I_i \supseteq \omega \\ 0 & \text{otherwise} \end{cases}$$

In pe-tables, $p(\mathbf{D} \supseteq \omega)$ is equivalent to marginal probability $p(X_1 \geq x_1, \dots, X_n \geq x_n)$. We call $p(\mathbf{D} \supseteq \omega)$ as the *marginal* of a pattern.

Pattern-based Summary. Denote the mapping \mathcal{E}_{max} from each pattern ω to its marginal:

$$\mathcal{E}_{max} = \{ \omega \rightarrow p(\mathbf{D} \supseteq \omega) \mid \omega \in \{0,1\}^n \}$$

A *pattern-based summary* \mathcal{E} is any such partial mapping $\mathcal{E} \subseteq \mathcal{E}_{max}$. We denote the marginal of pattern ω in summary \mathcal{E} by $\mathcal{E}[\omega]$ ($= p(\mathbf{D} \supseteq \omega)$). When it is clear from context, we abuse syntax and also use \mathcal{E} to denote the set of patterns it maps (i.e., $domain(\mathcal{E})$). Hence, $|\mathcal{E}|$ is the number of mapped patterns, which we call the summary's *verbosity*.

5.3 Marginal Estimation

A pattern-based summary not only delivers marginals that it contains, but also helps in estimating those that are not explicitly given.

Naive Summary Revisited. Before detailed explanation, we first consider an intuitive example. There is one specific family of pattern-based summaries that treat each tuple as being independent, which is equivalent to naive summaries discuss in Section 5.2.1. As a re-phrase, we call a pattern-based summary as a naive summary $\ddot{\mathcal{E}}$ if it is composed of all patterns that have exactly one possible tuple \mathbf{t}_i with marginal $p(X_i \geq 1)$:

$$domain(\ddot{\mathcal{E}}) = \{ (0, \dots, 0, x_i, 0, \dots, 0) \mid i \in [1, n], x_i = 1 \}$$

To estimate marginals not given in the summary, we apply multiplication

$$p(X_1 \geq x_1, \dots, X_n \geq x_n) = \prod_{i=1, n} p(X_i \geq x_i)$$

Patterns as Constraints. For an arbitrary summary \mathcal{E} , we characterize its process of marginal computation by defining a *space* (denoted by $\Omega_{\mathcal{E}}$) of distributions $\rho \in \Omega_{\mathcal{E}}$ allowed by a summary \mathcal{E} . This space is defined by constraints as follows: First, we have the general properties of probability distributions:

$$\forall D_i \in \{0, 1\}^n : \rho(D_i) \geq 0 \quad \sum_{i=1,2^n} \rho(D_i) = 1$$

Each pattern ω in the summary \mathcal{E} constrains relevant probabilities in distribution ρ to sum to the target marginal:

$$\forall \omega \in \text{domain}(\mathcal{E}) : \mathcal{E}[\omega] = \sum_{D \supseteq \omega} \rho(D)$$

The dual constraints $1 - \mathcal{E}[\omega] = \sum_{D \not\supseteq \omega} \rho(D)$ are redundant under the constraint $\sum_i \rho(D_i) = 1$.

The resulting space $\Omega_{\mathcal{E}}$ is the set of all permitted distributions, or equivalently all input *PDB*, that obey these constraints. As a result, the distribution $\rho \in \Omega_{\mathcal{E}}$ that the summary delivers is ambiguous: We model this ambiguity using a random variable $\mathcal{P}_{\mathcal{E}}$ with support $\Omega_{\mathcal{E}}$. The true distribution \mathbf{P} must appear in $\Omega_{\mathcal{E}}$, denoted as $\mathbf{P} \equiv \rho^* \in \Omega_{\mathcal{E}}$ (i.e., $p(\mathcal{P}_{\mathcal{E}} = \rho^*) > 0$). Of the remaining distributions ρ admitted by $\Omega_{\mathcal{E}}$, it is possible that some are more likely than others. This prior knowledge may be modeled as a prior on the distribution of $\mathcal{P}_{\mathcal{E}}$ or equivalently by an additional constraint. However, for the purposes of this paper, we take the uninformed prior by assuming that $\mathcal{P}_{\mathcal{E}}$ is uniformly distributed over $\Omega_{\mathcal{E}}$:

$$p(\mathcal{P}_{\mathcal{E}} = \rho) = \begin{cases} \frac{1}{|\Omega_{\mathcal{E}}|} & \text{if } \rho \in \Omega_{\mathcal{E}} \\ 0 & \text{otherwise} \end{cases}$$

The summary delivers marginals by randomly picking one permitted distribution from the entire space and aggregate probabilities accordingly.

5.4 Summary Accuracy

The accuracy of such delivery can be considered from two related, but subtly distinct perspectives: (1) *Ambiguity* measures how much room the summary leaves for

interpretation and (2) *Deviation* measures how reliably the summary approximates the target distribution ρ^* .

Ambiguity. We define the Ambiguity $\Gamma(\mathcal{E})$ of a summary as the entropy of the random variable $\mathcal{P}_{\mathcal{E}}$. The higher the entropy, the less precisely \mathcal{E} identifies a specific distribution.

$$\Gamma(\mathcal{E}) = \sum_{\rho} p(\mathcal{P}_{\mathcal{E}} = \rho) \log(p(\mathcal{P}_{\mathcal{E}} = \rho))$$

Deviation. The deviation from any permitted distribution ρ to the true distribution ρ^* can be measured by the Kullback-Leibler (K-L) divergence $\mathcal{D}_{KL}(\rho^* || \rho)$. We define the Deviation $\mathbb{D}(\mathcal{E})$ of a summary as the expectation of the K-L divergence over all permitted $\rho \in \Omega_{\mathcal{E}}$:

$$\mathbb{D}(\mathcal{E}) = \mathbb{E}_{\mathcal{P}_{\mathcal{E}}} [\mathcal{D}_{KL}(\rho^* || \mathcal{P}_{\mathcal{E}})] = \sum_{\rho \in \Omega_{\mathcal{E}}} p(\mathcal{P}_{\mathcal{E}} = \rho) \cdot \mathcal{D}_{KL}(\rho^* || \rho)$$

Limitations. Neither Deviation nor Ambiguity has a closed-form formula and it is not practical to enumerate the infinitely large space. Authors in [?] suggests that a single *representative* can be chosen from the space that approximates true distribution ρ^* .

Maximum Entropy Distribution. The representative distribution is chosen by applying the maximum entropy principle [23] commonly used in pattern-based encoders [35, 16]. That is, we select the distribution $\bar{\rho}_{\mathcal{E}}$ with maximum entropy:

$$\bar{\rho}_{\mathcal{E}} = \arg \max_{\rho \in \Omega_{\mathcal{E}}} \mathcal{H}(\rho) \quad \text{where } \mathcal{H}(\rho) = \sum_{i=[1,2^n]} -\rho(I_i) \log \rho(I_i)$$

The maximum entropy distribution best represents the current state of knowledge. That is, a distribution with lower entropy assumes additional constraints derived from patterns that we do not know, while one with higher entropy violates the constraints from patterns we do know.

Maximizing an objective function belonging to the exponential family (entropy in our case) under a mixture of linear equality/inequality constraints is a convex

optimization problem [9] which guarantees a *unique* solution and can be solved using the cvx toolkit [18], and/or by *iterative scaling* [35, 16].

Naive Summary revisited. For naive summaries $\ddot{\mathcal{E}}$ specifically, $\bar{\rho}_{\ddot{\mathcal{E}}}$ has the closed-form representation:

$$\bar{\rho}_{\ddot{\mathcal{E}}}(I) = \prod_i p(X_i = x_i) \quad \text{where } I = (x_1, \dots, x_n)$$

Reproduction Error. We define *Reproduction Error* $e(\mathcal{E})$ of summary \mathcal{E} as the entropy difference between the representative and true distribution:

$$e(\mathcal{E}) = \mathcal{H}(\bar{\rho}_{\mathcal{E}}) - \mathcal{H}(\rho^*) \quad \text{where } \bar{\rho}_{\mathcal{E}} = \arg \min_{\rho \in \Omega_{\mathcal{E}}} -\mathcal{H}(\rho)$$

It has been proved [?] that Reproduction Error closely parallels both Ambiguity and Deviation.

5.4.1 Query Accuracy

We then study the accuracy of pattern-based summaries when the corresponding probabilistic databases are involved in three basic relational operations: (1) *join*, (2) *selection* and (3) *project*. Specifically, we study how to modify existing pattern-based summaries in order to suit these operations and give upper-bounds on Reproduction Error of the resulting summaries.

5.4.1.1 Accuracy under join

Consider two uncertain tables, or equivalently possible worlds $\mathbf{P}_1, \mathbf{P}_2$ under *join* operation. Recall that possible worlds can be translated into joint distribution, that is, we have $\mathbf{P}_1 = p(\dots, X_i, \dots)$ and $\mathbf{P}_2 = p(\dots, Y_j, \dots)$ where X_i, Y_j indicate the presence of i, j th tuple in possible worlds $\mathbf{P}_1, \mathbf{P}_2$ respectively. The resulting uncertain table after *join*, or equivalently possible worlds $\mathbf{P}_{1 \times 2} = \mathbf{P}_1 \times \mathbf{P}_2$ can also be translated into joint distribution $\mathbf{P}_{1 \times 2} = p(\dots, Z_{i,j}, \dots)$ where $Z_{i,j} = 1$ iff $X_i = 1$ and $Y_j = 1$.

Join-derived summary. Consider pattern-based summaries $\mathcal{E}_1, \mathcal{E}_2$ for $\mathbf{P}_1, \mathbf{P}_2$ with corresponding maximum entropy distribution $\bar{\rho}_{\mathcal{E}_1}, \bar{\rho}_{\mathcal{E}_2}$. The resulting summary $\mathcal{E}_{1\times 2}$ can be obtained from *joining* patterns from both summaries and multiplying their corresponding marginals. Denote the resulting maximum entropy distribution after *join* as $\bar{\rho}_{\mathcal{E}_{1\times 2}}$.

Upper-bound of Reproduction Error. By definition, the Reproduction Error $e(\mathcal{E}_{1\times 2})$ of summary $\mathcal{E}_{1\times 2}$ measures the entropy difference $\mathcal{H}(\bar{\rho}_{\mathcal{E}_{1\times 2}}) - \mathcal{H}(\mathbf{P}_{1\times 2})$. We would like to approximate the two entropies to avoid populating mappings in distributions $\bar{\rho}_{\mathcal{E}_{1\times 2}}$ and $\mathbf{P}_{1\times 2}$. That is, we would like to provide an upper-bound $\tilde{e}(\mathcal{E}_{1\times 2})$ and such that $e(\mathcal{E}_{1\times 2}) \leq \tilde{e}(\mathcal{E}_{1\times 2})$.

Observe that there is a many-to-one mapping from valuations of jointly distributed variables $\mathbf{P}_{1,2} = p(\dots, X_i, \dots, Y_j, \dots)$ to $\mathbf{P}_{1\times 2}$. That is, by fixing the valuation of all other variables except for some pair X_i, Y_j , all three cases where $X_i = 0 \wedge Y_j = 0$, $X_i = 0 \wedge Y_j = 1$ or $X_i = 1 \wedge Y_j = 0$ map to the same valuation of $\mathbf{P}_{1\times 2}$ where $Z_{i,j} = 0$. Hence, the entropy $\mathcal{H}(\mathbf{P}_{1\times 2}) \leq \mathcal{H}(\mathbf{P}_{1,2})$. Trivially we also have $\mathcal{H}(\mathbf{P}_{1\times 2}) \geq \max(\mathcal{H}(\mathbf{P}_1), \mathcal{H}(\mathbf{P}_2))$. Since X_i are independent of Y_j in $\mathbf{P}_{1,2}$, we have $\mathcal{H}(\mathbf{P}_{1,2}) = \mathcal{H}(\mathbf{P}_1) + \mathcal{H}(\mathbf{P}_2)$. We have similar observation for maximum entropy distributions: $\mathcal{H}(\bar{\rho}_{\mathcal{E}_{1\times 2}}) \leq \mathcal{H}(\bar{\rho}_{\mathcal{E}_1}) + \mathcal{H}(\bar{\rho}_{\mathcal{E}_2})$. We thus compute the upper-bound $\tilde{e}(\mathcal{E}_{1\times 2}) = \mathcal{H}(\bar{\rho}_{\mathcal{E}_1}) + \mathcal{H}(\bar{\rho}_{\mathcal{E}_2}) - \max(\mathcal{H}(\mathbf{P}_1), \mathcal{H}(\mathbf{P}_2))$ or a looser but simpler bound $\mathcal{H}(\bar{\rho}_{\mathcal{E}_1}) + \mathcal{H}(\bar{\rho}_{\mathcal{E}_2})$.

5.4.1.2 Accuracy under selection predicate

Consider the case where only a subset of tuples in possible worlds $\mathbf{P} = p(\dots, X_i, \dots)$ are chosen, according to some arbitrary *selection* predicate σ . For convenience, we abuse notation and also represent the indices of its chosen tuples as σ . The resulting possible worlds, or equivalently the joint distribution \mathbf{P}_σ is the marginal distribution $p(\dots, X_{i \in \sigma}, \dots)$ by summing out variables $X_{i \notin \sigma}$.

Selection-derived summary. Consider a pattern-based summary \mathcal{E} for \mathbf{P} , the resulting summary \mathcal{E}_σ for \mathbf{P}_σ can be constructed in three steps: (1) remove tuples

that are not chosen by σ for each pattern in \mathcal{E} , (2) merge patterns that are the same and add their marginals and (3) remove empty patterns.

Upper-bound of Reproduction Error. Suppose we encode \mathbf{P} first by naive summary $\ddot{\mathcal{E}}$ whose maximum entropy distribution $\bar{\rho}_{\mathcal{E}}$ is equivalent to tuple-independent model. The resulting naive summary $\ddot{\mathcal{E}}_{\sigma} \subseteq \ddot{\mathcal{E}}$ after *select* operation is a subset of single-tuple patterns filtered by the predicate σ . Reproduction Error of the naive summary $\ddot{\mathcal{E}}_{\sigma}$ is thus computed as $e(\ddot{\mathcal{E}}_{\sigma}) = \sum_{i \in \sigma} \mathcal{H}(p(X_i)) - \mathcal{H}(\mathbf{P}_{\sigma})$. In cases when we do not want to populate mappings of distribution \mathbf{P}_{σ} for each given predicate, we can simply compute $\tilde{e}(\ddot{\mathcal{E}}_{\sigma}) = \sum_{i \in \sigma} \mathcal{H}(p(X_i))$ which is the upper-bound of $e(\ddot{\mathcal{E}}_{\sigma})$ for all possible σ . For any other summary \mathcal{E}_{σ} that contains naive summary $\ddot{\mathcal{E}}_{\sigma} \supseteq \mathcal{E}_{\sigma}$, it may improve naive summary by marginals that encode tuple-correlation. Hence the Reproduction Error $e(\mathcal{E}_{\sigma})$ is also upper-bounded by $\tilde{e}(\ddot{\mathcal{E}}_{\sigma})$. Trivially, we also have $e(\mathcal{E}_{\sigma}) \leq \mathcal{H}(\bar{\rho}_{\mathcal{E}})$. Hence the Reproduction Error of any summary $\mathcal{E}_{\sigma} \supseteq \ddot{\mathcal{E}}_{\sigma}$ is upper-bounded by $\min(\tilde{e}(\ddot{\mathcal{E}}_{\sigma}), \mathcal{H}(\bar{\rho}_{\mathcal{E}}))$.

5.4.1.3 Accuracy under project

Consider the case where tuples in uncertain table \mathbf{P} are *projected* onto an arbitrary set of attributes A . We group tuples whose values are the same on projected attributes A and create a new random variable Z for each group, where $Z = 0$ iff $X_i = 0$ for all tuples \mathbf{t}_i in the group, otherwise $Z = 1$. The possible worlds distribution \mathbf{P}_{π_A} after *project* operation, is thus the joint distribution $p(Z_1, \dots, Z_k)$ where k is the total number of tuple groups.

Project-derived summary. Suppose we have summary \mathcal{E} for \mathbf{P} . The resulting summary \mathcal{E}_{π_A} for \mathbf{P}_{π_A} can be constructed by (1) applying *project* on each pattern and (2) merge patterns that are the same after *project* and add their marginals accordingly.

Upper-bound of Reproduction Error. We fix the valuations for all Z_i except for the k th group Z_k . Equivalently we fix all valuations of X_i in \mathbf{P} except for those X_j that belong to k th group. Then we notice that there is a many-to-one mapping from all valuations of X_j in the group, where at least one $X_j = 1$, to

$Z_k = 1$. In other words, we have $\mathcal{H}(\mathbf{P}_{\pi_A}) \leq \mathcal{H}(\mathbf{P})$. Similarly we also have $\mathcal{H}(\bar{\rho}_{\mathcal{E}_{\pi_A}}) \leq \mathcal{H}(\bar{\rho}_{\mathcal{E}})$ for maximum entropy distributions of corresponding summaries. That is, the Reproduction Error of summary \mathcal{E}_{π_A} for any *project* operation π_A is upper-bounded by $\mathcal{H}(\bar{\rho}_{\mathcal{E}}) - \mathcal{H}(\mathbf{P}_{\pi_A})$ or a looser but simpler bound $\mathcal{H}(\bar{\rho}_{\mathcal{E}})$.

5.5 Constructing Summaries

5.6 Related Work

In this section, we walk through related work on encoding tuple correlation in probabilistic databases using various kinds of auxiliary information.

C-tables. The canonical representation of probabilistic databases is *c-table* [?]. A *c-table* is a table with each tuple annotated with an additional attribute: a propositional formula over random variables. Two tuples are correlated if random variables in their propositional formula depends on each other. Given a query over a c-table, the correlation among tuples in the query answer is encoded by the *lineage* (or provenance [?]) of each tuple, which is also a propositional formula.

Based on c-tables, Trio system [?] resolves tuple correlation by tracking lineage of tuples in query processing. By defining a probability space over the assignments of the random variables in lineage of tuples, a *confidence value* can be computed for each tuple in the query answer. The probability of some answer to be chosen as the "correct" one is proportional to the product of confidence values of all of its tuples. Similar idea can be found in MayBMS [?]. However, as [?] pointed out, there are cases where lineage is not obtainable. In our example application, tuples are correlated in a latent, complex way that populating propositional formula for tuples in the c-table is difficult.

Graphical models and Factors. Enlightened by Probabilistic Relational Model (PRM) [?], numerous works [?, ?] use *probabilistic graphical models* (PGM) to encode tuple correlation directly by probabilities of co-existing tuples. More specifically, considering any set of co-existing tuples, encoded as a binary vector $\mathbf{v} = (x_1, \dots, x_n)$ with $x_i = 1$ indicating the presence of the i th possible tuple in the co-existence

relationship. A *factor function* or simply *factor* $f : \mathbf{v} \in \{0, 1\}^n \rightarrow [0, 1]$ maps any set of co-existing tuples \mathbf{v} to its probability of being present in the probabilistic database. Since there are exponentially large (i.e., 2^n) number of mappings, instead of having a single factor that takes n variables as input, we can have multiple factors, each takes only a small subset of variables that are correlated. If two factors do not share variables, the variables in one are considered independent of those in the other. By multiplying the outputs of the two factors, we can efficiently compute the probability of any combination of variables in two factors.

Markov Logic Networks. Jha and Suciu [?] shows how to automatically generate factors and populate their content mappings using prior knowledge on tuple-correlation, expressed by the language of First Order Logic. Similar technique is implemented in DeepDive system [?].

As authors of MCDB [?] pointed out, most previously discussed approaches augment data with attribute-level or tuple-level annotations, which are loaded into the database along with the data itself. The mixture of data and uncertainty annotation is inflexible when changes have to be made to the underlying uncertainty model (e.g., lineage), as all related annotations typically need to be recomputed outside of the database and then loaded back in. To facilitate maintenance, MCDB proposes to separate uncertainty model from the data.

Monte Carlo approach. More specifically, MCDB allows a user to define arbitrary *variable generation* (VG) functions that encodes the uncertainty model. That is, treating a VG function as a black box, MCDB use it to pseudo-randomly generate samples of possible instances on the fly and run queries over the samples. Interestingly, sample instances are sets of co-existing tuples. By assigning a probability to each sample, the probability of any co-existing tuples that are contained in any sample is estimated naturally by iterating through the samples.

Chapter 6

Findings and Discussions

Chapter 7

Conclusions and Future Work

todo

Bibliography

- [1] D. J. Abadi, S. Madden, and M. Ferreira. Integrating compression and execution in column-oriented database systems. In S. Chaudhuri, V. Hristidis, and N. Polyzotis, editors, *Proceedings of the ACM SIGMOD International Conference on Management of Data, Chicago, Illinois, USA, June 27-29, 2006*, pages 671–682. ACM, 2006.
- [2] S. Agrawal, S. Chaudhuri, and V. R. Narasayya. Automated selection of materialized views and indexes in SQL databases. In A. E. Abbadi, M. L. Brodie, S. Chakravarthy, U. Dayal, N. Kamel, G. Schlageter, and K. Whang, editors, *VLDB 2000, Proceedings of 26th International Conference on Very Large Data Bases, September 10-14, 2000, Cairo, Egypt*, pages 496–505. Morgan Kaufmann, 2000.
- [3] J. Aligon, M. Golfarelli, P. Marcel, S. Rizzi, and E. Turricchia. Similarity measures for OLAP sessions. *Knowl. Inf. Syst.*, 39(2):463–489, 2014.
- [4] S. Amer-Yahia and T. Johnson. Optimizing queries on compressed bitmaps. In A. E. Abbadi, M. L. Brodie, S. Chakravarthy, U. Dayal, N. Kamel, G. Schlageter, and K. Whang, editors, *VLDB 2000, Proceedings of 26th International Conference on Very Large Data Bases, September 10-14, 2000, Cairo, Egypt*, pages 329–338. Morgan Kaufmann, 2000.
- [5] G. Antoshenkov. Byte-aligned bitmap compression. In *Proceedings of the Conference on Data Compression, DCC '95*, pages 476–, Washington, DC, USA, 1995. IEEE Computer Society.

- [6] G. Antoshenkov and M. Ziauddin. Query processing and optimization in oracle rdb. *VLDB J.*, 5(4):229–237, 1996.
- [7] K. Aouiche, P. Jouve, and J. Darmont. Clustering-based materialized view selection in data warehouses. In Y. Manolopoulos, J. Pokorný, and T. K. Sellis, editors, *Advances in Databases and Information Systems, 10th East European Conference, ADBIS 2006, Thessaloniki, Greece, September 3-7, 2006, Proceedings*, volume 4152 of *Lecture Notes in Computer Science*, pages 81–95. Springer, 2006.
- [8] D. M. Blei. Probabilistic topic models. *Commun. ACM*, 55(4):77–84, 2012.
- [9] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, New York, NY, USA, 2004.
- [10] N. Bruno and S. Chaudhuri. Automatic physical database tuning: A relaxation-based approach. In F. Özcan, editor, *Proceedings of the ACM SIGMOD International Conference on Management of Data, Baltimore, Maryland, USA, June 14-16, 2005*, pages 227–238. ACM, 2005.
- [11] N. Bruno, S. Chaudhuri, and L. Gravano. Stholes: A multidimensional workload-aware histogram. In S. Mehrotra and T. K. Sellis, editors, *Proceedings of the 2001 ACM SIGMOD international conference on Management of data, Santa Barbara, CA, USA, May 21-24, 2001*, pages 211–222. ACM, 2001.
- [12] G. Chatzopoulou, M. Eirinaki, S. Koshy, S. Mittal, N. Polyzotis, and J. S. V. Varman. The querie system for personalized query recommendations. *IEEE Data Eng. Bull.*, 34(2):55–60, 2011.
- [13] S. Chaudhuri and V. R. Narasayya. An efficient cost-driven index selection tool for microsoft SQL server. In M. Jarke, M. J. Carey, K. R. Dittrich, F. H. Lochovsky, P. Loucopoulos, and M. A. Jeusfeld, editors, *VLDB'97, Proceedings of 23rd International Conference on Very Large Data Bases, August 25-29, 1997, Athens, Greece*, pages 146–155. Morgan Kaufmann, 1997.

- [14] C. Dwork. Differential privacy. In M. Bugliesi, B. Preneel, V. Sassone, and I. Wegener, editors, *Automata, Languages and Programming, 33rd International Colloquium, ICALP 2006, Venice, Italy, July 10-14, 2006, Proceedings, Part II*, volume 4052 of *Lecture Notes in Computer Science*, pages 1–12. Springer, 2006.
- [15] S. J. Finkelstein, M. Schkolnick, and P. Tiberio. Physical database design for relational databases. *ACM Trans. Database Syst.*, 13(1):91–128, 1988.
- [16] K. E. Gebaly, P. Agrawal, L. Golab, F. Korn, and D. Srivastava. Interpretable and informative explanations of outcomes. *PVLDB*, 8(1):61–72, 2014.
- [17] A. Giacometti, P. Marcel, E. Negre, and A. Soulet. Query recommendations for OLAP discovery-driven analysis. *IJDWM*, 7(2):1–25, 2011.
- [18] M. Grant and S. Boyd. CVX: Matlab software for disciplined convex programming, version 2.1. <http://cvxr.com/cvx>, Mar. 2014.
- [19] J. Han, H. Cheng, D. Xin, and X. Yan. Frequent pattern mining: current status and future directions. *Data Min. Knowl. Discov.*, 15(1):55–86, 2007.
- [20] L. Hu, K. A. Ross, Y. Chang, C. A. Lang, and D. Zhang. Queryscope: visualizing queries for repeatable database tuning. *PVLDB*, 1(2):1488–1491, 2008.
- [21] D. A. Huffman. A method for the construction of minimum-redundancy codes. *Proceedings of the IRE*, 40(9):1098–1101, Sept 1952.
- [22] A. K. Jain. Data clustering: 50 years beyond k-means. *Pattern Recognition Letters*, 31(8):651–666, 2010.
- [23] E. T. Jaynes. Prior probabilities. *IEEE Trans. Systems Science and Cybernetics*, 4(3):227–241, 1968.
- [24] I. T. Jolliffe. Principal component analysis. In M. Lovric, editor, *International Encyclopedia of Statistical Science*, pages 1094–1096. Springer, 2011.

- [25] A. Kamra, E. Terzi, and E. Bertino. Detecting anomalous access patterns in relational databases. *VLDB J.*, 17(5):1063–1077, 2008.
- [26] R. Kannan, S. Vempala, and A. Vetta. On clusterings: Good, bad and spectral. *J. ACM*, 51(3):497–515, 2004.
- [27] O. Kennedy, J. A. Ajay, G. Challen, and L. Ziarek. Pocket data: The need for TPC-MOBILE. In R. Nambiar and M. Poess, editors, *Performance Evaluation and Benchmarking: Traditional to Big Data to Internet of Things - 7th TPC Technology Conference, TPCTC 2015, Kohala Coast, HI, USA, August 31 - September 4, 2015. Revised Selected Papers*, volume 9508 of *Lecture Notes in Computer Science*, pages 8–25. Springer, 2015.
- [28] N. Khoussainova, Y. Kwon, M. Balazinska, and D. Suciu. Snipsuggest: Context-aware autocomplete for SQL. *PVLDB*, 4(1):22–33, 2010.
- [29] K. Knight and D. Marcu. Summarization beyond sentence extraction: A probabilistic approach to sentence compression. *Artif. Intell.*, 139(1):91–107, 2002.
- [30] G. Kul, D. Luong, T. Xie, P. Coonan, V. Chandola, O. Kennedy, and S. J. Upadhyaya. Ettu: Analyzing query intents in corporate databases. In J. Bourdeau, J. Hendler, R. Nkambou, I. Horrocks, and B. Y. Zhao, editors, *Proceedings of the 25th International Conference on World Wide Web, WWW 2016, Montreal, Canada, April 11-15, 2016, Companion Volume*, pages 463–466. ACM, 2016.
- [31] G. Kul, D. T. A. Luong, T. Xie, V. Chandola, O. Kennedy, and S. Upadhyaya. Similarity metrics for sql query clustering. *IEEE Transactions on Knowledge and Data Engineering*, 30(12):2408–2420, Dec 2018.
- [32] G. Kul, S. J. Upadhyaya, and V. Chandola. Detecting data leakage from databases on android apps with concept drift. In *IEEE TrustCom*, pages 905–913, 2018.
- [33] D. D. Lee and H. S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788, 1999.

- [34] V. H. Makiyama, J. Raddick, and R. D. C. Santos. Text mining applied to SQL queries: A case study for the SDSS skyserver. In J. A. Lossio-Ventura and H. Alatrista-Salas, editors, *Proceedings of the 2nd Annual International Symposium on Information Management and Big Data - SIMBig 2015, Cusco, Peru, September 2-4, 2015.*, volume 1478 of *CEUR Workshop Proceedings*, pages 66–72. CEUR-WS.org, 2015.
- [35] M. Mampaey, J. Vreeken, and N. Tatti. Summarizing data succinctly with the most informative itemsets. *TKDD*, 6(4):16:1–16:42, 2012.
- [36] S. Mittal, J. S. V. Varman, G. Chatzopoulou, M. Eirinaki, and N. Polyzotis. Querie: A query recommender system supporting interactive database exploration. In W. Fan, W. Hsu, G. I. Webb, B. Liu, C. Zhang, D. Gunopoulos, and X. Wu, editors, *ICDMW 2010, The 10th IEEE International Conference on Data Mining Workshops, Sydney, Australia, 13 December 2010*, pages 1411–1414. IEEE Computer Society, 2010.
- [37] A. Moffat and J. Zobel. Compression and fast indexing for multi-gigabyte text databases. *Australian Computer Journal*, 26(1):1–9, 1994.
- [38] B. O’Donoghue, E. Chu, N. Parikh, and S. P. Boyd. Conic optimization via operator splitting and homogeneous self-dual embedding. *J. Optimization Theory and Applications*, 169(3):1042–1068, 2016.
- [39] A. Pavlo, G. Angulo, J. Arulraj, H. Lin, J. Lin, L. Ma, P. Menon, T. C. Mowry, M. Perron, I. Quah, S. Santurkar, A. Tomasic, S. Toor, D. V. Aken, Z. Wang, Y. Wu, R. Xian, and T. Zhang. Self-driving database management systems. In *CIDR*, 2017.
- [40] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

- [41] P. Skibinski and J. Swacha. Fast and efficient log file compression. In Y. E. Ioannidis, B. Novikov, and B. Rachev, editors, *Communications of the Eleventh East-European Conference on Advances in Databases and Information Systems, Varna, Bulgaria, September 29 - October 3, 2007*, volume 325 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2007.
- [42] D. Wang, S. Zhu, T. Li, and Y. Gong. Multi-document summarization using sentence-based topic models. In *ACL 2009, Proceedings of the 47th Annual Meeting of the Association for Computational Linguistics and the 4th International Joint Conference on Natural Language Processing of the AFNLP, 2-7 August 2009, Singapore, Short Papers*, pages 297–300. The Association for Computer Linguistics, 2009.
- [43] I. H. Witten, A. Moffat, and T. C. Bell. *Managing Gigabytes: Compressing and Indexing Documents and Images*. Van Nostrand Reinhold, 1994.
- [44] K. Wu, E. J. Otoo, and A. Shoshani. Compressing bitmap indexes for faster search operations. In *Proceedings of the 14th International Conference on Scientific and Statistical Database Management, July 24-26, 2002, Edinburgh, Scotland, UK*, pages 99–108. IEEE Computer Society, 2002.
- [45] T. Xie, O. Kennedy, and V. Chandola. Query log compression for workload analytics. *CoRR*, abs/1809.00405, 2018.
- [46] X. Yang, C. M. Procopiuc, and D. Srivastava. Recommending join queries via query log analysis. In Y. E. Ioannidis, D. L. Lee, and R. T. Ng, editors, *Proceedings of the 25th International Conference on Data Engineering, ICDE 2009, March 29 2009 - April 2 2009, Shanghai, China*, pages 964–975. IEEE Computer Society, 2009.
- [47] J. Ziv and A. Lempel. A universal algorithm for sequential data compression. *IEEE Trans. Information Theory*, 23(3):337–343, 1977.

- [48] J. Ziv and A. Lempel. Compression of individual sequences via variable-rate coding. *IEEE Trans. Information Theory*, 24(5):530–536, 1978.
- [49] J. Zobel and A. Moffat. Inverted files for text search engines. *ACM Comput. Surv.*, 38(2):6, 2006.

Appendix A

Appendix

todo