

A Formal Framework for Probabilistic Unclean Databases

Christopher De Sa,¹ Ihab F. Ilyas², Benny Kimelfeld³,
Christopher Ré⁴, and Theodoros Rekatsinas⁵

¹ Cornell, ² University of Waterloo, ³ Technion, ⁴ Stanford, ⁵ University of Wisconsin-Madison

Abstract

Most theoretical frameworks that focus on data errors and inconsistencies follow logic-based reasoning. Yet, practical data cleaning tools need to incorporate statistical reasoning to be effective in real-world data cleaning tasks. Motivated by these empirical successes, we propose a formal framework for unclean databases, where two types of statistical knowledge are incorporated: The first represents a belief of how intended (clean) data is generated, and the second represents a belief of how noise is introduced in the actual observed database. To capture this noisy channel model, we introduce the concept of a Probabilistic Unclean Database (PUD), a triple that consists of a probabilistic database that we call *the intention*, a probabilistic data transformer that we call *the realization* and captures how noise is introduced, and a dirty observed database that we call *the observation*. We define three computational problems in the PUD framework: cleaning (infer the most probable intended database given a PUD), probabilistic query answering (compute the probability of an answer tuple over the unclean observed database), and learning (estimate the most likely intention and realization models of a PUD, given examples as training data). We illustrate the PUD framework on concrete representations of the intention and realization, show that they generalize traditional concepts of repairs such as cardinality and value repairs, draw connection to consistent query answering, and prove tractability results. We further show that parameters can be learned in some practical instantiations, and in fact, prove that under certain conditions we can learn a PUD directly from a single dirty database without any need for clean examples.

2012 ACM Subject Classification General literature

Keywords and phrases Unclean databases, data cleaning, probabilistic databases

Digital Object Identifier 10.4230/LIPIcs.CVIT.2016.23

1 Introduction

Managing errors and inconsistency in databases is traditionally viewed as a challenge of a logical nature. It is typical that errors in a database are defined with respect to *integrity constraints* that capture normative aspects of downstream applications. The aim of integrity constraints is to guarantee the consistency of data used by these applications. Typically, an unclean database is defined as a database instance J that violates the underlying set of integrity constraints. In turn, a *repair* of database instance J is a clean database instance I wherein all integrity constraints hold, and is obtained from J by a set of operations (e.g., deletions of tuples or updates of tuple values) that feature some form of non-redundancy [2, 4].

Various computational problems around unclean databases have been investigated in prior work [4, 29, 30, 33]. Past theoretical research has established fundamental results that concentrate on tractability boundaries for repair-checking and consistent query answering [2, 15, 27]. In their majority, these works adopt a deterministic interpretation of data repairs and cast all repairs equally likely. These theoretical developments have inspired practical tools that aim to automate data cleaning [7, 12, 20, 43, 44]. The majority of proposed methods



© Christopher De Sa and Ihab F. Ilyas and Benny Kimelfeld and Christopher Ré and Theodoros Rekatsinas;

licensed under Creative Commons License CC-BY

42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:31

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

44 assume as input a set of integrity constraints and use those to identify possible repairs via
 45 search-based procedures. To prioritize across possible repairs during search, the proposed
 46 methods rely on the notion of *minimality* [11, 23, 31]. Informally, minimality states that given
 47 two candidate sets of repairs, the one with fewer changes with respect to the original database
 48 is preferable. The use of minimality as an operational principle to find data repairs is a
 49 practical artifact that is used to limit the search space. These approaches to data cleaning
 50 suffer from two major drawbacks: First, they do not permit concrete statements about the
 51 “likelihood” of possible repairs. Consequently, they categorize query answers to a limited set
 52 of validity labels (e.g., certain, possible, and impossible); these labels might be unsuitable for
 53 downstream applications. Second, combinatorial principles such as minimality, while desired,
 54 do not entail the richness of the arguments and evidences (e.g., statistical features of data)
 55 that are needed to reason about and generate correct repairs.

56 Effective data cleaning needs to incorporate statistical reasoning. Our recent work on
 57 HoloClean [35] casts data repairing as a statistical learning and inference problem and
 58 reasons about *the most probable repair* instead of the minimal repair. Our study shows that
 59 HoloClean obtains more accurate data cleaning results than competing minimality-based
 60 data cleaning tools for a diverse array of real-world data cleaning scenarios [35]. HoloClean
 61 uses training data to learn a probabilistic model for how clean data is generated and how
 62 data errors are injected. HoloClean’s model follows the *noisy channel model* [22], the de-facto
 63 probabilistic framework used in natural language tasks, such as spell checking and speech
 64 recognition, to reason about noisy data. Existing theoretical frameworks for data cleaning
 65 do not capture this type of probabilistic reasoning.

66 **Goals.** We aim to establish a formal framework for *probabilistic unclean databases* (PUD)
 67 that adopts a statistical view of database cleaning. We do so by following the aforementioned
 68 noisy channel paradigm of HoloClean. Within the PUD framework, we formalize fundamental
 69 computational problems: *cleaning*, *query answering*, and *learning*. With that, we aim to draw
 70 connections between theoretical database research and important aspects of practical systems.
 71 In particular, our goal is to open the way for analyses and algorithms with theoretical
 72 guarantees for such systems. We argue that our framework is basic enough to allow for
 73 nontrivial theoretical advances, as illustrated by our preliminary results that (a) draw
 74 connections to traditional deterministic concepts, and (b) devise algorithms for special cases.

75 **Probabilistic Unclean Databases.** We view an unclean database as if a clean database
 76 I had been “distorted” via a noisy channel into a dirty database J ; our goal is to establish a
 77 model of this channel. Given the observed unclean database J , we seek the true database
 78 I from which J is produced. This model adopts Bayesian inference: out of all possible I ,
 79 we seek the one for which the probability given J is highest. Following Bayes’ rule, our
 80 objective is to find $\arg \max_I \Pr(I) \cdot \Pr(J|I)$. This objective decomposes in two parts: (1) the
 81 *prior* model for a clean database captured by $\Pr(I)$, and (2) the *channel* or *error* model
 82 characterized by $\Pr(J|I)$. To capture that, we define a *Probabilistic Unclean Database (PUD)*
 83 as a triple $(\mathcal{I}, \mathcal{R}, J^*)$ where: (1) \mathcal{I} , referred to as the *intention model*, is a distribution that
 84 produces intended clean database instances; (2) \mathcal{R} , referred to as the *realization model*,
 85 is a function that maps each clean database instance I to a distribution \mathcal{R}_I that defines
 86 how noise is introduced into I ; and (3) J^* is an observed unclean database instance. The
 87 distribution \mathcal{I} defines the prior $\Pr(I)$ over clean instances, while the distribution \mathcal{R}_I defines
 88 the aforementioned noisy channel $\Pr(J|I)$.

89 **Computational problems.** We define and study three computational problems in the
 90 PUD framework: (1) *data cleaning*, where given a PUD $(\mathcal{I}, \mathcal{R}, J)$, we seek to compute a

91 database instance I that maximizes the probability $\mathcal{R}_I(J) \cdot \mathcal{I}(I)$; (2) *probabilistic query*
92 *answering*, i.e., the problem of evaluating a query Q over a PUD following the traditional
93 possible tuple semantics [13, 40]; and (3) *learning* a PUD, where we consider *parametric*
94 representations \mathcal{I}_Ξ and \mathcal{R}_Θ of the intention and realization models and seek to estimate the
95 parameter vectors Ξ^* and Θ^* that maximize the likelihood of training data.

96 **Preliminary analysis.** PUDs allow for different instantiations of the intention and real-
97 ization models. To establish preliminary complexity and convergence results, we focus on
98 specific instantiations of the intention and realization models. We study intention models
99 that can describe the distribution of tuple values as well as both soft and hard integrity
100 constraints. We also focus on simple noise models. We study (1) realizations that introduce
101 new tuples, hence, the clean database is a subset of the observed unclean database, and (2)
102 realizations that update table cells, hence, the clean database is obtained via value repairs
103 over the observed unclean database.

104 We present PUD instantiations for which solving the data cleaning problem has polynomial
105 complexity. For instance, we show that in the presence of only one key constraint, soft or
106 hard, data cleaning in PUDs can be solved in polynomial time. This result extends results for
107 deterministic repairs that focus on hard integrity constraints to *weak* (soft) key constraints
108 (e.g., two people are unlikely, but might, have the same first and last name). Here, the most
109 probable repair under the PUD framework may violate weak key constraints. We also draw
110 connections between data cleaning in the PUD framework and *minimal repairs*. We identify
111 conditions under which data cleaning in the PUD framework is equivalent to *cardinality*
112 *repairs* [33] and *optimal V-repairs* [23]. For PUD learning, we consider both *supervised* and
113 *unsupervised* learning. In the former case, we are given intention-realization pairs, and in the
114 former, we are given only realizations (i.e., dirty databases). Our results discuss convexity
115 and gradient computation for the optimization problem underlying the learning problem.

116 Our PUD model can be viewed as a generalization of the approach of Gribkoff et al. [19],
117 who view the dirty database as a tuple-independent probabilistic database [13], and seek the
118 *most-probable database* that satisfies a set of underlying integrity constraints (e.g., functional
119 dependencies). In contrast, our modeling allows for arbitrary distributions over the intention,
120 including ones with *weak* constraints that we discuss later on. Interestingly, our PUD model
121 goes in the reverse direction of the *operational* approach of Calautti et al. [10], who view
122 the dirty database as a deterministic object and its *cleaning* (rather than the *error*) as a
123 probabilistic process (namely a Markov chain of repairing operations).

124 **Vision.** This paper falls within the bigger vision of bridging database theory with learning
125 theory as outlined in a recent survey [1]. We aim to draw connections between the rich
126 theory on inconsistency management by the database community, and fundamentals of
127 statistical learning theory with emphasis on *structured prediction* [5]. Structured prediction
128 typically focuses on problems where, given a collection of observations, one seeks to predict
129 the most likely assignment of values to structured objects. In most practical structured
130 prediction problems, structure is encoded via logic-based constraints [18] in a way similar to
131 how consistency is enforced in data cleaning. It is our hope that this paper will commence a
132 line of work towards theoretical developments that take the benefit of both worlds, and will
133 lead to new techniques that are both practical and rooted in strong foundations.

134 **Organization.** We begin with preliminary definitions in Section 2. In Section 3 we present
135 the concept of PUDs. We present the three fundamental computational problems in Section 4,
136 and describe preliminary results in Sections 5 and 6. We conclude with a discussion in
137 Section 7. For space limitations, all proofs are in the Appendix.

138 **2 Preliminaries**

139 We first introduce concepts, definitions and notation that we need throughout the paper.

140 **Schemas and databases.** A *relation signature* is a sequence $\alpha = (A_1, \dots, A_k)$ of distinct
 141 *attributes* A_i , where k is the *arity* of α . A (*relational*) *schema* \mathbf{S} has a finite set of *relation*
 142 *symbols*, and it associates each relation symbol R with a signature that we denote by $\text{sig}_{\mathbf{S}}(R)$,
 143 or just $\text{sig}(R)$ if \mathbf{S} is clear from the context. We assume an infinite domain Const of *constants*.
 144 Let \mathbf{S} be a schema, and let R be a relation symbol of \mathbf{S} . A *tuple* t over R is a sequence
 145 (c_1, \dots, c_k) of constants, where k is the arity of $\text{sig}(R)$. If $t = (c_1, \dots, c_k)$ is a tuple over R
 146 and $\text{sig}(R) = (A_1, \dots, A_k)$, then we refer to the value c_j as $t.A_j$ (where $j = 1, \dots, k$). We
 147 denote by $\text{tuples}(R)$ the set of all tuples over R .

148 In our databases, tuples have unique record IDs. Formally, a table r over R is associated
 149 with a finite set $\text{ids}(r)$ of identifiers, and it maps each identifier i to a tuple $r[i]$ over R . A
 150 *database* I over \mathbf{S} consists of a table R^I over each relation symbol R of \mathbf{S} , such that no two
 151 occurrences of tuples have the same identifier; that is, if R_1 and R_2 are distinct relation
 152 symbols in \mathbf{S} , then $\text{ids}(R_1^I)$ and $\text{ids}(R_2^I)$ are disjoint sets. We denote by $\text{ids}(I)$ the union of
 153 the sets $\text{ids}(R^I)$ over all relation symbols R of \mathbf{S} . If $i \in \text{ids}(R^I)$, then we may refer to the
 154 tuple $R^I[i]$ simply as $I[i]$.

155 A *cell* of a database I is a pair (i, A) , where $i \in \text{ids}(R^I)$ for a relation symbol R , and A
 156 is an attribute in $\text{sig}(R)$. We denote the cell (i, A) also by $i.A$, and we denote by $\text{cells}(I)$ the
 157 set of all cells of I .

158 Let I and J be databases over the same schema \mathbf{S} . We say that I is a *subset* of J if I can
 159 be obtained from J by deleting tuples, that is, $\text{ids}(R^I) \subseteq \text{ids}(R^J)$ for all relation symbols R
 160 of \mathbf{S} (hence, $\text{ids}(I) \subseteq \text{ids}(J)$) and $I[i] = J[i]$ for all $i \in \text{ids}(I)$. We say that I is an *update* of
 161 J if I can be obtained from J by changing attribute values, that is, $\text{ids}(R^I) = \text{ids}(R^J)$ for
 162 all relation symbols R of \mathbf{S} .

163 A *query* Q over a schema \mathbf{S} is associated with fixed arity, and it maps every database D
 164 over \mathbf{S} into a finite set $Q(D)$ of tuples of constants over the fixed arity.

165 **Integrity constraints.** Various types of logical conditions are used for declaring integrity
 166 constraints, including *Functional Dependencies* (FDs), *conditional FDs* [7], *Denial Constraints*
 167 (DCs) [17], referential constraints [14], and so on. In this paper, by *integrity constraint* over
 168 a schema \mathbf{S} we refer to a general expression φ of the form $\forall x_1, \dots, x_m[\gamma(x_1, \dots, x_m)]$, where
 169 $\gamma(x_1, \dots, x_m)$ is a safe expression in Tuple Relational Calculus (TRC) over \mathbf{S} . For example,
 170 an FD $R : A \rightarrow B$ is expressed here as the integrity constraint

$$171 \quad \forall x, y [(x \in R \wedge y \in R) \rightarrow (x.A = y.A \rightarrow x.B = y.B)].$$

172 A *violation* of $\varphi = \forall x_1, \dots, x_m[\gamma(x_1, \dots, x_m)]$ in the database I is a sequence i_1, \dots, i_m of
 173 tuple identifiers in $\text{ids}(I)$ such that I violates $\gamma(I[i_1], \dots, I[i_m])$, and we denote by $V(\varphi, I)$
 174 the set of violations of φ in I . We say that I *satisfies* φ if I has no violations of φ , that
 175 is, $V(\varphi, I)$ is empty. Finally, I *satisfies* a set Φ of integrity constraints if I satisfies every
 176 integrity constraint φ in Φ .

177 **Minimum repairs.** Traditionally, database *repairs* are defined over inconsistent databases,
 178 where *inconsistencies* are manifested as violations of integrity constraints. A repair is a
 179 consistent database that is obtained from the inconsistent one by applying a *minimal* change,
 180 and we recall two types of repairs: *subset* (obtained by deleting tuples) and *update* (obtained
 181 by changing values). Moreover, the repairing operations may be weighted by tuple weights
 182 (in the first case) and cell weights (in the second case).

183 Formally, let \mathbf{S} be a schema, Φ a set of integrity constraints over \mathbf{S} , and J a database
 184 that does not necessarily satisfy Φ . A *consistent subset* (resp., *consistent update*) of J is a
 185 subset (resp., update) I of J such that I satisfies Φ . A *minimum subset repair* of J w.r.t. a
 186 weight function $w : ids(J) \rightarrow [0, \infty)$ is a consistent subset I of J that minimizes the sum
 187 $\sum_{i \in ids(J) \setminus ids(I)} w(i)$. As a special case, a *cardinality repair* of J is a minimum subset repair
 188 w.r.t. a constant weight (e.g., $w(i) = 1$), that is, a consistent subset with a maximal number of
 189 tuples. A *minimum update repair* of J w.r.t. a weight function $w : cells(J) \times \text{Const} \rightarrow [0, \infty)$
 190 is a consistent update I of J that minimizes the sum $\sum_{i.A \in cells(I)} w(i.A, I[i].A)$.

191 **Probabilistic Databases.** A *probabilistic database* is a probability distribution over
 192 ordinary databases. As a representation system, our model is a generalization of the *Tuple-*
 193 *Independent probabilistic Database* (TID) wherein each tuple might be either exist (with
 194 an associated probability) or not [13, 40]. In our model, each tuple comes from a general
 195 probability distribution over tuples (where inexistence is one of the options). This allows us
 196 to incorporate beliefs about the likelihood of tuples and cell values.

197 Formally, let \mathbf{S} be a schema. A *generalized tuple-independent database*, or just *generalized*
 198 *TID* for short, is a database \mathcal{K} that is defined similarly to an ordinary database over \mathbf{S} ,
 199 except that instead of a tuple, the entry $R^{\mathcal{K}}[i]$ is a discrete probability distribution over
 200 tuples $(R) \cup \{\perp\}$, where the special value \perp denotes that no tuple is generated. Hence, for
 201 every tuple t over R , the probability that $R^{\mathcal{K}}[i]$ produces t is given by $R^{\mathcal{K}}[i](t)$, or just
 202 $\mathcal{K}[i](t)$; moreover, the number $\mathcal{K}[i](\perp)$ is the probability that no tuple is generated for the
 203 identifier i . Therefore, \mathcal{K} defines a probabilistic distribution over databases I over \mathbf{S} such
 204 that $ids(I) \subseteq ids(\mathcal{K})$ and the probability $\mathcal{K}(I)$ of a database I is defined as follows:

$$205 \quad \mathcal{K}(I) \stackrel{\text{def}}{=} \prod_{i \in ids(I)} \mathcal{K}[i](I[i]) \times \prod_{i \in ids(\mathcal{K}) \setminus ids(I)} \mathcal{K}[i](\perp)$$

206 We incorporate weak integrity constraints by adopting the standard concept of *parametric*
 207 *factors* (or *parfactors* for short), which has been used in the *soft keys* of Jha et al. [21] and
 208 the *PrDB* model of Sen et al. [38], and which can be viewed as a special case of the *Markov*
 209 *Logic Network* (MLN) [36]. Under this concept, each constraint φ is associated with a weight
 210 $w(\varphi) > 0$ and each violation of φ contributes a factor of $\exp(-w(\varphi))$ to the probability of a
 211 random database I . Formally, a *parfactor database* over a schema \mathbf{S} is a triple $\mathcal{D} = (\mathcal{K}, \Phi, w)$,
 212 where \mathcal{K} is a generalized TID, Φ is a finite set of integrity constraints, both over \mathbf{S} , and
 213 $w : \Phi \rightarrow (0, \infty)$ is a weight function over Φ . The probability $\mathcal{D}(I)$ of a database I is defined
 214 as follows.

$$215 \quad \mathcal{D}(I) \stackrel{\text{def}}{=} \frac{1}{Z} \times \mathcal{K}(I) \times \exp \left(- \sum_{\varphi \in \Phi} w(\varphi) \times |V(\varphi, I)| \right)$$

216 Recall that $V(\varphi, I)$ the set of violations of φ in I . The number Z is a *normalization factor*
 217 (also called the *partition function*) that normalizes the sum of probabilities to one:

$$218 \quad Z \stackrel{\text{def}}{=} \sum_I \mathcal{K}(I) \times \exp \left(- \sum_{\varphi \in \Phi} w(\varphi) \times |V(\varphi, I)| \right)$$

219 Observe that the above sum is over a countable domain, since we assume that every $R^{\mathcal{K}}[i]$
 220 is discrete (hence, there are countably many random databases I). Since we normalize the
 221 probability, it is not really necessary for \mathcal{K} to be normalized, as \mathcal{D} would be a probability
 222 distribution even if \mathcal{K} is *not* normalized. In fact, in our analysis we will *not* make the
 223 assumption that \mathcal{K} is normalized.

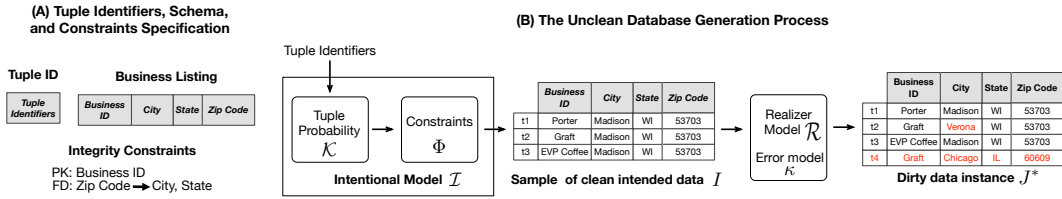


Figure 1 Overview of the PUD framework.

224 **3 Probabilistic Unclean Databases**

225 We introduce the Probabilistic Unclean Database (PUD) framework and describe example
 226 instantiations PUDs corresponding to data cleaning applications in the HoloClean system [35].
 227 In our framework, a PUD consists of three components following a noisy-channel model: (1)
 228 an *intention* model for generating clean database instances, (2) a noisy *realization* model
 229 that can distort the intended clean database instance, and (3) an observed unclean database.
 230 The formal definition follows.

231 ► **Definition 1.** Let \mathbf{S} be a schema. A PUD (over \mathbf{S}) is a triple $\mathcal{U} = (\mathcal{I}, \mathcal{R}, J^*)$ where:

- 232 1. \mathcal{I} is a probabilistic database, referred to as the *intention* model;
- 233 2. \mathcal{R} , referred to as the *realization* model, is a function that maps each database I to a
 234 probabilistic database \mathcal{R}_I ;
- 235 3. J^* is a database referred to as the *observed* or *unclean* database.

236 ► **Example 2.** Figure 1 illustrates a high-level example of the PUD framework. We use a
 237 running example from business listings. Figure 1(A) depicts the schema \mathbf{S} of the example.
 238 The constraints include a primary key and a functional dependency. Figure 1(B) depicts
 239 the unclean database generation process. Intention \mathcal{I} outputs a valid database I with three
 240 tuples. The realizer \mathcal{R} takes as input this instance I , injects the new tuple $t4$ and updates
 241 the *City* value of tuple $t2$ from “Madison” to “Verona.” ◀

242 A PUD $\mathcal{U} = (\mathcal{I}, \mathcal{R}, J^*)$ defines a probability distribution, denoted $\mathcal{R} \circ \mathcal{I}$, over pairs (I, J) .
 243 Conditioning on $J = J^*$, the PUD \mathcal{U} also defines a probability distribution, denoted \mathcal{U}^* , over
 244 intentions I (i.e., a probabilistic database). In the generative process of $\mathcal{R} \circ \mathcal{I}$, we sample the
 245 intention I from \mathcal{I} , and then we sample J from the realization \mathcal{R}_I . Hence, the probability of
 246 (I, J) is given by

247
$$\mathcal{R} \circ \mathcal{I}(I, J) \stackrel{\text{def}}{=} \mathcal{I}(I) \cdot \mathcal{R}_I(J).$$

248 In the probabilistic database \mathcal{U}^* , the probability of each candidate intention I' is given by

249
$$\mathcal{U}^*(I') \stackrel{\text{def}}{=} \Pr_{(I, J) \sim \mathcal{R} \circ \mathcal{I}}(I = I' \mid J = J^*) = \frac{\mathcal{R} \circ \mathcal{I}(I', J^*)}{\sum_I \mathcal{R} \circ \mathcal{I}(I, J^*)}$$

250 that is, the probability conditioned on the random J being J^* . For this distribution to be
 251 well defined, we require J^* to have nonzero probability; that is, there exists I such that
 252 $\mathcal{R} \circ \mathcal{I}(I, J^*) > 0$.

Business ID	City	State	Zip Code
Porter	Madison	WI	53703
Graft	Madison	WI	53703
EVP Coffee	Madison	WI	53703

Business ID	City	State	Zip Code
Porter	Madison	WI	53703
Graft	Madison	WI	53703
EVP Coffee	Madison	WI	53703
EVP Coffee	Madison	WI	53703

Business ID	City	State	Zip Code
Porter	Madison	WI	53703
Graft	Madison	WI	53704
EVP Coffee	adison	WI	53703

■ **Figure 2** Examples of a subset realizer and an update realizer.

3.1 Example Instantiations of PUDs

Our definition of a PUD is abstract, and not associated with any specific representation model. We now present concrete instantiations of PUD representations. These instantiations are probabilistic generalizations of the (deterministic) concepts of *subset repairs* [2,33] and *update repair* [24,31], respectively. More precisely, in both instantiations, the PUD $\mathcal{U} = (\mathcal{I}, \mathcal{R}, J^*)$ is such that \mathcal{I} is represented as a parfactor database \mathcal{D} (as defined in Section 2) and J^* is an ordinary database (as expected); the two differ in the representation of the realization model \mathcal{R} . In the first instantiation, \mathcal{R} is allowed to introduce new random tuples (hence, the intended database is a *subset* of the unclean one) and in the second, \mathcal{R} is allowed to randomly change tuples (hence, the intended database is an *update* of the unclean one). Formally, let \mathbf{S} be a schema.

- A *parfactor/subset* PUD is a triple (\mathcal{D}, τ, J^*) where \mathcal{D} is a parfactor database, τ maps every identifier $i \in \text{ids}(R^{\mathcal{D}})$, where $R \in \mathbf{S}$, to a discrete distribution $\tau[i]$ over $\text{tuples}(R) \cup \{\perp\}$, and J^* is an ordinary database. As usual, \perp means that no tuple is generated.
- A *parfactor/update* PUD is a triple $(\mathcal{D}, \kappa, J^*)$ where \mathcal{D} is a parfactor database, κ maps every identifier $i \in \text{ids}(R^{\mathcal{D}})$ and tuple $t \in \text{tuples}(R)$, where $R \in \mathbf{S}$, to a discrete distribution $\kappa[i, t]$ over $\text{tuples}(R)$, and J^* is an ordinary database.

In a parfactor/subset PUD $\mathcal{U} = (\mathcal{D}, \tau, J^*)$, the probability $\mathcal{R} \circ \mathcal{I}(I, J)$ is then defined as follows. If I is not a subset of J , then $\mathcal{R} \circ \mathcal{I}(I, J) = 0$; otherwise:

$$\mathcal{R} \circ \mathcal{I}(I, J) \stackrel{\text{def}}{=} \mathcal{D}(I) \times \prod_{\substack{i \in \text{ids}(J) \setminus \\ \text{ids}(I)}} \tau[i](J[i]) \times \prod_{\substack{i \in \text{ids}(\mathcal{D}) \setminus \\ \text{ids}(J)}} \tau[i](\perp)$$

That is, $\mathcal{R} \circ \mathcal{I}(I, J)$ is the probability of I (i.e., $\mathcal{D}(I)$), multiplied by the probability that each new tuple of J is produced by τ (i.e., $\tau[i](J[i])$), multiplied by the probability that each tuple identifier i missing in J is indeed not produced (i.e., $\tau[i](\perp)$).

In a parfactor/update PUD $\mathcal{U} = (\mathcal{D}, \kappa, J^*)$, the probability $\mathcal{R} \circ \mathcal{I}(I, J)$ is then defined as follows. If I is not an update of J , then $\mathcal{R} \circ \mathcal{I}(I, J) = 0$; otherwise:

$$\mathcal{R} \circ \mathcal{I}(I, J) \stackrel{\text{def}}{=} \mathcal{D}(I) \times \prod_{i \in \text{ids}(I)} \kappa[i, I[i]](J[i])$$

That is, $\mathcal{R} \circ \mathcal{I}(I, J)$ is the probability of I (i.e., $\mathcal{D}(I)$), multiplied by the probability that κ changes each tuple $I[i]$ to $J[i]$ (i.e., $\kappa[i, I[i]](J[i])$).

► **Example 3.** Figure 2 shows the intended database from Example 2 and two dirty versions obtained by a subset realizer and an update realizer. The update realizer introduces a duplicate, while the update realizer introduces two typos. These correspond to two types of common errors in relational data. Our PUD framework can naturally model such cases. ◀

287 In Section 5, we discuss connections between these PUD instantiations and the determin-
 288 istic models. Finally, in Section 6, we provide more concrete cases of PUD instantiations.

289 4 Computational Problems

290 We define three computational problems over PUDs that are motivated by the need to clean
 291 and query unclean data, and learn the intention and realization models from observed data.

292 **Data Cleaning.** Given a PUD $(\mathcal{I}, \mathcal{R}, J^*)$, we wish to compute a Most Likely Intention
 293 (MLI) database I , given the observed unclean database J^* . We refer to this problem as *data*
 294 *cleaning* in PUDs.

295 ► **Definition 4 (Cleaning).** Let \mathbf{S} be a schema and \mathbf{R} a representation system for PUDs. The
 296 problem (\mathbf{S}, \mathbf{R}) -*cleaning* is that of computing an *MLI* of a given PUD $\mathcal{U} = (\mathcal{I}, \mathcal{R}, J^*)$, that
 297 is, computing a database I such that the probability $\mathcal{U}^*(I)$ is maximal (or, equivalently, the
 298 probability $\mathcal{R} \circ \mathcal{I}(I, J^*)$ is maximal).

299 **Probabilistic query answering.** A PUD defines a probabilistic database—a probability
 300 space over the intensions I . The problem of *Probabilistic Query Answering (PQA)* is that
 301 of evaluating a query over this probabilistic database. We adopt the standard semantics
 302 of query evaluation over probabilistic databases [13, 40], where the confidence in an answer
 303 tuple is its marginal probability.

304 ► **Definition 5 (PQA).** Let \mathbf{S} be a schema, Q a query over \mathbf{S} , and \mathbf{R} a representation system
 305 for PUDs. The problem $(\mathbf{S}, Q, \mathbf{R})$ -*PQA* is the following. Given a PUD \mathcal{U} and a tuple \mathbf{a} ,
 306 compute the *confidence* of \mathbf{a} , that is, the probability $\Pr_{I \sim \mathcal{U}^*}(\mathbf{a} \in Q(I))$.

307 For now, we assume that both \mathcal{I} and \mathcal{R} are fully specified. We next define the problem
 308 of learning models \mathcal{I} and \mathcal{R} using training (potentially labeled) data.

309 **PUD learning.** For a PUD $(\mathcal{I}, \mathcal{R}, J^*)$, the models \mathcal{I} and \mathcal{R} are typically represented using
 310 numeric *parameters*. For example, the parameters of a parfactor/subset PUD (\mathcal{D}, τ, J^*) are
 311 those needed to represent \mathcal{D} (e.g., the weights of the constraints), and the parameters that
 312 define the distributions over the tuples in both \mathcal{D} and τ . By a *parametric intention* we refer
 313 to an intension model \mathcal{I}_{Ξ} with a vector Ξ of uninitialized parameters, and by $\mathcal{I}_{\Xi/\mathbf{c}}$ we denote
 314 the actual intention model where Ξ is assigned the values in the vector \mathbf{c} . Similarly, by a
 315 *parametric realization* we refer to a realization model \mathcal{R}_{Θ} with a vector Θ of uninitialized
 316 parameters, and by $\mathcal{R}_{\Theta/\mathbf{d}}$ we denote the actual realization model where Θ is set to \mathbf{d} .

317 Following the concept of *maximum likelihood estimation*, the goal in learning is to find
 318 the parameters that best explain (i.e., maximize the probability) of the training examples.
 319 In the *supervised* variant, we are given examples of both unclean databases and their clean
 320 versions; in the *unsupervised* variant, we are given only unclean databases.

321 ► **Definition 6 (Learning).** Let \mathbf{S} be a schema, and \mathbf{R} a representation system for parametric
 322 intensions and realizations. In the following problems, we are given as part of the input
 323 parametric intention and realization models \mathcal{I}_{Ξ} and \mathcal{R}_{Θ} , respectively.

- 324 ■ In the *supervised* (\mathbf{S}, \mathbf{R}) -*learning* problem we are also given a collection $(I_j, J_j)_{j=1}^n$ of
 325 database pairs (intention-realization examples), and the goal is to find parameter values
 326 \mathbf{c} and \mathbf{d} that maximize $\prod_{j=1}^n \mathcal{R} \circ \mathcal{I}(I_j, J_j)$ for $\mathcal{I} = \mathcal{I}_{\Xi/\mathbf{c}}$ and $\mathcal{R} = \mathcal{R}_{\Theta/\mathbf{d}}$.
- 327 ■ In the *unsupervised* (\mathbf{S}, \mathbf{R}) -*learning* problem we are also given a collection $(J_j)_{j=1}^n$ of
 328 databases (realization examples), and the goal is to find parameter values \mathbf{c} and \mathbf{d} that

329 maximize $\prod_{j=1}^n \mathcal{R} \circ \mathcal{I}(J_j)$ for $\mathcal{I} = \mathcal{I}_{\Xi/c}$ and $\mathcal{R} = \mathcal{R}_{\Theta/d}$, where $\mathcal{R} \circ \mathcal{I}(J_j)$ is the marginal
 330 probability of J_j , that is, $\sum_I \mathcal{R} \circ \mathcal{I}(I, J_j)$.

331 Note that the summation in the unsupervised variant is over the sample space of the
 332 intention model \mathcal{I} . While the reader might be concerned about the source of many examples
 333 (I_j, J_j) and J_j in the phrasing of the learning problems, it is oftentimes the case that a single
 334 large example (I, J) (or just J in the unsupervised variant) can be decomposed into many
 335 smaller examples. This depends on the independence assumptions in the parametric models
 336 \mathcal{I}_{Ξ} and \mathcal{R}_{Θ} as we discuss in Section 6.1. In the next sections, we give preliminary results on
 337 the above problems, focusing on parfactor/subset and parfactor/update PUDs.

338 5 Cleaning and Querying Unclean Data

339 We draw connections between data cleaning in the PUD framework (MLIs) and traditional
 340 *minimum repairs*. We also give preliminary results on the complexity of cleaning. Finally,
 341 we draw a connection between probabilistic query answering and certain answers.

342 5.1 Generalizing Minimum Repairs

343 We now show that the concept of an MLI in parfactor/subset PUDs generalizes the concept of
 344 a minimum subset repair, and the concept of an MLI in parfactor/update PUDs generalizes
 345 the concept of an optimal update repair. Minimum subset repairs correspond to MLIs of
 346 PUDs with hard (or heavy) constraints. Minimum update repairs correspond to MLIs over
 347 PUDs that assume both hard (or heavy) constraints, and assumptions of independence among
 348 the attributes. From the viewpoint of computational complexity, this means that finding
 349 an exact MLI is not easier than finding a minimum repair, which is often computationally
 350 hard [31]. Therefore, we should aim for *approximation* guarantees (which have clear semantics
 351 in the probabilistic setting) if we wish to avoid restricting the generality of the input.

352 **Subset repairs and parfactor/subset PUDs.** Recall that in parfactor/subset PUDs
 353 (as defined in Section 3.1), every intention I with a nonzero probability is a subset of the
 354 observed unclean database J^* . In particular, every MLI is subset of J^* . Our first result
 355 relates cleaning in parfactor/subset PUDs to the traditional minimum subset (or *cardinality*)
 356 repairs. This result states, intuitively, that the notion of an MLI in a parfactor/subset PUD
 357 converges with the notion of a minimum subset repair if the weight of the formulas is high
 358 enough and the probability of introducing error is small enough.

359 ► **Theorem 7.** *Let (\mathcal{D}, τ, J^*) be a parfactor/subset PUD with $\mathcal{D} = (\mathcal{K}, \Phi, w)$. For $i \in \text{ids}(J^*)$,
 360 assume that $\mathcal{K}[i](\perp) > 0$ and $\tau[i](J^*[i]) > 0$, let $q(i) = \mathcal{K}[i](J^*[i]) / (\mathcal{K}[i](\perp) \cdot \tau[i](J^*[i]))$, and
 361 assume that $q(i) \geq 1$. There is a number M such that if $w(\varphi) > M$ for all $\varphi \in \Phi$ then the
 362 following are equivalent for all $I \subseteq J^*$:*

- 363 1. I is an MLI.
- 364 2. I is a minimum subset repair of J^* w.r.t. the weight function $w(i) = \log(q(i))$.

365 Note that in the theorem, $q(i)$ is the ratio between $\mathcal{K}[i](J^*[i])$, namely the probability that
 366 \mathcal{K} produces the i th tuple of J^* , and $\mathcal{K}[i](\perp) \cdot \tau[i](J^*[i])$, namely the probability that \mathcal{K} does
 367 not generate the i th tuple of J^* but τ does.

368 Next, we draw a similar connection between minimum update repairs and MLIs of
 369 parfactor/update PUDs.

370 **Update repairs and parfactor/update PUDs.** We now turn our attention to update
 371 repairs. Recall that in a parfactor/update PUD (defined in Section 3.1), the intended clean

372 database I is assumed to be an update of the observed unclean database J^* . We establish a
 373 result analogous to Theorem 7, stating conditions under which MLIs for parfactor/update
 374 PUDs converge to traditional minimum update repairs.

375 Let \mathbf{S} be a schema, and let $\mathcal{U} = (\mathcal{D}, \kappa, J^*)$ be a parfactor/update PUD with $\mathcal{D} = (\mathcal{K}, \Phi, w)$.
 376 We say that \mathcal{U} is *attribute independent* if there is probabilistic independence between the
 377 attributes in each $\mathcal{K}[i]$ and $\kappa[i, t]$. More precisely, if $i \in R^{\mathcal{D}}$ for $R \in \mathbf{S}$ with $\text{sig}(R) =$
 378 (A_1, \dots, A_k) , then we assume that $\mathcal{K}[i](a_1, \dots, a_k)$ can be written as

$$379 \quad \mathcal{K}[i](a_1, \dots, a_k) = \prod_{j=1}^k \mathcal{K}_{A_j}[i](a_j)$$

380 and, for $t = (b_1, \dots, b_k)$, that $\kappa[i, t](a_1, \dots, a_k)$ can be written as

$$381 \quad \kappa[i, t](a_1, \dots, a_k) = \prod_{j=1}^k \kappa_{A_j}[i, b_j](a_j).$$

382 In particular, the choice of the value a_j depends only on b_j and not on other values $b_{j'}$. The
 383 following theorem states that the concept of an MLI of a parfactor/update PUD converges
 384 with the concepts of a minimum update repair when the PUD is attribute independent and,
 385 moreover, the weight of the integrity constraints is high.

386 ► **Theorem 8.** *Let $\mathcal{U} = (\mathcal{D}, \tau, J^*)$ be an attribute-independent parfactor/update PUD with
 387 $\mathcal{D} = (\mathcal{K}, \Phi, w)$ such that $\mathcal{U}^*(I) > 0$ for at least one consistent update I of J^* . There is a
 388 number M such that if $w(\varphi) > M$ for all $\varphi \in \Phi$, then these are equivalent for all $I \subseteq J^*$:*

- 389 1. I is an MLI.
- 390 2. I is a minimum update repair w.r.t. the weight function

$$391 \quad w(i.A, a) = -\log(\mathcal{K}_A[i](a) \cdot \kappa_A[i, a](J^*[i].A)).$$

392 Note that $\mathcal{K}_A[i](a) \cdot \kappa_A[i, a](J^*[i].A)$ is the probability that a is produced by \mathcal{K} for the cell
 393 $i.A$, and that a is then changed to $J^*[i].A$ via κ . Note that in the case where this product is
 394 zero, we slightly abuse the notation by assuming that the weight is infinity.

395 5.2 Complexity of Cleaning with Key Constraints

396 We now present a complexity result on computing an MLI of a parfactor/subset PUD in
 397 the presence of key constraints. The following theorem states that in the case of a single
 398 key constraint per relation (which is the common setup, e.g., for the analysis of certain
 399 query answering [3, 25, 27]), an MLI can be found in polynomial time. Note that we do not
 400 make any assumption about the parameters; in particular, it may be the case that an MLI
 401 violates the key constraints since the constraints are weak. Regarding the representation of
 402 the probability spaces \mathcal{K} and τ , the only assumption we make is that, given a tuple t , the
 403 probabilities $\mathcal{K}[i](t)$ and $\tau[i](t)$ can be computed in polynomial time.

404 ► **Theorem 9.** *Let (\mathcal{D}, τ, J^*) be a parfactor/subset PUD with $\mathcal{D} = (\mathcal{K}, \Phi, w)$. If Φ consists
 405 of (at most) one key constraint per relation, and no relation of J^* has duplicate tuples, then
 406 an MLI can be computed in polynomial time.*

407 It is left for future investigation to seek additional constraints (e.g., functional depen-
 408 dencies) for which an MLI can be found in polynomial time. Note that Theorem 7 implies
 409 that (under conventional complexity assumptions) we cannot generalize the polynomial-time
 410 result to all sets of functional dependencies, since finding a minimum subset repair might be
 411 computationally hard [31, 33].

5.3 Probabilistic Query Answering

For probabilistic query answering, we again focus on the parametric/subset PUDs, and now we draw a connection to *consistent query answering* over the cardinality repairs. Recall that a *consistent answer* for a query Q over an inconsistent database J is a tuple t that belongs to $Q(I)$ for every cardinality repair I of J .

Let \mathbf{S} be a schema, J^* a database, and Φ a set of integrity constraints. Let $M = |\text{ids}(J^*)|$. The *uniform* parfactor/subset PUD for J^* and Φ with the parameters p and u , denoted $\mathcal{U}_{p,u}(J^*, \Phi)$ or just $\mathcal{U}_{p,u}$ if J^* and Φ are clear from the context, is the parfactor/subset PUD (\mathcal{D}, τ, J^*) with $\mathcal{D} = (\mathcal{K}, \Phi, w)$ such that the following hold.

- For all $R \in \mathbf{S}$ we have $\text{ids}(R^{\mathcal{K}}) = \text{ids}(R^{J^*})$ and $\mathcal{K}[i](x) = 1/M$ for every tuple identifier i and argument x in $R^{\mathcal{K}} \cup \{\perp\}$. The remaining mass (required to reaching 1) is given to an arbitrary tuple outside of J^* .
- $w(\varphi) = u$ for every $\varphi \in \Phi$.
- $\tau[i](\perp) = p$, and $\tau[i](t) = (1-p)/M$ for every identifier i and tuple t .

Observe that \mathcal{D} is defined in such a way that every subset I of J^* has the same prior probability $\mathcal{K}(I)$, namely $1/M^{|J^*|}$.

The following theorem states that, for $\mathcal{U}_{p,u}$, the consistent answers are precisely the answers whose probability approaches one when all of the following hold: (1) the probability of introducing error (i.e., $1-p$) approaches zero; and (2) the weight of the weak constraints (i.e., u) approaches infinity, that is, the constraints strengthen towards hardness.

► **Theorem 10.** *Let J^* be a database, Φ a set of integrity constraints, Q a query, and t a tuple. The following are equivalent:*

1. t is a consistent answer over the cardinality repairs.
2. $\lim_{p \rightarrow 1} \lim_{w \rightarrow \infty} \Pr_{I \sim \mathcal{U}_{p,u}^*}(t \in Q(I)) = 1$.

Therefore, Theorem 10 sheds light on the role that the consistent answers have in probabilistic query answering over parfactor/subset PUDs.

6 Learning Probabilistic Unclean Databases

We now give preliminary results on PUD learning, focusing on parfactor/update PUDs. We begin by describing the setup we consider in this section and the representation system \mathbf{R} we use to describe the parametric intention and realization models \mathcal{I}_{Ξ} and \mathcal{R}_{Θ} .

6.1 Setup

To discuss the learning of parameters, we need to specify the actual parametric model we assume. Let \mathbf{S} be a schema, and let $\mathcal{U} = (\mathcal{D}, \kappa, J^*)$ be a parfactor/update PUD with $\mathcal{D} = (\mathcal{K}, \Phi, w)$. Since we restrict the discussion to parfactor/update PUDs, we assume (without loss of generality) that the identifiers in \mathcal{D} are exactly those in J^* , that is, $\text{ids}(\mathcal{D}) = \text{ids}(J^*)$.

In our setup, \mathcal{K} of \mathcal{D} and κ of \mathcal{U} are expressed in a parametric form that allows us to define the parametric intention and realization models \mathcal{I}_{Ξ} and \mathcal{R}_{Θ} as *Gibbs* distributions. We refer to these as *Gibbs parfactor/update PUD models*, and define them as follows.

Parametric intention. To specify \mathcal{K} , we assume that for each relation symbol $R \in \mathbf{S}$, the probability $\mathcal{K}[i](t)$, with $i \in \text{ids}(R^{\mathcal{D}})$, is expressed in the form of an exponential distribution $\mathcal{K}[i](t) = \exp\left(-\sum_{f \in F} w_f f(t)\right)$ where each $f \in F$ is an arbitrary function (*feature*) over t , and each weight w_f is a real number.

454 For example, a feature $f \in F$ may be a function that takes as input a tuple and returns
 455 a value in $\{-1, 1\}$. An example feature f can state that $f(t) = 1$ if $t[\text{gender}] = \text{female}$ and,
 456 otherwise, $f(t) = -1$. Another example is $f[t] = 1$ if $t[\text{zip}]$ starts with 53 and $t[\text{state}] = \text{WI}$,
 457 and otherwise $f[t] = -1$. Additional examples of such features include the ones used in our
 458 prior work on HoloClean [35] to capture the co-occurrence probability of attribute value
 459 pairs. Each weight w_f corresponds to a parameter of the model. An assignment to these
 460 weights gives as a probability distribution, similarly to probabilistic graphical models [26].

461 The parameter vector Ξ of the parametric intention model \mathcal{I}_Ξ consists of two sets of
 462 parameters: the weights w_f for each feature $f \in F$, and the weights $w(\varphi)$, which we write as
 463 w_φ for uniformity of presentation, for each constraint $\varphi \in \Phi$. Thus, the overall parametric
 464 intention model \mathcal{I}_Ξ is expressed as a parametric Gibbs distribution.

465 We will take a special interest in the case where the integrity constraints are *unary*, which
 466 means that they have the form $\forall x[\gamma(x)]$, where γ is quantifier free; hence, a unary constraint
 467 is a statement about a single tuple. Examples of unary constraints are restricted cases of
 468 conditional functional dependencies [7, 16]. An example of such a constraint can be “age
 469 smaller than 10 cannot co-occur with a salary greater than \$100k.”

470 **Parametric realization.** We consider a parametric realization model that is similar
 471 to the parametric intention model presented above, which is again a parametric Gibbs
 472 distribution. For each relation symbol $R \in \mathbf{S}$ and every pair $(t, t') \in \text{tuples}(R) \times \text{tuples}(R)$,
 473 the probability $\kappa[i, t](t')$, with $i \in \text{ids}(\mathcal{D})$, is expressed in the form of the Gibbs distribution
 474 $\kappa[i, t](t') = \frac{1}{Z_\kappa(t)} \exp\left(\sum_{g \in G} w_g g(t, t')\right)$, where G is the set of features, and each g is an
 475 arbitrary function (*feature*) over (t, t') , each weight w_g is a real number, and $Z_\kappa(t)$ is a
 476 normalization constant defined as $Z_\kappa(t) = \sum_{t' \in \text{tuples}(R)} \exp\left(\sum_{g \in G} w_g g(t, t')\right)$. Hence, we
 477 get a parametric model for probability distribution κ .

478 As an example, a feature function $g(t, t')$ may capture spelling errors: $g(t, t') = 1$ if $t'[\text{city}]$
 479 can be obtained by deleting one character from $t[\text{city}]$, and otherwise, $g(t, t') = -1$. The
 480 parameter vector Θ of the parametric realization model \mathcal{R}_Θ consists of the set of weights w_g
 481 for each feature $g \in G$.

482 **Assumptions.** We make two assumption here. First, we assume that the attributes of all
 483 relation symbols in \mathbf{S} take values over a finite and given set. This means that for each relation
 484 symbol $R \in \mathbf{S}$, $\text{tuples}(R)$ is also finite and given as input. Second, all features describing \mathcal{K}
 485 and κ can be computed efficiently, that is, in polynomial time in the size of the input.

486 **Obtaining examples for learning.** For supervised learning, we require a training collection
 487 $(I_j, J_j)_{j=1}^n$ and for unsupervised learning, a training collection $(J_j)_{j=1}^n$. We would like to
 488 make the case that, oftentimes, a single large example can be broken down into many small
 489 examples. Recall that in our setup (Gibbs parfactor/update PUDs), cross-tuple correlations
 490 can be introduced only by the integrity constraints in Φ . Consider, for instance, the case
 491 where all constraints in Φ are unary. Then, we get *tuple-independent* parfactor/update PUDs,
 492 and each tuple identifier i in can become an example database: $(I[i], J[i])$ in the supervised
 493 case, and $J[i]$ in the unsupervised case. For general constraints, cross-tuple correlations exist,
 494 and each example (I_j, J_j) and (J_j) can be obtained by taking correlated groups of tuples
 495 from J^* by considering different values for the attributes participating in each constraint
 496 $\varphi \in \Phi$. The number of tuples contained in each example depends on the constraints. This is
 497 a standard practice with parameterized probabilistic models as the ones we consider here [32].

6.2 Supervised Learning

We begin by considering supervised (\mathbf{S}, \mathbf{R}) -learning. All results presented in this section build upon standard tools from statistical learning. We are given a collection $(I_j, J_j)_{j=1}^n$ of intention-realization examples, and the goal is to find parameter values \mathbf{c} and \mathbf{d} that maximize the *likelihood* of pairs $(I_j, J_j)_{j=1}^n$, that is, $\prod_{j=1}^n \mathcal{R} \circ \mathcal{I}(I_j, J_j)$ for $\mathcal{I} = \mathcal{I}_{\Xi/\mathbf{c}}$ and $\mathcal{R} = \mathcal{R}_{\Theta/\mathbf{d}}$. To facilitate the analysis, we write this objective function as a sum over terms by considering the negative log-likelihood $l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (I_j, J_j)_{j=1}^n) = -\log \prod_{j=1}^n \mathcal{R} \circ \mathcal{I}(I_j, J_j) = -\sum_{j=1}^n \log \mathcal{R} \circ \mathcal{I}(I_j, J_j)$, and seek parameter values \mathbf{c} and \mathbf{d} that minimize it. For parfactor/update PUDs we have that $\mathcal{R} \circ \mathcal{I}(I, J) = \mathcal{D}(I) \times \prod_{i \in \text{ids}(I)} \kappa[i, I[i]](J[i])$ where $\mathcal{D}(I) = \mathcal{K}(I) \times 1/Z \times \exp(-\sum_{\varphi \in \Phi} w(\varphi) \times |V(\varphi, I)|)$. For the Gibbs parfactor/update PUD, \mathbf{c} corresponds to an assignment of parameters w_f describing \mathcal{K} and parameters $w_\varphi = w(\varphi)$ for constraints Φ . Similarly, \mathbf{d} corresponds to an assignment of parameters w_g describing κ . We use $Z(\mathbf{c})$ to denote the partition function Z under the parameters \mathbf{c} . We have:

$$l(\mathbf{c}, \mathbf{d}; (I_j, J_j)_{j=1}^n) = -\sum_{j=1}^n \log \left(\mathcal{K}(I_j; \mathbf{c}) \times \exp \left(-\sum_{\varphi \in \Phi} w_\varphi \times |V(\varphi, I_j)| \right) \right) + n \log Z(\mathbf{c}) - \sum_{j=1}^n \sum_{i \in \text{ids}(I_j)} \log(\kappa[i, I[i]; \mathbf{d}](J_j[i])) \quad (1)$$

where $\mathcal{K}(I_j; \mathbf{c})$ denotes that \mathcal{K} is parametrized by \mathbf{c} and $\kappa[i, I[i]; \mathbf{d}]$ denotes that κ is parametrized by \mathbf{d} .

Our goal becomes to minimize the expression in (1). It is well-known from the ML literature that there is no analytical solution to such minimization problems, and one needs to use iterative *gradient-based* methods [26]. We investigate whether gradient-based methods can indeed find a global minimum, and whether computing the gradient of this objective during each iteration is tractable. For the first question, the answer is positive.

► **Proposition 11.** $l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (I_j, J_j)_{j=1}^n)$ is a convex function of (Ξ, Θ) .

This proposition implies that the optimization objective for supervised learning has only global optima. Hence, it is guaranteed that any gradient-based optimization method will converge to a global optimum. Next, we study when the gradient of $l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (I_j, J_j)_{j=1}^n)$ with respect to \mathbf{c} and \mathbf{d} can be computed efficiently.

To compute the gradient of $l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (I_j, J_j)_{j=1}^n)$ one has to compute $\frac{\partial}{\partial c_i} \log Z(\mathbf{c})$. To compute this derivative a full inference step is required [26]. This is because computing this gradient is equivalent to computing the expected value for each feature (corresponding to each parameter c_i) according to the distribution defined by \mathbf{c} (Proposition 20.2 in [26]). However, marginal inference is often #P-hard [18]. In our setup, the constraints in Φ correspond to features, and it is not clear whether the gradient can be efficiently computable. In general, one can still estimate the aforementioned gradient by using approximate inference methods such as *Markov chain Monte Carlo* (MCMC) methods [9, 39] or *belief propagation* [41]. While effective in practice, these methods do not come with guarantees on the quality of the obtained solution. Next, we focus on an instance of PUD learning where exact inference is tractable (linear on tuples(R)), hence, we can compute the exact gradients of the aforementioned optimization objective efficiently.

Tuple independence. We focus on Gibbs parfactor/update PUD models where all constraints in Φ are unary. Here, $(I_j, J_j)_{j=1}^n$ corresponds to a collection of n examples, each of which having one tuple identifier. We use $I_j[0]$ and $J_j[0]$ to denote the tuples in the example

(I_j, J_j). In the appendix, we show that the negative log-likelihood $l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (I_j, J_j)_{j=1}^n)$ factorizes as $l(\mathbf{c}, \mathbf{d}; (I_j, J_j)_{j=1}^n) = \sum_{j=1}^n l'(\mathbf{c}, \mathbf{d}; I_j[0], J_j[0])$, where each $l'(\mathbf{c}, \mathbf{d}; I_j[0], J_j[0])$ is a convex function of Ξ and Θ (see Proposition 15 in the appendix). Moreover we show that the gradient of each $l'(\mathbf{c}, \mathbf{d}; I_j[0], J_j[0])$ can be evaluated in time linear to $\text{tuples}(R)$ where R is the relation corresponding to the tuple identifier associated with example (I_j, J_j) (see Proposition 16 in the appendix). Hence, we get the following.

► **Theorem 12.** *Given a training collection $(I_j, J_j)_{j=1}^n$ and a Gibbs parfactor/update PUD model with unary constraints, the exact gradient of $l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (I_j, J_j)_{j=1}^n)$ can be evaluated in $O(n \cdot \max_{R \in \mathcal{S}} |\text{tuples}(R)|)$ time.*

The above theorem implies that convex-optimization techniques such as *stochastic gradient descent* (SGD) [8] can be used to scale to large PUD learning instances (i.e., for large n).

A question that arises is about the number of examples (i.e., n) required to learn a PUD model. To answer this question, we study the convergence of *supervised (S, R)-learning* for Gibbs parfactor/update PUD models with unary constraints. For that, we view the collection $(I_j, J_j)_{j=1}^n$ as independent and identically distributed (i.i.d.) examples, drawn from a distribution that corresponds to a Gibbs parfactor/update PUD model with unary constraints and true parameters \mathbf{c}^* and \mathbf{d}^* . By the law of large numbers, the maximum likelihood estimates \mathbf{c} and \mathbf{d} are guaranteed to converge to \mathbf{c}^* and \mathbf{d}^* in probability. This means that for arbitrarily small $\epsilon > 0$ we have that $P(|\mathbf{c} - \mathbf{c}^*| > \epsilon) \rightarrow 0$ as $n \rightarrow \infty$. Same for \mathbf{d} . Moreover, we show that the MLE estimates \mathbf{c} and \mathbf{d} satisfy the property of *asymptotic normality* [28]. Intuitively, asymptotic normality states that the estimator not only converges to the unknown parameter, but it converges fast enough at a rate of $1/\sqrt{n}$. This implies that to achieve error ϵ for \mathbf{c} and \mathbf{d} one only needs $n = O(\epsilon^{-2})$ training examples.

► **Theorem 13.** *Consider a training collection $(I_j, J_j)_{j=1}^n$ drawn i.i.d. from a Gibbs parfactor/update PUD model with unary constraints and true parameters \mathbf{c}^* and \mathbf{d}^* . The estimates \mathbf{c} and \mathbf{d} obtained by solving MLE satisfy asymptotic normality, that is,*

$$\sqrt{n}(\mathbf{c} - \mathbf{c}^*) \rightarrow \mathcal{N}(0, \Sigma_{\mathbf{c}^*}^2) \text{ as } n \rightarrow \infty \text{ and } \sqrt{n}(\mathbf{d} - \mathbf{d}^*) \rightarrow \mathcal{N}(0, \Sigma_{\mathbf{d}^*}^2) \text{ as } n \rightarrow \infty$$

where $\Sigma_{\mathbf{c}^*}^2$ and $\Sigma_{\mathbf{d}^*}^2$ are the asymptotic variance of the estimates \mathbf{c} and \mathbf{d} .

Note that both the *multivariate Gaussian* distribution $\mathcal{N}(\mu, \Sigma^2)$ and the *asymptotic variance* are defined in classic statistics literature [28].

6.3 Unsupervised Learning

We now present preliminary results for unsupervised (S, R)-learning. We are given a training collection $(J_j)_{j=1}^n$ and seek to find \mathbf{c} and \mathbf{d} that minimize the negative log-likelihood $l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (J_j)_{j=1}^n) = -\sum_{j=1}^n \log \sum_I \mathcal{R} \circ \mathcal{I}(I, J_j)$. Again, there is no analytical solution for finding optimal \mathbf{c} and \mathbf{d} . Hence, one needs to use iterative gradient-based approaches, and again the questions of convexity and gradient computation arise. In the general case, this function is *not* necessarily convex. Hence, gradient-based methods are not guaranteed to converge to a global optimum. However, one can still solve the corresponding optimization problem using non-convex optimization methods [6]. Nonetheless, we show next that when realizers do not introduce *too much error*, we can establish guarantees.

Low-noise condition. Consider a Gibbs parfactor/update PUD model. We say that a PUD defined by $\mathcal{I}_{\Xi/\mathbf{c}}$ and $\mathcal{R}_{\Theta/\mathbf{d}}$ satisfies the *low-noise condition* with probability p if the realizer introduces an error with probability at most p . That is, for all intensions I and identifiers $i \in \text{ids}(I)$, it is the case that $\Pr(J[i] = I[i]|I) \geq 1 - p$. We have the following.

585 ► **Theorem 14.** *Consider a Gibbs parfactor/update PUD model where Ξ takes values from*
 586 *a compact convex set. Given a training collection $(J_j)_{j=1}^n$, there exists a fixed probability*
 587 *$p > 0$ such that, under the low-noise condition with probability p , the negative log-likelihood*
 588 *$l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (J_j)_{j=1}^n)$ is a convex function of Ξ .*

589 Hence, in certain cases, it is possible to find a global optimum of the overall negative
 590 log-likelihood. For that, the low-noise condition should hold with probability p that is also
 591 bounded, that is, it cannot be arbitrarily large. We can show that, if the low-noise condition
 592 holds with probability p , it is indeed bounded for Gibbs parfactor/update PUDs with unary
 593 constraints (see Proposition 19 in the appendix). We then find a global optimum as follows.
 594 We assume a simple parametric realization \mathcal{R}_Θ , that is, a model for which we can efficiently
 595 perform *grid search* over the space of parameter values \mathbf{d} . To find the global optimum for the
 596 negative log-likelihood, we solve a series of convex optimization problems over Ξ for different
 597 fixed $\Theta = \mathbf{d}$. For each of these problems, we are guaranteed to find a corresponding global
 598 optimum \mathbf{c} , and by performing a grid search we are guaranteed to find the overall global
 599 optimum \mathbf{c} and \mathbf{d} . This approach has been shown to converge for similar simple non-convex
 600 problems [34, 37].

601 Finally, similarly to supervised learning, the negative log-likelihood for fixed $\Theta = \mathbf{d}$
 602 decomposes into a sum of convex losses over $(J_j)_{i=1}^n$ where each example J_j contains a
 603 single tuple. We use $J_j[0]$ to denote that tuple. We have that $l(\Xi = \mathbf{c}, \Theta = \mathbf{d}, (J_j)_{j=1}^n) =$
 604 $\sum_{j=1}^n l'(\mathbf{c}, \mathbf{d}; J_j[0])$ where \mathbf{d} is fixed. We show that the gradient of each $l'(\mathbf{c}, \mathbf{d}; J[0])$ can be
 605 evaluated in polynomial time to tuples(R) where R is the relation corresponding to the tuple
 606 identifier associated with the example (J_j) .

607 It is left for future work to find sufficient conditions for p to be bounded for PUD models
 608 with more general constraints, as well as the complexity and convergence aspects.

609 7 Concluding Remarks

610 Taking inspiration from our experience with the HoloClean system [35], we introduced the
 611 concept of Probabilistic Unclean Databases (PUDs), a framework for unclean data that
 612 follows a noisy channel approach to model how errors are introduced in data. We defined
 613 three fundamental problems in the framework: cleaning, probabilistic query answering, and
 614 PUD learning (parameter estimation). We introduced PUD instantiations that generalize
 615 the deterministic concepts of subset repairs and update repairs, presented computational
 616 complexity results for probabilistic data cleaning, and discussed under which conditions we
 617 can learn a PUD from a single noisy database without the use of any labeled training data.

618 This paper opens up many research directions for future exploration. One is to investigate
 619 the complexity of cleaning in more general configurations than the ones covered here.
 620 Moreover, in cases where probabilistic cleaning is computationally hard, it is of natural
 621 interest to find *approximate* repairs that have a probability (provably) close to the maximum.
 622 Another direction is the complexity of probabilistic query answering and approximation
 623 thereof, starting with the most basic constraints (e.g., primary keys) and queries (e.g.,
 624 determine the marginal probability of a certain fact). Finally, an important direction is to
 625 devise learning algorithms for more general cases than those covered here. In particular, it
 626 is of high importance to understand when we can learn parameters without training data,
 627 based only on the given dirty database, under more generic noisy realization models than
 628 the ones discussed in this paper.

629 — References

- 630 1 S. Abiteboul, M. Arenas, P. Barceló, M. Bienvenu, D. Calvanese, C. David, R. Hull,
631 E. Hüllermeier, B. Kimelfeld, L. Libkin, W. Martens, T. Milo, F. Murlak, F. Neven, M. Or-
632 tiz, T. Schwentick, J. Stoyanovich, J. Su, D. Suciu, V. Vianu, and K. Yi. Research directions
633 for principles of data management (abridged). *SIGMOD Record*, 45(4):5–17, 2016.
- 634 2 F. N. Afrati and P. G. Kolaitis. Repair checking in inconsistent databases: algorithms and
635 complexity. In *ICDT*, pages 31–41. ACM, 2009.
- 636 3 P. Andritsos, A. Fuxman, and R. J. Miller. Clean answers over dirty databases: A proba-
637 bilistic approach. In *ICDE*, page 30. IEEE Computer Society, 2006.
- 638 4 M. Arenas, L. E. Bertossi, and J. Chomicki. Consistent query answers in inconsistent
639 databases. In *PODS*, pages 68–79. ACM, 1999.
- 640 5 G. H. Bakir, T. Hofmann, B. Schölkopf, A. J. Smola, B. Taskar, and S. V. N. Vishwanathan.
641 *Predicting Structured Data (Neural Information Processing)*. The MIT Press, 2007.
- 642 6 D. P. Bertsekas. *Nonlinear programming*. Athena scientific Belmont, 1999.
- 643 7 P. Bohannon, W. Fan, F. Geerts, X. Jia, and A. Kementsietsidis. Conditional functional
644 dependencies for data cleaning. In *ICDE*, pages 746–755. IEEE, 2007.
- 645 8 S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- 646 9 N. E. Breslow and D. G. Clayton. Approximate inference in generalized linear mixed models.
647 *Journal of the American Statistical Association*, 88(421):9–25, 1993.
- 648 10 M. Calautti, L. Libkin, and A. Pieris. An operational approach to consistent query answer-
649 ing. In *Proceedings of the 37th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles
650 of Database Systems, Houston, TX, USA, June 10-15, 2018*, pages 239–251. ACM, 2018.
- 651 11 J. Chomicki and J. Marcinkowski. Minimal-change integrity maintenance using tuple dele-
652 tions. *Information and Computation*, 197(1):90–121, 2005.
- 653 12 X. Chu, I. F. Ilyas, and P. Papotti. Holistic data cleaning: Putting violations into context.
654 In *ICDE*, pages 458–469, 2013.
- 655 13 N. N. Dalvi and D. Suciu. Efficient query evaluation on probabilistic databases. In *VLDB*,
656 pages 864–875. Morgan Kaufmann, 2004.
- 657 14 C. J. Date. Referential integrity. In *VLDB*, pages 2–12. VLDB Endowment, 1981.
- 658 15 R. Fagin, B. Kimelfeld, and P. G. Kolaitis. Dichotomies in the complexity of preferred
659 repairs. In *PODS*, pages 3–15, New York, NY, USA, 2015. ACM.
- 660 16 W. Fan, F. Geerts, X. Jia, and A. Kementsietsidis. Conditional functional dependencies
661 for capturing data inconsistencies. *ACM Trans. Database Syst.*, 33(2):6:1–6:48, June 2008.
- 662 17 T. Gaasterland, P. Godfrey, and J. Minker. An overview of cooperative answering. *J. Intell.*
663 *Inf. Syst.*, 1(2):123–157, 1992.
- 664 18 A. Globerson, T. Roughgarden, D. Sontag, and C. Yildirim. How hard is inference for
665 structured prediction? In *ICML*, pages 2181–2190. JMLR.org, 2015.
- 666 19 E. Gribkoff, G. V. den Broeck, and D. Suciu. The most probable database problem. In
667 *BUDA*, 2014.
- 668 20 I. F. Ilyas. Effective data cleaning with continuous evaluation. *IEEE Data Eng. Bull.*,
669 39:38–46, 2016.
- 670 21 A. K. Jha, V. Rastogi, and D. Suciu. Query evaluation with soft-key constraints. In *PODS*,
671 pages 119–128, 2008.
- 672 22 D. Jurafsky and J. H. Martin. *Speech and Language Processing (2Nd Edition)*. Prentice-
673 Hall, Inc., Upper Saddle River, NJ, USA, 2009.
- 674 23 S. Kolahi and L. V. S. Lakshmanan. On approximating optimum repairs for functional
675 dependency violations. In *ICDT*, volume 361, pages 53–62. ACM, 2009.
- 676 24 S. Kolahi and L. V. S. Lakshmanan. On approximating optimum repairs for functional
677 dependency violations. In *ICDT*, pages 53–62, 2009.

- 678 25 P. G. Kolaitis and E. Pema. A dichotomy in the complexity of consistent query answering
679 for queries with two atoms. *Inf. Process. Lett.*, 112(3):77–85, 2012.
- 680 26 D. Koller and N. Friedman. *Probabilistic Graphical Models: Principles and Techniques -*
681 *Adaptive Computation and Machine Learning*. The MIT Press, 2009.
- 682 27 P. Koutris and J. Wijsen. Consistent query answering for self-join-free conjunctive queries
683 under primary key constraints. *ACM Trans. Database Syst.*, 42(2):9:1–9:45, 2017.
- 684 28 S. Kullback. *Information theory and statistics*. Courier Corporation, 1997.
- 685 29 M. Lenzerini. Data integration: A theoretical perspective. In *PODS*, pages 233–246, New
686 York, NY, USA, 2002. ACM.
- 687 30 L. Libkin. Incomplete data: What went wrong, and how to fix it. In *PODS*, pages 1–13,
688 New York, NY, USA, 2014. ACM.
- 689 31 E. Livshits, B. Kimelfeld, and S. Roy. Computing optimal repairs for functional dependen-
690 cies. In *Proceedings of the 37th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles*
691 *of Database Systems, Houston, TX, USA, June 10-15, 2018*, pages 225–237. ACM, 2018.
- 692 32 B. London, B. Huang, B. Taskar, and L. Getoor. Collective stability in structured predic-
693 tion: Generalization from one example. In *Proceedings of the 30th International Conference*
694 *on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 828–
695 836, Atlanta, Georgia, USA, 17–19 Jun 2013.
- 696 33 A. Lopatenko and L. E. Bertossi. Complexity of consistent query answering in databases
697 under cardinality-based and incremental repair semantics. In *ICDT*, pages 179–193, 2007.
- 698 34 C. Ma, K. Wang, Y. Chi, and Y. Chen. Implicit regularization in nonconvex statistical
699 estimation: Gradient descent converges linearly for phase retrieval, matrix completion and
700 blind deconvolution. *arXiv preprint arXiv:1711.10467*, 2017.
- 701 35 T. Rekatsinas, X. Chu, I. F. Ilyas, and C. Ré. Holoclean: Holistic data repairs with
702 probabilistic inference. *PVLDB*, 10(11), 2017.
- 703 36 M. Richardson and P. Domingos. Markov logic networks. *Mach. Learn.*, 62(1-2):107–136,
704 Feb. 2006.
- 705 37 C. D. Sa, C. Ré, and K. Olukotun. Global convergence of stochastic gradient descent
706 for some non-convex matrix problems. In *ICML*, volume 37 of *JMLR Proceedings*, pages
707 2332–2341. JMLR.org, 2015.
- 708 38 P. Sen, A. Deshpande, and L. Getoor. Prdb: managing and exploiting rich correlations in
709 probabilistic databases. *VLDB J.*, 18(5):1065–1090, 2009.
- 710 39 S. Singh, M. Wick, and A. McCallum. Monte carlo mcmc: Efficient inference by ap-
711 proximate sampling. In *MNLP-CoNLL*, pages 1104–1113. Association for Computational
712 Linguistics, 2012.
- 713 40 D. Suciu, D. Olteanu, C. Ré, and C. Koch. *Probabilistic Databases*. Morgan & Claypool
714 Publishers, 1st edition, 2011.
- 715 41 M. J. Wainwright, T. S. Jaakkola, and A. S. Willsky. Tree-reweighted belief propagation
716 algorithms and approximate ML estimation via pseudo-moment matching. In *AISTATS*,
717 Jan. 2003.
- 718 42 M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational
719 inference. *Found. Trends Mach. Learn.*, 1(1-2):1–305, Jan. 2008.
- 720 43 J. Wang and N. Tang. Towards dependable data repairing with fixing rules. In *SIGMOD*,
721 pages 457–468. ACM, 2014.
- 722 44 M. Yakout, A. K. Elmagarmid, J. Neville, M. Ouzzani, and I. F. Ilyas. Guided data repair.
723 *PVLDB*, 4(5):279–289, 2011.

■ **Table 1** Main symbols

\mathbf{S}	A schema.
\mathcal{U}	A PUD $(\mathcal{I}, \mathcal{R}, J^*)$.
\mathcal{I}	An intention model (probabilistic database).
\mathcal{R}	A realization model, maps every I into a probabilistic database \mathcal{R}_I .
J^*	An observed unclean database.
$\mathcal{R} \circ \mathcal{I}$	Distribution over pairs (I, J) given by $\mathcal{R} \circ \mathcal{I}(I, J) \stackrel{\text{def}}{=} \mathcal{I}(I) \cdot \mathcal{R}_I(J)$.
\mathcal{U}^*	A probabilistic database given by $\mathcal{U}^*(I') = \Pr_{(I, J) \sim \mathcal{R} \circ \mathcal{I}}(I = I' \mid J = J^*)$.
(\mathcal{D}, τ, J^*)	A parfactor/subset PUD.
$(\mathcal{D}, \kappa, J^*)$	A parfactor/update PUD.
\mathcal{D}	A parfactor database (\mathcal{K}, Φ, w) with $w : \Phi \rightarrow (0, \infty)$.
\mathcal{K}	A generalized tuple-independent database (generalized TID).
Φ	A set of integrity constraints φ .
τ	Maps $i \in \text{ids}(R^{\mathcal{D}})$ to a discrete distribution $\tau[i]$ over $\text{tuples}(R) \cup \{\perp\}$.
κ	Maps $i \in \text{ids}(R^{\mathcal{D}})$ and $t \in \text{tuples}(R)$ to a discrete distribution $\kappa[i, t]$ over $\text{tuples}(R)$.

724 **A** Symbol Table

725 Table 1 lists the main symbols in the framework, along with their meaning.

726 **B** Proofs

727 In this section, we provide proofs that are missing from the body of the paper.

728 **B.1** Proof of Theorem 7

729 **THEOREM 7.** *Let (\mathcal{D}, τ, J^*) be a parfactor/subset PUD with $\mathcal{D} = (\mathcal{K}, \Phi, w)$. For $i \in \text{ids}(J^*)$,
 730 assume that $\mathcal{K}[i](\perp) > 0$ and $\tau[i](J^*[i]) > 0$, let $q(i) = \mathcal{K}[i](J^*[i]) / (\mathcal{K}[i](\perp) \cdot \tau[i](J^*[i]))$, and
 731 assume that $q(i) \geq 1$. There is a number M such that if $w(\varphi) > M$ for all $\varphi \in \Phi$ then the
 732 following are equivalent for all $I \subseteq J^*$:*

- 733 1. I is an MLI.
- 734 2. I is a minimum subset repair of J^* w.r.t. the weight function $w(i) = \log(q(i))$.

735 **Proof.** We analyze the probability of an MLI I . We have the following for all subsets I of
 736 J^* .

$$737 \quad \mathcal{R} \circ \mathcal{I}(I, J^*) = \mathcal{D}(I) \times \prod_{\substack{i \in \text{ids}(J^*) \setminus \\ \text{ids}(I)}} \tau[i](J^*[i]) \times \prod_{\substack{i \in \text{ids}(\mathcal{D}) \setminus \\ \text{ids}(J^*)}} \tau[i](\perp) \sim \mathcal{D}(I) \times \prod_{\substack{i \in \text{ids}(J^*) \setminus \\ \text{ids}(I)}} \tau[i](J^*[i]) \quad (2)$$

738 The proportionality is due to the fact that the third factor in the first multiplication is the
 739 same for all I . If w_φ is large enough for all $\varphi \in \Phi$, then we can assume that every I that
 740 satisfies Φ has a higher probability than every I that violates Φ . Observe that our integrity
 741 constraints and assumption that $\mathcal{K}[i](\perp) > 0$ and $\tau[i](J^*[i]) > 0$ implies that at least one
 742 consistent intention has a nonzero probability—the *empty* database. Hence, we can assume
 743 to begin with that the MLI I is selected from the subsets of J^* that satisfy Φ , and therefore

744 $\mathcal{D}(I) \sim \mathcal{K}(I)$. Moreover, we have the following.

$$745 \quad \mathcal{D}(I) \sim \mathcal{K}(I) = \prod_{i \in \text{ids}(I)} \mathcal{K}[i](I[i]) \times \prod_{\substack{i \in \text{ids}(\mathcal{K}) \setminus \\ \text{ids}(I)}} \mathcal{K}[i](\perp) \sim \prod_{i \in I} \mathcal{K}[i](I[i]) \times \prod_{\substack{i \in \text{ids}(J^*) \setminus \\ \text{ids}(I)}} \mathcal{K}[i](\perp) \quad (3)$$

746 Here, the proportionality is due to the fact that we divide the probability by the same factor
747 for all I . From Equations (2) and (3), we conclude the following.

$$748 \quad \mathcal{R} \circ \mathcal{I}(I, J^*) \sim \prod_{i \in I} \mathcal{K}[i](I[i]) \times \prod_{\substack{i \in \text{ids}(J^*) \setminus \\ \text{ids}(I)}} \mathcal{K}[i](\perp) \cdot \tau[i](J^*[i]) \\ 749 \quad = \prod_{i \in I} \frac{\mathcal{K}[i](I[i])}{\mathcal{K}[i](\perp) \cdot \tau[i](J^*[i])} \times \prod_{i \in \text{ids}(J^*)} \mathcal{K}[i](\perp) \cdot \tau[i](J^*[i]) \quad (4)$$

$$750 \quad \sim \prod_{i \in I} \frac{\mathcal{K}[i](I[i])}{\mathcal{K}[i](\perp) \cdot \tau[i](J^*[i])} = \prod_{i \in I} \frac{\mathcal{K}[i](J^*[i])}{\mathcal{K}[i](\perp) \cdot \tau[i](J^*[i])} = \prod_{i \in I} q(i) \quad (5)$$

752 Again, the last proportionality is due to the fact that the right factor in (4) is the same
753 for all I . The second equality in (5) is due to the fact that I is a subset of J^* , and so,
754 $I[i] = J^*[i]$ for all $i \in \text{ids}(I)$. Hence, I is a most likely repair if and only if I satisfies Φ and
755 maximizes (5), which is the same as maximizing the sum $\sum_{i \in \text{ids}(I)} \log q(i)$, as claimed. \blacktriangleleft

756 B.2 Proof of Theorem 8

757 **THEOREM 8.** *Let $\mathcal{U} = (\mathcal{D}, \tau, J^*)$ be an attribute-independent parfactor/update PUD with*
758 *$\mathcal{D} = (\mathcal{K}, \Phi, w)$ such that $\mathcal{U}^*(I) > 0$ for at least one consistent update I of J^* . There is a*
759 *number M such that if $w(\varphi) > M$ for all $\varphi \in \Phi$, then these are equivalent for all $I \subseteq J^*$:*

- 760 1. I is an MLI.
- 761 2. I is a minimum update repair w.r.t. the weight function

$$762 \quad w(i.A, a) = -\log(\mathcal{K}_A[i](a) \cdot \kappa_A[i, a](J^*[i].A)).$$

763 **Proof.** Since we consider only updates I of J^* , we can ignore all of the tuples $i \in \text{ids}(\mathcal{D}) \setminus J^*$.
764 Hence, we have the following.

$$765 \quad \mathcal{R} \circ \mathcal{I}(I, J^*) = \mathcal{D}(I) \times \prod_{i.A \in \text{cells}(J^*)} \kappa_A[i, I[i].A](J^*[i].A) \quad (6)$$

766 If w_φ is large enough for all $\varphi \in \Phi$, then we can assume that every I that satisfies Φ has
767 a higher probability than every I that violates Φ . Hence, due to our assumption that at
768 least one consistent update has a nonzero probability, we can assume to begin with that the
769 MLI I is selected from the subsets of J^* that satisfy Φ , and therefore $\mathcal{D}(I) \sim \mathcal{K}(I)$. Hence,
770 from (6) we conclude the following.

$$771 \quad \mathcal{R} \circ \mathcal{I}(I, J^*) \sim \prod_{i \in \text{ids}(I)} \mathcal{K}[i](I[i]) \times \prod_{i.A \in \text{cells}(J^*)} \kappa_A[i, I[i].A](J^*[i].A) \\ 772 \quad = \prod_{i.A \in \text{cells}(J^*)} \mathcal{K}_A[i](I[i].A) \cdot \kappa_A[i, I[i].A](J^*[i].A) \quad (7)$$

774 The last equation is due to the fact that I is an update of J^* and that \mathcal{U} is an attribute-
775 independent parfactor/update PUD. Hence, maximizing the probability of I amounts to
776 minimizing the weight function $w(i.A, a) = -\log(\mathcal{K}_A[i](a) \cdot \kappa_A[i, a](J^*[i].A))$. \blacktriangleleft

777 **B.3 Proof of Theorem 9**

778 **THEOREM 9.** *Let (\mathcal{D}, τ, J^*) be a parfactor/subset PUD with $\mathcal{D} = (\mathcal{K}, \Phi, w)$. If Φ consists of*
 779 *(at most) one key constraint per relation, and no relation of J^* has duplicate tuples, then an*
 780 *MLI can be computed in polynomial time.*

781 **Proof.** Our goal is to compute a subset I of J^* that maximizes $\mathcal{R} \circ \mathcal{I}(I, J^*)$. As argued in
 782 the proof of Theorem 7, we can ignore tuple identifiers not in $\text{ids}(J^*)$. By a *block* of J^* we
 783 refer to the (maximal) set of tuples of J^* that share a key value in a single relation. By
 784 our assumption, every block of more than a single tuple is a violation of the corresponding
 785 key constraint. With a straightforward spelling out of $\mathcal{R} \circ \mathcal{I}(I, J^*)$ we can observe that this
 786 probability factorizes across the blocks of J^* and their corresponding subsets in I , due to
 787 probabilistic independence between blocks. In particular, it suffices to find an MLI for each
 788 block separately, and take the union of the MLIs as our solution I . Therefore, we can assume
 789 that J^* has a single block; that is, all the tuples in J^* belong to the same relation, and all of
 790 them share the same key.

791 We denote our single key constraint by φ . Recall that $w(\varphi)$ is the weight of φ in \mathcal{D} ,
 792 and that $V(\varphi, I)$ the set of violations of φ in I . From our assumptions it follows that
 793 $|V(\varphi, I)| = |I| \cdot (|I| - 1)$. We conclude the following.

$$\begin{aligned}
 794 \quad \mathcal{R} \circ \mathcal{I}(I, J^*) &= \mathcal{D}(I) \times \prod_{\substack{i \in \text{ids}(J^*) \setminus \\ \text{ids}(I)}} \tau(J^*[i]) \\
 795 \quad &= \prod_{i \in \text{ids}(I)} \mathcal{K}(I[i]) \times \exp(-w(\varphi) \cdot |I| \cdot (|I| - 1)) \times \prod_{\substack{i \in \text{ids}(J^*) \setminus \\ \text{ids}(I)}} \tau(J^*[i]) \\
 796 \quad &= \prod_{i \in \text{ids}(I)} \frac{\mathcal{K}(I[i])}{\tau(J^*[i])} \times \exp(-w(\varphi) \cdot |I| \cdot (|I| - 1)) \times \prod_{i \in \text{ids}(J^*)} \tau(J^*[i]) \\
 797 \quad &\sim \prod_{i \in \text{ids}(I)} \frac{\mathcal{K}(I[i])}{\tau(J^*[i])} \times \exp(-w(\varphi) \cdot |I| \cdot (|I| - 1)) \\
 798 \quad &= \prod_{i \in \text{ids}(I)} \frac{\mathcal{K}(J^*[i])}{\tau(J^*[i])} \times \exp(-w(\varphi) \cdot |I| \cdot (|I| - 1)) \\
 799 \quad &
 \end{aligned}$$

800 The last equality is due to the fact that I is a subset of J^* , and therefore, $I[i] = J^*[i]$. Denote
 801 $q(i) = \mathcal{K}(J^*[i]) / \tau(J^*[i])$. In order to find a subset I of J^* that maximizes $\mathcal{R} \circ \mathcal{I}(I, J^*)$, we
 802 can consider every size $\ell = 0, \dots, |J^*|$ of I and find a subset I_ℓ with $|I_\ell| = \ell$ that maximizes
 803 the product $\prod_{i \in I_\ell} q(i)$. Then, we take the I_ℓ with the maximal $\prod_{i \in I_\ell} q(i) \exp(-w(\varphi)\ell(\ell - 1))$.
 804 In turn, to find I_ℓ it suffices to select ℓ tuple identifiers i with the maximal $q(i)$. ◀

805 **B.4 Proof of Theorem 10**

806 **THEOREM 10.** *Let J^* be a database, Φ a set of integrity constraints, Q a query, and t a*
 807 *tuple. The following are equivalent:*

- 808 1. t is a consistent answer over the cardinality repairs.
- 809 2. $\lim_{p \rightarrow 1} \lim_{w \rightarrow \infty} \Pr_{I \sim \mathcal{U}_{p,u}^*}(t \in Q(I)) = 1$.

810 **Proof.** Recall the definition of the unnormalized probability $\mathcal{R}\circ\mathcal{I}(I, J^*)$ of I .

$$811 \quad \mathcal{R}\circ\mathcal{I}(I, J^*) = \mathcal{D}(I) \times \prod_{\substack{i \in \text{ids}(J^*) \setminus \\ \text{ids}(I)}}} \tau[i](J^*[i]) \times \prod_{\substack{i \in \text{ids}(\mathcal{D}) \setminus \\ \text{ids}(J^*)}} \tau[i](\perp) \sim \mathcal{D}(I) \times \prod_{\substack{i \in \text{ids}(J^*) \setminus \\ \text{ids}(I)}}} \tau[i](J^*[i])$$

813 We make the following observations.

- 814 ■ For any constant p , as $u \rightarrow \infty$ the mass of \mathcal{D} concentrates on the databases consistent
- 815 subsets of J^* (i.e., the ones that satisfy Φ), or, in other words, as $u \rightarrow \infty$ the probability
- 816 that I satisfies Φ approaches 1.
- 817 ■ Since τ assigns the same probability to every tuple, we can write $\mathcal{R}\circ\mathcal{I}(I, J^*) \sim \mathcal{D}(I) \cdot$
- 818 $C^{|J^* \setminus I|}$ for some $C \in (0, 1)$. Moreover, as p approaches 1, the constant C approaches 0.
- 819 ■ From our definition of \mathcal{D} and \mathcal{K} it follows that every consistent subset has the same
- 820 probability in \mathcal{D} (i.e., $\mathcal{D}(I) = \mathcal{D}(I')$ whenever I and I' are consistent subsets of J^*).

821 We conclude that if I is cardinality repair and I' is any subset that is not a cardinality repair

822 (i.e., inconsistent or not maximal in cardinality), then the ratio between $\mathcal{R}\circ\mathcal{I}(I, J^*)$ and

823 $\mathcal{R}\circ\mathcal{I}(I', J^*)$, and so the ratio between the probabilities $\mathcal{U}^*(I)$ and $\mathcal{U}^*(I')$, approaches infinity

824 as p approaches 1.

825 We conclude that the total probability of the cardinality repairs approaches 1 as p

826 approaches 1, and all cardinality repairs have the same probability. In particular, if t is

827 a consistent answer, then its probability approaches 1. Conversely, if t is *not* a consistent

828 answer, then there is a constant portion of the probability that is missing for every p —this is

829 the probability of a cardinality repair I in which $t \notin Q(I)$. ◀

830 B.5 Proof of Proposition 11

831 PROPOSITION 11. $l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (I_j, J_j)_{j=1}^n)$ is a convex function of (Ξ, Θ) .

832 **Proof.** The negative log-likelihood is:

$$833 \quad l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (I_j, J_j)_{j=1}^n) = - \sum_{j=1}^n \log \left(\mathcal{K}(I_j; c) \times \exp \left(- \sum_{\varphi \in \Phi} w_\varphi \times |V(\varphi, I_j)| \right) \right)$$

$$834 \quad + n \log Z(\mathbf{c})$$

$$835 \quad - \sum_{j=1}^n \sum_{i \in \text{ids}(I_j)} \log(\kappa[i, I_j[i]; d](J_j[i]))$$

837 Additionally we have that:

$$838 \quad \mathcal{K}(I) = \prod_{i \in \text{ids}(I)} \mathcal{K}[i](I[i]) = \prod_{i \in \text{ids}(I)} \exp \left(\sum_{f \in F} w_f f(I[i]) \right)$$

839 and

$$840 \quad \kappa[i, I[i]](J[i]) = \frac{1}{Z_\kappa(I[i])} \exp \left(\sum_{g \in G} w_g g(I[i], J[i]) \right)$$

841 .

842 Replacing these to $l(\mathbf{c}, \mathbf{d}; (I_j, J_j)_{j=1}^n)$ we have:

$$\begin{aligned}
 843 \quad l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (I_j, J_j)_{j=1}^n) &= - \sum_{j=1}^n \log \left(\prod_{i \in \text{ids}(I_j)} \exp \left(\sum_{f \in F} w_f f(I_j[i]) \right) \right) \\
 844 &+ \sum_{j=1}^n \sum_{\varphi \in \Phi} w_\varphi |V(\varphi, I_j)| \\
 845 &+ n \log Z(\mathbf{c}) \\
 846 &- \sum_{j=1}^n \sum_{i \in \text{ids}(I_j)} \log \left(\frac{1}{Z_\kappa(I_j[i])} \exp \left(\sum_{g \in G} w_g g(I_j[i], J_j[i]) \right) \right) \\
 847 &
 \end{aligned}$$

848 or

$$\begin{aligned}
 849 \quad l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (I_j, J_j)_{j=1}^n) &= - \sum_{j=1}^n \sum_{i \in \text{ids} I_j} \left(\sum_{f \in F} w_f f(I_j[i]) \right) \\
 850 &+ \sum_{j=1}^n \sum_{\varphi \in \Phi} w_\varphi |V(\varphi, I_j)| \\
 851 &+ n \log Z(\mathbf{c}) \\
 852 &- \sum_{j=1}^n \sum_{i \in \text{ids}(I_j)} \left(\sum_{g \in G} w_g g(I_j[i], J_j[i]) - \log Z_\kappa(I_j[i]) \right) \\
 853 &
 \end{aligned}$$

854 or

$$\begin{aligned}
 855 \quad l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (I_j, J_j)_{j=1}^n) &= - \sum_{j=1}^n \sum_{i \in \text{ids} I_j} \sum_{f \in F} w_f f(I_j[i]) \\
 856 &+ \sum_{j=1}^n \sum_{\varphi \in \Phi} w_\varphi |V(\varphi, I_j)| \\
 857 &+ n \log \left(\sum_I \exp \left(\sum_{i \in \text{ids}(I)} \sum_{f \in F} w_f f(I[i]) \right) \right) \\
 858 &+ n \log \left(\sum_I \exp \left(- \sum_{\varphi \in \Phi} w_\varphi |V(\varphi, I)| \right) \right) \\
 859 &- \sum_{j=1}^n \sum_{i \in \text{ids}(I_j)} \sum_{g \in G} w_g g(I_j[i], J_j[i]) \\
 860 &+ \sum_{j=1}^n \sum_{i \in \text{ids}(I_j)} \log \left(\sum_{t' \in \text{tuples}(R^i)} \exp \left(\sum_{g \in G} w_g g(I_j[i], t') \right) \right) \\
 861 &
 \end{aligned}$$

862 where R^i denotes the relation symbol $R \in \mathbf{S}$ associated with tuple identifier i . Now, recall
 863 that vector Ξ contains all weights w_f and w_φ while vector Θ contains all weights w_g . We
 864 have that the negative log-likelihood is a convex function of Ξ as it corresponds to the sum of
 865 LogSumExp functions—those components corresponding to partition functions—with affine
 866 functions, which is well-known to be convex [8]. Similarly for Θ . ◀

867 **B.6 Proof of Theorem 12**

868 Before we prove Theorem 12, we prove two necessary propositions. Recall that we consider
 869 Gibbs parfactor/update PUD models with unary constraints. Here, $(I_j, J_j)_{j=1}^n$ is a collection
 870 of n i.i.d. examples each of which is associated with one tuple identifier. We use $I_j[0]$ and
 871 $J_j[0]$ to denote the tuples in example (I_j, J_j) .

872 First, we show that for Gibbs parfactor/update PUD models with unary constraints the
 873 negative log-likelihood $l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (I_j, J_j)_{j=1}^n)$ can be written as $l(\mathbf{c}, \mathbf{d}; (I_j, J_j)_{j=1}^n) =$
 874 $\sum_{j=1}^n l'(\mathbf{c}, \mathbf{d}; I_j[0], J_j[0])$, where each $l'(\mathbf{c}, \mathbf{d}; I_j[0], J_j[0])$ is a convex function of Ξ and Θ .
 875 Recall that We have:

876 **► Proposition 15.** *Given a collection $(I_j, J_j)_{j=1}^n$ of intention-realization examples and a*
 877 *Gibbs parfactor/update PUD model with unary constraints we have that $l(\mathbf{c}, \mathbf{d}; (I_j, J_j)_{j=1}^n) =$*
 878 *$\sum_{j=1}^n l'(\mathbf{c}, \mathbf{d}; I_j[0], J_j[0])$, where each $l'(\mathbf{c}, \mathbf{d}; I_j[0], J_j[0])$ is a convex function of Ξ and Θ .*

879 **Proof.** From the proof of Proposition 11 we have that negative log-likelihood is:

$$\begin{aligned}
 880 \quad l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (I_j, J_j)_{j=1}^n) &= - \sum_{j=1}^n \sum_{i \in \text{ids}(I_j)} \sum_{f \in F} w_f f(I_j[i]) \\
 &+ \sum_{j=1}^n \sum_{\varphi \in \Phi} w_\varphi |V(\varphi, I_j)| \\
 881 &+ n \log \left(\sum_I \exp \left(\sum_{i \in \text{ids}(I)} \sum_{f \in F} w_f f(I[i]) \right) \right) \\
 882 &+ n \log \left(\sum_I \exp \left(- \sum_{\varphi \in \Phi} w_\varphi |V(\varphi, I)| \right) \right) \\
 883 &- \sum_{j=1}^n \sum_{i \in \text{ids}(I_j)} \sum_{g \in G} w_g g(I_j[i], J_j[i]) \\
 884 &+ \sum_{j=1}^n \sum_{i \in \text{ids}(I_j)} \log \left(\sum_{t' \in \text{tuples}(R^i)} \exp \left(\sum_{g \in G} w_g g(I_j[i], t') \right) \right) \\
 885 & \\
 886 &
 \end{aligned}$$

887 where R^i denotes the relation symbol $R \in \mathbf{S}$ associated with tuple identifier i . Given that

each example (I_j, J_j) corresponds to one tuple identifier we have that:

$$\begin{aligned}
 l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (I_j, J_j)_{j=1}^n) &= - \sum_{j=1}^n \sum_{f \in F} w_f f(I_j[0]) \\
 &+ \sum_{j=1}^n \sum_{\varphi \in \Phi} w_\varphi |V(\varphi, I_j[0])| \\
 &+ n \log \left(\sum_I \exp \left(\sum_{i \in ids(I)} \sum_{f \in F} w_f f(I[i]) \right) \right) \\
 &+ n \log \left(\sum_I \exp \left(- \sum_{\varphi \in \Phi} w_\varphi |V(\varphi, I)| \right) \right) \\
 &- \sum_{j=1}^n \sum_{g \in G} w_g g(I_j[0], J_j[0]) \\
 &+ \sum_{j=1}^n \log \left(\sum_{t' \in tuples(R^i)} \exp \left(\sum_{g \in G} w_g g(I_j[0], t') \right) \right)
 \end{aligned}$$

The only components that do not immediately decompose over individual tuple identifiers are $\log \left(\sum_I \exp \left(\sum_{i \in ids(I)} \sum_{f \in F} w_f f(I[i]) \right) \right)$ and $\log \left(\sum_I \exp \left(- \sum_{\varphi \in \Phi} w_\varphi |V(\varphi, I)| \right) \right)$. We next show that for Gibbs parfactor/update PUD models with unary constraints both can be decomposed into a sum over terms that are defined over individual identifiers.

For the first term we have: By sum separation we have that:

$$\sum_I \exp \left(\sum_{i \in ids(I)} \sum_{f \in F} w_f f(I[i]) \right) = \prod_{i \in ids(\mathcal{D})} \sum_{t \in tuples(R^i)} \exp \sum_{f \in F} w_f f(t)$$

Hence

$$\begin{aligned}
 \log \left(\sum_I \exp \left(\sum_{i \in ids(I)} \sum_{f \in F} w_f f(I[i]) \right) \right) &= \log \left(\prod_{i \in ids(\mathcal{D})} \sum_{t \in tuples(R^i)} \exp \sum_{f \in F} w_f f(t) \right) \\
 &= \sum_{i \in ids(\mathcal{D})} \log \left(\sum_{t \in tuples(R^i)} \exp \sum_{f \in F} w_f f(t) \right)
 \end{aligned}$$

For the second term we have: Since we focus on unary constraints, each formula $\varphi \in \Phi$ has only one free variable. Therefore, we can let $grd(i, \varphi)$ denote that formula grounded with tuple identifier $i \in ids(\mathcal{D})$, and rewrite our expression as $\log Z(\mathbf{c}) = \log \left(\sum_I \exp \left(- \sum_{\varphi \in \Phi} \sum_{\substack{i \in ids(I) \\ I \models grd(i, \varphi)}} w_\varphi \right) \right)$. However, it holds that $I \models grd(i, \varphi)$ if and only if $I[i] \models grd(i, \varphi)$. From this we have that $\log \left(\sum_I \exp \left(- \sum_{\varphi \in \Phi} \sum_{\substack{i \in ids(I) \\ I[i] \models grd(i, \varphi)}} w_\varphi \right) \right)$.

Now, if we look at the argument of the above logarithm more closely, we notice that the body of the sum can be factored in terms of expressions that only depend on $I[i]$. It follows

914 by sum separation that:

$$915 \quad \sum_I \exp \left(- \sum_{\varphi \in \Phi} \sum_{\substack{i \in ids(I) \\ I[i] = grd(i, \varphi)}} w_\varphi \right) = \prod_{i \in ids(\mathcal{D})} \sum_{t \in tuples(R^i)} \exp \sum_{\varphi \in \Phi} \sum_{I[i] = grd(i, \varphi)} -w_\varphi$$

916 From this we have that:

$$917 \quad \log \left(\prod_{i \in ids(\mathcal{D})} \sum_{t \in tuples(R^i)} \exp \sum_{\varphi \in \Phi} \sum_{I[i] = grd(i, \varphi)} -w_\varphi \right)$$

$$918 \quad = \sum_{i \in ids(\mathcal{D})} \log \left(\sum_{t \in tuples(R^i)} \exp \sum_{\varphi \in \Phi} \sum_{I[i] = grd(i, \varphi)} -w_\varphi \right)$$

919

920 Based on our discussion in Section 6.1, we have that for Gibbs parfactor/update PUD
 921 models with unary constraints one can assume without loss of generality that the set of tuple
 922 identifiers in $ids(\mathcal{D})$ is exactly those present in the training examples $(I_j, J_j)_{j=1}^n$. Given all
 923 the above we can rewrite the negative log-likelihood as:

$$924 \quad l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (I_j, J_j)_{j=1}^n) = - \sum_{j=1}^n \sum_{f \in F} w_f f(I_j[0])$$

$$925 \quad + \sum_{j=1}^n \sum_{\varphi \in \Phi} w_\varphi |V(\varphi, I_j[0])|$$

$$926 \quad + \sum_{j=1}^n \log \left(\sum_{t \in tuples(R^i)} \exp \sum_{f \in F} w_f f(t) \right)$$

$$927 \quad + \sum_{j=1}^n \log \left(\sum_{t \in tuples(R^i)} \exp \sum_{\varphi \in \Phi} \sum_{I[i] = grd(i, \varphi)} -w_\varphi \right)$$

$$928 \quad - \sum_{j=1}^n \sum_{g \in G} w_g g(I_j[0], J_j[0])$$

$$929 \quad + \sum_{j=1}^n \log \left(\sum_{t' \in tuples(R^i)} \exp \left(\sum_{g \in G} w_g g(I_j[0], t') \right) \right)$$

930

931 Eventually we have that:

$$\begin{aligned}
932 \quad l'(\Xi = \mathbf{c}, \Theta = \mathbf{d}; I_j[0], J_j[0]) &= - \sum_{f \in F} w_f f(I_j[0]) + \sum_{\varphi \in \Phi} w_\varphi |V(\varphi, I_j[0])| \\
933 \quad &\quad - \sum_{g \in G} w_g g(I_j[0], J_j[0]) \\
934 \quad &\quad + \log \left(\sum_{t \in \text{tuples}(R^i)} \exp \sum_{f \in F} w_f f(t) \right) \\
935 \quad &\quad + \log \left(\sum_{t \in \text{tuples}(R)} \exp \sum_{\varphi \in \Phi} \sum_{I[i]=\text{grad}(i, \varphi)} -w_\varphi \right) \\
936 \quad &\quad + \log \left(\sum_{t' \in \text{tuples}(R)} \exp \left(\sum_{g \in G} w_g g(I_j[0], t') \right) \right) \\
937
\end{aligned}$$

938 where R denotes the relation symbol $R \in \mathbf{S}$ associated with the tuple identifier in $I_j[0]$.

939 Function $l'(\Xi = \mathbf{c}, \Theta = \mathbf{d}; I_j[0], J_j[0])$ is convex over Ξ and Θ as it corresponds to the sum
940 of LogSumExp functions with affine functions. ◀

941 Second we show that:

942 ▶ **Proposition 16.** *The gradient of $l'(\mathbf{c}, \mathbf{d}; I_j[0], J_j[0])$ with respect to Ξ and Θ can be*
943 *evaluated in time linear to $\text{tuples}(R)$ where R is the relation corresponding to the tuple*
944 *identifier associated with example (I_j, J_j) .*

945 **Proof.** To compute the gradient of $l'(\Xi = \mathbf{c}, \Theta = \mathbf{d}; I_j[0], J_j[0])$ we need to compute the
946 partial derivatives with respect to each parameter w_f , w_φ and w_g . It is trivial to see that
947 to compute these partial derivatives for the first three terms of $l'(\Xi = \mathbf{c}, \Theta = \mathbf{d}; I_j[0], J_j[0])$
948 one must evaluate each feature function f , g and $V(\varphi, \cdot)$. All of these functions are assumed
949 to be efficiently computable. It is also trivial to see that to compute the partial derivative of
950 each LogSumExp term of $l'(\Xi = \mathbf{c}, \Theta = \mathbf{d}; I_j[0], J_j[0])$ one needs to iterate over all tuples
951 $t \in \text{tuples}(R)$ as the expression inside the logarithm appear in the denominator of the each
952 partial derivative. Hence, the time required to compute the gradient for $l'(\Xi = \mathbf{c}, \Theta =$
953 $\mathbf{d}; I_j[0], J_j[0])$ is $O(\text{tuples}(R))$. ◀

954 **THEOREM 12.** *Given a training collection $(I_j, J_j)_{j=1}^n$ and a Gibbs parfactor/update PUD*
955 *model with unary constraints, the exact gradient of $l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (I_j, J_j)_{j=1}^n)$ can be*
956 *evaluated in $O(n \cdot \max_{R \in \mathbf{S}} |\text{tuples}(R)|)$ time.*

957 **Proof.** We have that $l(\mathbf{c}, \mathbf{d}; (I_j, J_j)_{j=1}^n) = \sum_{j=1}^n l'(\mathbf{c}, \mathbf{d}; I_j[0], J_j[0])$. To compute the overall
958 gradient of l we need to compute the gradient of each function l' . Hence, the overall
959 complexity is $O(n \cdot \max_{R \in \mathbf{S}} |\text{tuples}(R)|)$. ◀

960 B.7 Proof of Theorem 13

961 Before we present the proof for this theorem we discuss some notation that we use for
962 convenience and we also discuss the tools used to prove asymptotic normality.

963 First, we switch to matrix notation to denote the Gibbs parametric intention and
964 realization models \mathcal{I}_Ξ and \mathcal{R}_Θ . It is a simple exercise to show that that parametric models
965 introduced in Section 6.1 can be written as $\mathcal{I}_\Xi(I) = \frac{1}{Z_I} \exp(\Xi^T \cdot u_I(I))$ and $\mathcal{R}_\Theta(J|I) =$

966 $\frac{1}{Z_{\mathcal{R}}(I)} \exp(\Theta^T \cdot u_{\mathcal{R}}(I, J))$ where $u_{\mathcal{I}}(\cdot)$ is a vector function that corresponds to the features
 967 characterizing \mathcal{I}_{Ξ} and $u_{\mathcal{R}}(\cdot)$ is a vector function that corresponds to the features characterizing
 968 \mathcal{R}_{Θ} . These models correspond to the standard *exponential family* [42].

969 Vector functions $u_{\mathcal{I}}(\cdot)$ and $u_{\mathcal{R}}(\cdot)$ can also be represented as matrices $u_{\mathcal{I}} : \mathcal{R}^{|\Omega| \times d_{\mathcal{I}}}$ and
 970 $u_{\mathcal{R}} : \mathcal{R}^{|\Omega| \times d_{\mathcal{R}}}$ where Ω is the sample space of our PUD (e.g., if we had a single relation R
 971 that would be $\text{tuples}(R)$), $d_{\mathcal{I}}$ is the number of features describing \mathcal{I}_{Ξ} , and $d_{\mathcal{R}}$ is the number
 972 of features describing \mathcal{R}_{Θ} .

973 Second, we introduce the notion of *Fisher information* of the available training data [28].
 974 The Fisher information determines the amount of information that observed database
 975 instances carry about the unknown parameters Ξ and Θ . Intuitively, Fisher information can
 976 be interpreted as a measure of how quickly the distribution density will change when we
 977 slightly change a parameter in θ near the optimal θ^* .

978 Next, we define the Fisher information of a Gibbs parfactor/update PUD:

979 ► **Definition 17.** The *intention's Fisher information* of a Markov parametric PUD is:

$$980 \quad \mathcal{I}(\Xi) = \mathbf{Cov}_{I \sim \mathcal{I}_{\Xi}} [u_{\mathcal{I}}(I)].$$

981 Similarly, the *realizer's Fisher information* of a parametric PUD is:

$$982 \quad \mathcal{I}(\Theta) = \mathbf{E}_{I \sim \mathcal{I}_{\Xi}} [\mathbf{Cov}_{J \sim \mathcal{R}_{\Theta}(\cdot|I)} [u_{\mathcal{R}}(J|I)]] .$$

983 For general parameter learning, the Fisher information matrices can be *singular*, i.e., our
 984 observations carry no information about the parameters in some direction. Two conditions
 985 that lead to singular Fisher information matrices are: (1) the parameters of our PUD model
 986 are *redundant* (e.g., we can have the same formula listed twice with different weights) or (2)
 987 there is a parameter in our PUD model that has no effect on the distribution (e.g., we have
 988 a parameter associated with a formula that always evaluates to true). Notice that both cases
 989 described above correspond to misspecified parametric PUD models.

990 We now show that for any Gibbs parfactor/update PUD model the Fisher information
 991 matrices are positive definite, thus, not singular. Given that the Fisher information matrices
 992 are not singular, and thus invertible, we show the *asymptotic normality* of the MLE estimates
 993 Ξ and Θ for tuple independent MLD/update PUDs.

994 ► **Lemma 18.** For any PUD, if matrix $u_{\mathcal{I}}$ is always full-rank, and similarly for $u_{\mathcal{R}}$, and all
 995 parameters (Ξ, Θ) are finite, then

$$996 \quad \mathcal{I}(\Xi) \succ 0 \text{ and } \mathcal{I}(\Theta) \succ 0.$$

997 That is that $\mathcal{I}(\Xi)$ and $\mathcal{I}(\Theta)$ are positive definite and thus invertible.

998 **Proof.** For the intention model, the Fisher information is the covariance of $u_{\mathcal{I}}(I)$. The only
 999 way this matrix could be singular is if there is some unit vector ϕ such that

$$1000 \quad \mathbf{Var}_{I \sim \mathcal{I}_{\Xi}} [\phi^T u_{\mathcal{I}}(I)] = 0.$$

1001 This, in turn, will only happen if $\phi^T u_{\mathcal{I}}(I)$ is constant across all I on which \mathcal{I}_{Ξ} is supported.
 1002 Since Ξ is finite, \mathcal{I}_{Ξ} is supported everywhere on Ω . This means that if we define $\bar{u}_{\mathcal{I}}$ as

$$1003 \quad \bar{u}_{\mathcal{I}}(I) = u_{\mathcal{I}}(I) - \phi \phi^T u_{\mathcal{I}}(I),$$

1004 then $\bar{u}_{\mathcal{I}}(I)$ will also be a sufficient statistics function for the same PUD. But, $\bar{u}_{\mathcal{I}}(I)$ is rank
 1005 deficient, because $\phi^T \bar{u}_{\mathcal{I}} = 0$. But this cannot happen, since we supposed that any sufficient
 1006 statistics matrix would be full rank. Therefore, the Fisher information is positive definite,
 1007 which is what we wanted to prove. ◀

1008 Given that the Fisher information matrices are invertible we can now proceed to show
 1009 asymptotic normality. Before we proceed with our theorem, recall that since we consider
 1010 Gibbs parfactor/update PUDs with unary constraints each example (I_j, J_j) in the collection
 1011 of training examples $(I_j, J_j)_{j=1}^n$ corresponds to an independent single-tuple database example.
 1012 In the theorem below, we use quantities $\mathcal{I}_1(\Xi^*)$ and $\mathcal{I}_1(\Theta^*)$ to denote the Fisher information
 1013 of a Gibbs parfactor/update PUD model over a single-tuple database. We have for asymptotic
 1014 normality:

1015 **THEOREM 13.** *Consider a training collection $(I_j, J_j)_{j=1}^n$ drawn i.i.d. from a Gibbs*
 1016 *parfactor/update PUD model with unary constraints and true parameters \mathbf{c}^* and \mathbf{d}^* . The*
 1017 *estimates \mathbf{c} and \mathbf{d} obtained by solving MLE satisfy asymptotic normality, that is,*

$$1018 \quad \sqrt{n}(\mathbf{c} - \mathbf{c}^*) \rightarrow \mathcal{N}(0, \Sigma_{\mathbf{c}^*}^2) \text{ as } n \rightarrow \infty \text{ and } \sqrt{n}(\mathbf{d} - \mathbf{d}^*) \rightarrow \mathcal{N}(0, \Sigma_{\mathbf{d}^*}^2) \text{ as } n \rightarrow \infty$$

1019 where $\Sigma_{\mathbf{c}^*}^2$ and $\Sigma_{\mathbf{d}^*}^2$ are the asymptotic variance of the estimates \mathbf{c} and \mathbf{d} .

1020 **Proof.** We will prove this by the standard proof technique that is used to prove asymptotic
 1021 normality. The gradient of the negative log-likelihood of an exponential family model with
 1022 features u and parameters θ is:

$$1023 \quad \nabla f(\theta) = -\nabla_{\theta} \log \left(\frac{\exp(\theta^T u(x))}{\sum_{y \in \Omega} \exp(\theta^T u(y))} \right)$$

$$1024 \quad = -u(x) + \frac{\sum_{y \in \Omega} \exp(\theta^T u(y)) u(y)}{\sum_{y \in \Omega} \exp(\theta^T u(y))} = -u(x) + \mathbf{E}_{y \sim \pi_{\theta}} [u(y)]$$

1026 and the Hessian (the matrix that corresponds to the second-order partial derivatives) is

$$1027 \quad \nabla^2 f(\theta) = \frac{\sum_{y \in \Omega} \exp(\theta^T u(y)) u(y) u(y)^T}{\sum_{y \in \Omega} \exp(\theta^T u(y))} -$$

$$1028 \quad \frac{\sum_{y \in \Omega} \exp(\theta^T u(y)) u(y)}{\sum_{y \in \Omega} \exp(\theta^T u(y))} \cdot \frac{\sum_{y \in \Omega} \exp(\theta^T u(y)) u(y)^T}{\sum_{y \in \Omega} \exp(\theta^T u(y))}$$

$$1029 \quad = \mathbf{E}_{x \sim \pi_{\theta}} [u(x) u(x)^T] - \mathbf{E}_{x \sim \pi_{\theta}} [u(x)] \mathbf{E}_{x \sim \pi_{\theta}} [u(x)]^T = \mathbf{Cov}_{x \sim \pi_{\theta}} [u(x)].$$

1031 In the limit, we have that the gradient at the true parameter values is, for training examples
 1032 x_1, \dots, x_n ,

$$1033 \quad \nabla f(\theta) = \mathbf{E}_{y \sim \pi_{\theta}} [u(y)] - \frac{1}{n} \sum_{i=1}^n u(x_i).$$

1035 By a Taylor expansion, we expand the negative log-likelihood about the true parameters θ^*
 1036 and obtain that:

$$1037 \quad \nabla f(\theta) = H(\theta - \theta^*).$$

1038 From the above and Lemma 18, which states that the Fisher information is positive definite
 1039 and thus invertible, it follows that

$$1040 \quad \theta - \theta^* = \mathcal{I}(\theta)^{-1} \left(\mathbf{E}_{y \sim \pi_{\theta}} [u(y)] - \frac{1}{n} \sum_{i=1}^n u(x_i) \right)$$

$$1041 \quad = \left(\mathbf{Cov}_{x \sim \pi_{\theta}} [u(x)] \right)^{-1} \left(\mathbf{E}_{y \sim \pi_{\theta}} [u(y)] - \frac{1}{n} \sum_{i=1}^n u(x_i) \right)$$

1043 where $\mathcal{I}(\theta) = \mathbf{Cov}_{x \sim \pi_\theta} [u(x)]$ is the Fisher information of the model for x . This has expected
1044 value 0, and covariance

$$1045 \quad \mathbf{Cov}[\theta - \theta^*] = \frac{1}{n} \mathbf{Cov}_{x \sim \pi_\theta} [u(x)]^{-1}.$$

1046 If we plug in our Gibbs parfactor/update PUD model to the above we have for the intention
1047 model:

$$1048 \quad \mathbf{Cov}[\Xi - \Xi^*] = \frac{1}{n} \mathbf{Cov}_{I \sim \mathcal{I}_\Xi} [u_{\mathcal{I}}(I)]^{-1}$$

1049 Therefore:

$$1050 \quad \mathbf{Cov}[\Xi - \Xi^*] = \frac{1}{n} \mathcal{I}_1(\Xi^*)^{-1} \text{ as } n \rightarrow \infty.$$

1051 For a conditional exponential family distribution, we have:

$$1052 \quad \pi(X|I) = \frac{\exp(\theta^T u(X|I))}{\sum_Y \exp(\theta^T u(Y|I))}.$$

1054 The gradient of the negative log-likelihood is

$$1055 \quad \nabla_\theta - \log \pi(X|I) = -u(X|I) + \frac{\sum_Y \exp(\theta^T u(Y|I)) u(Y|I)}{\sum_Y \exp(\theta^T u(Y|I))}$$

$$1056 \quad = -u(X|I) + \mathbf{E}_{Y \sim \pi(\cdot|I)} [u(Y|I)].$$

1058 The Hessian is

$$1059 \quad \nabla_\theta^2 - \log \pi(X|I) = \frac{\sum_Y \exp(\theta^T u(Y|I)) u(Y|I) u(Y|I)^T}{\sum_Y \exp(\theta^T u(Y|I))}$$

$$1060 \quad - \left(\frac{\sum_Y \exp(\theta^T u(Y|I)) u(Y|I)}{\sum_Y \exp(\theta^T u(Y|I))} \right) \left(\frac{\sum_Y \exp(\theta^T u(Y|I)) u(Y|I)}{\sum_Y \exp(\theta^T u(Y|I))} \right)^T$$

$$1061 \quad = \mathbf{Cov}_{Y \sim \pi(\cdot|I)} [u(Y|I)].$$

1063 In expectation over I , this will be

$$1064 \quad \nabla^2 f(\theta) = \mathbf{E}_I [\mathbf{Cov}_{Y \sim \pi(\cdot|I)} [u(Y|I)]] .$$

1066 This leaves us with a typical gradient of using samples

$$1067 \quad \nabla f(\theta) = \frac{1}{n} \sum_{k=1}^n (\mathbf{E}_{Y \sim \pi(\cdot|I_k)} [u(Y|I_k)] - u(X_k|I_k)) .$$

1069 This will have expected value 0, and covariance

$$1070 \quad \mathbf{Cov}[\nabla f(\theta)] = \mathbf{Cov} \left[\frac{1}{n} \sum_{k=1}^n (\mathbf{E}_{Y \sim \pi(\cdot|I_k)} [u(Y|I_k)] - u(X_k|I_k)) \right]$$

$$1071 \quad = \frac{1}{n} \mathbf{Cov} [\mathbf{E}_{Y \sim \pi(\cdot|I_1)} [u(Y|I_1)] - u(X_1|I_1)]$$

$$1072 \quad = \frac{1}{n} \mathbf{E}_{I_1} [\mathbf{Cov}_{X_1} [\mathbf{E}_{Y \sim \pi(\cdot|I_1)} [u(Y|I_1)] - u(X_1|I_1)]]$$

$$1073 \quad = \frac{1}{n} \mathbf{E}_{I_1} [\mathbf{Cov}_{X_1} [u(X_1|I_1)]] .$$

1074

1075 So in this setting,

$$1076 \quad \mathbf{Cov}[\theta - \theta^*] = \frac{1}{n} \mathbf{E}_z \left[\mathbf{Cov}_{x \sim \pi_\theta} [u(x|z)]^{-1} \right].$$

1077 If we plug in our Gibbs parfactor/update PUD model to the above we have for the realizer
1078 model:

$$1079 \quad \mathbf{Cov}[\Theta - \Theta^*] = \frac{1}{n} \mathbf{E}_{I \sim \mathcal{I}_\Xi} \mathbf{Cov}_{J \sim \mathcal{R}_\Theta(\cdot|I)} [u_{\mathcal{R}}(J|I)]^{-1}$$

1080 Therefore:

$$1081 \quad \mathbf{Cov}[\Theta - \Theta^*] = \frac{1}{n} \mathcal{I}_1(\Theta^*)^{-1} \text{ as } n \rightarrow \infty.$$

1082 This concludes the proof. ◀

1083 B.8 Proof of Theorem 14

1084 For convenience we use the notation introduced at the beginning of Section B.7. We continue
1085 with our proof.

1086 **THEOREM 14.** *Consider a Gibbs parfactor/update PUD model where Ξ takes values from*
1087 *a compact convex set. Given a training collection $(J_j)_{j=1}^n$, there exists a fixed probability*
1088 *$p > 0$ such that, under the low-noise condition with probability p , the negative log-likelihood*
1089 *$l(\Xi = \mathbf{c}, \Theta = \mathbf{d}; (J_j)_{j=1}^n)$ is a convex function of Ξ .*

1090 **Proof.** Based on the discussion in Section 6.1, we have that $(J_j)_{j=1}^n = J^*$. Given this, the
1091 negative log-likelihood with respect to Ξ is:

$$1092 \quad l(\Xi) = -\log \left(\sum_{I \in \Omega} \mathcal{I}_\Xi(I) \cdot \mathcal{R}_\Theta(J^*|I) \right) = -\log \left(\sum_{I \in \Omega} \frac{\exp(\Xi^T u_{\mathcal{I}}(I))}{\sum_{J \in \Omega} \exp(\Theta^T u_{\mathcal{I}}(J))} \cdot \mathcal{R}_\Theta(J^*|I) \right).$$

1093 The gradient of this is

$$1094 \quad \nabla f(\Xi) = \frac{\sum_{J \in \Omega} \exp(\Xi^T u_{\mathcal{I}}(J)) \cdot u_{\mathcal{I}}(J)}{\sum_{J \in \Omega} \exp(\Xi^T u_{\mathcal{I}}(J))} - \frac{\sum_{I \in \Omega} \exp(\Xi^T u_{\mathcal{I}}(I)) \cdot \mathcal{R}_\Theta(J^*|I) u_{\mathcal{I}}(I)}{\sum_{I \in \Omega} \exp(\Xi^T u_{\mathcal{I}}(I)) \cdot \mathcal{R}_\Theta(J^*|I)}$$

1096 and the Hessian is

$$1097 \quad \nabla^2 f(\Xi) = \frac{\sum_{J \in \Omega} \exp(\Xi^T u_{\mathcal{I}}(J)) \cdot u_{\mathcal{I}}(J) \cdot u_{\mathcal{I}}(J)}{\sum_{J \in \Omega} \exp(\Xi^T u_{\mathcal{I}}(J))}$$

$$1098 \quad - \left(\frac{\sum_{J \in \Omega} \exp(\Xi^T u_{\mathcal{I}}(J)) \cdot u_{\mathcal{I}}(J)}{\sum_{J \in \Omega} \exp(\Xi^T u_{\mathcal{I}}(J))} \right) \left(\frac{\sum_{J \in \Omega} \exp(\Xi^T u_{\mathcal{I}}(J)) \cdot u_{\mathcal{I}}(J)}{\sum_{J \in \Omega} \exp(\Xi^T u_{\mathcal{I}}(J))} \right)^T$$

$$1099 \quad - \frac{\sum_{I \in \Omega} \exp(\Xi^T u_{\mathcal{I}}(I)) \cdot \mathcal{R}_\Theta(J^*|I) \cdot u_{\mathcal{I}}(I) \cdot u_{\mathcal{I}}(I)}{\sum_{I \in \Omega} \exp(\Xi^T u_{\mathcal{I}}(I)) \cdot \mathcal{R}_\Theta(J^*|I)}$$

$$1100 \quad + \left(\frac{\sum_{I \in \Omega} \exp(\Xi^T u_{\mathcal{I}}(I)) \cdot \mathcal{R}_\Theta(J^*|I) \cdot u_{\mathcal{I}}(I) \cdot u_{\mathcal{I}}(I)}{\sum_{I \in \Omega} \exp(\Xi^T u_{\mathcal{I}}(I)) \cdot \mathcal{R}_\Theta(J^*|I)} \right)$$

$$1101 \quad \left(\frac{\sum_{I \in \Omega} \exp(\Xi^T u_{\mathcal{I}}(I)) \cdot \mathcal{R}_\Theta(J^*|I) \cdot u_{\mathcal{I}}(I) \cdot u_{\mathcal{I}}(I)}{\sum_{I \in \Omega} \exp(\Xi^T u_{\mathcal{I}}(I)) \cdot \mathcal{R}_\Theta(J^*|I)} \right)^T$$

$$1102 \quad = \mathbf{Cov}[u_{\mathcal{I}}(I)] - \mathbf{Cov}[u_{\mathcal{I}}(I)|J^*].$$

1104 By the argument in Lemma 18, we know that $\mathbf{Cov}[u_{\mathcal{I}}(I)] \succ 0$ on all Ξ . It follows by
 1105 continuity that there exists a δ such that on Ξ , $\mathbf{Cov}[u_{\mathcal{I}}(I)] \succ \delta \mathbf{I}$ where \mathbf{I} is the identity
 1106 matrix. On the other hand, if we have a $\mathcal{R}_{\Theta}(J^*|I)$ such that for each tuple identifier $i \in I$
 1107 $I[i] = J^*[i]$ with probability at least $1 - p$, then

$$1108 \quad \mathbf{Cov}[u_{\mathcal{I}}(I)|J^*] \preceq p \max_{J \in \Omega} \|u_{\mathcal{I}}(J)\|^2.$$

1109 It follows that we can choose a p small enough that the Hessian is always positive definite on
 1110 Ξ , which means f is convex. This completes the proof. \blacktriangleleft

1111 Notice that for the Hessian to be positive definite it must be that $p \cdot \max_{J \in \Omega} \|u_{\mathcal{I}}(J)\|^2 \prec \delta I$
 1112 which in turn means that the value of p depends on the maximum value of features $u_{\mathcal{I}}$ when
 1113 computed over J . This means that probability p might not be bounded. In Section 6.3,
 1114 we discuss methods for solving the unsupervised version of PUD learning. These methods
 1115 require that p is bounded.

1116 Based on the analysis for the above theorem, we show that p is bounded for Gibbs
 1117 parfactor/update PUD models with unary constraints. We have the following proposition:

1118 **► Proposition 19.** *Given a Gibbs parfactor/update PUD model with unary constraints for*
 1119 *which the low noise condition holds, then we have for probability p that $p \cdot \|C\|^2 < \delta$ where*
 1120 $C = \max_{R \in \mathbf{S}, t \in \text{tuples}(R), f \in F, \phi \in \Phi} (f(t), V(\phi, t))$ *and δ is a constant with $\delta > 0$.*

1121 **Proof.** From the proof of Theorem 14 we have that when the noise condition holds it must
 1122 be that:

$$1123 \quad p \cdot \max_{J \in \Omega} \|u_{\mathcal{I}}(J)\|^2 \prec \delta \mathbf{I}$$

1124 It is easy to see that in the case of Gibbs parfactor/update PUD models with unary constraints
 1125 the maximum value for $\|u_{\mathcal{I}}(J)\|^2$ for any J scales independently of the number of tuple
 1126 identifiers in $\text{ids}(\mathcal{D})$ and depends only on the features of distribution \mathcal{K} for the intention
 1127 model \mathcal{I}_{Ξ} . \blacktriangleleft