

Learning Without Cleaning Over Content-Heterogeneous Databases

Jose Picado
Oregon State University
picadolj@oregonstate.edu

Arash Termehchy
Oregon State University
termehca@oregonstate.edu

ABSTRACT

Given a relational database and training examples for a target relation, relational learning algorithms learn a Datalog program that defines the target relation in terms of the existing relations in the database. The information in a domain is usually spread across several databases, which often represent the same entities under different names. Therefore, learning over multiple databases may result in inaccurate definitions. This type of heterogeneities is also common within a single database due to low data quality. Currently, users have to spend a great deal of time and effort to resolve these heterogeneities and create a unified and clean database instance to be used for learning. As it is *not* usually clear which values represent the same real-world entity or data item, e.g., whether *J. Smith* represents *John Smith* or *Jeremy Smith*, users may end up with numerous database instances after cleaning and resolving the data values. We propose HDLearn, a relational learning system that learns directly over heterogeneous databases without any pre-processing step. The user specifies the attributes across relations that contain values that may refer to the same real-world entity or data item through a set of declarative constraints called matching dependencies. HDLearn leverages these dependencies to learn accurate definitions over the heterogeneous data. Our empirical study on real-world datasets indicate that HDLearn learns accurate definitions over large databases efficiently.

PVLDB Reference Format:

Jose Picado and Arash Termehchy. Learning Over Heterogeneous Relational Databases. *PVLDB*, 12(xxx): xxxx-yyyy, 2019. DOI: <https://doi.org/10.14778/xxxxxxx.xxxxxxx>

1. INTRODUCTION

Users often use machine learning methods to discover interesting and novel relations over relational databases [12, 21, 41]. For instance, consider the IMDb database (*imdb.com*), which contains information about movies, for which schema fragments are shown in Table 1 (top). Given this database,

a user may want to find the definition for the new relation *highGrossing(title)*, which indicates that the movie with a given title is high grossing. Given a relational database and training examples for a new target relation, *relational machine learning* (relational learning) algorithms attempt to learn (approximate) relational definitions of the target relation in terms of existing relations in the database [12, 21]. Definitions are usually restricted to (non-recursive) Datalog programs for efficiency reasons. For instance, the user who wants to learn a definition for the new target relation *highGrossing* using the IMDb database may provide a set of high grossing movies as positive examples and a set of low grossing movies as negative ones to a relational learning algorithm. Given the IMDb database and these examples, the algorithm may learn the following definition:

$$\text{highGrossing}(x) \leftarrow \text{movies}(y, x, z), \text{mov2genres}(y, \text{'comedy'}), \\ \text{mov2releasedate}(y, \text{'May'}, u),$$

which indicates that high grossing movies are often released in May and their genre is *comedy*. Since relational learning algorithms leverage the structure of the database directly to learn new relations and their results are interpretable, they have been widely used in database management, analytics, and machine learning applications, such as query learning and information extraction [12, 21, 27, 2, 29, 41].

Databases, however, often contain heterogeneities in representing data values [14, 20, 8, 17], which may prevent the learning algorithms from finding an accurate definition. In particular, the information in a domain is usually spread across several databases. For example, IMDb does *not* contain the information about the budget or total grossing of movies. This information is available in another database called Box Office Mojo (BOM) (*boxofficemojo.com*), for which schema fragments are shown in Table 1 (bottom). Thus, to learn an accurate definition for the *highGrossing* relation, the user has to collect (more) data from the BOM database. For instance, the information about the budget of movies may be helpful in learning an effective definition for the *highGrossing* relation as high grossing movies may have high

This work is licensed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 12, No. xxx
ISSN 2150-8097.

DOI: <https://doi.org/10.14778/xxxxxxx.xxxxxxx>

IMDb	
movies(id,title,year)	mov2countries(id,name)
mov2genres(id,name)	mov2releasedate(id,month,year)
Box Office Mojo	
mov2totalGross(title,gross)	highBudgetMovies(title)

Table 1: Schema fragments for the IMDb and Box Office Mojo databases.

budgets. Thus, users may integrate these two databases under a unified schema over which to learn the target relation. However, the same entity and/or value may be represented in various forms in the integrated database. For instance, the *titles* of the same movie in IMDb and BOM have different formats., e.g., the title of the movie *Star Wars: Episode IV* is represented in IMDb as *Star Wars: Episode IV- 1977* and in BOM as *Star Wars - IV*. In this case, the learning algorithm *cannot* learn the aforementioned accurate definition for the *highGrossing* relation as it will *not* be able to connect the information about the same movie from IMDb and BOM movie relations. Such heterogeneities in representing names of entities and data values appear also in a single data source due to the inconsistencies in data entry and collection among others [8, 17].

Currently, users have to resolve the heterogeneities in representing data values manually and then learn over the clean database. However, it is difficult, time-consuming, and labor-intensive to resolve these heterogeneities over large databases [14, 20]. First, the user has to develop and/or train a matching function that distinguishes and unifies different values that refer to the same entity, apply that function on the database, and materialize the produced instance. Second, a unification may lead to new inconsistencies and opportunities to do more cleaning. For example, after unifying the titles of movies in a database, the user may notice that the names of directors of these movies are different in the database. To keep the database consistent, the user has to unify the names of the directors whose movies have been unified in the recently produced database. This process will be repeated for entities related to directors, e.g., production companies with which a director has worked, and other entities related to the movies until there are no more values to reconcile. Thus, it will take a long time to produce and materialize a clean database instance. Third, the process of unifying values may produce numerous clean database instances as one value may match multiple distinct values [8, 17, 9]. For example, title *Star Wars* may match both titles *Star Wars: Episode IV- 1977* and *Star Wars: Episode III- 2005*. Since we know that the *Star Wars: Episode IV- 1977* and *Star Wars: Episode III- 2005* refer to two different movies, the title *Star Wars* must be unified with only one of them. For each choice, one ends up with a distinct database instance. Since a large database has often many potentially matches, the number of clean database instances may be enormous. It is *not* clear which clean instance one should use to learn an accurate definition for a target relation.

The aforementioned tasks substantially increase the time needed to get insights. In fact, most data scientists spend about 80% of their time on data preparation tasks, such as dealing with data heterogeneities [28]. Since users *cannot* wait for a long time and/or allocate enormous resources to prepare large datasets, researchers have proposed on-demand approaches to data preparation [23, 39, 28, 24]. These methods may reduce the time and effort of data cleaning by preparing only a subset of the data that may be relevant to the task at hand. More importantly, the process of data analysis is inherently iterative [24]. Users may start by collecting the data items which they deem relevant to the target relation, clean and prepare them, and learn a definition over the data. If the learned definition is *not* sufficiently accurate, users may look for other data items and repeat this process. An on-demand approach enables users to evaluate

the relevance of a dataset to the target relation faster.

In this paper, we build upon the idea of on-demand data analysis and propose a novel learning method that directly learns over the original data without finding and materializing appropriate clean instance(s) of the data. Our approach advances the on-demand approaches as users do *not* have to perform any cleaning and preparation on the underlying data. Instead, users provide only a set of declarative (logical) constraints on the database schema, which determine values of what attributes can meaningfully match and be unified [17, 5, 9, 7, 38, 26, 3, 19, 25]. For instance, the values of the attribute *title* in relation *mov2totalGross* in BOM and the ones in the attribute *title* in relation *movies* in IMDb databases may match and be unified. On the other hand, it is *not* meaningful for the values of attributes *gross* in relation *mov2totalGross* and *year* in relation *movies* to be unified as they refer to different types of entities and values. Hence, our method substantially reduces the effort needed to produce clean databases instance(s) for data analysis. Our contributions are as follows:

- Since the representational conflicts in a heterogeneous database are not resolved, the properties of a desirable learned definition over a heterogeneous database are *not* clear. In particular, one heterogeneous database may have several clean instances. We introduce and formalize the problem of learning over a content-heterogeneous relational database (Section 3).
- We propose a novel relational learning algorithm called *HDLearn* that extends current learning methods over databases without heterogeneity to learn over a database with heterogeneous values (Section 4). It leverages the set of declarative constraints on attributes whose values can meaningfully match and unify during learning to learn an effective definition.
- Every learning algorithm chooses the final result based on its coverage of the training data, i.e., roughly speaking, the more (fewer) positive (negative) examples a definition covers the more accurate it is. It is challenging to perform these steps effectively over a database without cleaning it, as resolving certain entities, e.g., movies, may lead to unifying other entities of the same or other types, e.g., directors. Furthermore, one database may have numerous possible clean versions. We propose an efficient method to compute the coverage of a definition directly over the heterogeneous database without cleaning it (Section 4.2).
- We provide an implementation of HDLearn over a main-memory relational database management system. We use sampling techniques to scale learning over large databases while keeping it effective (Section 5).
- We perform an extensive empirical study over real-world datasets and show that HDLearn scales to large databases and learns efficiently over heterogeneous databases.

2. BACKGROUND

2.1 Basic Definitions

We fix two mutually exclusive sets of relation and attribute symbols. A database schema \mathcal{S} is a finite set of relation symbols R_i , $1 \leq i \leq n$. Each relation R_i is associated with a set of attribute symbols denoted as $R_i(A_1, \dots, A_m)$. We denote the domain of values for attribute A as $dom(A)$. Each database instance I of schema \mathcal{S} maps a finite set of

tuples to every relation R_i in \mathcal{S} . Each tuple t is a function that maps each attribute symbol in R_i to a value from its domain. We denote the value of the attribute X of tuple t in the database I by $t^I[X]$. An *atom* is a formula in the form of $R(u_1, \dots, u_n)$, where R is a relation symbol and u_1, \dots, u_n are *terms*. Each term is either a variable or a constant, i.e., value. A *ground atom* is an atom that only contains constants. A *literal* is an atom, or the negation of an atom. A *Horn clause* (clause for short) is a finite set of literals that contains exactly one positive literal. A *ground clause* is a clause that only contains ground atoms. Horn clauses are also called Datalog rules (without negation) or conjunctive queries. A *Horn definition* is a set of Horn clauses with the same positive literal, i.e., non-recursive Datalog program or union of conjunctive queries. A literal L_i in a clause $H \leftarrow L_1, \dots, L_n$ is *head-connected* if and only if at least one variable in L_i appears either in H or in a body literal L_j , where $1 \leq j < i$.

2.2 Relational Learning

In this section, we present the basic concepts of relational learning over databases without any heterogeneity and matching dependencies [12, 21]. Relational learning algorithms learn first-order logic definitions from an input relational database and training examples. Training examples E are usually tuples of a single *target relation*, and express positive (E^+) or negative (E^-) examples. The input relational database is also called *background knowledge*. The *hypothesis space* is the set of all possible first-order logic definitions that the algorithm can explore. It is usually restricted to Horn definitions to keep learning efficient. Each member of the hypothesis space is a *hypothesis*. Clause C covers an example e if $I \wedge C \models e$, where \models is the entailment operator, i.e., if I and C are true, then e is true. Definition H covers an example e if any of the clauses in H covers e . The goal of a learning algorithm is to find the definition, i.e., hypothesis, in the hypothesis space that covers all positive and the fewest negative examples as possible.

EXAMPLE 2.1. *The IMDb database contains the tuples movie(10, ‘Star Wars: Episode IV- 1977’, 1977), mov2genres(10, ‘comedy’), and mov2releasedate(10, ‘May’, 1977). Therefore, the definition that indicates that high grossing movies are often released in May and their genre is comedy shown in Section 1 covers the positive example highGrossing(‘Star Wars: Episode IV- 1977’).*

Most relational learning algorithms follow a covering approach, depicted in Algorithm 1 [34, 32, 30, 41, 33]. They construct one clause at a time using the *LearnClause* function. If the clause satisfies a criterion, e.g., covers at least certain fraction of the positive examples and does *not* cover certain fraction of negative ones, they add the clause to the learned definition and discards the positive examples covered by the clause. They stop when all positive examples are covered by the learned definition.

2.3 Matching & Cleaning

2.3.1 Matching Dependencies

Learning over databases with heterogeneity in representing values may deliver inaccurate answers as same entities and values may be represented under different names. Thus,

Algorithm 1: Covering approach algorithm.

Input : Database instance I , examples E
Output: Horn definition H

```

1  $H = \{\}$ 
2  $U = E^+$ 
3 while  $U$  is not empty do
4    $C = \text{LearnClause}(I, U, E^-)$ 
5   if  $C$  satisfies minimum criterion then
6      $H = H \cup C$ 
7      $U = U - \{e \in U \mid H \wedge I \models e\}$ 
8 return  $H$ 

```

one must resolve these representational differences to produce a high-quality database to learn an effective definition. If one wants to match and resolve values in a couple of attributes, one may use a supervised or unsupervised matching function to identify and unify their potential matches according to the domain of those attributes [14, 20]. For example, given that the domain of attributes is a set of strings, one may use string similarity functions, such as edit distance to find potential matches. Users develop and use matching and resolution rules and functions based on their domain knowledge. For example, consider relation *Employee(id, name, home-phone, address)*. According to her domain knowledge, the user knows that if phone numbers of two tuples are sufficiently similar, e.g., *001-333-1020* and *333-1020*, their addresses must be equal. Knowing this rule, the user may manipulate the database to ensure that all tuples with a similar phone number have equal addresses.

There may be a large number of matching rules over several sets of attributes in a large relational database [17, 5, 9, 7, 38, 26, 3, 19, 25]. Furthermore, these rules may interact with each other. For example, consider again the relation *Employee(id, name, home-phone, address)*. Assume that the user also knows that if two tuples have sufficiently similar addresses, e.g., *1 Main St., NY* and *1 Main Street, New York*, and names, they must have the same values for attribute *id*. Now, consider two tuples whose names and home phone numbers are sufficiently similar but their addresses are *not*. According to this rule, one *cannot* unify the ids of these tuples. But, after applying the rule mentioned in the preceding paragraph on the phone number and address, their addresses become sufficiently similar. Then, one can apply the second rule to unify the values of *id* in these tuples. The abundance of rules to match values and their interactions and dependencies call for a more formal approach to data cleaning where one can guarantee that all matching and resolution rules have been applied to the original data.

The database community has proposed declarative matching and resolution rules to express the domain knowledge about matching and resolution [17, 5, 9, 7, 38, 26, 3, 19, 25]. *Matching dependencies (MD)* is a popular type of such declarative rules, which provide a simple yet powerful method of expressing domain knowledge on matching values [17, 8, 6]. Let \mathcal{S} be the schema of the original database and R_1 and R_2 two relations in \mathcal{S} . Attributes A_1 and A_2 from relations R_1 and R_2 , respectively, are comparable if $\text{dom}(A_1) = \text{dom}(A_2)$. Given two pairs of pairwise comparable attributes A_1, A_2 and B_1, B_2 from relations R_1 and R_2 , respectively, an MD ϕ in \mathcal{S} is a sentence of the form $R_1[A_1] \approx R_2[A_2] \rightarrow R_1[B_1] \doteq R_2[B_2]$, where \approx is a similarity operator and

$R_1[B_1] \equiv R_2[B_2]$ indicates that the values of $R_1[B_1]$ and $R_2[B_2]$ refer to the same value, i.e., are interchangeable. Intuitively, the aforementioned MD says that given the values of $R_1[A_1]$ and $R_2[A_2]$ are sufficiently similar, the values of $R_1[B_1]$ and $R_2[B_2]$ are essentially different representations of the same value. For example, consider again the database that contains relations from *IMDb* and *BOM* whose schema fragments are shown in Table 1. According to our discussion in Section 1, one can define the following MD:

$$\begin{aligned} \text{movies}[title] &\approx \text{highBudgetMovies}[title] \rightarrow \\ &\text{movies}[title] \equiv \text{highBudgetMovies}[title]. \end{aligned}$$

The exact implementation of the similarity operator depends on the underlying domains of attributes. Our results are orthogonal to the implementation details of the similarity operator. The definition of MDs extend for the case where A_1, A_2 are list of attributes. For the sake of simplicity, we assume that MDs have only a single attribute in their left-hand sides. Our results extend to MDs with a list of attributes on their left-hand side. We also assume that all attributes share the same domain dom in the rest of this paper. Our results generalize to other cases.

2.3.2 Stable Instances

Given an input database with MDs on its relations, one must enforce the MDs to generate a high-quality database. Let t_1 and t_2 belong to R_1 and R_2 in database I of schema \mathcal{S} , respectively, such that $t_1^I[A_1] \approx t_2^I[A_2]$. To enforce the MD $\phi : R_1[A_1] \approx R_2[A_2] \rightarrow R_1[B_1] \equiv R_2[B_2]$ on I , one makes the values of $t_1^I[B_1]$ and $t_2^I[B_2]$ identical as they actually refer to the same value [17, 8]. For example, if attributes B_1 and B_2 contain titles of movies, one unifies both values *Star Wars - 1977* and *Star Wars - IV* to *Star Wars Episode IV - 1977* as it deems this value as the one to which $t_1^I[B_1]$ and $t_2^I[B_2]$ refer. The following definition formalizes the concept of applying an MD to the tuples t_1 and t_2 on I .

DEFINITION 2.2. *Database I' of schema \mathcal{S} is the immediate result of enforcing MD ϕ on t_1 and t_2 in I , denoted by $(I, I')_{[t_1, t_2]} \models \phi$ if*

1. $t_1^I[A_1] \approx t_2^I[A_2]$, but $t_1^I[B_1] \neq t_2^I[B_2]$;
2. $t_1^{I'}[B_1] = t_2^{I'}[B_2]$; and
3. I and I' agree on every other tuple and attribute value.

One may define a matching function over some domains to map the values that refer to the same value/ entity to the correct value/ entity in the stable instance. It may, however, be difficult to define such a function in many cases due to the lack of information. For example, let B_1 and B_2 in Definition 2.2 contain information about names of people and $t_1^I[B_1]$ and $t_2^I[B_2]$ have values *J. Bryan Smith* and *John B. Smith*, respectively, which according to an MD refer to the same actual name. It is *not* clear how to compute this name using the values of $t_1^I[B_1]$ and $t_2^I[B_2]$. Thus, we know that after enforcing ϕ , the values of $t_1^{I'}[B_1]$ and $t_2^{I'}[B_2]$ will be identical, but we do *not* know their exact values. Because we aim at developing learning algorithms that are efficient and effective over databases from various domains, we do *not* fix any matching method in this paper. We assume that matching every pair of values a and b in the database creates a fresh and new value denoted as $\langle a, b \rangle$.

Given the database I of schema \mathcal{S} with set of MDs Σ , I is *stable* iff $(I, I)_{[t_1, t_2]} \models \phi$ for all $\phi \in \Sigma$ and all tuples

$t_1, t_2 \in I$. In a stable database instance, all values that represent the same data item according to the database MDs are assigned equal values. Thus, it does *not* have the heterogeneities that may reduce the effectiveness of learning and analyzing databases. Given a database I with set of MDs Σ , one can produce a stable instance for I by starting from I and iteratively applying each MD in Σ according to Definition 2.2 to the output of the previous application finitely many times. In other words, let I, I_1, \dots, I_k denote the sequence of databases produced by applying MDs according to Definition 2.2 starting from I such that I_k is stable. We say that (I, I_k) satisfy Σ and denote it as $(I, I_k) \models \Sigma$. Each database may have (finitely) many stable instances depending on the order of applications of MDs [17, 8].

EXAMPLE 2.3. *Let $(10, \text{'Star Wars: Episode IV- 1977', 1977})$ and $(40, \text{'Star Wars: Episode III- 2005', 2005})$ be tuples in relation *movies* and ('Star Wars') be a tuple in relation *highBudgetMovies* whose schemas are shown in Table 1. Consider MD $\text{movies}[title] \approx \text{highBudgetMovies}[title] \rightarrow \text{movies}[title] \equiv \text{highBudgetMovies}[title]$. Let $\text{'Star Wars: Episode IV- 1977'} \approx \text{'Star Wars'}$ and $\text{'Star Wars: Episode III- 2005'} \approx \text{'Star Wars'}$ be true. Since the movies with titles $\text{'Star Wars: Episode IV- 1977'}$ and $\text{'Star Wars: Episode III- 2005'}$ are different movies with distinct titles, one can unify the title in the tuple ('Star Wars') in *highBudgetMovies* with only one of them in each stable instance. Each unification alternative leads to a distinct stable instance.*

The number of stable instances of a database generally grows exponentially with the number of matches and unification applications in the database. For example, in Example 2.3, if relation *highBudgetMovies* has a hundred titles that each matches and unifies with at least two titles in relation *movies*, the database will have around 2^{100} stable instances. Since there are often many matches and unification applications in a large database, one may have to generate numerous stable and clean databases. To produce only one stable instance, one may assume that all distinct unified values for a given value in the database represent the same value. For example, in Example 2.3, if one assumes that all movies with prefix *Black Sheep* represent the same movie, one can unify the titles of all movies to the same value. Nevertheless, this approach often violates the semantics of the underlying database and adds misleading information to it. For example, in Example 2.3, this approach leads to assuming that movies with titles $\text{'Star Wars: Episode IV- 1977'}$ and $\text{'Star Wars: Episode III- 2005'}$ have the same titles, which is obviously incorrect and misleading. This approach may also violate the integrity constraints in the database. For example, if matches and unifications are between key attributes, this approach will violate the key or functional dependency constraint by adding duplicate keys to a relation and makes the database inconsistent.

3. LEARNING A DEFINITION OVER HETEROGENEOUS DATA

3.1 Approaches to Learning Over Heterogeneous Data

Given the database I of schema \mathcal{S} with MDs Σ , and a set of training examples E , we wish to learn a Horn definition for a target relation T in terms of the relations in schema

S. Obviously, one may *not* learn an accurate definition by applying current learning algorithms over the set of input databases as it may consider different occurrences of the same entity/value to be distinct because of their different representations. Let $StableInstances(I, \Sigma)$ be the set of stable instances of I . One can learn definitions by generating all stable instances, i.e., $\mathbf{J} = StableInstances(I, \Sigma)$, learning a definition over each stable instance $J \in \mathbf{J}$ separately, and computing a union (disjunction) of all learned definitions. Since the discrepancies in value representations are resolved in stable instances, this approach may learn accurate definitions.

However, this method is *not* desirable and/or feasible for learning over large databases. As a large database may have numerous stable instances, it takes a great deal of time and storage to compute and materialize all of them. Moreover, we have to run the learning algorithm once for each stable instance, which may take an extremely long time. More importantly, as the learning has been done separately over each stable instance, it is *not* clear whether the final definition is sufficiently effective considering the information of all stable instances. For example, let database I have two stable instances I_1^s and I_2^s over which the aforementioned approach learns definitions H_1 and H_2 , respectively. H_1 and H_2 must cover a relatively small number of negative examples over I_1^s and I_2^s , respectively. However, H_1 and H_2 may cover a lot of negative examples over I_2^s and I_1^s , respectively. Thus, the disjunction of H_1 and H_2 will *not* be effective considering the information in I_1^s and I_2^s . Hence, it is *not* clear whether the disjunction of H_1 and H_2 is the definition that covers the most positive and the least negative examples over I_2^s and I_1^s . Finally, it is *not* clear how to encode and represent useably the final result as we may end up with numerous different definitions each of which is accurate over one stable instance. One may end up with hundreds or thousands of learned definitions. Thus, learning over each stable instance may take a very long time, *not* return an effective definition that considers information of all stable instances, and/or will be hard to use.

Another approach is to consider only the information shared among all stable instances for learning. The resulting definition will cover the most positive and the least negative examples considering the information common among all clean instances. This idea has been used in the context of query answering over inconsistent data [4, 8]. However, this approach may lead to ignoring many positive and negative examples as their connections to other relations in the database may *not* be present in all stable instances. For example, consider the tuples in relations *movies* and *highBudgetMovies* in Example 2.3. The training example (*Star Wars*) has different values in different stable instances of the database, therefore, it will be ignored. It will also be connected to two distinct movies with vastly different properties in each instance. Nevertheless, the training examples are very important in delivering an effective result and are usually costly to attain. In a sufficiently heterogeneous database, most or all positive and negative examples will *not* have the same information among all stable instances. The learning algorithm may simply learn an empty definition in these cases.

Thus, we hit a middle-ground. We follow the approach of learning directly over the original database. But, we also give the language of definitions and semantic of learning

enough flexibility to take advantage of as much (training) information as possible. Each definition will be a compact representation of a set of definitions, each of which is sufficiently accurate over some stable instances. If one increases the expressivity of the language, learning and checking coverage for each clause may become inefficient, e.g., disjunctive Datalog [15]. We ensure that the added capability to the language of definitions is minimal so learning remains efficient. In this section, we present our modifications to the hypothesis space and semantic of learning.

3.2 Representing Heterogeneity in Definitions

We represent the heterogeneity of the underlying data in the our language for definitions. Each new definition encapsulates the varieties of definitions learned over the stable instances of the underlying database. To achieve this, we increase the hypothesis space of relational learning algorithms over schema \mathcal{S} by introducing a fresh and unique (builtin) relation symbol per each MD in \mathcal{S} . These symbols do *not* belong to \mathcal{S} , and are called *matching relations*. More precisely, given schema \mathcal{S} and MD ϕ in \mathcal{S} , we add relation symbol m_ϕ with arity two to the set of relation symbols used by the Datalog definitions in the hypothesis space. A literal with a matching relation symbol is a *matching literal*. Consider again the database created by integrating IMDb and BOM datasets, whose schema fragments are in Table 1, with MD $\phi : movies[title] \approx highBudgetMovies[title] \rightarrow movies[title] \equiv highBudgetMovies[title]$. We may learn the following Datalog definition for the target relation *highGrossing*.

$$highGrossing(x) \leftarrow movies(y, t, z), mov2genres(y, 'comedy'), highBudgetMovies(x), m_\phi(x, t).$$

The matching literal $m_\phi(x, t)$ represents the relationship between titles in relations *highGrossing* and *movies*. According to the definition of MDs, for every matching literal m_ϕ , we have $m_\phi(x, y)$ iff $m_\phi(y, x)$.

We call a clause (definition) *stable* if it does *not* have any matching literal. Each clause (definition) with matching literals represents a set of stable clauses (definitions). We convert a clause with matching literals to a set of stable clauses by iteratively applying the MDs that correspond to each matching literal and eliminating the matching literal from the clause, similar to the process of applying MDs to a database instance described in Section 2.3. Given a clause C and an MD ϕ , to apply ϕ to C , we do the following. For each matching literal $m_\phi(x, y)$ in C , we create a new variable $\langle x, y \rangle$. We then replace every occurrence of x and y in C with $\langle x, y \rangle$. Finally, we remove the original matching literal $m_\phi(x, y)$ from C . The algorithm progressively applies MDs to the created clause until no matching literal is left in the clause. Similar to the one explained in Section 2.3, this algorithm is guaranteed to terminate. The resulting set is called the *stable clauses* of the input clause.

EXAMPLE 3.1. Consider the following clause over the movie database of IMDb and BOM.

$$highGrossing(x) \leftarrow movies(y, t, z), mov2genres(y, 'comedy'), highBudgetMovies(x), m_\phi(x, t).$$

The application of MD $\phi : highGrossing[title] \approx movies[title] \rightarrow highGrossing[title] \equiv movies[title]$, unifies variables t and x

to variable $\langle x, t \rangle$ and generates the following clause:

$highCrossing(\langle x, t \rangle) \leftarrow movies(y, \langle x, t \rangle, z),$
 $mov2genres(y, 'comedy'), highBudgetMovies(\langle x, t \rangle).$

After each application of an MD to a clause, the modified variables will *not* be similar to any term in the clause, which is equivalent to the same way we treat the result of application of MDs on databases in Section 2.3. This is because we do *not* have enough information to declare the variable/constant created after matching equal to any other variable/constant. Similar to the application of MDs to a database instance, the application of MDs to a clause may generate multiple stable clauses.

EXAMPLE 3.2. Consider a target relation $T(A)$, an input database with schema $\{R(B), S(C)\}$, and MDs $\phi_1 : T[A] \approx R[B] \rightarrow T[A] \equiv R[B]$ and $\phi_2 : T[A] \approx S[C] \rightarrow T[A] \equiv S[C]$. The definition $H : T(x) \leftarrow R(y), m_{\phi_1}(x, y), S(z), m_{\phi_2}(x, z)$ over this schema has two stable definitions: $H'_1 : T(\langle x, y \rangle) \leftarrow R(\langle x, y \rangle), S(z)$ and $H'_2 : T(\langle x, z \rangle) \leftarrow R(y), S(\langle x, z \rangle).$

As Example 3.2 illustrates, using matching literals enables us to provide a compact representation of multiple learned clauses/definitions where each may explain the patterns in the training data in some stable instances of the input database. Given an input definition H , the *stable definitions* of H are a set of definitions where each definition contains exactly one stable clause per each clause in H .

3.3 Coverage Over Heterogeneous Data

A learning algorithm evaluates the accuracy of a clause/definition according to the number of its covered positive and negative examples. For example, the *LearnClause* function in Algorithm 1 usually computes the difference of the number of positive and negative examples as the score and picks the clause(s) with the highest score. Each definition may have multiple stable definitions each of which may cover a different number of positive and negative examples on stable instances of the underlying database. Thus, it is *not* clear how to compute the score of a definition over a database.

One approach is to consider that a definition covers a positive example if at least one of its stable definitions covers it in some stable instance(s). Given all other conditions are the same, this approach may lead to learning a definition with numerous stable definitions where each may *not* cover sufficiently many or none positive examples. Thus, it is *not* clear whether each stable definition is accurate. A more restrictive approach is to define that a definition covers a positive example if all its stable definitions cover the example. This method will deliver a definition whose stable definitions have high positive coverage over stable instances. There are similar alternative for defining coverage of negative examples. One may consider that a definition covers a negative example if all its stable definitions cover it. Thus, if at least one stable definition does *not* cover the negative example, the definition will *not* cover it. This may lead to learning numerous stable definitions, which cover many negative examples. On the other hand, a restrictive approach may define a negative example covered by a definition if at least one of its stable definitions covers it. In this case, generally speaking, each learned stable definition will *not* cover too many negative examples.

We follow the more restrictive approach for definition of coverage of positive and negative examples.

DEFINITION 3.3. A definition H covers a positive example e with regard to database I iff every stable definition of H covers e in some stable instance of I .

EXAMPLE 3.4. As an example consider again the schema, MDs, and definition H in Examples 3.2 and the database of this schema with training example $T(a)$ and tuples $\{R(b), S(c)\}$. Assume that $a \approx b$ and $a \approx c$ are true. The database has two stable instances $I'_1 : \{T(\langle a, b \rangle), R(\langle a, b \rangle), S(c)\}$ and $I'_2 : \{T(\langle a, c \rangle), R(b), S(\langle a, c \rangle)\}$. Definition H covers the single training example in the original database according to Definition 3.3 as its stable definitions H'_1 and H'_2 cover the training example in stable instances I'_1 and I'_2 , respectively.

Definition 3.3 provides a more flexible semantic than considering only the common information between all stable instances as described in Section 3.1. The latter semantic considers that the definition H covers a positive example if it covers the example in all stable instances of a database. As explained in Section 3.1, this approach may lead to ignoring many if not all examples for learning.

DEFINITION 3.5. A definition H covers a negative example e with regard to database I if at least one of the stable definitions of H covers e in some stable instance of I .

To compute the coverage of a clause C for a set of positive examples, learning algorithms simply check the coverage of the clause for each example in the set and sum the covered examples. This approach can be safely applied for the case where C is stable. However, if C is *not* stable, this method deliver a clause that does *not* cover sufficiently many positive examples in any stable instance. For example, consider a database with one hundred stable instance and a set of 100 positive examples. Let each stable clause of a clause C cover only one example over each stable instance of the database. Obviously, C is *not* effective on any of these instances as its stable clauses cover an extremely small number of positive examples over each stable instance. This method, however, declares that C covers all the positive examples and is effective. Thus, a more reasonable approach is to ensure that each stable clause of C also covers sufficiently many positive examples over a stable instance. Given database I , stable clause C , and example e , let $\mathcal{I}_{C,e}$ be the set of stable instances of I on which C covers e .

DEFINITION 3.6. Given a database I , clause C , and set of positive examples E^+ , C covers E^+ in some stable instance if C covers each $e \in E^+$ and for every stable clause C_s of C and every pair of examples $e, g \in E^+$, we have $\mathcal{I}_{C_s,e} \cap \mathcal{I}_{C_s,g} \neq \emptyset$.

The definition naturally extends to (Horn) definitions and the set coverage for negative examples.

As explained in Section 3.1, our semantic of coverage based on Definitions 3.3 and 3.5 provides a middle-ground between the approaches of using only the information common between all stable instances and learning over each stable instance separately and unioning the results. It avoids ignoring too many examples and also ensures that the learned definition does *not* cover too many negatives example overall. It also guarantees that each stable definition of the learned definition covers sufficiently many positive examples and does *not* cover too many negative examples if the learning algorithm checks these restrictions on the learned definition.

4. HDLEARN: A LEARNING ALGORITHM FOR HETEROGENEOUS DATA

There are, generally speaking, two types of relational learning algorithms. In the top-down approach [34, 36, 41], the *LearnClause* function in Algorithm 1, searches the hypothesis space from general to specific, by using a refinement operator that is generally adding a new literal to the current clause. It starts with an empty clause, iteratively adds some literals to it, evaluates the resulting clause(s) over the training examples and database, and stops as soon as the clause satisfies certain criteria, e.g., covers the most positive and least negative examples. This method is challenging to use for learning over heterogeneous data, as it has to evaluate numerous clauses, i.e., queries, over the original (dirty) dataset. Since the dataset may have many stable solutions, query answering is shown to be generally difficult and inefficient over a relatively large and heterogeneous dataset [8].

The second approach for relational learning algorithms is bottom-up, in which the algorithm starts by exploring the patterns and clauses available in the data and then generalizes them to cover the training examples [30, 32]. More specifically, in these algorithms, the *LearnClause* function has two steps. It first builds the most specific clause in the hypothesis space that covers a given positive example, called a *bottom-clause*. Then, it generalizes the bottom-clause to cover as most positive and as fewest negative examples as possible. In this section, we propose a bottom-up algorithm called HDLearn for learning over heterogeneous data efficiently. To do so, HDLearn integrates the input MDs into the learning process. It also uses novel techniques to pick the most accurate clauses (and definitions) without computing the stable instances of the data. Because it is guided by the training data, HDLearn automatically determines whether and which databases one has to integrate to learn an accurate definition for the target concept. Next, we explain the two main steps of the algorithm in detail.

4.1 Bottom-clause Construction

A *bottom-clause* C_e associated with an example e is the most specific clause in the hypothesis space that covers e . Let I be the input database of schema \mathcal{S} , and Σ be the input matching dependencies. The bottom-clause construction algorithm consists of two phases. First, HDLearn finds all the information in I relevant to e . The information relevant to example e is the set of tuples $I_e \subseteq I$ that are connected to e . A tuple t is connected to e if we can reach t using a sequence of (similarity) search operations, starting from e . Given the information relevant to e , HDLearn creates the bottom-clause C_e .

EXAMPLE 4.1. *Given example $highGrossing(Superbad)$, the database in Table 2, and MD*

$$\begin{aligned} \phi : highGrossing[title] \approx movies[title] \rightarrow \\ highGrossing[title] = movies[title], \end{aligned}$$

*HDLearn finds the relevant tuples $movies(m1, Superbad(2007), 2007)$, $mov2genres(m1, comedy)$, $mov2countries(m1, c1)$, $englishMovies(m1)$, $mov2releasedate(m1, August, 2007)$, and $countries(c1, USA)$. As the movie title in the training example, e.g., *Superbad*, does not match with the movie title in the movies relation, e.g., *Superbad(2007)*, the tuple $movies(m1, Superbad(2007), 2007)$ is obtained through a similarity search according to MD ϕ . The other tuples are obtained through exact search.*

$movies(m1, Superbad(2007), 2007)$	$mov2genres(m1, comedy)$
$movies(m2, Zoolander(2001), 2001)$	$mov2genres(m2, comedy)$
$movies(m3, Orphanage(2007), 2007)$	$mov2genres(m3, drama)$
$mov2countries(m1, c1)$	$countries(c1, USA)$
$mov2countries(m2, c1)$	$countries(c2, Spain)$
$mov2countries(m3, c2)$	$englishMovies(m1)$
$mov2releasedate(m1, August, 2007)$	$englishMovies(m2)$
$mov2releasedate(m2, September, 2001)$	$spanishMovies(m3)$

Table 2: Example database.

To find the information relevant to e , HDLearn uses Algorithm 2. HDLearn maintains a set M that contains all seen constants. Let $e = T(a_1, \dots, a_n)$ be a training example. First, HDLearn adds a_1, \dots, a_n to M (line 3). These constants are values that appear in tuples in I . Then, HDLearn searches all tuples in I that contain at least one constant in M and adds them to I_e (lines 5-11). For exact search, HDLearn uses simple *select* operations (line 7). For similarity search, HDLearn uses MDs in Σ . If M contains constants in some relation R_i and given an MD $\phi \in \Sigma$, $\phi : R_i[A] \approx R_j[B] \rightarrow R_i[A] = R_j[B]$, HDLearn performs a similarity search over $R_j[B]$ to find relevant tuples in R_j , denoted by $\psi_{B \approx M}(R_j)$ (lines 8-9). For each new tuple in I_e , the algorithm extracts new constants and adds them to M (lines 10-11). The algorithm repeats this process for a fixed number of iterations d (line 4).

Algorithm 2: Bottom-clause construction.

Input : example e , # of iterations d , sample size s
Output: bottom-clause C_e

- 1 $I_e = \{\}$
- 2 $M = \{\}$ // M stores known constants
- 3 add constants in e to M
- 4 **for** $i = 1$ **to** d **do**
- 5 **foreach** *relation* $R \in I$ **do**
- 6 **foreach** *attribute* A **in** R **do**
- 7 $I_R = \sigma_{A \in M}(R)$
- 8 **if** \exists MD $\phi \in \Sigma$,
- 9 $\phi : S[B] \approx R[A] \rightarrow S[B] = R[A]$ **then**
- 10 $I_R = I_R \cup \psi_{A \approx M}(R)$
- 11 **foreach** *tuple* $t \in I_R$ **do**
- 12 add t to I_e and constants in t to M
- 13 $C_e =$ create clause from e and I_e
- 14 **return** C_e

To create the bottom-clause C_e from I_e , HDLearn first maps each constant in M to a new variable. It creates the head of the clause by creating a literal for e and replacing the constants in e with their assigned variables. Then, for each tuple $t \in I_e$, HDLearn creates a literal and adds it to the body of the clause, replacing each constant in t with its assigned variable. If there is an input MD $\phi : R_i[A] \approx R_j[B] \rightarrow R_i[A] = R_j[B]$, and there is a tuple $t' \in R_i$, then I_e may contain a tuple t obtained through similarity search. In this case, we add a *matching literal* $m_\phi(v_1, v_2)$, where v_1 and v_2 are the variables assigned to $t'[A]$ and $t[B]$, respectively. A matching literal indicates that the two values represented by the variables in the literal are interchangeable, according to the MD.

EXAMPLE 4.2. Given the relevant tuples found in Example 4.1, HDLearn creates the following bottom-clause:

```
highGrossing(x) ← movies(y, t, z), mφ(x, t),
mov2genres(y, 'comedy'), mov2countries(y, v),
countries(v, 'USA'), englishMovies(y),
mov2releasedate(y, 'August', u).
```

4.2 Generalization

After creating the bottom-clause C_e for example e , HDLearn generalizes C_e iteratively. To generalize C_e , HDLearn randomly picks a subset $E^{+s} \subseteq E^+$ of positive examples. For each example e' in E^{+s} , HDLearn generalizes C_e to generate a candidate clause C' , which is more general than C_e and covers e' . It does so by removing the *blocking literals*. Let $C_e = T \leftarrow L_1, \dots, L_n$ be the bottom-clause. L_i is a *blocking literal* iff i is the least value such that for all substitutions θ where $e' = T\theta$, the clause $(T \leftarrow L_1, \dots, L_i)\theta$ does not cover e' [32]. From the set of generalized clauses, HDLearn selects the highest scoring candidate clause. The score of a clause is the number of positive examples covered by the clause minus the number of negative examples covered by the clause. It then repeats this process with the selected clause until the clause cannot be improved.

EXAMPLE 4.3. Consider the bottom-clause C_e in Example 4.2 and positive example $e' = \text{highGrossing}(\text{Zoolander})$. To generalize C_e to cover e' , HDLearn drops the literal $\text{mov2releasedates}(y, \text{'August'}, u)$ because the movie *Zoolander* was not released in August.

A matching literal can also be a blocking literal.

EXAMPLE 4.4. Consider again the bottom-clause C_e in Example 4.2 and positive example $e'' = \text{highGrossing}(\text{Inception})$. The movie with title 'Inception' does not appear in the *movies* relation. Therefore, the first blocking literal in C_e is $m_\phi(x, t)$, which is in the join path between $\text{highGrossing}(x)$ and $\text{movies}(y, t, z)$.

HDLearn checks whether a candidate clause covers training examples in order to find blocking literals in a clause. It also computes the score of a clause by computing the number of training examples covered by the clause. Coverage tests dominate the time for learning [12]. One approach to perform a coverage test is to transform the clause into a SQL query and evaluate it over the input database to determine the training examples covered by the clause. However, since bottom-clauses over large databases normally have many literals, e.g., hundreds, the SQL query will involve long joins, making the evaluation extremely slow. Furthermore, it is challenging to evaluate clauses using this approach over heterogeneous data [8]. Moreover, it is *not* clear how to evaluate clauses that contain matching literals.

Therefore, HDLearn uses an approach called θ -subsumption. Assume that clause C does not contain matching literals, i.e., C is a stable clause. Clause C θ -subsumes clause G , denoted by $C \subseteq_\theta G$, iff there is some substitution θ such that $C\theta \subseteq G$ [12, 1]. $C\theta \subseteq G$ means that the result of applying substitution θ to clause C is a subset of clause G . The relation of θ -subsumption is both sound and complete for Horn definitions without functions, which are the focus of this paper, i.e., $C\theta \subseteq G$ iff $C \models G$ [1]. To evaluate whether a clause C covers an example e over database I , we first build a ground bottom-clause G_e for e in I , i.e., a bottom-clause

for e without replacing constants with variables. Then, we check whether $C \subseteq_\theta G_e$, which in turn indicates whether $C \wedge I \models e$. This approach is used by other learning algorithms [32, 33]. Ideally, G_e must have one literal per each tuple in the database that is connected to e through some joins, otherwise, the θ -subsumption test may declare that C does *not* cover e when C actually covers e . Ground bottom-clauses may be very large over large databases. The θ -subsumption is expensive over large bottom-clauses. We will discuss how to handle these cases in Section 5.

Clause C and ground bottom-clause G_e may *not* be stable. This means that there are possibly multiple stable clauses that can be derived from C and G_e . Following Definitions 3.3 and 3.5, we define θ -subsumption between potentially non-stable clauses. Let G_e^+ and G_e^- denote the ground bottom-clauses for the cases where e is a positive and negative example, respectively.

DEFINITION 4.5. Given clauses C , G_e^+ , and G_e^- over schema S , clause C θ -subsumes G_e^+ if each stable clause of C θ -subsumes at least one stable clause of G_e^+ . C θ -subsumes G_e^- if at least one stable clause of C θ -subsumes at least one stable clause of G_e^- .

According to Definition 4.5, given training example e , to check whether C θ -subsumes G_e^+ (G_e^-), one has to enumerate and check θ -subsumption of almost every pair of stable clauses of C and G_e^+ (G_e^-) in the worst case. Since both (ground) bottom-clauses normally contain many literals and θ -subsumption is NP-hard [1], this method is time-consuming. More importantly, because the learning algorithm performs numerous coverage tests, learning a definition may be extremely inefficient.

Actually, we can check whether C θ -subsumes G_e^+ (G_e^-) in a much more efficient way and without enumerating all stable clauses of C and G_e^+ (G_e^-). We consider matching literals in C and G_e^+ (G_e^-) that correspond to the same MD as the same relation symbol. For positive coverage, to find a substitution θ for C such that $C\theta \subseteq G_e^+$, we treat matching literals in C and G_e^+ as normal literals. If such substitution exists, then C θ -subsumes G_e^+ , hence C covers e .

EXAMPLE 4.6. Consider clause C

```
highGrossing(x) ← movies(y, t, z), mφ(x, t),
mov2genres(y, comedy).
```

Consider positive example $e' = \text{highGrossing}(\text{Zoolander})$ and ground bottom-clause G_e^+

```
highGrossing(Zoolander) ←
movies(m2, Zoolander(2001), 2001),
mφ(Zoolander, Zoolander(2001)),
mov2genres(m2, comedy), mov2countries(m2, c1),
countries(c1, USA), englishMovies(m2),
mov2releasedate(m2, September, 2001).
```

C θ -subsumes G_e^+ , using substitution $\theta = \{x = \text{Zoolander}, t = \text{Zoolander}(2001), y = m2, z = 2001\}$. Thus, C covers e' .

We follow a quite similar approach for the negative examples. Let G_e^- be the ground bottom-clause for the negative example e . We generate all stable clauses of C as described in Section 3 and check whether the generated stable clause θ -subsumes G_e^- exactly the same way as checking θ -subsumption for C and a ground bottom-clause for a positive example explained in the preceding paragraph. We

declare the C θ -subsumes G_e^- as soon as one stable clause of C θ -subsumes G_e^- . The following theorem shows that the aforementioned algorithms correctly detects θ -subsumption between clauses.

THEOREM 4.7. *Given a clause C and G_e^+ over schema \mathcal{S} , clause C θ -subsumes G_e^+ if and only if the aforementioned algorithm for checking the θ -subsumption of C and a positive ground bottom-clause finds a substitution θ such that $C\theta \subseteq G_e^+$. Similarly, given a clause C and G_e^- over \mathcal{S} , clause C θ -subsumes G_e^- if and only if the aforementioned algorithm for checking the θ -subsumption of C and a negative ground bottom-clause finds a substitution θ for one stable clause of C , C_s such that $C_s\theta \subseteq G_e^-$.*

PROOF. We prove the theorem for ground bottom clauses of positive examples. Assume that there is a substitution θ such that $C\theta \subseteq C'$. Thus, each matching literal of C , such as $m_\phi(x, y)$ is mapped and unified with a matching literal $m_\phi(x', y')$ in C' . Let C_ϕ and C'_ϕ be the results of applying $m_\phi(x, y)$ and $m_\phi(x', y')$ to C and C' , respectively. We show that the applications of mapped matching literals in these clauses preserve the substitution, i.e., $C_\phi\theta \subseteq C'_\phi$. Applying a matching literal does *not* eliminate any non-matching literal from a clause and only replaces some of their variables with new ones. Let C_ϕ and C'_ϕ be created by applying $m_\phi(x, y)$ and $m_\phi(x', y')$ to $R(\bar{x}), S(\bar{y})$ and $R(\bar{x}'), S(\bar{y}')$, respectively. The application of $m_\phi(x, y)$ replaces x and y in $R(\bar{x})$ and $S(\bar{y})$ with a new variable $\langle x, y \rangle$. Similarly, by applying $m_\phi(x', y')$, we replace x' and y' in $R(\bar{x}')$ and $S(\bar{y}')$ with $\langle x', y' \rangle$. The applications of $m_\phi(x, y)$ and $m_\phi(x', y')$ do *not* modify the literals that do *not* share any variable with them. Thus, we have $C_\phi\theta \subseteq C'_\phi$. Furthermore, as we assume that matching functions are *not* similarity preserving, i.e., each application of a matching function creates a fresh constant, subsequent applications of the matching literals other than $m_\phi(x, y)$ and $m_\phi(x', y')$ do *not* modify the variables introduced by applying these matching literals. Thus, the substitution θ is preserved over the application of each matching literal in C and its corresponding one in C' . After exhausting the application of all matching literals in C and their corresponding ones in C' , C' may still contain matching literals. However, applying these literals to C' retains the substitution between the stable clause of C and subsequent results of application of produced stable clauses of C' . The necessity of the theorem is proved by reversing the aforementioned steps. Now, we consider the case that C' is a negative ground bottom-clause. In this case, the algorithm checks the subsumption using every stable clause of the bottom-clause. Let C be one of these stable clauses. If C subsumes C' , according to the first part of this proof, C subsumes every stable clause of C' . The proof for the necessity of this part is also similar to the one of the part of the theorem on positive ground bottom-clauses.

□

Given bottom-clause C and a set of positive examples E^+ , according to Definition 3.6, it seems that one *cannot* use the simple method of computing the coverage of C for each example in E^+ and summing up the results to calculate the coverage of C over E^+ . However, the following theorem shows that if one uses the aforementioned subsumption method to check the coverage of C for each example, this simple method returns the desired coverage value for

C . Therefore, our algorithm is able to compute the coverage of a clause over a set of examples without any extra overhead. In the following theorem, we assume that no MD is defined over relation that contains the training examples, i.e., all MDs are defined on (other) relations in the background knowledge.

THEOREM 4.8. *Given a database I , set of positive examples E^+ , and clause C , let C θ -subsumes $G_{e_1}^+$ and $G_{e_2}^+$, which are ground bottom-clauses for $e_1, e_2 \in E^+$, respectively. For each stable clause C_s of C , we have $\mathcal{I}_{C_s, e_1} \cap \mathcal{I}_{C_s, e_2} \neq \emptyset$.*

PROOF. For brevity, we denote $G_{e_1}^+$ and $G_{e_2}^+$ as G_{e_1} and G_{e_2} , respectively. Our goal is to show that the stable instances in which C_s θ -subsumes some stable clauses of both G_{e_1} , $G_{e_1}^s$, and G_{e_2} , $G_{e_2}^s$, is not empty. This set becomes empty if and only if a tuple in the original database has different values in the stable instance(s) in which C θ -subsumes $G_{e_1}^s$ than all stable instance(s) in which C θ -subsumes $G_{e_2}^s$. If this holds, there will not be any stable instance in which C_s θ -subsumes both $G_{e_1}^s$ and $G_{e_2}^s$. Otherwise, one can construct at least one stable instance over which C_s θ -subsumes both $G_{e_1}^s$ and $G_{e_2}^s$.

First, assume that G_{e_1} and G_{e_2} do *not* share any normal, i.e., non-matching, literal. Since G_{e_1} and G_{e_2} do *not* have any literal in common, $G_{e_1}^s$ and $G_{e_2}^s$ are produced by applying matching literals on two completely different sets of literals. Let M_1 and M_2 be the sequence of MDs that correspond to the matching literals whose applications on G_{e_1} and G_{e_2} create the stable clauses $G_{e_1}^s$ of $G_{e_2}^s$, respectively. Since C_s θ -subsumes both $G_{e_1}^s$ of $G_{e_2}^s$, \mathcal{I}_{C_s, e_1} and \mathcal{I}_{C_s, e_2} contain stable instances produced by applying the MDs that correspond to M_1 and M_2 , respectively. There is *not* any tuple in the original database on which an MD from the MDs of M_1 and an MD from the set of MDs that correspond to M_2 are applied. Hence, the intersection of \mathcal{I}_{C_s, e_1} and \mathcal{I}_{C_s, e_2} is not empty and contain all the stable instances in which C_s θ -subsumes both $G_{e_1}^s$ and $G_{e_2}^s$.

Now, assume that G_{e_1} and G_{e_2} share a literal $R(\bar{a})$. \mathcal{I}_{C_s, e_1} and \mathcal{I}_{C_s, e_2} are disjoint if and only if 1) G_{e_1} and G_{e_2} have matching literals $m_1(b, c)$ and $m_2(b, d)$, $c \neq d$, respectively, that are used to create both $G_{e_1}^s$ and $G_{e_2}^s$ and 2) $m_1(b, c)$ is applied to literals $(R(\bar{a}), S(\bar{f}))$ in G_{e_1} and $m_2(b, d)$ is applied to literals $(R(\bar{d}), U(\bar{g}))$ in G_{e_2} . This way the tuple in the original database that corresponds to the literal $R(\bar{a})$ will have different constants in the stable instance(s) in which C θ -subsumes $G_{e_1}^s$ than the ones in which C θ -subsumes $G_{e_2}^s$.

For the sake of simplicity, let $R(\bar{a})$ be the only literal shared between the bodies of G_{e_1} and G_{e_2} . Our proof extends to other cases. According to the algorithm of creating ground bottom-clauses, all literals in G_{e_1} that have a constant in common with $R(\bar{a})$, i.e., their corresponding tuples in the database are connected to the tuple of $R(\bar{a})$ in the data through some join paths, will also appear in G_{e_2} . The same claim is true about the literals constructed and added to the ground bottom-clauses based on the tuples that are connected to $R(\bar{a})$ via the similarity predicate of some MD. Let us call this set of literals and their corresponding matching literals $L_{R(\bar{a})}$. According to the algorithm of ground bottom-clause construction, we may extend $L_{R(\bar{a})}$ to another set L , which contains all and only the literals that are connected to at least one literal in L via a join path, i.e., common constant, or a matching literal. We also add the matching literals that connect the (normal) literals of L to

L . The literals in L appear in the bodies of both G_{e_1} and G_{e_2} .

Let θ_{e_1} and θ_{e_2} be the substitution mappings from constants in the literals of G_{e_1} and G_{e_2} , respectively, to the variables/constants in the literals of C . Members of θ_{e_1} are pairs of (a, x) where a is a constant in G_{e_1} and x is a variable or constant in C . Let $\theta_{e_1}^L \subseteq \theta_{e_1}$ and $\theta_{e_2}^L \subseteq \theta_{e_2}$ be the substitution mappings from G_{e_1} and G_{e_2} to C , respectively. Let's replace $\theta_{e_2}^L$ with $\theta_{e_1}^L$ in θ_{e_2} and get a substitution θ_{e_2}' . Since the literals in L do not have any constant or matching literal in common with the rest of literals in G_{e_2} , C still θ -subsumes G_{e_2} using θ_{e_2}' . θ_{e_1} and θ_{e_2}' unify every (matching) literal in L with the same (matching) literal in C .

Matching literals that are not in L do not share any constant with the literals in L . Hence, applying them do not change the constants of any literal in L . Thus, without loss of generality we assume that G_{e_1} and G_{e_2} contain only literals that appear in L . Let M_C be the sequence of matching literals in C whose application creates C_s . Let M_L be the sequence of matching literals in L that are unified with M_C using θ_{e_1} and θ_{e_2}' . One can generate $G_{e_1}^s$ and $G_{e_2}^s$ by first applying M_L and then the rest of matching literals in L in an arbitrary order. C_s θ -subsumes both $G_{e_1}^s$ and $G_{e_2}^s$. Since $G_{e_1}^s$ and $G_{e_2}^s$ are created by applying the same sequence of matching literals on the same set of literals, $R(\bar{a})$ has the same value in both stable clauses. \square

Theorem 4.8 extends for negative examples.

After learning, one may use these coverage testing methods to use the learned definition and apply the learned definition to do prediction over the original database.

4.3 Commutativity Property

As discussed in Section 3, our proposed semantic for coverage over heterogeneous data as described in Definitions 3.3 and 3.5 is more restrictive than the method of learning definitions over every stable instance and then providing a union of all definitions as the final result. Since the method of learning over every stable instance separately, contains all possible definitions learned over the data, it is important to understand whether our algorithm works with the same set of (stable) definitions as this method to get an understanding of the coverage of the definitions produced by our algorithm in terms of the learned definitions over the individual stable instances. We show that our algorithm starts with and generalizes the same set of stable definitions used to learn over each stable instance. In other words, assume that we apply the bottom-clause construction part of our learning algorithm, as described in 4.1 directly to each stable instance of database I and output a set of stable bottom-clauses B . As we explained in the beginning of this section, the algorithm of HDLearn over non-heterogeneous, i.e., stable, databases is essentially existing state-of-the-art relational learning algorithms [32, 33]. Now, assume that we apply the bottom-clause construction part of our proposed algorithm directly over the original heterogeneous database and get a (unstable) definition H . We show that B is exactly the set of stable instances of H . Let $StableInstances(I, \Sigma)$ denote the set of stable instances of I with MDs Σ . Also, assume that $StableInstances(H)$ returns the set of all stable definitions of H . Let $BottomClause(I, \Sigma, e)$ denote the bottom-clause generated by applying the bottom-clause construction algorithm using example e over (heterogeneous) database I with

the set of MDs Σ . Also, let $BottomClause_s(StableInstances(I, \Sigma), e)$ be the set of stable clauses generated by applying the bottom-clause construction algorithm to every stable instance of I using example e .

THEOREM 4.9. *Given database I with MDs Σ and set of positive examples E^+ , for every positive example $e \in E^+$*

$$BottomClause_s(StableInstances(I, \Sigma), e) = StableDefinitions(BottomClause(I, \Sigma, e))$$

PROOF. Without loss of generality, assume that all learned definitions contain a single clause. Let $StableInstances(I, \Sigma) = \mathbf{J} = \{J_1, \dots, J_n\}$. We show that $BottomClause(I, e, \Sigma) = C$ is a compact representation of $BottomClause_s(\mathbf{J}, e) = \{C_1, \dots, C_n\}$.

Let $StableDefinitions(C) = \{C'_1, \dots, C'_m\}$. We remove the literals that are not head-connected in each clause in $\{C'_1, \dots, C'_m\}$. Let $\{J^{C'_1}, \dots, J^{C'_m}\}$ be the canonical database instances of $\{C'_1, \dots, C'_m\}$ [1]. This set can also be generated by running $StableInstances(I^C, \Sigma)$, where I^C is the canonical database instance of C .

Let $\{J^{C_1}, \dots, J^{C_n}\}$ be the canonical database instances of $\{C_1, \dots, C_n\}$. By definition, I^C contains all tuples that are related to e , either by exact or similarity matching (according to MDs in Σ). Because $StableInstances(I^C, \Sigma) = \{J^{C'_1}, \dots, J^{C'_m}\}$, all tuples that may appear in an instance in $\{J^{C_1}, \dots, J^{C_n}\}$ must also appear in an instance in $\{J^{C'_1}, \dots, J^{C'_m}\}$.

A tuple t may appear in an instance in $J^{C'_j} \in \{J^{C'_1}, \dots, J^{C'_m}\}$, but not appear in the corresponding instance $J^{C_i} \in \{J^{C_1}, \dots, J^{C_n}\}$. This is because t became disconnected of training example e when generating the stable instance J_i , which is a superset of J^{C_i} . Then, when building bottom-clause C_i from J_i , a literal was not created for t . However, the same tuple would also become disconnected of training example e in $J^{C'_j}$. Because we remove literals that are not head-connected in each clause in $\{C'_1, \dots, C'_m\}$, we would remove t from C'_j .

The sets of canonical database instances $\{J^{C'_1}, \dots, J^{C'_m}\}$ and $\{J^{C_1}, \dots, J^{C_n}\}$ are both generated using the function $StableInstances$ with the same matching dependencies Σ , and only contain tuples related to e . Therefore, $StableDefinitions(C) = \{C'_1, \dots, C'_m\}$ is equal to $\{C_1, \dots, C_n\}$. \square

This theorem shows that our algorithm considers exactly the same set of candidate clauses and definitions as the one that learns over each stable instance. We further show that this set is intact during generalization. Similar to above, assume that $Generalize(C, I, e', \Sigma)$ denotes the clause generated by generalizing clause C to cover example e' over (heterogeneous) database I with the set of MDs Σ . Also, let \mathbf{C} be a set of stable clauses and let $Generalize_s(\mathbf{C}, StableInstances(I, \Sigma), e')$ be the set of stable clauses generated by generalizing every stable clause in \mathbf{C} to cover example e' over every stable instance of I .

THEOREM 4.10. *Given database I with MDs Σ and set of examples E*

$$Generalize_s(\mathbf{C}, StableInstances(I, \Sigma), e') = StableDefinitions(Generalize(C, I, e', \Sigma))$$

PROOF. We show that $Generalize(C, I, e', \Sigma) = C^*$ is a compact representation of $Generalize_s(\mathbf{C}, \mathbf{J}, e') = \{C_1^*, \dots, C_n^*\}$, i.e., $StableDefinitions(C^*) = \{C_1^*, \dots, C_n^*\}$.

Assume that the schema is $\mathcal{R} = \{R_1(A, B), R_2(B, C)\}$ and we have MD $\phi : R_1[B] \approx R_2[B] \rightarrow R_1[B] \doteq R_2[B]$. Assume that database instance I contains tuples $R_1(a, b), R_2(b', c), R_2(b'', c)$, and that $b \approx b'$ and $b \approx b''$. Then, bottom-clause C has the form

$$\begin{aligned} T(u) \leftarrow & L'_1, \dots, L'_{l-1}, \\ & R_1(a, b), R_2(b', c), m_\phi(b, b'), R_2(b'', d), m_\phi(b, b''), \\ & L'_l, \dots, L'_n, \end{aligned}$$

where $L'_k, 1 \leq k \leq n$, is a literal.

Now consider two stable instances generated in $\text{StableInstances}(I, \Sigma)$: J_1 , which contains tuples $R_1(a, < b, b' >)$, $R_2(< b, b' >, c)$, $R_2(b'', c)$; and J_2 , which contains tuples $R_1(a, < b, b'' >)$, $R_2(b', c)$, $R_2(< b, b'' >, c)$. The bottom-clause C_1 over instance J_1 has the form

$$\begin{aligned} T(u) \leftarrow & L_1, \dots, L_{l-1}, \\ & R_1(a, < b, b' >), R_2(< b, b' >, c), R_2(b'', c), \\ & L_l, \dots, L_n, \end{aligned}$$

and the bottom-clause C_2 over instance J_2 has the form

$$\begin{aligned} T(u) \leftarrow & L_1, \dots, L_{l-1}, \\ & R_1(a, < b, b'' >), R_2(b', c), R_2(< b, b'' >, c), \\ & L_l, \dots, L_n, \end{aligned}$$

where $L_k, 1 \leq k \leq n$, is a literal.

We want to generalize C_1 to cover another training example e' . Let $G_{e'}$ be the ground bottom-clause for e' and $G'_{e'}$ be a stable clause of $G_{e'}$. The literals in C_1 that are blocking will depend on the content of the ground bottom-clause $G'_{e'}$. Assume that the sets of literals $\{L'_1, \dots, L'_n\}$ and $\{L_1, \dots, L_n\}$ are equal. We consider the following cases for the literals that are not equal. The same cases apply when we want to generalize any other clause generated from a stable instance, e.g., C_2 .

Case 1: $G'_{e'}$ contains the literals $R_1(a, < b, b' >)$ and $R_2(< b, b' >, c)$. In this case, $R_1(a, < b, b' >)$ and $R_2(< b, b' >, c)$ are *not* blocking literals, i.e., they are not removed from C_1 . $G_{e'}$ also contains literals $R_1(a, b), R_2(b', c), m_\phi(b, b')$. Therefore, the same literals are *not* blocking literals in C either.

Case 2: $G'_{e'}$ contains literals with same relation names but not the same pattern. Assume that $G'_{e'}$ contains the literals $R_1(a, b)$ and $R_2(d, c)$, i.e., they do not join. In this case, literal $R_2(< b, b' >, c)$ in C_1 is a blocking literal because it joins with a literal that appears previously in the clause, $R_1(a, < b, b' >)$. Hence, it is removed. $G_{e'}$ also contains literals $R_1(a, b)$ and $R_2(d, c)$. Because in clause $G'_{e'}$, created from the stable instance, these literals do not join, in $G_{e'}$ they do not join either. In this case, the blocking literal in C is $m_\phi(b, b')$. After literal $m_\phi(b, b')$ is removed, literal $R_2(b', c)$ is not head-connected. Therefore, it is also removed.

Case 3: $G'_{e'}$ contains $R_1(a, < b, b' >)$, but not $R_2(< b, b' >, c)$. In this case, literal $R_2(< b, b' >, c)$ is a blocking literal in C_1 . Therefore, it is removed. $G_{e'}$ also contains literals $R_1(a, b)$ and $m_\phi(b, b')$, but not $R_2(b', c)$. Therefore, literal $R_2(b', c)$ in C is also blocking and it is removed.

Case 4: $G'_{e'}$ contains $R_2(< b, b' >, c)$, but not $R_1(a, < b, b' >)$. This is similar to the previous case.

Case 5: $G'_{e'}$ does not contain either $R_1(a, < b, b' >)$ nor $R_2(< b, b' >, c)$. In this case, both $R_1(a, < b, b' >)$ and $R_2(< b, b' >, c)$ are blocking; hence they are removed. $G_{e'}$ does not contain literals $R_1(a, b), R_2(b', c), m_\phi(b, b')$. Hence,

these literals are also blocking literals in C and are removed.

The generalization operations $\text{Generalize}(C, I, e', \Sigma)$ and $\text{Generalize}_s(\mathbf{C}, \mathbf{J}, e')$ consist of removing blocking literals from C and \mathbf{C} respectively. We have shown that the same literals are blocking over both the clauses. Therefore, $\text{StableDefinitions}(C^*) = \{C_1^*, \dots, C_n^*\}$. \square

This property indicates that HDLearn leverages all information available in the heterogeneous database and learns all interesting patterns in its stable instances. Of course, as we follow a more restrictive semantic to provide more accurate definitions overall, our output may *not* be exactly as the one by the method that learns over each stable instance.

5. IMPLEMENTATION DETAILS

HDLearn is based on a relational learning system called Castor [33]. As Castor, HDLearn is implemented on top of a main-memory relational database management system. In this section, we explain implementation details of HDLearn.

5.1 Efficient Learning Through Sampling

The set of tuples I_e that HDLearn gathers to build a bottom-clause may be large if many tuples in I are relevant to e , particularly when learning over a large database. Thus, bottom-clause C_e would be very large, making the learning process extremely time-consuming. This problem is exacerbated when using similarity search operators to find relevant tuples, as many entities have similar names. To overcome this problem, HDLearn samples from the tuples in I_e to obtain a smaller tuple set $I_e^s \subseteq I_e$. To do so, HDLearn restricts the number of literals added to the bottom-clause per relation through a parameter called *sample size*. Let C_e be a bottom-clause associated with example e . A *sampled clause* C_e^s of clause C_e is the clause obtained the following way. Let I_R be the set of tuples in relation R that can be added to I_e (lines 10-11 in Algorithm 2). We obtain a *uniform* sample I_R^s of I_R and only add tuples in I_R^s to I_e . In a uniform sample, every tuple in I_R is sampled independently with the same probability. This sampling technique is used by existing relational learning algorithms [32, 33].

5.1.1 Approximate Coverage Testing

As mentioned in Section 4.2, HDLearn uses θ -subsumption to compute the coverage of candidate clauses. Ideally, a ground bottom-clause G_e for example e must contain one literal per each tuple in the database that is connected to e through some (similarity) joins. Otherwise, the θ -subsumption test may declare that C does *not* cover e when C actually covers e . However, it is expensive to check θ -subsumption for large clauses. Therefore, to improve efficiency, we use sampling to build ground bottom-clauses. Because we use a sampled bottom-clause G_e^s , checking whether $C \subseteq_\theta G_e^s$ is an approximation of checking whether $I \wedge C \models e$. Our empirical results in Section 6.4 indicate that this issue does *not* significantly impact the accuracy of the learned definitions. The learning algorithm involves many (thousands) coverage tests. Because HDLearn reuses ground bottom-clauses, it can run efficiently over large databases.

5.2 Matching Dependencies

To implement matching dependencies, HDLearn uses the similarity operator defined as the average of the *Smith-Waterman-Gotoh* and the *Length* similarity functions. The *Smith-Waterman-Gotoh* function [22] measures the similarity of two strings

Name	#R	#T	#P	#N
IMDB	9	3.3M	100	200
OMDB	15	4.8M		
DBLP	4	15K	500	1000
Google Scholar	4	328K		
JMDB	43	9.2M	1000	2000
BOM	8	92K		

Table 3: Numbers of relations (#R), tuples (#T), positive examples (#P), and negative examples (#N) for each dataset.

based on their local sequence alignments. The *Length* function computes the similarity of the length of two strings by dividing the length of the smaller string by the length of the larger string. Given a string s_1 , HDLearn considers a string s_2 similar to s_1 if their similarity score is greater than or equal to 0.65 and s_2 is within the top- k_m similar strings to s_1 . In our empirical evaluation in Section 6, we vary the value of k_m . To improve efficiency, we pre-compute the pairs of similar values according to the similarity operator. We use the pre-computed pairs as an index for similarity search.

6. EXPERIMENTS

We empirically evaluate HDLearn to answer the following questions:

1. Can HDLearn learn over heterogeneous databases effectively and efficiently? (Section 6.2)
2. What is the benefit of using matching dependencies during learning? (Section 6.2)
3. How does the number of training examples affect HDLearn’s effectiveness and efficiency? (Section 6.3)
4. How does sampling affect HDLearn’s effectiveness and efficiency? (Section 6.4)

6.1 Experimental Settings

6.1.1 Datasets

We use three pairs of databases whose statistics are shown in Table 3.

1. **IMDB + OMDB:** The IMDB and OMDB databases contain information about movies [11]. We learn the target relation *dramaRestrictedMovies(imdbId)*, which contains the *imdbId* of movies that are of *drama* genre and are rated R. The *imdbId* is only contained in the IMDB database, the genre information is contained in both databases, and the rating information is only contained in the OMDB database. Therefore, in order to learn a definition for this relation, one needs information from both databases. We specify the MD

$$IMDB.movies[title] \approx OMDB.movies[title] \rightarrow$$

$$IMDB.movies[title] \Rightarrow OMDB.movies[title].$$

The original database [11] contains ground truth that matches IMDB and OMDB movies. We use these to generate positive examples. To generate negative examples, we generate movies that are not of *drama* genre or rated R and sample them to obtain twice the number of positive examples. We refer to this dataset with one MD as **IMDB + OMDB (one MD)**. We also create MDs that match cast members and writer names between the two databases. We refer to the dataset that contains the three MDs as **IMDB + OMDB (three MDs)**.

2. **DBLP + Google Scholar:** The DBLP and Google Scholar databases contain information about academic

papers [11]. We learn the target relation *gsPaperYear(gsId, year)*, which contains the Google Scholar id *gsId* and the *year* of publication of the paper as indicated in the DBLP database. We specify the MDs

$$DBLP.papers[title] \approx GoogleScholar.papers[title] \rightarrow$$

$$DBLP.papers[title] \Rightarrow GoogleScholar.papers[title]$$

$$DBLP.papers[venue] \approx GoogleScholar.papers[venue] \rightarrow$$

$$DBLP.papers[venue] \Rightarrow GoogleScholar.papers[venue].$$

The original database contains ground truth that matches DBLP and Google Scholar papers. We use these to generate positive examples. To generate negative examples, we generate pairs of Google Scholar ids and years of publication that are not correct, and sample them to obtain twice the number of positive examples.

3. **JMDB + BOM:** We scrape the Box Office Mojo (BOM) website to obtain a list of movies and their total grossing. We use the JMDB database (*jmdb.de*), which contains information from the IMDb website in relational format. We use these databases to learn a definition for the target relation *highGrossing(title)*, which indicates that the movie with title *title* is high grossing. From BOM, we obtain the top 1K grossing movies and use them as positive examples, and obtain the lowest 2K grossing movies and use them as negative examples. We specify the MD

$$BOM.highGrossing[title] \approx JMDB.movies[title] \rightarrow$$

$$BOM.highGrossing[title] \Rightarrow JMDB.movies[title]$$

Notice that the learning time of relational learning algorithms is not only affected by the number of training examples, but also by the number of tuples in the underlying database. HDLearn builds bottom-clauses for training examples. Learning over a large database often results in large bottom-clauses, which translates to long learning times. Therefore, learning over relational databases is generally time-consuming [33, 41]. In Section 6.3 we evaluate the impact of the number of training examples on HDLearn.

6.1.2 Systems and Environment

We compare HDLearn against two baseline systems.

1. **Castor-NoMDs:** We use the relational learning system Castor [33] to learn over the original databases. Castor does *not* use MDs. We use Castor because it scales to large databases.
2. **Castor-ExactMatch:** We use Castor [33], but allow the attributes that appear in an MD to be joined through exact joins. Therefore, this system uses information from MDs but only considers exact matches between values.
3. **HDLearn:** We use HDLearn, as described in Section 4. HDLearn uses the parameter *sample size* to restrict the size of (ground) bottom-clauses. We fix *sample size* to 10. In Section 6.4, we evaluate the impact of this parameter on HDLearn’s effectiveness and efficiency.

We perform 5-fold cross validation over all datasets. We evaluate the F1-score and learning time, showing the average over the cross validation. F1-score, which is the harmonic average of the precision and recall, measures the effectiveness of the learned definitions.

All systems use 16 threads to parallelize coverage testing. All experiments were run on a server with 30 2.3GHz Intel Xeon E5-2670 processors, running CentOS Linux with 500GB of main memory.

6.2 Effectiveness and Efficiency

Table 4 shows the results of learning over all datasets using HDLearn and the baseline systems. Over all datasets, HDLearn obtains a better F1-score than the baselines. This is expected as HDLearn uses information in the input MDs to find relevant information from all databases. Castor-NoMDs does not learn any definition in the DBLP + Google Scholar dataset. This is because Castor cannot access information from the DBLP database. Therefore, it is not able to find a reasonable definition. We see a similar result in the JMDB + BOM dataset. Castor-ExactMatch is able to learn a definition over all datasets. However, as it relies on exact matches, the learned definitions are not as effective as the definitions learned by HDLearn. In particular, we see that Castor-ExactMatch obtains a low F1-score over the JMDB + BOM dataset. This is because the same entities over these databases are consistently represented by different names. On the other hand, Castor-ExactMatch obtains a competitive F1-score in the IMDB + BOM dataset with three MDs. The MDs that match cast members and writer names between the two databases contain many exact matches. Therefore, Castor-ExactMatch is able to find paths that connect movies in IMDB with movies in OMDB.

HDLearn is able to learn effective definitions over heterogeneous databases efficiently. In general, as k_m , the number of matches considered in MDs, increases, HDLearn’s effectiveness increases. This shows that, even though the number of incorrect matches by the similarity function may increase, HDLearn is able to ignore these false matches during learning. In other words, using the training data, HDLearn figures that some patterns are useful and some are not. However, this is not always the case. When using multiple MDs or when learning a difficult concept, a high value for k_m affects HDLearn’s effectiveness. In these cases, incorrect matches represent noise that affects HDLearn ability to learn an effective definition. This is reflected in the IMDB + OMDB (three MDs) dataset and in the JMDB + BOM dataset. For this reason, it is *not* always better to use multiple MDs. HDLearn obtains an effective definition over the IMDB + OMDB with one MD and is significantly more efficient than using three MDs.

HDLearn performs best when it uses the right value for k_m . Therefore, one should find the right value for k_m by experimenting with different values. In general, HDLearn performs well with a sufficiently small k_m , with the benefit that it can learn directly from the dirty data.

One extreme is to use $k_m = 1$, i.e., keep only the top similar match. This can be seen as generating only one stable instance and learning over it. Although efficient, the results are not always effective. Therefore, there is a benefit of looking at multiple stable instances for learning. As k_m increases, the learning time also increases. Therefore, there is a trade-off between the looseness of the similarity function used in MDs and HDLearn’s effectiveness and efficiency.

The other extreme is to generate all stable instances, and then learn over each of them. However, there may exist an extremely large number of stable instances. For instance, the movie *American Splendor - 2003* in IMDB matches 145 movies in the OMDB database. Just from this movie, we would get 145 distinct stable instances. Integrating and cleaning databases would take significant amount of time and manual labor. Materializing all stable instances would take a huge amount of space. Further, learning over each

stable instance would be time-consuming. Therefore, the biggest benefit of HDLearn is that it avoids these problems by learning directly from dirty data.

The definitions learned by HDLearn and Castor are interpretable. Therefore, we can analyze the learned definitions. For instance, over the JMDB + BOM dataset, HDLearn learns the following definition:

$$\begin{aligned} \text{highGrossing}(x) \leftarrow & \text{movies}(y, x, z), \text{mov2genres}(y, \text{'Fantasy'}), \\ & \text{mov2countries}(y, \text{'USA'}), \text{releaseseason}(y, \text{'Winter'}), \\ & \text{technical}(y, \text{'PFM : 35mm'}). \end{aligned}$$

On the other hand, Castor-NoMDs is *not* able to learn any definition. Castor-ExactMatch can only access information in JMDB for movies whose titles match exactly with titles in BOM, and it learns the following definition:

$$\text{highGrossing}(x) \leftarrow \text{movies}(y, x, z), \text{mov2genres}(y, \text{'Fantasy'}).$$

6.3 Scalability of HDLearn

One advantage of relational learning algorithms is that they are data-efficient, i.e., they can learn effective definitions from a small number of training examples [16]. However, for a given learning task, it is not clear how many examples are needed. In this section, we evaluate the impact of the number of training examples in both HDLearn’s effectiveness and efficiency. We use the IMDB + OMDB (three MDs) dataset and fix $k_m = 2$. We generate 2100 positive and 4200 negative examples. From these sets, we use 100 positive and 200 negative examples for testing. From the remaining examples, we generate training sets containing 100, 500, 1000, and 2000 positive examples, and double the amount of negative examples. For each training set, we use HDLearn to learn a definition. Figure 1 (left) shows the F1-scores and learning times for each training set. With 100 positive and 200 negative examples, HDLearn obtains an F1-score of 0.80. With 500 positive and 1000 negative examples, the F1-score increases to 0.91. More training examples do not impact the F1-score significantly. On the other hand, the learning time consistently increases with the number of training examples. Nevertheless, HDLearn is able to learn efficiently even with the largest training set.

6.4 Impact of Sampling

In this section, we evaluate the impact of sampling (Section 5.1) on HDLearn’s effectiveness and efficiency. We use the IMDB + OMDB (three MDs) dataset and fix $k_m = 2$ and $k_m = 5$. We use 800 positive and 1600 negative examples for training, and 200 positive and 400 negative examples for testing. Figure 1 (middle and right) show the F1-score and learning time of HDLearn with $k_m = 2$ and $k_m = 5$, respectively, when varying the sample size. For both values of k_m , the F1-score does not change significantly with different sampling sizes. With $k_m = 2$, the learning time remains almost the same with different sampling sizes. However, with $k_m = 5$, the learning time increases significantly. Therefore, using a small sample size is enough for learning an effective definition efficiently.

7. RELATED WORK

Surveys on relational learning algorithms and applications can be found in [31, 21]. We have witnessed a surge in building efficient relational learning algorithms for large databases [41, 40, 18, 33]. Our work complements this line

Dataset	Metric	Castor- NoMDs	Castor- ExactMatch	HDLearn			
				$k_m = 1$	$k_m = 2$	$k_m = 5$	$k_m = 10$
IMDB + OMDB (one MD)	F1-score	0.47	0.59	0.86	0.90	0.92	0.92
	Learning time (m)	0.12	0.13	0.18	0.26	0.42	0.87
IMDB + OMDB (three MDs)	F1-score	0.47	0.82	0.86	0.90	0.93	0.89
	Learning time (m)	0.12	0.48	0.21	0.30	25.87	285.39
DBLP + Google Scholar	F1-score	0	0.54	0.61	0.67	0.71	0.82
	Learning time (m)	2.5	2.5	3.1	2.7	2.7	2.7
JMDB + Box Office Mojo	F1-score	0	0.01	0.83	0.85	0.66	0.67
	Learning time (m)	0.3	0.8	5.2	8.5	24.2	33.7

Table 4: Results of learning over all datasets. Number of top similar matches denoted by k_m .

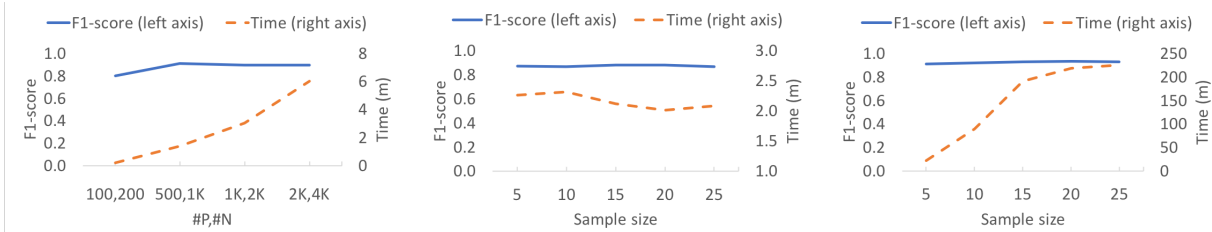


Figure 1: Results of learning over the IMDB+OMDB (three MDs) dataset while increasing the number of positive and negative ($\#P, \#N$) training examples (left) and while increasing sample size for $k_m = 2$ (middle) and $k_m = 5$ (right).

of work by creating an efficient relational learning system over heterogeneous databases.

Entity resolution is an active and important area of research in data management whose aim is to find whether two data values or tuples refer to the same real-world entity [20]. One may use these efforts in implementing precise similarity operation in MDs.

There have been multiple efforts on modeling value and entity heterogeneity in databases using declarative constraints and defining the properties of clean instances under them [17, 8, 5, 9, 7, 38, 26, 3, 19, 25]. Our effort extends the same approach of dealing with heterogeneous data by efficiently learning over the original database.

A different approach to querying and analyzing multiple databases is to use the available constraints, including MDs, define concrete similarity and matching functions for the databases, and create and materialize a unified database instance [10, 35, 13]. One may represent the inherent uncertainties in similarities between producing matched values as probabilities and store the final integrated database in a probabilistic database [37]. This approach is preferable if one is able to define a similarity function that accurately quantifies the degree of similarity between values in the domain and a matching function that is able to generate the correct unified value from its input. Such a matching function is hard and resource-intensive to define and usually requires a great deal of expert attention and manual labor. Also, one is *not* usually able to define a general one that can be used across various domains. Our goal is to create a learning algorithm that can be used over databases from various domains. Moreover, if the input databases are relatively large, the materialized database has to accommodate tuples from numerous stable instance causing the database to become extremely large. This in turn makes efficient learning over such a database challenging. Our proposed method is able to learn over the original database, which contains significantly fewer tuples and avoids this challenge.

ActiveClean gradually cleans a dirty dataset to learn a convex-loss model, such as Logistic Regression [28]. Its goal

is to clean the underlying dataset such that the learned model becomes more effective as it receives more cleaned records. Cleaning may be done manually by an expert. Our objective, however, is to learn a model over dirty data without cleaning and transforming it. Furthermore, ActiveClean does *not* address the problem of having multiple cleaned versions of the database and assumes that the dataset has a unique clean instance.

8. CONCLUSION

In this paper, we investigated the problem of learning directly over heterogeneous data. We leveraged the concept of declarative constraint for data cleaning and defined a novel semantic for learning over such databases. We have proposed a new algorithm to learn over unclean data. Our empirical results show that our proposed algorithm is efficient and effective over large datasets.

9. REFERENCES

- [1] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases: The Logical Level*. Addison-Wesley, 1994.
- [2] A. Abouzied, D. Angluin, C. Papadimitriou, J. Hellerstein, and A. Silberschatz. Learning and verifying quantified boolean queries by example. In *PODS*, 2013.
- [3] A. Arasu, S. Chaudhuri, and R. Kaushik. Transformation-based framework for record matching. In *2008 IEEE 24th International Conference on Data Engineering*, pages 40–49, April 2008.
- [4] M. Arenas, L. Bertossi, and J. Chomicki. Consistent query answers in inconsistent databases. In *PODS*, 1999.
- [5] Z. Bahmani, L. E. Bertossi, S. Kolahi, and L. V. S. Lakshmanan. Declarative entity resolution via matching dependencies and answer set programs. In *Principles of Knowledge Representation and Reasoning: Proceedings of the Thirteenth International*

- Conference, KR 2012, Rome, Italy, June 10-14, 2012, 2012.*
- [6] Z. Bahmani, L. E. Bertossi, and N. Vasiloglou. Erblox: Combining matching dependencies with machine learning for entity resolution. In *SUM*, 2015.
 - [7] O. Benjelloun, H. Garcia-Molina, D. Menestrina, Q. Su, S. E. Whang, and J. Widom. Swoosh: A generic approach to entity resolution. *The VLDB Journal*, 18(1):255–276, Jan. 2009.
 - [8] L. E. Bertossi, S. Kolahi, and L. V. S. Lakshmanan. Data cleaning and query answering with matching dependencies and matching functions. *Theory of Computing Systems*, 52:441–482, 2011.
 - [9] D. Burdick, R. Fagin, P. G. Kolaitis, L. Popa, and W.-C. Tan. A declarative framework for linking entities. *ACM Trans. Database Syst.*, 41(3):17:1–17:38, July 2016.
 - [10] X. Chu, I. F. Ilyas, S. Krishnan, and J. Wang. Data cleaning: Overview and emerging challenges. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD Conference 2016, San Francisco, CA, USA, June 26 - July 01, 2016*, pages 2201–2206, 2016.
 - [11] S. Das, A. Doan, P. S. G. C., C. Gokhale, and P. Konda. The magellan data repository. <https://sites.google.com/site/anahidgroup/projects/data>.
 - [12] L. De Raedt. *Logical and Relational Learning*. Springer Publishing Company, Incorporated, 1st edition, 2010.
 - [13] D. Deng, R. C. Fernandez, Z. Abedjan, S. Wang, M. Stonebraker, A. K. Elmagarmid, I. F. Ilyas, S. Madden, M. Ouzzani, and N. Tang. The data civilizer system. In *CIDR 2017, 8th Biennial Conference on Innovative Data Systems Research, Chaminade, CA, USA, January 8-11, 2017, Online Proceedings*, 2017.
 - [14] A. Doan, A. Halevy, and Z. Ives. *Principles of Data Integration*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1st edition, 2012.
 - [15] T. Eiter, G. Gottlob, and H. Mannila. Disjunctive datalog. *ACM Trans. Database Syst.*, 22(3):364–418, Sept. 1997.
 - [16] R. Evans and E. Grefenstette. Learning explanatory rules from noisy data. *J. Artif. Intell. Res.*, 61:1–64, 2018.
 - [17] W. Fan, X. Jia, J. Li, and S. Ma. Reasoning about record matching rules. *Proc. VLDB Endow.*, 2(1):407–418, Aug. 2009.
 - [18] L. Galárraga, C. Teflioudi, K. Hose, and F. Suchanek. Fast Rule Mining in Ontological Knowledge Bases with AMIE+. In *VLDB Journal*, 2015.
 - [19] H. Galhardas, D. Florescu, D. Shasha, E. Simon, and C.-A. Saita. Declarative data cleaning: Language, model, and algorithms. In *Proceedings of the 27th International Conference on Very Large Data Bases, VLDB '01*, pages 371–380, San Francisco, CA, USA, 2001. Morgan Kaufmann Publishers Inc.
 - [20] L. Getoor and A. Machanavajjhala. Entity resolution for big data. In *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '13*, pages 1527–1527, New York, NY, USA, 2013. ACM.
 - [21] L. Getoor and B. Taskar. *Introduction to Statistical Relational Learning*. MIT Press, 2007.
 - [22] O. Gotoh. An improved algorithm for matching biological sequences. *Journal of molecular biology*, 162 3:705–8, 1982.
 - [23] A. Halevy, M. Franklin, and D. Maier. Principles of dataspace systems. In *Proceedings of the Twenty-fifth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS '06*, pages 1–9, New York, NY, USA, 2006. ACM.
 - [24] J. M. Hellerstein. People, computers, and the hot mess of real data. In *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16*, pages 7–7, New York, NY, USA, 2016. ACM.
 - [25] M. Hernández, G. Koutrika, R. Krishnamurthy, L. Popa, and R. Wisnesky. Hil: A high-level scripting language for entity integration. In *Proceedings of the 16th International Conference on Extending Database Technology, EDBT '13*, pages 549–560, New York, NY, USA, 2013. ACM.
 - [26] M. A. Hernández and S. J. Stolfo. The merge/purge problem for large databases. In *Proceedings of the 1995 ACM SIGMOD International Conference on Management of Data, SIGMOD '95*, pages 127–138, New York, NY, USA, 1995. ACM.
 - [27] D. V. Kalashnikov, L. V. Lakshmanan, and D. Srivastava. Fastqre: Fast query reverse engineering. In *Proceedings of the 2018 International Conference on Management of Data, SIGMOD '18*, pages 337–350, New York, NY, USA, 2018. ACM.
 - [28] S. Krishnan, J. Wang, E. Wu, M. J. Franklin, and K. Goldberg. Activeclean: Interactive data cleaning for statistical modeling. 9(12), 2016.
 - [29] H. Li, C.-Y. Chan, and D. Maier. Query from examples: An iterative, data-driven approach to query construction. *PVLDB*, 8(13), 2015.
 - [30] L. Mihalkova and R. Mooney. Bottom-Up Learning of Markov Logic Network Structure. In *ICML*, 2007.
 - [31] S. Muggleton, L. Raedt, D. Poole, I. Bratko, P. Flach, K. Inoue, and A. Srinivasan. Ilp turns 20. *Mach. Learn.*, 86(1):3–23, Jan. 2012.
 - [32] S. Muggleton, J. C. A. Santos, and A. Tamaddoni-Nezhad. Progolem: A system based on relative minimal generalisation. In *ILP*, volume 5989, 2009.
 - [33] J. Picado, A. Termehchy, A. Fern, and P. Ataei. Schema Independent Relational Learning. In *SIGMOD*, 2017.
 - [34] J. R. Quinlan. Learning Logical Definitions From Relations. *Machine Learning*, 5, 1990.
 - [35] T. Rekatsinas, X. Chu, I. F. Ilyas, and C. Ré. Holoclean: Holistic data repairs with probabilistic inference. *Proc. VLDB Endow.*, 10(11):1190–1201, Aug. 2017.
 - [36] M. Richardson and P. Domingos. Markov logic networks. *Machine Learning*, 62(1-2):107–136, Feb. 2006.
 - [37] D. Suciu, D. Olteanu, C. Ré, and C. Koch. *Probabilistic Databases*. Synthesis Lectures on Data Management. Morgan & Claypool Publishers, 2011.

- [38] M. Weis, F. Naumann, U. Jehle, J. Lufter, and H. Schuster. Industry-scale duplicate detection. *Proc. VLDB Endow.*, 1(2):1253–1264, Aug. 2008.
- [39] Z. Yan, N. Zheng, Z. G. Ives, P. P. Talukdar, and C. Yu. Actively soliciting feedback for query answers in keyword search-based data integration. *PVLDB*, 6(3):205–216, 2013.
- [40] X. Yin, J. Han, J. Yang, and P. S. Yu. CrossMine: Efficient Classification Across Multiple Database Relations. In *ICDE*, 2004.
- [41] Q. Zeng, J. M. Patel, and D. Page. Quickfoil: Scalable inductive logic programming. *Proc. VLDB Endow.*, 8(3):197–208, Nov. 2014.