

Learning Concepts Definable in First-Order Logic with Counting

Steffen van Bergerem
RWTH Aachen University
vanbergerem@informatik.rwth-aachen.de

Abstract—We study classification problems over relational background structures for hypotheses that are defined using logics with counting. The aim of this paper is to find learning algorithms running in time sublinear in the size of the background structure. We show that hypotheses defined by FOCN(P)-formulas over structures of polylogarithmic degree can be learned in sublinear time. Furthermore, we prove that for structures of unbounded degree there is no sublinear learning algorithm for first-order formulas.

I. INTRODUCTION

In this paper, we study Boolean classification problems over relational structures. We consider the relational structure, also called *background structure*, as fixed. For fixed $k \in \mathbb{N}$, the goal is to learn a classification function, called a *hypothesis*, that maps k -ary tuples from the relational structure to Booleans. Given a sequence of training examples, each of which consists of a k -ary tuple from the background structure and a Boolean, we aim to find a hypothesis that is consistent with the examples. In other words, our goal is to learn a description of a k -ary relation that is consistent with a given sequence of positive and negative examples, i.e., the relation contains all positive and no negative examples.

Example I.1. Consider a relational database containing data from an online encyclopedia. The universe of the structure consists of all pages of the encyclopedia. We have a binary relation representing hyperlinks between pages and a unary relation representing category pages. Pages that are not category pages are article pages. Our task is to learn a binary relation containing tuples of pages, where the first element of the tuple is a category page and the second element is a page that belongs to the category. For this we are given a training sequence of classified tuples, e.g. tuples that have been classified by experts beforehand. The goal is to learn the description of a relation that is consistent with the training sequence.

Consider the background structure and the training examples given in Figure 1. A description of a consistent relation would be the following. The relation contains all tuples, where the first element is a category page and the second page is linked from the category page or there is another page linked from the category page and both pages have at least two common linked pages. The corresponding relation can be seen in Figure 2. For example the tuple $(1, 8)$ is contained in the relation since Page 1 is a category, Page 2 is linked from the category, and there

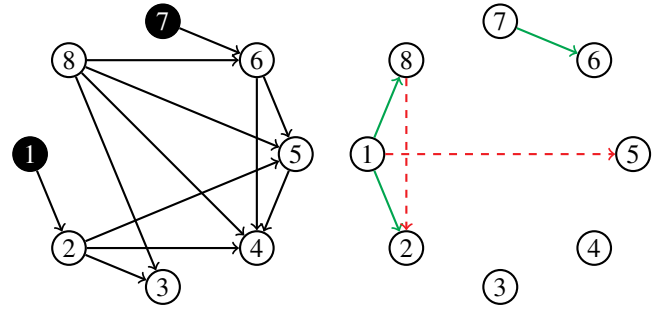


Figure 1. (Left) A background structure viewed as a directed graph. Vertices represent pages in the online encyclopedia, category pages are black and edges represent hyperlinks. (Right) Training examples. Green edges denote positive examples, i.e., the tuple is contained in the relation. Red edges denote negative examples.

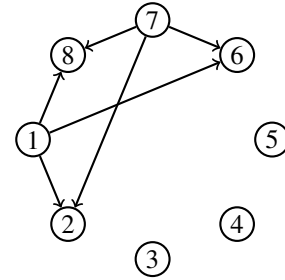


Figure 2. The learned relation from Example I.1.

are at least two pages (Pages 3, 4, and 5) that both Page 2 and Page 8 link to. Note that the relation is consistent with the training examples.

Since the data are contained in a relational database, it would be convenient to learn an SQL query that defines the relation. Figure 3 shows an SQL query for the learned relation. ┘

This paper studies the classification problems in the declarative framework that has been introduced in [1], [2], where logics are used to describe the hypotheses. Grohe and Ritzert showed in [2] that learning hypotheses in first-order logic is possible in time polynomial in the number of training examples and the degree of the background structure. For structures of polylogarithmic degree and training sequences of polylogarithmic length, measured in the size of the background structure, this yields a learning algorithm that runs in polylogarithmic time. We are interested in learning hypotheses that can be expressed in SQL. While first-order logic can be seen as

```

SELECT C.'page', CatLink.'to'
FROM Categories C, Links CatLink
WHERE CatLink.'from' = C.'page'
UNION
SELECT C.'page', L1.'from'
FROM Categories C, Links CatLink,
      Links L1, Links L2
WHERE CatLink.'from' = C.'page'
AND CatLink.'to' = L2.'from'
AND L1.'to' = L2.'to'
GROUP BY C.'page', L1.'from'
HAVING count(*) >= 2;

```

Figure 3. An SQL query that defines the relation learned in Example I.1.

the “logical core” of SQL, there are aggregating operators in SQL, namely COUNT, AVG, SUM, MIN, and MAX, that do not have corresponding expressions in first-order logic. Motivated by this, we study the logic FOCN(P), which Kuske and Schweikardt introduced in [3] and which extends first-order logic by cardinality conditions similar to the COUNT operator in SQL. The logic depends on a collection of numerical predicates P , i.e., functions $P: \mathbb{Z}^m \rightarrow \{0, 1\}$, that we use in FOCN(P)-formulas to express restrictions on the results of counting terms.

Let $\mathcal{B} = (U(\mathcal{B}), L, C)$ be the background structure from Example I.1, where $U(\mathcal{B})$ is the set of all pages, L is the binary relation of links between pages and C is the unary relation of category pages. The SQL query from Figure 3 can be expressed as the FOCN(P)-formula

$$\varphi(c, p) = Cc \wedge \left(Lcp \vee \exists x (Lcx \wedge \#(y).(Lxy \wedge Lpy) \geq 2) \right),$$

where $\geq \in P$. The counting term $\#(y).(Lxy \wedge Lpy)$ counts the number of pages y such that both x and p link to y . The formula $\#(y).(Lxy \wedge Lpy) \geq 2$ checks whether this number is at least 2. In a more general approach, we may use the formula

$$\varphi'(c, p; \kappa) = Cc \wedge \left(Lcp \vee \exists x (Lcx \wedge \#(y).(Lxy \wedge Lpy) \geq \kappa) \right)$$

with the free number variable κ . When viewed as a parameter, for every assignment of κ we obtain a new hypothesis.

A hypothesis consists of the following components: a formula $\varphi(\bar{x}; \bar{y}, \bar{\kappa})$ as well as two parameter tuples \bar{v} in $U(\mathcal{B})^\ell$ and $\bar{\lambda}$ in \mathbb{Z}^m . Together, they describe a function $\llbracket \varphi(\bar{x}; \bar{v}, \bar{\lambda}) \rrbracket^{\mathcal{B}}: U(\mathcal{B})^k \rightarrow \{0, 1\}$ with

$$\llbracket \varphi(\bar{x}; \bar{v}, \bar{\lambda}) \rrbracket^{\mathcal{B}}(\bar{u}) := \begin{cases} 1 & \text{if } \mathcal{B} \models \varphi(\bar{u}; \bar{v}, \bar{\lambda}) \\ 0 & \text{otherwise.} \end{cases}$$

Given a finite training sequence T of tuples (\bar{u}, c) with $\bar{u} \in U(\mathcal{B})^k$ and $c \in \{0, 1\}$, our goal is to learn a hypothesis $(\varphi(\bar{x}; \bar{y}, \bar{\kappa}), \bar{v}, \bar{\lambda})$ that is consistent with the training examples, i.e., $\llbracket \varphi(\bar{x}; \bar{v}, \bar{\lambda}) \rrbracket^{\mathcal{B}}(\bar{u}) = c$ for all training examples (\bar{u}, c) . In the context of machine learning, the hypothesis is a *parametric model* and the described learning problem is called *model learning*.

A. Our Results

We study the model-learning problem for FOCN(P) over relational background structures. Instead of random access to the background structure, the algorithms we consider are granted only *local access* (see Section II for details). We measure the complexity of FOCN(P)-formulas in terms of their *binding width* and *binding rank* [3]. The binding width bounds the number of variables that occur in quantifier and counting terms. The binding rank bounds the nesting depth of these constructs. We bound the complexity of first-order formulas in Hanf normal form in terms of their *locality radius*. This allows us to consider only local neighborhoods in the model-checking problem for the given formula. We give a more detailed description of these parameters in Section II.

We give an algorithm for the model-learning problem for FOCN(P) that runs in time polylogarithmic in the size of the background structure, polynomial in the length of the training sequence, and quasipolynomial in the degree of the background structure. The behavior of the algorithm depends on the existence of a *target hypothesis* that is consistent with the training sequence and whose FOCN(P)-formula has binding rank at most r and binding width at most w . Both r and w are parameters of the algorithm. If there is no such hypothesis, then the algorithm may reject the input. Otherwise it always returns a hypothesis. Although the FOCN(P)-formula of the target hypothesis might contain counting terms, the algorithm will only return first-order formulas. This is due to the surprising fact that on a fixed background structure first-order formulas are as expressive as FOCN(P)-formulas. The following theorem is the main result of this paper.

Theorem I.1. *Let $k, \ell, r, w \in \mathbb{N}$. Then there is a learning algorithm \mathcal{L}_{con} for the k -ary learning problem over some finite relational structure \mathcal{B} , that receives a training sequence T and the degree of the structure $\Delta \mathcal{B}$ as input, with the following properties:*

- (1) *If there are an FOCN(P)-formula $\varphi(\bar{x}; \bar{y}, \bar{\kappa})$ of binding rank at most r and binding width at most w and parameter tuples $\bar{v} \in U(\mathcal{B})^\ell$ and $\bar{\lambda} \in \{0, \dots, |\mathcal{B}|\}^{|\bar{\kappa}|}$ such that $\llbracket \varphi(\bar{x}; \bar{v}, \bar{\lambda}) \rrbracket^{\mathcal{B}}$ is consistent with T , then \mathcal{L}_{con} always returns a hypothesis.*
- (2) *If the algorithm returns a hypothesis H , then H is of the form $(\varphi^*(\bar{x}; \bar{y}), \bar{v}^*)$ for some first-order formula $\varphi^*(\bar{x}; \bar{y})$ in Hanf normal form with locality radius smaller than $(2w+1)^r$ and $\bar{v}^* \in U(\mathcal{B})^\ell$, and $\llbracket \varphi^*(\bar{x}; \bar{v}^*) \rrbracket^{\mathcal{B}}$ is consistent with the input sequence T of training examples.*
- (3) *The algorithm runs in time $(\log n + t)^{\mathcal{O}(1)} \cdot d^{\text{polylog}(d)}$ with only local access to \mathcal{B} , where $n := |\mathcal{B}|$, $d := \Delta \mathcal{B}$ and $t := |T|$.*
- (4) *The hypothesis returned by the algorithm can be evaluated in time $(\log n + t)^{\mathcal{O}(1)} \cdot d^{\text{polylog}(d)}$ with only local access to \mathcal{B} .*

When the degree of the background structure and the length of the training sequence are bounded by $\text{polylog}(|\mathcal{B}|)$, then

the running time of the algorithm is sublinear in the size of the background structure. Thus we obtain the following result.

Corollary I.2. *There is a consistent model-learning algorithm for FOCN(P)-formulas with only local access on background structures with polylogarithmic degree that runs in sublinear time on training sequences of polylogarithmic length, measured in the size of the background structure.*

Theorem I.1 is a direct generalization of the corresponding theorem for first-order logic due to Grohe and Ritzert [2], albeit with a slightly worse running time that is quasipolynomial in the degree. This generalization is well-motivated by the fact that typical machine-learning models have numerical parameters; our results may be seen as a first step towards including numerical aspects in the declarative framework. While the generalization may seem straightforward at first sight, at least for background structures of small (say, logarithmic) but unbounded degree, it is not obvious that an extension of the first-order result to FOCN(P) holds at all. The reason is that FOCN(P) loses its strong locality properties on structures of unbounded degree. For example, by comparing the degree sequences of the neighbors of nodes one can establish quite complex relations that may range over long distances. Indeed, as shown by Grohe and Schweikardt [4], various algorithmic meta theorems, whose proofs are also based on locality properties, fail when extended from first-order logic to first-order logic with counting. We show that it suffices to consider only FO-formulas to find a consistent hypothesis. However, the quantifier rank of the FO-formulas is polynomial in the degree of the background structure. Hence, a direct application of Grohe’s and Ritzert’s results [2] would not yield a sublinear-time learning algorithm, since the running time of their algorithm is non-elementary in the quantifier rank.

Thus it is not surprising that, even though our theorem looks similar to the corresponding result for first-order logic, there are significant differences in the proofs. The proof of the first-order result in [2] is based on Gaifman’s theorem, but there is no analogue of Gaifman’s theorem for the counting logic FOCN(P). Instead, our proof is based on Hanf’s theorem. But this causes the technical difficulty that we have to deal with isomorphism types of local neighborhoods in our structures. To be able to do this within the desired time bounds, we apply a recent new graph isomorphism test running in time $n^{\text{polylog}(d)}$ for n -vertex graphs of maximum degree d [5].

In addition to the results for structures of polylogarithmic degree, in Section IV we obtain a *probably approximately correct* (PAC) learning algorithm for structures of bounded degree, i.e., the algorithm returns on most of the training sets (probably) a hypothesis that has a small expected error on new examples (approximately). One could also say that the hypotheses returned by the algorithm *generalize well* on unseen examples with high probability.

We also investigated learnability on background structures without a degree restriction and obtained the following negative result.

Theorem I.3. *There is no consistent model-learning algorithm for first-order formulas with only local access on background structures with no degree restriction whose running time is sublinear in the size of the background structure.*

B. Related Work

The descriptive framework has been considered in [1], [2], [6]. Grohe and Ritzert [2] showed that first-order formulas are PAC-learnable over background structures of polylogarithmic degree. We did not obtain an analogous result for FOCN(P)-formulas, but we prove PAC-learnability over structures of bounded degree and learnability of consistent hypotheses over structures of polylogarithmic degree. Grohe, Löding and Ritzert [6] obtained learning algorithms for monadic second-order logic over string data.

The framework of inductive logic programming (ILP) is closely related to the framework we consider (see, for example, [7], [8], [9], [10], [11]). One of the main differences is that we encode the background knowledge in a relational structure, whereas in ILP it is represented in a background theory. Furthermore, ILP focuses on first-order logic, whereas in our framework different logics have been considered. Other related learning frameworks in the context of databases can be found in [12], [13].

In [3] Kuske and Schweikardt introduced FOCN(P), which extends first-order logic by counting quantifiers and numerical predicates. The logic generalizes logics like FO(Cnt) from [14] and FO+C from [15]. Our results rely on the fact that Hanf normal forms for FOCN(P) always exist. We use the structure of the normal form to argue that considering FO-formulas for our hypotheses suffices for fixed background structures. Hella et al. studied other aggregating operators in [16].

II. BACKGROUND FROM LOGICS

A. Structures

We only consider finite structures. A *signature* is a finite set σ of *relation symbols* and *constant symbols*. Every relation symbol $R \in \sigma$ has an arity $\text{ar}(R) \in \mathbb{N}$. A signature is called *relational* if it does not contain any constant symbol. A (σ) -*structure* \mathcal{A} consists of a finite set $U(\mathcal{A})$, called the *universe* of \mathcal{A} , a relation $R(\mathcal{A}) \subseteq (U(\mathcal{A}))^{\text{ar}(R)}$ for every relation symbol $R \in \sigma$, and an element $c(\mathcal{A})$ for every constant symbol $c \in \sigma$. The *order* of \mathcal{A} is $|\mathcal{A}| := |U(\mathcal{A})|$.

A structure \mathcal{B} is a *substructure* of \mathcal{A} , denoted by $\mathcal{B} \subseteq \mathcal{A}$, if $U(\mathcal{B}) \subseteq U(\mathcal{A})$, $R(\mathcal{B}) \subseteq R(\mathcal{A})$ for every relation symbol $R \in \sigma$ and $c(\mathcal{B}) = c(\mathcal{A})$ for every constant symbol $c \in \sigma$. For a relational structure \mathcal{A} and a set $V \subseteq U(\mathcal{A})$, the structure *induced* by \mathcal{A} on V is the structure $\mathcal{A}[V]$ with universe V and $R(\mathcal{A}[V]) = R(\mathcal{A}) \cap V^{\text{ar}(R)}$ for every relation symbol $R \in \sigma$. The *union* of two relational structures \mathcal{A} and \mathcal{B} is the structure $\mathcal{A} \cup \mathcal{B}$ with universe $U(\mathcal{A}) \cup U(\mathcal{B})$ and $R(\mathcal{A} \cup \mathcal{B}) = R(\mathcal{A}) \cup R(\mathcal{B})$ for all $R \in \sigma$. The *intersection* is defined analogously.

The *Gaifman graph* of a σ -structure \mathcal{A} is the graph $G_{\mathcal{A}}$ with vertex set $U(\mathcal{A})$ and an edge between two vertices $u, v \in U(\mathcal{A})$ if there is a relation symbol $R \in \sigma$ and a tuple $(u_1, \dots, u_{\text{ar}(R)}) \in R(\mathcal{A})$ with $u, v \in \{u_1, \dots, u_{\text{ar}(R)}\}$. The

degree of \mathcal{A} is the maximum degree of its Gaifman graph, i.e., the maximum number of neighbors of a vertex in $G_{\mathcal{A}}$.

The distance $\text{dist}^{\mathcal{A}}(u, v)$ between two vertices $u, v \in U(\mathcal{A})$ is the length of a shortest path between u and v in $G_{\mathcal{A}}$, and $\text{dist}^{\mathcal{A}}(u, v) = \infty$ if there is no path between u and v . For $r \in \mathbb{N}$, the r -neighborhood of a vertex $u \in U(\mathcal{A})$ is the set $N_r^{\mathcal{A}}(u) = \{v \in U(\mathcal{A}) \mid \text{dist}^{\mathcal{A}}(u, v) \leq r\}$ and the r -neighborhood of a tuple $\bar{u} = (u_1, \dots, u_k) \in U(\mathcal{A})^k$ is the set $N_r^{\mathcal{A}}(\bar{u}) = \bigcup_{i=1}^k N_r^{\mathcal{A}}(u_i)$.

We say that an algorithm has *local access* to a σ -structure \mathcal{A} if it may use queries such as “Is $(u_1, \dots, u_{\text{ar}(R)}) \in R(\mathcal{A})$?” and “Return a list of all neighbors of $u \in U(\mathcal{A})$ ”.

B. First-Order Logic with Counting

We assume that the reader is familiar with first-order logic. Let σ be a relational signature. Let P a countable set of predicate names, $\text{ar}: P \rightarrow \mathbb{N}_{\geq 1}$ an arity function, and $\llbracket P \rrbracket \subseteq \mathbb{Z}^{\text{ar}(P)}$ the semantics of the predicate name $P \in P$. Then we call the tuple $(P, \text{ar}, \llbracket \cdot \rrbracket)$ a numerical predicate collection.

Definition II.1 (FOCN(P) [3]). Let $(P, \text{ar}, \llbracket \cdot \rrbracket)$ be a numerical predicate collection, and let vars and nvars be disjoint sets of structure variables and number variables, respectively. The set of *formulas* for FOCN(P) is built according to the following rules.

- (F1) $x_1 = x_2$ and $R(x_1, \dots, x_{\text{ar}(R)})$ are formulas for $R \in \sigma$ and structure variables $x_1, \dots, x_{\text{ar}(R)} \in \text{vars}$.
- (F2) If φ and ψ are formulas, then $\neg\varphi$, $(\varphi \wedge \psi)$, and $(\varphi \vee \psi)$ are also formulas.
- (F3) $\exists x \varphi$ and $\forall x \varphi$ are formulas for $x \in \text{vars}$ and a formula φ .
- (F4) If $t_1, \dots, t_{\text{ar}(P)}$ are counting terms and $P \in P$, then $P(t_1, \dots, t_{\text{ar}(P)})$ is a formula.
- (F5) $\exists \kappa \varphi$ is a formula for every number variable $\kappa \in \text{nvars}$ and every formula φ .

The set of *counting terms* for FOCN(P) is built according to the following rules.

- (C1) $\# \bar{x}.\varphi$ is a counting term for $s \in \mathbb{N}$, $\bar{x} = (x_1, \dots, x_s) \in \text{vars}^s$ pairwise distinct structure variables and a formula φ .
- (C2) Every $i \in \mathbb{Z}$ is a counting term.
- (C3) If t_1 and t_2 are counting terms, then $(t_1 + t_2)$ and $(t_1 \cdot t_2)$ are also counting terms.
- (C4) Every $\kappa \in \text{nvars}$ is a counting term.

Let $\mathcal{I} = (\mathcal{A}, \beta)$ be an interpretation, where \mathcal{A} is a relational structure with universe $U(\mathcal{A})$, and $\beta: \text{vars} \cup \text{nvars} \rightarrow U(\mathcal{A}) \cup \mathbb{Z}$ with $\beta(x) \in U(\mathcal{A})$ for $x \in \text{vars}$ and $\beta(\kappa) \in \mathbb{Z}$ for $\kappa \in \text{nvars}$.

For $k, \ell \in \mathbb{N}$, $x_1, \dots, x_k \in \text{vars}$, $\kappa_1, \dots, \kappa_\ell \in \text{nvars}$, $a_1, \dots, a_k \in U(\mathcal{A})$, and $\lambda_1, \dots, \lambda_\ell \in \mathbb{Z}$, let $\mathcal{I} \frac{a_1, \dots, a_k}{x_1, \dots, x_k} \frac{\lambda_1, \dots, \lambda_\ell}{\kappa_1, \dots, \kappa_\ell} := (\mathcal{A}, \beta')$ with $\beta'(x_i) := a_i$ for all $i \in [k]$, $\beta'(\kappa_i) := \lambda_i$ for all $i \in [\ell]$, and $\beta'(z) := \beta(z)$ for all $z \in (\text{vars} \cup \text{nvars}) \setminus \{x_1, \dots, x_k, \kappa_1, \dots, \kappa_\ell\}$. Then the semantics for a formula or a counting term is defined as follows.

- (F1) $\llbracket x_1 = x_2 \rrbracket^{\mathcal{I}} := 1$ if $\beta(x_1) = \beta(x_2)$, and $\llbracket x_1 = x_2 \rrbracket^{\mathcal{I}} := 0$ otherwise.

- $\llbracket R(x_1, \dots, x_{\text{ar}(R)}) \rrbracket^{\mathcal{I}} := 1$ if $(\beta(x_1), \dots, \beta(x_{\text{ar}(R)})) \in R^{\mathcal{A}}$, and $\llbracket R(x_1, \dots, x_{\text{ar}(R)}) \rrbracket^{\mathcal{I}} := 0$ otherwise.
- (F2) $\llbracket \neg\varphi \rrbracket^{\mathcal{I}} := 1 - \llbracket \varphi \rrbracket^{\mathcal{I}}$, $\llbracket \varphi \wedge \psi \rrbracket^{\mathcal{I}} := \min \{ \llbracket \varphi \rrbracket^{\mathcal{I}}, \llbracket \psi \rrbracket^{\mathcal{I}} \}$ and $\llbracket \varphi \vee \psi \rrbracket^{\mathcal{I}} := \max \{ \llbracket \varphi \rrbracket^{\mathcal{I}}, \llbracket \psi \rrbracket^{\mathcal{I}} \}$.
- (F3) $\llbracket \exists x \varphi \rrbracket^{\mathcal{I}} := \max \{ \llbracket \varphi \rrbracket^{\mathcal{I} \frac{a}{x}} \mid a \in U(\mathcal{A}) \}$ and $\llbracket \forall x \varphi \rrbracket^{\mathcal{I}} := \min \{ \llbracket \varphi \rrbracket^{\mathcal{I} \frac{a}{x}} \mid a \in U(\mathcal{A}) \}$.
- (F4) $\llbracket P(t_1, \dots, t_{\text{ar}(P)}) \rrbracket^{\mathcal{I}} := 1$ if $(\llbracket t_1 \rrbracket^{\mathcal{I}}, \dots, \llbracket t_{\text{ar}(P)} \rrbracket^{\mathcal{I}}) \in \llbracket P \rrbracket$ and otherwise $\llbracket P(t_1, \dots, t_{\text{ar}(P)}) \rrbracket^{\mathcal{I}} := 0$.
- (F5) $\llbracket \exists \kappa \varphi \rrbracket^{\mathcal{I}} := \max \{ \llbracket \varphi \rrbracket^{\mathcal{I} \frac{k}{\kappa}} \mid k \in \{0, \dots, |U(\mathcal{A})|\} \}$.
- (C1) $\llbracket \# \bar{x}.\varphi \rrbracket^{\mathcal{I}} := \left| \{ \bar{a} \in (U(\mathcal{A}))^s \mid \llbracket \varphi \rrbracket^{\mathcal{I} \frac{\bar{a}}{x_1, \dots, x_s}} = 1 \} \right|$.
- (C2) $\llbracket i \rrbracket^{\mathcal{I}} := i$.
- (C3) $\llbracket (t_1 + t_2) \rrbracket^{\mathcal{I}} := \llbracket t_1 \rrbracket^{\mathcal{I}} + \llbracket t_2 \rrbracket^{\mathcal{I}}$ and $\llbracket (t_1 \cdot t_2) \rrbracket^{\mathcal{I}} := \llbracket t_1 \rrbracket^{\mathcal{I}} \cdot \llbracket t_2 \rrbracket^{\mathcal{I}}$.
- (C4) $\llbracket \kappa \rrbracket^{\mathcal{I}} := \beta(\kappa)$. J

Let φ be an FOCN(P)-formula. The *binding rank* $\text{br}(\varphi)$ of φ is the maximal nesting depth of constructs of the form $\exists x$ and $\forall x$ with $x \in \text{vars}$ and $\# \bar{x}$, where \bar{x} is a tuple in vars . The *binding width* $\text{bw}(\varphi)$ of φ is the maximal arity of an \bar{x} of a term $\# \bar{x}.\psi$ in φ . If φ contains no such term, then $\text{bw}(\varphi) = 1$ if φ contains a quantifier $\exists x$ or $\forall x$ with $x \in \text{vars}$, and $\text{bw}(\varphi) = 0$ otherwise.

C. Types, spheres and sphere formulas

Let $r \geq 0$ and $n \geq 1$, and let c_1, c_2, \dots be a sequence of pairwise distinct constant symbols. An r -type (with n centers) is a structure $\tau = (\mathcal{A}, a_1, \dots, a_n)$ over the signature $\sigma \uplus \{c_1, \dots, c_n\}$, where \mathcal{A} is a σ -structure with $a_1, \dots, a_n \in U(\mathcal{A})$ and $U(\mathcal{A}) = N_r^{\mathcal{A}}(a_1, \dots, a_n)$. The elements a_1, \dots, a_n are the *centers* of τ . For every tuple $\bar{a} \in (U(\mathcal{A}))^n$, the r -sphere of \bar{a} in \mathcal{A} is the r -type with n centers

$$\mathcal{N}_r^{\mathcal{A}}(\bar{a}) = (\mathcal{A}[N_r^{\mathcal{A}}(\bar{a})], \bar{a}).$$

Let τ be an r -type. A first-order formula $\text{sph}_\tau(\bar{x})$ is a *sphere-formula* if for every σ -structure \mathcal{B} and every tuple $\bar{b} \in U(\mathcal{B})^n$ we have

$$\mathcal{B} \models \text{sph}_\tau(\bar{b}) \iff \mathcal{N}_r^{\mathcal{B}}(\bar{b}) \cong \tau.$$

The *locality radius* of $\text{sph}_\tau(\bar{x})$ is r .

D. Numerical oc-type conditions and hnf-formulas

A *basic counting term* is a counting term of the form $\#(x).\text{sph}_\tau(x)$, where x is a structure variable in vars and τ is an r -type with a single center. The radius r is called the *locality radius* of the basic counting term. In a σ -structure \mathcal{A} , a basic counting term evaluates to the number of elements $a \in U(\mathcal{A})$ with $\mathcal{N}_r^{\mathcal{A}}(a) \cong \tau$, i.e., the number of r -neighborhoods with one center in \mathcal{A} that are isomorphic to τ .

A *numerical condition on occurrences of types with one center* (or *numerical oc-type condition*) is an FOCN(P)-formula that is built from basic counting terms and rules (F2), (F4), (F5), (C2), (C3), and (C4), i.e., using number variables and integers, and combining them by addition, multiplication, numerical predicates from $P \cup \{P_\exists\}$, Boolean combinations, and quantification of number variables. Its locality radius is the maximal locality radius of the involved basic counting terms.

Numerical oc-type conditions do not have any free structure variables.

A formula is in *Hanf normal form for FOCN(P)* or an *hnf-formula for FOCN(P)* if it is a Boolean combination of numerical oc-type conditions and sphere formulas. The locality radius of an hnf-formula is the maximal locality radius of the involved conditions and formulas.

E. Local hnf-formulas

Let \mathcal{A} be a relational structure over σ , $\bar{u} \in U(\mathcal{A})^k$ for some $k \in \mathbb{N}$ and $r \in \mathbb{N}$. Then the *local hnf-formulas (for FOCN(P)) of \bar{u} with locality radius smaller than r in \mathcal{A}* are

$$\text{lhfr}_r(\mathcal{A}, \bar{u}) = \{\varphi(\bar{x}) \text{ hnf-formula} \mid \mathcal{A} \models \varphi(\bar{u}), \text{ locality radius of } \varphi \text{ is smaller than } r\}.$$

The following results help us to show that reduced formula and parameter spaces suffice to find consistent hypotheses. The first lemma states that two tuples satisfy the same local hnf-formulas if their neighborhoods are isomorphic.

Lemma II.2. *Let \mathcal{A} be a structure over a relational signature σ , $k, r \in \mathbb{N}$ and $\bar{u}, \bar{u}' \in U(\mathcal{A})^k$. If $\mathcal{N}_r^{\mathcal{A}}(\bar{u}) \cong \mathcal{N}_r^{\mathcal{A}}(\bar{u}')$, then $\text{lhfr}_{r+1}(\mathcal{A}, \bar{u}) = \text{lhfr}_{r+1}(\mathcal{A}, \bar{u}')$.*

Proof: Let $\varphi(\bar{x})$ be an hnf-formula with locality radius at most r . Then φ is a Boolean combination of numerical oc-type conditions and sphere formulas with locality radius at most r . The numerical oc-type conditions do not have any free structure variables and are thus independent from the assignment for \bar{x} . The free variables of the sphere formulas are a subset of $\text{free}(\varphi)$. Let $\text{sph}_\tau(x_{i_1}, \dots, x_{i_\ell})$ be a sphere-formula used in φ with $x_{i_1}, \dots, x_{i_\ell} \in \text{free}(\varphi)$ and τ an r' -type with ℓ centers for $r' \leq r$. Then $\mathcal{N}_r^{\mathcal{A}}(u_{i_1}, \dots, u_{i_\ell}) \cong \mathcal{N}_r^{\mathcal{A}}(u'_{i_1}, \dots, u'_{i_\ell})$ and $\mathcal{N}_{r'}^{\mathcal{A}}(u_{i_1}, \dots, u_{i_\ell}) \cong \mathcal{N}_{r'}^{\mathcal{A}}(u'_{i_1}, \dots, u'_{i_\ell})$ and hence

$$\begin{aligned} \mathcal{A} \models \text{sph}_\tau(u_{i_1}, \dots, u_{i_\ell}) \\ \iff \mathcal{N}_{r'}^{\mathcal{A}}(\bar{u}) \models \text{sph}_\tau(u_{i_1}, \dots, u_{i_\ell}) \\ \iff \mathcal{N}_{r'}^{\mathcal{A}}(\bar{u}') \models \text{sph}_\tau(u'_{i_1}, \dots, u'_{i_\ell}) \\ \iff \mathcal{A} \models \text{sph}_\tau(u'_{i_1}, \dots, u'_{i_\ell}). \end{aligned}$$

This holds for all sphere formulas in φ . Thus $\mathcal{A} \models \varphi(\bar{u})$ if and only if $\mathcal{A} \models \varphi(\bar{u}')$. \square

The following result is a local variant of the Feferman-Vaught Theorem [17] translated to our context. It allows us to analyze the parameters we choose by splitting them into two parts with disjoint neighborhoods.

Lemma II.3 (Local Composition Lemma). *Let $\mathcal{A}, \mathcal{A}'$ be structures over a relational signature σ , $\bar{u} \in U(\mathcal{A})^k$, $\bar{v} \in U(\mathcal{A})^\ell$, $\bar{u}' \in U(\mathcal{A}')^k$, $\bar{v}' \in U(\mathcal{A}')^\ell$, and $r \in \mathbb{N}$, such that $\mathcal{N}_r^{\mathcal{A}}(\bar{u}) \cap \mathcal{N}_r^{\mathcal{A}}(\bar{v}) = \mathcal{N}_r^{\mathcal{A}'}(\bar{u}') \cap \mathcal{N}_r^{\mathcal{A}'}(\bar{v}') = \emptyset$, $\text{lhfr}_{r+1}(\mathcal{A}, \bar{u}) = \text{lhfr}_{r+1}(\mathcal{A}', \bar{u}')$ and $\text{lhfr}_{r+1}(\mathcal{A}, \bar{v}) = \text{lhfr}_{r+1}(\mathcal{A}', \bar{v}')$. Then*

$$\text{lhfr}_{r+1}(\mathcal{A}, \bar{u}\bar{v}) = \text{lhfr}_{r+1}(\mathcal{A}', \bar{u}'\bar{v}').$$

Proof: All hnf-sentences φ with locality radius at most r that satisfy $\mathcal{A} \models \varphi$ are contained in $\text{lhfr}_{r+1}(\mathcal{A}, \bar{u})$ and thus

also in $\text{lhfr}_{r+1}(\mathcal{A}', \bar{u}')$. Hence \mathcal{A} and \mathcal{A}' model the same hnf-sentences with locality radius at most r .

We know that the sphere-formula $\text{sph}_{\mathcal{N}_r^{\mathcal{A}}(\bar{u})}(\bar{x})$ is contained in $\text{lhfr}_{r+1}(\mathcal{A}, \bar{u}) = \text{lhfr}_{r+1}(\mathcal{A}', \bar{u}')$, so $\mathcal{N}_r^{\mathcal{A}}(\bar{u}) \cong \mathcal{N}_r^{\mathcal{A}'}(\bar{u}')$. Analogously it follows that $\mathcal{N}_r^{\mathcal{A}}(\bar{v}) \cong \mathcal{N}_r^{\mathcal{A}'}(\bar{v}')$. Since $\mathcal{N}_r^{\mathcal{A}}(\bar{u}) \cap \mathcal{N}_r^{\mathcal{A}}(\bar{v}) = \emptyset$ and $\mathcal{N}_r^{\mathcal{A}'}(\bar{u}') \cap \mathcal{N}_r^{\mathcal{A}'}(\bar{v}') = \emptyset$, we also obtain $\mathcal{N}_r^{\mathcal{A}}(\bar{u}\bar{v}) \cong \mathcal{N}_r^{\mathcal{A}'}(\bar{u}'\bar{v}')$. Thus for all sphere formulas $\psi(x_1, \dots, x_m)$ with locality radius at most r , we know that $\mathcal{A} \models \psi(w_1, \dots, w_m)$ for $w_1, \dots, w_m \in \bar{u} \cup \bar{v}$ if and only if $\mathcal{A}' \models \psi(w'_1, \dots, w'_m)$.

Let $\varphi \in \text{lhfr}_{r+1}(\mathcal{A}, \bar{u}\bar{v})$. Then φ is a Boolean combination of hnf-sentences and sphere formulas. These hnf-sentences and sphere formulas hold for $(\mathcal{A}, \bar{u}\bar{v})$ if and only if they hold for $(\mathcal{A}', \bar{u}'\bar{v}')$. Thus $\varphi \in \text{lhfr}_{r+1}(\mathcal{A}', \bar{u}'\bar{v}')$ and $\text{lhfr}_{r+1}(\mathcal{A}, \bar{u}\bar{v}) \subseteq \text{lhfr}_{r+1}(\mathcal{A}', \bar{u}'\bar{v}')$. Analogously we can show that $\text{lhfr}_{r+1}(\mathcal{A}', \bar{u}'\bar{v}') \subseteq \text{lhfr}_{r+1}(\mathcal{A}, \bar{u}\bar{v})$. \square

For $d \in \mathbb{N}$, two formulas φ, φ' are called *d-equivalent* if $\llbracket \varphi \rrbracket^{\mathcal{I}} = \llbracket \varphi' \rrbracket^{\mathcal{I}}$ for all interpretations $\mathcal{I} = (\mathcal{A}, \beta)$ for all structures \mathcal{A} of degree at most d . The following result is due to Kuske and Schweikardt [3].

Theorem II.4. *Let $(P, \text{ar}, \llbracket \cdot \rrbracket)$ be a numerical predicate collection. For any relational signature σ , any degree bound $d \in \mathbb{N}$, and any FOCN(P)[σ]-formula φ , there exists a d-equivalent hnf-formula ψ for FOCN(P)[σ] of locality radius smaller than $(2 \text{bw}(\varphi) + 1)^{\text{br}(\varphi)}$ with $\text{free}(\psi) = \text{free}(\varphi)$.*

Using this result, we can show that a formula behaves the same for two different assignments of structure variables if the local hnf-formulas of the assigned tuples are the same.

Lemma II.5. *Let \mathcal{A} be a structure over a relational signature σ , $\varphi(\bar{x}, \bar{y})$ an FOCN(P)-formula, $\bar{u}, \bar{u}' \in U(\mathcal{A})^{|\bar{x}|}$ and $\bar{\lambda} \in \{0, \dots, |\mathcal{A}|\}^{|\bar{y}|}$. If*

$$\text{lhfr}_{(2 \text{bw}(\varphi)+1)^{\text{br}(\varphi)}}(\mathcal{A}, \bar{u}) = \text{lhfr}_{(2 \text{bw}(\varphi)+1)^{\text{br}(\varphi)}}(\mathcal{A}, \bar{u}'),$$

then

$$\mathcal{A} \models \varphi(\bar{u}, \bar{\lambda}) \iff \mathcal{A} \models \varphi(\bar{u}', \bar{\lambda}).$$

Proof: Let $\varphi'(\bar{x}) := \varphi(\bar{x}, \bar{\lambda})$, $r := \text{br}(\varphi) = \text{br}(\varphi')$ and $w := \text{bw}(\varphi) = \text{bw}(\varphi')$. Using Theorem II.4 we obtain an hnf-formula ψ with locality radius smaller than $(2w+1)^r$ that is $\Delta\mathcal{A}$ -equivalent to φ' and, just like φ' , doesn't have any free number variables. Then

$$\begin{aligned} \mathcal{A} \models \varphi(\bar{u}, \bar{\lambda}) &\iff \mathcal{A} \models \psi(\bar{u}) \\ &\iff \psi \in \text{lhfr}_{(2w+1)^r}(\mathcal{A}, \bar{u}) \\ &\iff \psi \in \text{lhfr}_{(2w+1)^r}(\mathcal{A}, \bar{u}') \\ &\iff \mathcal{A} \models \psi(\bar{u}') \\ &\iff \mathcal{A} \models \varphi(\bar{u}', \bar{\lambda}). \end{aligned}$$

\square

III. LEARNING FOCN(P)-DEFINABLE CONCEPTS OVER STRUCTURES OF POLYLOGARITHMIC DEGREE

In this section, we prove Theorem I.1 and a variant of it, as well as negative results, including Theorem I.3. Let $k, \ell, t, r, w \in \mathbb{N}$ be fixed, $\bar{x} := (x_1, \dots, x_k)$ instance

variables, $\bar{y} := (y_1, \dots, y_\ell)$ parameter variables, \mathcal{B} a background structure over a relational signature σ and $\mathcal{T} := (U(\mathcal{B})^k \times \{0, 1\})^t$ training sequences. For $s \in \mathbb{N}$ and $T \in \mathcal{T}$, $T = ((\bar{u}_1, c_1), \dots, (\bar{u}_t, c_t))$, let $N_s^{\mathcal{B}}(T) := \bigcup_{i=1}^t N_s^{\mathcal{B}}(\bar{u}_i)$. Let Φ be the set of FOCN(P)-formulas $\varphi(\bar{x}; \bar{y}, \bar{\kappa})$ with binding width at most w , binding rank at most r and free number variables $\bar{\kappa}$, and

$$\mathcal{C} := \{ \llbracket \varphi(\bar{x}; \bar{v}, \bar{\lambda}) \rrbracket^{\mathcal{B}} \mid \varphi(\bar{x}; \bar{y}, \bar{\kappa}) \in \Phi, \bar{v} \in U(\mathcal{B})^\ell, \bar{\lambda} \in \{0, \dots, |\mathcal{B}|^{|\bar{\kappa}|}\} \}.$$

To prove Theorem I.1, we give an algorithm that uses brute force. Given a training sequence that is consistent with some hypothesis in \mathcal{C} , the algorithm finds a consistent hypothesis consisting of a first-order formula and a parameter tuple $\bar{v} \in U(\mathcal{B})^\ell$. The following lemma states that it suffices to search in a reduced parameter space and to check a single formula per parameter. This enables us to find a consistent hypothesis with a sufficient running time.

Lemma III.1. *Let $T = ((\bar{u}_1, c_1), \dots, (\bar{u}_t, c_t)) \in \mathcal{T}$ be consistent with some $C \in \mathcal{C}$. Then there is a tuple $\bar{v}^* = (v_1, \dots, v_\ell) \in N_{2\ell[(2w+1)^{r-1}]}^{\mathcal{B}}(T)^\ell$ and some $m \leq \ell$ such that $\llbracket \varphi^*(\bar{x}, \bar{v}^*) \rrbracket^{\mathcal{B}}$ is consistent with T for*

$$\varphi^*(\bar{x}; \bar{y}) := \bigvee_{i \in [t], c_i = 1} \vartheta_i(\bar{x}; \bar{y}^\circ),$$

$\vartheta_i(\bar{x}; \bar{y}^\circ) := \text{sph}_{\mathcal{N}_{(2w+1)^{r-1}}^{\mathcal{B}}(\bar{u}_i \bar{v}^\circ)}(\bar{x} \bar{y}^\circ)$, $\bar{v}^\circ := (v_1, \dots, v_m)$, and $\bar{y}^\circ := (y_1, \dots, y_m)$.

Proof: Let $\varphi(\bar{x}; \bar{y}, \bar{\kappa}) \in \Phi$ and $\bar{v} = (v_1, \dots, v_\ell) \in U(\mathcal{B})^\ell$, $\bar{\lambda} \in \{0, \dots, |\mathcal{B}|^{|\bar{\kappa}|}\}$ such that $C = \llbracket \varphi(\bar{x}; \bar{v}, \bar{\lambda}) \rrbracket^{\mathcal{B}}$ is consistent with T . Now define for some $m \leq \ell$ the elements $v^{(1)}, \dots, v^{(m)} \in \{v_1, \dots, v_\ell\}$ and the sets $N^{(0)}, \dots, N^{(m)} \subseteq U(\mathcal{B})$ by setting $N^{(0)} := N_{(2w+1)^{r-1}}^{\mathcal{B}}(T)$ and defining the rest as follows: Given $N^{(i)}$, if there is some $v \in \{v_1, \dots, v_\ell\} \setminus \{v^{(1)}, \dots, v^{(i)}\}$ such that $\text{dist}^{\mathcal{B}}(v, N^{(i)}) \leq (2w+1)^r$, then let $v^{(i+1)} := v$. Choose arbitrarily if there is more than one. Let $N^{(i+1)} := N^{(i)} \cup N_{(2w+1)^{r-1}}^{\mathcal{B}}(v^{(i+1)})$. If there is no such v , then set $m := i$ and stop.

Let $N^\circ := N^{(m)}$ and w.l.o.g. let $v^{(i)} = v_i$ for $i \in [m]$. Let $\bar{v}^\circ := (v_1, \dots, v_m)$ and $\bar{v}^\bullet := (v_{m+1}, \dots, v_\ell)$. Then $\bar{v}^\circ \in N_{2\ell[(2w+1)^{r-1}]}^{\mathcal{B}}(T)^m$ and

$$N^\circ = \bigcup_{i=1}^t N_{(2w+1)^{r-1}}^{\mathcal{B}}(\bar{u}_i) \cup \bigcup_{i=1}^m N_{(2w+1)^{r-1}}^{\mathcal{B}}(v_i) \quad (1)$$

and

$$N_{(2w+1)^{r-1}}^{\mathcal{B}}(\bar{v}^\bullet) \cap N^\circ = \emptyset. \quad (2)$$

Claim 1. *Let $i, j \in [t]$ such that $\mathcal{B} \models \vartheta_i(\bar{u}_j; \bar{v}^\circ)$. Then $c_i = c_j$.*

Proof: From $\mathcal{B} \models \vartheta_i(\bar{u}_j; \bar{v}^\circ)$ it follows that $\mathcal{N}_{(2w+1)^{r-1}}^{\mathcal{B}}(\bar{u}_i \bar{v}^\circ) \cong \mathcal{N}_{(2w+1)^{r-1}}^{\mathcal{B}}(\bar{u}_j \bar{v}^\circ)$. Using Lemma II.2 we obtain $\text{lh}_{(2w+1)^r}(\mathcal{B}, \bar{u}_i) = \text{lh}_{(2w+1)^r}(\mathcal{B}, \bar{u}_j)$. Furthermore, from Equation (1) and Equation (2) it follows that $N_{(2w+1)^{r-1}}^{\mathcal{B}}(\bar{u}_i \bar{v}^\circ) \cap N_{(2w+1)^{r-1}}^{\mathcal{B}}(\bar{v}^\bullet) = \emptyset$ and

Algorithm \mathcal{L}_{con}

Input: Training sequence $T = ((\bar{u}_1, c_1), \dots, (\bar{u}_t, c_t)) \in \mathcal{T}$,
 $d = \Delta \mathcal{B}$, local access to background structure \mathcal{B}

- 1: $N \leftarrow N_{2\ell[(2w+1)^{r-1}]}^{\mathcal{B}}(T)$
- 2: **for all** $\bar{v}^* \in N^m$, $m \leq \ell$ **do**
- 3: **for** $i = 1, \dots, t$ **do**
- 4: $\mathcal{N}_i \leftarrow \mathcal{N}_{(2w+1)^{r-1}}^{\mathcal{B}}(\bar{u}_i \bar{v}^*)$
- 5: $\vartheta_i(\bar{x}; \bar{y}) \leftarrow \text{sph}_{\mathcal{N}_i}(\bar{x} \bar{y})$ \triangleright a d -bounded
- 6: $[(2w+1)^r - 1]$ -type with $k + m$ centers
- 7: $\varphi^*(\bar{x}; \bar{y}) \leftarrow \bigvee_{i \in [t], c_i = 1} \vartheta_i(\bar{x}; \bar{y})$
- 8: $\text{consistent} \leftarrow \text{true}$
- 9: **for** $i \in [t]$ with $c_i = 0$ **do**
- 10: **for** $j \in [t]$ with $c_j = 1$ **do**
- 11: **if** $\mathcal{N}_i \cong \mathcal{N}_j$ **then**
- 12: $\text{consistent} \leftarrow \text{false}$
- 13: **if consistent then**
- 14: **return** $(\varphi^*(\bar{x}; \bar{y}), \bar{v}^*)$
- 15: **reject**

Figure 4. Learning algorithm \mathcal{L}_{con} of Theorem I.1

$N_{(2w+1)^{r-1}}^{\mathcal{B}}(\bar{u}_j \bar{v}^\circ) \cap N_{(2w+1)^{r-1}}^{\mathcal{B}}(\bar{v}^\bullet) = \emptyset$. Thus with Lemma II.3 we obtain

$$\text{lh}_{(2w+1)^r}(\mathcal{B}, \bar{u}_i \bar{v}) = \text{lh}_{(2w+1)^r}(\mathcal{B}, \bar{u}_j \bar{v}).$$

and hence with Lemma II.5 we get

$$\mathcal{B} \models \varphi(\bar{u}_i; \bar{v}, \bar{\lambda}) \iff \mathcal{B} \models \varphi(\bar{u}_j; \bar{v}, \bar{\lambda}).$$

This implies $c_i = c_j$. ┐

Set

$$\bar{v}^* := (v_1, \dots, v_m, \underbrace{v, \dots, v}_{\ell-m \text{ times}})$$

for some arbitrary $v \in N_{2\ell[(2w+1)^{r-1}]}^{\mathcal{B}}(T)$. Then $\bar{v}^* \in N_{2\ell[(2w+1)^{r-1}]}^{\mathcal{B}}(T)^\ell$. It remains to show that $\llbracket \varphi^*(\bar{x}; \bar{v}^*) \rrbracket^{\mathcal{B}}$ is consistent with T . If $\mathcal{B} \models \varphi^*(\bar{u}_i; \bar{v}^*)$, then there is some $p \in [t]$ with $c_p = 1$ and $\mathcal{B} \models \vartheta_p(\bar{u}_i; \bar{v}^*)$. Using the claim it follows that $c_i = c_p = 1$. On the other hand, if $c_i = 1$, then $\mathcal{B} \models \varphi^*(\bar{u}_i; \bar{v}^*)$. Thus $\llbracket \varphi^*(\bar{x}; \bar{v}^*) \rrbracket^{\mathcal{B}}$ is consistent with T . \square

In our algorithm, we have to compare isomorphism types of local neighborhoods. To do this within the desired time bounds, we apply the following result due to Grohe, Neuen and Schweitzer [5].

Theorem III.2. *Let \mathcal{A}_1 and \mathcal{A}_2 be two σ -structures with $n := \max\{|\mathcal{A}_1|, |\mathcal{A}_2|\}$, $d := \max\{\Delta \mathcal{A}_1, \Delta \mathcal{A}_2\}$ and $m := \max_{R \in \sigma} \text{ar}(R)$. One can check whether $\mathcal{A}_1 \cong \mathcal{A}_2$ in time $n^{\mathcal{O}(m \cdot (\log d)^c)}$ for some constant c .*

A. Consistent Hypotheses

We now prove the main theorem.

Proof of Theorem I.1: The pseudocode for our algorithm is shown in Figure 4. Because of Lemma III.1 the algorithm satisfies Condition (1). Let $(\varphi^*(\bar{x}; \bar{y}), \bar{v}^*)$ be the hypothesis returned by the algorithm. Note that $\mathcal{N}_{(2w+1)^{r-1}}^{\mathcal{B}}(\bar{u} \bar{v}^*) \models \varphi^*(\bar{u}, \bar{v}^*)$ iff $\mathcal{B} \models \varphi^*(\bar{u}, \bar{v}^*)$ for all $(\bar{u}, c) \in T$, because φ^* is a disjunction

of sphere formulas of locality radius smaller than $(2w+1)^r$. For $(\bar{u}_i, 1) \in T$ we know by construction that $\mathcal{N}_i \models \varphi^*(\bar{u}_i; \bar{v}^*)$. For all $(\bar{u}_i, 0) \in T$ the algorithm checks that there is no $(\bar{u}_j, 1) \in T$ with $\mathcal{N}_i \cong \mathcal{N}_j$ and thus $\mathcal{N}_i \not\models \varphi^*(\bar{u}_i; \bar{v}^*)$. Hence the hypothesis is consistent with the input sequence and the algorithm satisfies Condition (2) of the theorem.

Algorithm \mathcal{L}_{con} computes $\mathcal{N}_i = \mathcal{N}_{(2w+1)^{r-1}}^{\mathcal{B}}(\bar{u}_i \bar{v}^*)$ for all $i \in [t]$. $|\mathcal{N}_i| \leq (k + \ell) \cdot d^{(2w+1)^r}$ and the representation size is in $(\log n + d)^{\mathcal{O}(1)}$, because k, ℓ, r and w are constant. Each sphere-formula can be computed in time $(\log n \cdot (k + \ell) \cdot d^{(2w+1)^r})^{\mathcal{O}(\|\sigma\|)} = (\log n + d)^{\mathcal{O}(1)}$. ([3]) Thus φ^* can be computed in time $(\log n + d + t)^{\mathcal{O}(1)}$.

When checking whether φ^* is consistent with T , we check whether \mathcal{N}_i and \mathcal{N}_j are isomorphic. Using Theorem III.2 a single isomorphism test takes time

$$\begin{aligned} & (\max\{|\mathcal{N}_i|, |\mathcal{N}_j|\})^{\mathcal{O}(\max_{R \in \sigma} \text{ar}(R) \cdot (\log \max\{\Delta \mathcal{N}_i, \Delta \mathcal{N}_j\})^c)} \\ & \leq \left((k + \ell) \cdot d^{(2w+1)^r} \right)^{\mathcal{O}((\log d)^c)} \leq d^{\mathcal{O}((\log d)^c)} \end{aligned}$$

Before checking whether two substructures are isomorphic, we rename the vertices in the substructure such that the representation size no longer depends on $\log n$, but on the size of the substructures. This can be done in time $(\log n)^{\mathcal{O}(1)}$. Hence the consistency of the formula can be checked in time $(\log n + t)^{\mathcal{O}(1)} \cdot d^{\mathcal{O}((\log d)^c)}$.

The maximum number of iterations in the outer loop is $\sum_{m=1}^{\ell} |N|^m \leq \ell \cdot (2tkd^{2\ell(2w+1)^r})^{\ell} \in (t + d)^{\mathcal{O}(1)}$ and N can be computed with only local access to \mathcal{B} in time $(t + d)^{\mathcal{O}(1)}$. All in all the running time of the algorithm is in $(\log n + t)^{\mathcal{O}(1)} \cdot d^{\mathcal{O}((\log d)^c)}$. This shows that the algorithm satisfies Condition (3).

To evaluate the formula returned by the algorithm for \bar{u} , we only have to evaluate it using the structure $\mathcal{N}_{(2w+1)^{r-1}}^{\mathcal{B}}(\bar{u})$ with only local access to \mathcal{B} , so the running time is in $(\log n + t)^{\mathcal{O}(1)} \cdot d^{\mathcal{O}((\log d)^c)}$ and Condition (4) is satisfied. \square

If we consider k, ℓ, r and w as part of the input, then a more thorough analysis shows that the running time of the algorithm is exponential in $\text{polylog}(d) \log(k) \ell^2 (2w+1)^r + \ell \cdot \log(t)$ and polylogarithmic in n .

Theorem I.3 stated that there is no consistent model-learning algorithm for first-order formulas on background structures with no degree restriction that runs in sublinear time with only local access. In the proof we use the fact that an algorithm is unable to see all vertices of the background structure in sublinear time.

Proof of Theorem I.3: Let $k, \ell = 1$ and consider undirected graphs G_{ij} of size n for $i \in \{1, 2\}$ and $1 \leq j \leq t$, where t is an even integer. Let G_{11} and G_{21} be graphs with no edges and G_{12} and G_{22} be graphs with a single edge. For $j \geq 3$ let G_{ij} be the empty graph iff $(j - i)$ is even and let G_{ij} be the graph with a single edge else.

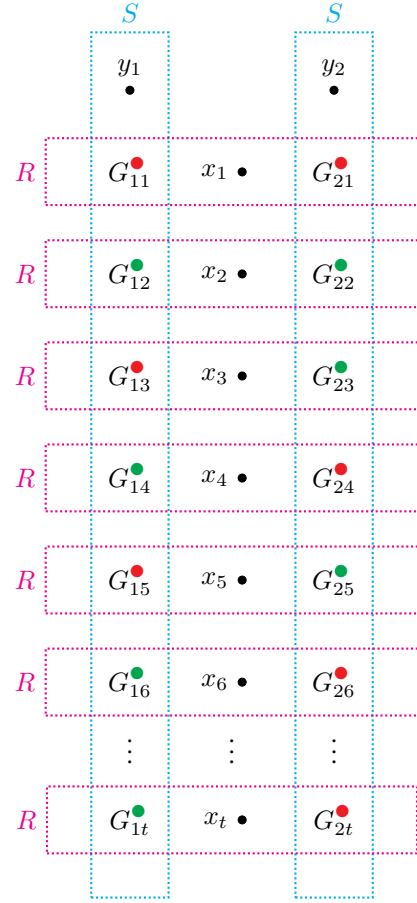


Figure 5. The background structure \mathcal{B} in the proof of Theorem I.3. The red dot marks graphs with no edges. The green dot marks graphs with a single edge.

Define the background structure \mathcal{B} for the relational signature $\sigma = \{E, R, S\}$ by

$$U(\mathcal{B}) := \{x_1, \dots, x_t, y_1, y_2\} \bigcup_{\substack{i \in \{1, 2\}, \\ 1 \leq j \leq t}} V(G_{ij}),$$

$$E^{\mathcal{B}} := \bigcup_{\substack{i \in \{1, 2\}, \\ 1 \leq j \leq t}} E(G_{ij}),$$

$$R^{\mathcal{B}} := \{(x_j, v) \mid 1 \leq j \leq t, v \in G_{ij} \text{ for some } i \in \{1, 2\}\},$$

and

$$S^{\mathcal{B}} := \{(y_i, v) \mid i \in \{1, 2\}, v \in G_{ij} \text{ for some } 1 \leq j \leq t\}.$$

The background structure can be seen in Figure 5. The size of the background structure is $|\mathcal{B}| = t \cdot (2|G_{11}| + 1) + 2$ and thus linear in the size of the graphs $|G_{ij}|$. Let

$$\varphi(x; y) := \exists v_1 \exists v_2 \ Rxv_1 \wedge Rxv_2 \wedge Syv_1 \wedge Syv_2 \wedge Ev_1v_2$$

and consider training examples $(x_1, c_{11}), \dots, (x_t, c_{it})$ with $c_{ij} = \llbracket \varphi(x_j; y_i) \rrbracket^{\mathcal{B}}$ for $i \in \{1, 2\}$ and $1 \leq j \leq t$. The formula

$\varphi(x_j; y_i)$ is evaluated to 1 if and only if G_{ij} contains an edge. Define two training sequences

$$\begin{aligned} T_1 &:= ((x_1, c_{11}), \dots, (x_t, c_{1t})) \\ &= ((x_1, 0), (x_2, 1), (x_3, 0), (x_4, 1), \\ &\quad (x_5, 0), (x_6, 1), \dots, (x_{t-1}, 0), (x_t, 1)) \text{ and} \\ T_2 &:= ((x_1, c_{21}), (x_2, c_{22}), (x_4, c_{24}), (x_3, c_{23}), \\ &\quad (x_6, c_{26}), (x_5, c_{25}), \dots, (x_t, c_{2t}), (x_{t-1}, c_{2(t-1)})) \\ &= ((x_1, 0), (x_2, 1), (x_4, 0), (x_3, 1), \\ &\quad (x_6, 0), (x_5, 1), \dots, (x_t, 0), (x_{t-1}, 1)). \end{aligned}$$

Let \mathcal{L} be a sublinear deterministic model-learning algorithm and run it on the training sequences T_1 and T_2 . Then there is some vertex order such that \mathcal{L} is unable to find even a single edge from E . In both sequences the algorithm receives the same sequence of Booleans and all visited vertices have a single R -edge, a single S -edge and no E -edge. Thus both training sequences are indistinguishable for the learner and hence it has to return the same formula $\psi(x; y) \in \text{FO}$ for both sequences T_1 and T_2 .

The learner \mathcal{L} can distinguish two different vertex sets $\{x_j\} \cup V(G_{1j}) \cup V(G_{2j})$ and $\{x_{j'}\} \cup V(G_{1j'}) \cup V(G_{2j'})$ only by the values of c_{ij} and $c_{ij'}$. Hence, by choosing a suitable vertex ordering, we can assume that the algorithm chooses the parameters only from the set $\{y_1, y_2, x_1, x_2\} \cup V(G_{11}) \cup V(G_{12}) \cup V(G_{21}) \cup V(G_{22})$. Since T_1 and T_2 are indistinguishable, it has to choose the same parameter for T_1 and T_2 , so $v_1 = v_2$.

But then \mathcal{L} can't be consistent with both T_1 and T_2 , because $c_{13} = 0 \neq 1 = c_{23}$ and both values would have to be equal to $\llbracket \psi(x_3, v_1) \rrbracket^{\mathcal{B}} = \llbracket \psi(x_3, v_2) \rrbracket^{\mathcal{B}}$. \square

Although we do not give a formal introduction into *probably approximately correct* (PAC) learning, we would like to mention that one can easily extend Theorem I.3 analogously to Grohe, Löding and Ritzert [6] to show that there is no sublinear model-learning algorithm for first-order formulas, that is a PAC-learning algorithm, on background structures with no degree restriction and only local access.

In addition to model learning, Grohe and Ritzert [2] also considered *parameter learning*, where we assume a fixed formula and we only want to find parameters such that the resulting hypothesis is consistent with the training examples.

Grohe and Ritzert showed that parameter learning is not possible in sublinear time with only local access. Here we prove a stronger result that shows that parameter learning is at least as hard as solving q -CLIQUE.

Theorem III.3. *If the exponential-time hypothesis (ETH) holds, then there is no consistent parameter-learning algorithm for first-order formulas φ of quantifier rank at most q on background structures \mathcal{B} with no degree restriction running in time $f(q) \cdot |\mathcal{B}|^{o(q)}$ for some function f , i.e. that, given φ and a sequence of training examples T , returns a tuple \bar{v} such that $\llbracket \varphi(\bar{x}; \bar{v}) \rrbracket^{\mathcal{B}}$ is consistent with all training examples.*

Proof: Let $q \in \mathbb{N}$. For the background structure let G be a graph of size $|G| \gg q$, G' a copy of G and H_+ and H_-

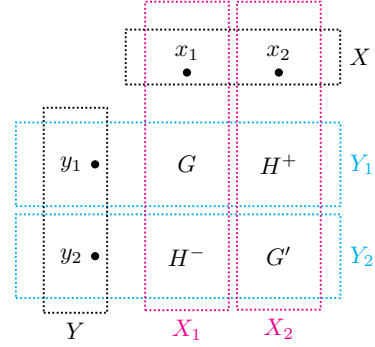


Figure 6. The background structure \mathcal{B} in the proof of Theorem III.3.

graphs of size $|G|$ such that H_+ has a q -clique and H_- does not, e.g. choose $H_+ = K_q \uplus \overline{K}_{|G|-q}$ and $H_- = \overline{K}_{|G|}$ where K_s and \overline{K}_s are the complete graph and the empty graph on s vertices for $s \in \mathbb{N}$.

Define the background structure \mathcal{B} for the relational signature $\sigma := \{E, X, Y, X_1, X_2, Y_1, Y_2\}$ by

$$\begin{aligned} U(\mathcal{B}) &:= \{x_1, x_2, y_1, y_2\} \uplus V(G) \uplus V(G') \\ &\quad \uplus V(H_+) \uplus V(H_-), \\ E(\mathcal{B}) &:= E(G) \uplus E(G') \uplus E(H_+) \uplus E(H_-), \\ X(\mathcal{B}) &:= \{x_1, x_2\}, \\ Y(\mathcal{B}) &:= \{y_1, y_2\}, \\ X_1(\mathcal{B}) &:= \{x_1\} \uplus V(G) \uplus V(H_+), \\ X_2(\mathcal{B}) &:= \{x_2\} \uplus V(G') \uplus V(H_-), \\ Y_1(\mathcal{B}) &:= \{y_1\} \uplus V(G) \uplus V(H_-), \text{ and} \\ Y_2(\mathcal{B}) &:= \{y_2\} \uplus V(G') \uplus V(H_+). \end{aligned}$$

The background structure can be seen in Figure 6. Let

$$\begin{aligned} \psi_{ij}(x; y) &:= X_i x \wedge Y_j y \wedge \exists v_1 \dots \exists v_q \\ &\quad \left(\bigwedge_{s=1}^q (X_i v_s \wedge Y_j v_s) \wedge \bigwedge_{1 \leq s_1 < s_2 \leq q} E v_{s_1} v_{s_2} \right) \end{aligned}$$

and

$$\varphi(x; y) := Xx \wedge Yy \wedge \bigvee_{i=1}^2 \bigvee_{j=1}^2 \psi_{ij}(x; y).$$

For a learned parameter v the hypothesis $\llbracket \varphi(x; v) \rrbracket^{\mathcal{B}}$ is only consistent with the training sequence $((x_1, \text{true}), (x_2, \text{false}))$, if $v \in \{y_1, y_2\}$. The parameter v is equal to y_1 if and only if G has a q -clique. Assuming ETH, there is no algorithm that checks whether a graph G has a q -clique in time $f(q) \cdot |G|^{o(q)}$ [18]. One can check whether G has a q -clique by computing \mathcal{B} in time quadratic in $|G|$ (or linear in $|V(G)| + |E(G)|$), learning the parameter v and then checking whether v is equal to y_1 . Thus learning the parameter is not possible in time $f(q) \cdot |G|^{o(q)} = f(q) \cdot |\mathcal{B}|^{o(q)}$. \square

B. Minimizing the training error

We continue the analysis of model learning. In Theorem I.1 we allow the algorithm to reject a training sequence T if

there is no consistent hypothesis within the given binding rank and binding width bounds. Instead of requiring a consistent hypothesis, we now try to find a hypothesis $H: U(\mathcal{B})^k \rightarrow \{0, 1\}$ that minimizes the number of errors, i.e., that minimizes the *training error*

$$\text{err}_T(H) = \frac{1}{|T|} |\{(\bar{u}, c) \in T \mid H(\bar{u}) \neq c\}|.$$

In the next theorem we generalize the results of Theorem I.1.

Theorem III.4. *Let $k, \ell, r, w \in \mathbb{N}$. Then there is a learning algorithm \mathcal{L}_{\min} for the k -ary learning problem over some finite relational structure \mathcal{B} , that receives a training sequence T and the degree of the structure $\Delta\mathcal{B}$ as input, with the following properties:*

- (1) *The algorithm always returns a hypothesis H of the form $(\varphi^*(\bar{x}; \bar{y}), \bar{v}^*)$ for some first-order formula $\varphi^*(\bar{x}; \bar{y})$ that is a Boolean combination of sphere formulas with locality radius smaller than $(2w+1)^r$ and $\bar{v}^* \in U(\mathcal{B})^\ell$.*
- (2) *If there are an FOCN(P)-formula $\varphi(\bar{x}; \bar{y}, \bar{\kappa})$ of binding rank at most r and binding width at most w and parameter tuples $\bar{v} \in U(\mathcal{B})^\ell$ and $\bar{\lambda} \in \{0, \dots, |\mathcal{B}|\}^{|\bar{\kappa}|}$ such that $\text{err}_T(\llbracket \varphi(\bar{x}; \bar{v}, \bar{\lambda}) \rrbracket^{\mathcal{B}}) \leq \varepsilon$ then \mathcal{L}_{\min} returns a hypothesis $(\varphi^*(\bar{x}; \bar{y}), \bar{v}^*)$ with $\text{err}_T(\llbracket \varphi^*(\bar{x}; \bar{v}^*) \rrbracket^{\mathcal{B}}) \leq \varepsilon$.*
- (3) *The algorithm runs in time $(\log n + t)^{\mathcal{O}(1)} \cdot d^{\text{polylog}(d)}$ with only local access to \mathcal{B} , where $n := |\mathcal{B}|$, $d := \Delta\mathcal{B}$ and $t := |T|$.*
- (4) *The hypothesis returned by the algorithm can be evaluated in time $(\log n + t)^{\mathcal{O}(1)} \cdot d^{\text{polylog}(d)}$ with only local access to \mathcal{B} .*

Proof: The pseudocode for our algorithm is shown in Figure 7. Let $(\varphi^*(\bar{x}; \bar{y}), \bar{v}^*)$ be the hypothesis returned by the algorithm. By construction we know that the hypothesis satisfies (1).

Let $\varphi(\bar{x}; \bar{y}, \bar{\kappa})$ be a FOCN(P)-formula of binding width at most w and binding rank at most r and some tuples $\bar{v} \in U(\mathcal{B})^\ell$ and $\bar{\lambda} \in \{0, \dots, |\mathcal{B}|\}^{|\bar{\kappa}|}$ of parameters such that $\text{err}_T(\llbracket \varphi(\bar{x}; \bar{v}, \bar{\lambda}) \rrbracket^{\mathcal{B}})$ is minimal, especially $\text{err}_T(\llbracket \varphi(\bar{x}; \bar{v}, \bar{\lambda}) \rrbracket^{\mathcal{B}}) \leq \varepsilon$ for the input sequence T of training examples.

Note that there is a $\Delta\mathcal{B}$ -equivalent hnf-formula for φ of locality radius smaller than $(2w+1)^r - 1$ and thus $\mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}}(\bar{u}\bar{v}) \models \varphi(\bar{u}, \bar{v}, \bar{\lambda})$ iff $\mathcal{B} \models \varphi(\bar{u}, \bar{v}, \bar{\lambda})$ for all $(\bar{u}, c) \in T$. Furthermore note that $\mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}}(\bar{u}\bar{v}^*) \models \varphi^*(\bar{u}, \bar{v}^*)$ iff $\mathcal{B} \models \varphi^*(\bar{u}, \bar{v}^*)$ for all $(\bar{u}, c) \in T$, because φ^* is a Boolean combination of sphere formulas of locality radius smaller than $(2w+1)^r$. Let $T' \subseteq T$ be the subsequence of examples that are consistent with $\llbracket \varphi(\bar{x}; \bar{v}, \bar{\lambda}) \rrbracket^{\mathcal{B}}$. For every (u_i, c_i) in T let

$$\text{pos}_i := \{j \in [t] \mid c_j = 1 \text{ and } \mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}}(\bar{u}_i\bar{v}) \cong \mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}}(\bar{u}_j\bar{v})\}$$

and

$$\text{neg}_i := \{j \in [t] \mid c_j = 0 \text{ and } \mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}}(\bar{u}_i\bar{v}) \cong \mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}}(\bar{u}_j\bar{v})\}.$$

Algorithm \mathcal{L}_{\min}

Input: Training sequence $T = ((\bar{u}_1, c_1), \dots, (\bar{u}_t, c_t)) \in \mathcal{T}$,
 $d = \Delta\mathcal{B}$, local access to background structure \mathcal{B}

- 1: $N \leftarrow N_{2\ell[(2w+1)^r-1]}^{\mathcal{B}}(T)$
- 2: $\text{consistent}_{\max} \leftarrow -1$ \triangleright maximal number of examples consistent with a chosen formula
- 3: **for all** $\bar{v} \in N^m$, $m \leq \ell$ **do**
- 4: **for** $i = 1, \dots, t$ **do**
- 5: $\mathcal{N}_i \leftarrow \mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}}(\bar{u}_i\bar{v})$
- 6: $\vartheta_i(\bar{x}; \bar{y}) \leftarrow \text{sph}_{\mathcal{N}_i}(\bar{x}\bar{y})$ \triangleright a d -bounded $[(2w+1)^r-1]$ -type with $k+m$ centers
- 7: **for** $i = 1, \dots, t$ **do**
- 8: $\text{neg}_i \leftarrow 0$
- 9: $\text{pos}_i \leftarrow 0$
- 10: **for all** $j = 1, \dots, t$ **do**
- 11: **if** $\mathcal{N}_i \cong \mathcal{N}_j$ **then**
- 12: **if** $c_j = 1$ **then**
- 13: $\text{pos}_i \leftarrow \text{pos}_i + 1$
- 14: **else**
- 15: $\text{neg}_i \leftarrow \text{neg}_i + 1$
- 16: $\varphi(\bar{x}; \bar{y}) \leftarrow \bigvee_{i \in [t], \text{pos}_i \geq \text{neg}_i} \vartheta_i(\bar{x}; \bar{y})$
- 17: $\text{consistent}_{\text{cur}} \leftarrow |\{i \in [t] \mid (\text{pos}_i \geq \text{neg}_i \text{ and } c_i = 1) \text{ or } (\text{pos}_i < \text{neg}_i \text{ and } c_i = 0)\}|$
- 18: **if** $\text{consistent}_{\text{cur}} > \text{consistent}_{\max}$ **then**
- 19: $\varphi^* \leftarrow \varphi$
- 20: $\bar{v}^* \leftarrow \bar{v}$
- 21: $\text{consistent}_{\max} \leftarrow \text{consistent}_{\text{cur}}$
- 22: **return** $(\varphi^*(\bar{x}; \bar{y}), \bar{v}^*)$

Figure 7. Learning algorithm \mathcal{L}_{\min} of Theorem III.4

Claim 1. *If $(u_i, 1)$ in T' , then $|\text{pos}_i| \geq |\text{neg}_i|$.*

Proof: Consider the formula

$$\varphi'(\bar{x}; \bar{y}, \bar{\kappa}) := \varphi(\bar{x}; \bar{y}, \bar{\kappa}) \wedge \neg \text{sph}_{\mathcal{N}_i}(\bar{x}\bar{y})$$

with $\mathcal{N}_i = \mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}}(\bar{u}_i\bar{v})$. We know that $c_i = 1$ and thus $\{(u_j, c_j) \mid j \in \text{pos}_i\} = \{(u_j, 1) \mid j \in \text{pos}_i\} \subseteq T'$ and $\{(u_j, c_j) \mid j \in \text{neg}_i\} \subseteq T \setminus T'$. The hypothesis $\llbracket \varphi'(\bar{x}; \bar{v}, \bar{\lambda}) \rrbracket^{\mathcal{B}}$ is consistent with the examples

$$T' \setminus \{(u_j, c_j) \mid j \in \text{pos}_i\} \cup \{(u_j, c_j) \mid j \in \text{neg}_i\}.$$

The cardinality of this set is $|T'| - |\text{pos}_i| + |\text{neg}_i|$. The claim follows from the optimality of $(\varphi, \bar{v}, \bar{\lambda})$. \lrcorner

Claim 2. *If $(u_i, 0)$ in T' , then $|\text{neg}_i| \geq |\text{pos}_i|$.*

Proof: The proof is analogous to the proof of the first claim. Here we consider the formula

$$\varphi'(\bar{x}; \bar{y}, \bar{\kappa}) := \varphi(\bar{x}; \bar{y}, \bar{\kappa}) \vee \text{sph}_{\mathcal{N}_i}(\bar{x}\bar{y})$$

and again use the optimality of $(\varphi, \bar{v}, \bar{\lambda})$. \lrcorner

When using $\bar{v}^* = \bar{v}$, our algorithm constructs a formula that is consistent with all examples $(u_i, 0)$ where $|\text{pos}_i| < |\text{neg}_i|$ and all examples $(u_i, 1)$ where $|\text{pos}_i| \geq |\text{neg}_i|$. Using both claims we can follow that the hypothesis returned by the

algorithm is consistent with at least as many examples as the optimal solution. This proves (2).

Analogous to Theorem I.1 there are at most $(t+d)^{\mathcal{O}(1)}$ tuples in $\bigcup_{m=1}^{\ell} N^m$ and N can be computed with only local access to \mathcal{B} in time $(t+d)^{\mathcal{O}(1)}$. The algorithm in Figure 7 can compute all \mathcal{N}_i and sphere formulas ϑ_i in time $(\log n + d + t)^{\mathcal{O}(1)}$. The isomorphism tests can be done in time $(\log n + t)^{\mathcal{O}(1)} \cdot d^{\text{polylog}(d)}$. All in all the running time of the algorithm is in $(\log n + t)^{\mathcal{O}(1)} \cdot d^{\text{polylog}(d)}$. This shows that the algorithm satisfies Condition (3).

To evaluate the formula returned by the algorithm for \bar{u} , we only have to evaluate it using the structure $\mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}}(\bar{u})$ with only local access to \mathcal{B} , so the running time is in $(\log n + t)^{\mathcal{O}(1)} \cdot d^{\text{polylog}(d)}$ and Condition (4) is satisfied. \square

IV. STRUCTURES OF BOUNDED DEGREE

In this section we consider structures of bounded degree. We prove results that improve the running times we obtain for structures of polylogarithmic degree and we also extend these results to prove that there is a sublinear-time PAC-learning algorithm for hypotheses using FOCN(P)-formulas on structures of bounded degree. Let $d \in \mathbb{N}$ be fixed and let \mathcal{B} be a background structure with maximum degree at most d over a relational signature σ .

As in Section III, let Φ be the set of FOCN(P)-formulas $\varphi(\bar{x}; \bar{y}, \bar{\kappa})$ with binding width at most w , binding rank at most r and free number variables $\bar{\kappa}$, and

$$\mathcal{C} := \{ \llbracket \varphi(\bar{x}; \bar{v}, \bar{\lambda}) \rrbracket^{\mathcal{B}} \mid \varphi(\bar{x}; \bar{y}, \bar{\kappa}) \in \Phi, \bar{v} \in U(\mathcal{B})^{\ell}, \bar{\lambda} \in \{0, \dots, |\mathcal{B}|\}^{|\bar{\kappa}|} \}.$$

Let Φ^* be the set of normalized formulas $\varphi^*(\bar{x}; \bar{y}) = \bigwedge_i \bigvee_j \psi_{ij}(\bar{x}, \bar{y})$ where ψ_{ij} are (possibly negated) sphere formulas (modulo equivalence) with locality radius smaller than $(2w+1)^r$ for structures of degree at most d and

$$\mathcal{C}^* := \{ \llbracket \varphi^*(\bar{x}; \bar{v}) \rrbracket^{\mathcal{B}} \mid \varphi^*(\bar{x}; \bar{y}) \in \Phi^*, \bar{v} \in U(\mathcal{B})^{\ell} \}.$$

Lemma IV.1. $|\Phi^*|$ is finite and does not depend on $|\mathcal{B}|$.

Proof: There are at most $\nu_d(r) := 1 + d \cdot \sum_{i=0}^{r-1} (d-1)^i \leq d^{r+1} + 2r + 1$ elements in a $((2w+1)^r - 1)$ -type with a single center [3]. Hence there are at most

$$\begin{aligned} E(d, k, \ell, r, w) &:= (k + \ell) \cdot \nu_d((2w+1)^r - 1) \\ &\leq (k + \ell) \cdot (d^{(2w+1)^r} + 2(2w+1)^r) \end{aligned}$$

elements in a $((2w+1)^r - 1)$ -type with $k + \ell$ centers [3]. Thus there are at most

$$F(d, k, \ell, r, w, \sigma) := \prod_{R \in \sigma} 2^{(E(d, k, \ell, r, w)^{\text{ar}(R)})}$$

non-isomorphic $((2w+1)^r - 1)$ -types with $k + \ell$ centers.

$$\begin{aligned} \varphi^* &= \bigwedge_i \bigvee_j \underbrace{(\neg) \psi_{ij}}_{\leq 2^{F(\dots) \cdot 2^k k!}} \\ &\quad \underbrace{\leq 2^{2F(\dots) \cdot 2^k k!}}_{\leq 2^{2^{2F(\dots) \cdot 2^k k!}}} \end{aligned}$$

Hence the number of normalized formulas in Φ^* is at most $2^{2^{2F(d, k, \ell, r, w, \sigma)} \cdot 2^k k!}$. \square

Lemma IV.2. $|\mathcal{C}|$ is finite and independent from $|\mathcal{B}|$.

Proof: The set Φ^* is finite and thus also \mathcal{C}^* . For all formulas in Φ there is a d -equivalent formula in Φ^* and thus $\mathcal{C} \subseteq \mathcal{C}^*$. Hence \mathcal{C} is finite as well. \square

A. Consistent Hypotheses

As in Theorem I.1, we use brute force to learn a consistent hypothesis. Instead of building a formula from the training examples, we show that it suffices to use a formula from Φ^* . This makes the complexity of the formula independent from the number of training examples and thus we obtain a running time for the evaluation of the hypothesis that is independent from the number of training examples.

Corollary IV.3. Let $T \in \mathcal{T}$ be consistent with some $C \in \mathcal{C}$. Then there is a formula $\varphi^*(\bar{x}; \bar{y}) \in \Phi^*$ and a tuple $\bar{v}^* \in N_{2\ell[(2w+1)^r-1]}^{\mathcal{B}}(T)^{\ell}$ such that $\llbracket \varphi^*(\bar{x}; \bar{v}^*) \rrbracket^{\mathcal{B}}$ is consistent with T .

Proof: Let $\varphi^*(\bar{x}; \bar{y})$ and \bar{v}^* be as in Lemma III.1. Then, after some normalization for φ^* , $\varphi^*(\bar{x}; \bar{y}) \in \Phi^*$, $\bar{v}^* \in N_{2\ell[(2w+1)^r-1]}^{\mathcal{B}}(T)^{\ell}$ and $\llbracket \varphi^*(\bar{x}; \bar{v}^*) \rrbracket^{\mathcal{B}}$ is consistent with T . \square

To check the consistency of a hypothesis, we use the following model-checking result due to Grohe [19].

Definition IV.4. We say that a class C of structures has *low degree* if for every $\varepsilon > 0$ there is an integer N_{ε} such that for all $\mathcal{A} \in C$ with $|\mathcal{A}| \geq N_{\varepsilon}$ we have $\Delta \mathcal{A} \leq |\mathcal{A}|^{\varepsilon}$. \lrcorner

Theorem IV.5. There is an algorithm \mathfrak{A} for FO-MODEL-CHECKING and a function f such that for every class C of structures that has low degree and for every $\varepsilon > 0$ the running time of \mathfrak{A} on an input $(\mathcal{A}, \varphi) \in C \times \text{FO}$ is in $\mathcal{O}(f(\|\varphi\|) \cdot |\mathcal{A}|^{1+\varepsilon})$.

Using this result we obtain a model-learning algorithm with an improved running time for hypothesis evaluation.

Theorem IV.6. Let $d, k, \ell, r, w \in \mathbb{N}$. Then there is a learning algorithm $\mathfrak{L}_{\text{con}}^d$ for the k -ary learning problem over some finite relational structure \mathcal{B} of degree at most d , that receives a training sequence T as input, with the following properties:

- (1) If the algorithm returns a hypothesis H , then H is of the form $(\varphi^*(\bar{x}; \bar{y}), \bar{v}^*)$ for some first-order formula $\varphi^*(\bar{x}; \bar{y})$ that is a Boolean combination of sphere formulas with locality radius smaller than $(2w+1)^r$ and $\bar{v}^* \in U(\mathcal{B})^{\ell}$, and $\llbracket \varphi^*(\bar{x}; \bar{v}^*) \rrbracket^{\mathcal{B}}$ is consistent with T .
- (2) If there are an FOCN(P)-formula $\varphi(\bar{x}; \bar{y}, \bar{\kappa})$ of binding rank at most r and binding width at most w and parameter tuples $\bar{v} \in U(\mathcal{B})^{\ell}$ and $\bar{\lambda} \in \{0, \dots, |\mathcal{B}|\}^{|\bar{\kappa}|}$ such that $\llbracket \varphi(\bar{x}; \bar{y}, \bar{\lambda}) \rrbracket^{\mathcal{B}}$ is consistent with the input sequence T , then $\mathfrak{L}_{\text{con}}^d$ always returns a hypothesis.
- (3) The algorithm runs in time $(\log n + t)^{\mathcal{O}(1)}$ with only local access to \mathcal{B} , where $n := |\mathcal{B}|$ and $t := |T|$.

Algorithm \mathcal{L}_{con}^d

Input: Training sequence $T \in \mathcal{T}$, local access to background structure \mathcal{B}

```

1:  $N \leftarrow N_{2\ell[(2w+1)^r-1]}^{\mathcal{B}}(T)$ 
2: for all  $\bar{v}^* \in N^\ell$  do
3:   for all  $\varphi^*(\bar{x}; \bar{y}) \in \Phi^*$  do
4:      $consistent \leftarrow \text{true}$ 
5:     for all  $(\bar{u}, c) \in T$  do
6:       if  $(\mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}}(\bar{u}\bar{v}^*) \models \varphi^*(\bar{u}; \bar{v}^*) \text{ and } c = 0)$ 
or  $(\mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}}(\bar{u}\bar{v}^*) \not\models \varphi^*(\bar{u}; \bar{v}^*) \text{ and } c = 1)$  then
7:          $consistent \leftarrow \text{false}$ 
8:       if consistent then
9:         return  $(\varphi^*(\bar{x}; \bar{y}), \bar{v}^*)$ 
10: reject
```

Figure 8. Learning algorithm \mathcal{L}_{con}^d of Theorem IV.6

(4) The hypothesis returned by the algorithm can be evaluated in time $(\log n)^{\mathcal{O}(1)}$ with only local access to \mathcal{B} .

Proof: The pseudocode for our algorithm is shown in Figure 8. The algorithm goes through all tuples $\bar{v}^* \in N_{2\ell[(2w+1)^r-1]}^{\mathcal{B}}(T)^\ell$ and all formulas $\varphi^*(\bar{x}; \bar{y}) \in \Phi^*$ and checks, whether $\llbracket \varphi^*(\bar{x}; \bar{v}^*) \rrbracket^{\mathcal{B}}$ is consistent with T . The algorithm returns the first consistent (φ^*, \bar{v}^*) and rejects if there is none. Note that $\mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}} \models \varphi^*(\bar{u}; \bar{v}^*)$ iff $\mathcal{B} \models \varphi^*(\bar{u}; \bar{v}^*)$ for all $(\bar{u}, c) \in T$, because φ^* is a Boolean combination of sphere formulas of locality radius smaller than $(2w+1)^r$. Thus the algorithm satisfies Condition (1) of Theorem IV.6 and because of Corollary IV.3 it satisfies (2) as well.

Let $n := |\mathcal{B}|$ and $t := |T|$. Then for all $\bar{u} \in U(\mathcal{B})^k$ and $\bar{v}^* \in U(\mathcal{B})^\ell$ we have $|\mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}}(\bar{u}\bar{v}^*)| \leq (k+\ell) \cdot 2d^{(2w+1)^r}$ and the representation size is $(\log n)^{\mathcal{O}(1)}$, because d, k, ℓ, w and r are constant. Lemma IV.1 tells us that the number of formulas in Φ^* to check is constant. Thus for every real valued function f it follows that $\max_{\psi \in \Phi^*} f(\|\psi\|)$ is finite. Hence according to Theorem IV.5 it takes time polynomial in the size of $\mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}}(\bar{u}\bar{v}^*)$ to check whether the structure satisfies a formula $\varphi^*(\bar{u}; \bar{v}^*)$. The number of tuples in N^ℓ is $|N|^\ell \leq (2tkd^{2\ell(2w+1)^r})^\ell \in t^{\mathcal{O}(1)}$ and N can be computed with only local access to \mathcal{B} in time $t^{\mathcal{O}(1)}$. All in all the running time of the algorithm is in $t^{\mathcal{O}(1)} \cdot t \cdot (\log n)^{\mathcal{O}(1)} \leq (\log n + t)^{\mathcal{O}(1)}$. This shows that the algorithm satisfies Condition (3).

To evaluate the formula returned by the algorithm for (\bar{u}, \bar{v}) , we only have to evaluate it using the structure $\mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}}(\bar{u}\bar{v})$ with only local access to \mathcal{B} , so the running time is in $(\log n)^{\mathcal{O}(1)}$ and Condition (4) is satisfied. \square

B. Minimizing the training error

Corollary IV.7. Let $T \in \mathcal{T}$ be such that $\text{err}_T(C) \leq \varepsilon$ for some $C \in \mathcal{C}$. Then there is a formula $\varphi^*(\bar{x}; \bar{y}) \in \Phi^*$ and a tuple $\bar{v}^* \in N_{2\ell[(2w+1)^r-1]}^{\mathcal{B}}(T)^\ell$ such that $\text{err}_T(\llbracket \varphi^*(\bar{x}, \bar{v}^*) \rrbracket^{\mathcal{B}}) \leq \varepsilon$.

Proof: If $\text{err}_T(C) \leq \varepsilon$, then there is some $S \subseteq T$ with $|S| \geq (1 - \varepsilon) \cdot |T|$ such that C is consistent with S .

Algorithm \mathcal{L}_{min}^d

Input: Training sequence $T \in \mathcal{T}$, local access to background structure \mathcal{B}

```

1:  $N \leftarrow N_{2\ell[(2w+1)^r-1]}^{\mathcal{B}}(T)$ 
2:  $\text{err}_{min} \leftarrow |T| + 1$ 
3: for all  $\bar{v}^* \in N^\ell$  do
4:   for all  $\varphi^*(\bar{x}; \bar{y}) \in \Phi^*$  do
5:      $\text{err}_{cur} \leftarrow 0$ 
6:     for all  $(\bar{u}, c) \in T$  do
7:       if  $(\mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}}(\bar{u}\bar{v}^*) \models \varphi^*(\bar{u}; \bar{v}^*) \text{ and } c = 0)$ 
or  $(\mathcal{N}_{(2w+1)^r-1}^{\mathcal{B}}(\bar{u}\bar{v}^*) \not\models \varphi^*(\bar{u}; \bar{v}^*) \text{ and } c = 1)$  then
8:          $\text{err}_{cur} \leftarrow \text{err}_{cur} + 1$ 
9:     if  $\text{err}_{cur} < \text{err}_{min}$  then
10:       $\text{err}_{min} \leftarrow \text{err}_{cur}$ 
11:       $\varphi_{min}^* \leftarrow \varphi^*$ 
12:       $\bar{v}_{min}^* \leftarrow \bar{v}^*$ 
13: return  $(\varphi_{min}^*, \bar{v}_{min}^*)$ 
```

Figure 9. Learning algorithm \mathcal{L}_{min}^d of Theorem IV.8 and Theorem IV.10

Using Corollary IV.3 on S we obtain a formula $\varphi^*(\bar{x}; \bar{y})$ and a tuple $\bar{v}^* \in N_{2\ell[(2w+1)^r-1]}^{\mathcal{B}}(S)^\ell \subseteq N_{2\ell[(2w+1)^r-1]}^{\mathcal{B}}(T)^\ell$ such that $\llbracket \varphi^*(\bar{x}, \bar{v}^*) \rrbracket^{\mathcal{B}}$ is consistent with S and thus $\text{err}_T(\llbracket \varphi^*(\bar{x}, \bar{v}^*) \rrbracket^{\mathcal{B}}) \leq \varepsilon$. \square

Theorem IV.8. Let $d, k, \ell, r, w \in \mathbb{N}$. Then there is a learning algorithm \mathcal{L}_{min}^d for the k -ary learning problem over some finite relational structure \mathcal{B} of degree at most d , that receives a training sequence T as input, with the following properties:

- (1) The algorithm always returns a hypothesis H of the form $(\varphi^*(\bar{x}; \bar{y}), \bar{v}^*)$ for some first-order formula $\varphi^*(\bar{x}; \bar{y})$ that is a Boolean combination of sphere formulas with locality radius smaller than $(2w+1)^r$ and $\bar{v}^* \in U(\mathcal{B})^\ell$.
- (2) If there are an FOCN(P)-formula $\varphi(\bar{x}; \bar{y}, \bar{\kappa})$ of binding rank at most r and binding width at most w and parameter tuples $\bar{v} \in U(\mathcal{B})^\ell$ and $\bar{\lambda} \in \{0, \dots, |\mathcal{B}|\}^{|\bar{\kappa}|}$ such that $\text{err}_T(\llbracket \varphi(\bar{x}; \bar{v}, \bar{\lambda}) \rrbracket^{\mathcal{B}}) \leq \varepsilon$, then \mathcal{L}_{min}^d returns a hypothesis $(\varphi^*(\bar{x}; \bar{y}), \bar{v}^*)$ with $\text{err}_T(\llbracket \varphi^*(\bar{x}; \bar{v}^*) \rrbracket^{\mathcal{B}}) \leq \varepsilon$.
- (3) The algorithm runs in time $(\log n + t)^{\mathcal{O}(1)}$ with only local access to \mathcal{B} , where $n := |\mathcal{B}|$ and $t := |T|$.
- (4) The hypothesis returned by the algorithm can be evaluated in time $(\log n)^{\mathcal{O}(1)}$ with only local access to \mathcal{B} .

Proof: The pseudocode for our algorithm is shown in Figure 9. The algorithm goes through all tuples $\bar{v}^* \in N_{2\ell[(2w+1)^r-1]}^{\mathcal{B}}(T)^\ell$ and all formulas $\varphi^*(\bar{x}; \bar{y}) \in \Phi^*$ and counts the number of errors that $\llbracket \varphi^*(\bar{x}; \bar{v}^*) \rrbracket^{\mathcal{B}}$ makes on T . The algorithm returns the hypothesis with minimum error. Using Corollary IV.7 one can show analogously to Theorem IV.6 that (1) and (2) hold. The running-time analysis for (3) and (4) is also analogous to the proof of Theorem IV.6. \square

C. Agnostic PAC Learning

We give a short introduction to *probably approximately correct* (PAC) learning. For more background, we refer to [20]. Instead of focussing on the training error, we are now interested

in hypotheses that generalize well. We assume that there is an (unknown) probability distribution \mathcal{D} on $U(\mathcal{B})^k \times \{0, 1\}$ and that training examples are drawn independently from this distribution. Our goal is to find an algorithm that with high probability on training examples drawn from \mathcal{D} , returns a hypothesis that has a small expected error on instances drawn from the same distribution \mathcal{D} .

The *generalization error* of a hypothesis $H: U(\mathcal{B})^k \rightarrow \{0, 1\}$ for a probability distribution \mathcal{D} on $U(\mathcal{B})^k \times \{0, 1\}$ is

$$\text{err}_{\mathcal{D}}(H) := \mathcal{P}_{(\bar{u}, c) \sim \mathcal{D}}(H(\bar{u}) \neq c).$$

A hypothesis class \mathcal{H} is *agnostic PAC-learnable* if there is some function $t_{\mathcal{H}}: (0, 1)^2 \rightarrow \mathbb{N}$ and a learning algorithm \mathcal{L} such that for every $\varepsilon, \delta \in (0, 1)$ and for every distribution \mathcal{D} over $U(\mathcal{B})^k \times \{0, 1\}$, when running the algorithm on a sequence T of $t_{\mathcal{H}}(\varepsilon, \delta)$ examples drawn i.i.d. from \mathcal{D} , it satisfies

$$\mathcal{P}_{T \sim \mathcal{D}^t} \left[\text{err}_{\mathcal{D}}(\mathcal{L}(T)) \leq \inf_{H \in \mathcal{H}} \text{err}_{\mathcal{D}}(H) + \varepsilon \right] \geq 1 - \delta.$$

We use the following lemma from [20] to bound the generalization error.

Lemma IV.9 (Uniform Convergence). *Let \mathcal{H} be a finite hypothesis class of hypotheses $H: U(\mathcal{B})^k \rightarrow \{0, 1\}$ and consider training sequences T of length*

$$t \geq t_{\mathcal{H}}^{UC}(\varepsilon, \delta) := \left\lceil \frac{\log(2|\mathcal{H}|/\delta)}{2\varepsilon^2} \right\rceil$$

where the examples are drawn i.i.d. from a probability distribution \mathcal{D} over $U(\mathcal{B})^k \times \{0, 1\}$. Then

$$\mathcal{P}_{T \sim \mathcal{D}^t} [|\text{err}_{\mathcal{D}}(H) - \text{err}_T(H)| \leq \varepsilon \text{ for all } H \in \mathcal{H}] \geq 1 - \delta.$$

Now we show that there is an agnostic PAC-learning algorithm for hypotheses using FOCN(P)-formulas on background structures of bounded degree.

Theorem IV.10. *Let $d, k, \ell, r, w \in \mathbb{N}$. Then there is some $s \in \mathbb{N}$ such that the learning algorithm \mathcal{L}_{min}^d for the k -ary learning problem over some finite relational structure \mathcal{B} of degree at most d has the following properties:*

- (1) *The algorithm always returns a hypothesis H of the form $(\varphi^*(\bar{x}; \bar{y}^*))$ for some first-order formula $\varphi^*(\bar{x}; \bar{y})$ that is a Boolean combination of sphere formulas with locality radius smaller than $(2w + 1)^r$ and $\bar{v}^* \in U(\mathcal{B})^\ell$ such that for an input sequence T of at least $|T| =: t = s \left\lceil \frac{\log(|\mathcal{B}|/\delta)}{\varepsilon^2} \right\rceil$ training examples it holds that*

$$\mathcal{P}_{T \sim \mathcal{D}^t} \left[\text{err}_{\mathcal{D}}(\llbracket \varphi^*(\bar{x}, \bar{v}^*)^B \rrbracket) - \min_{C \in \mathcal{C}} \text{err}_{\mathcal{D}}(C) \leq \varepsilon \right] \geq 1 - \delta.$$

- (2) *If $|T| = s \left\lceil \frac{\log(|\mathcal{B}|/\delta)}{\varepsilon^2} \right\rceil$, then the algorithm runs in time $(\log |\mathcal{B}| + 1/\varepsilon + \log 1/\delta)^{\mathcal{O}(1)}$ with only local access to \mathcal{B} .*
- (3) *The hypothesis returned by the algorithm can be evaluated in time $(\log n)^{\mathcal{O}(1)}$ with only local access to \mathcal{B} .*

Proof: By Lemma IV.2, the hypothesis class $\mathcal{H}^* := \Phi^* \times U(\mathcal{B})^\ell$ is finite. Thus we can bound the generalization error of a

hypothesis H returned by \mathcal{L}_{min}^d using the uniform convergence Lemma IV.9 by

$$\mathcal{P}_{T \sim \mathcal{D}^t} [|\text{err}_{\mathcal{D}}(H) - \text{err}_T(H)| \leq \varepsilon/2] \geq 1 - \delta/2$$

for input sequences T of length at least $t = \left\lceil \frac{\log(4|\mathcal{H}^*|/\delta)}{2(\varepsilon/2)^2} \right\rceil = \left\lceil \frac{4 \log |\Phi^*| \cdot \ell \cdot \log(|\mathcal{B}|/\delta)}{\varepsilon^2} \right\rceil$. Furthermore for all $C \in \mathcal{C}$,

$$\mathcal{P}_{T \sim \mathcal{D}^t} [|\text{err}_T(C) - \text{err}_{\mathcal{D}}(C)| \leq \varepsilon/2] \geq 1 - \delta/2$$

for input sequences T of length at least $t = \left\lceil \frac{\log(4|\mathcal{C}|/\delta)}{2(\varepsilon/2)^2} \right\rceil \leq \left\lceil \frac{4 \log |\Phi^*| \cdot \ell \cdot \log(|\mathcal{B}|/\delta)}{\varepsilon^2} \right\rceil$, because $|\mathcal{C}^*| \leq |\Phi^*|$ and $\mathcal{C} \subseteq \mathcal{C}^*$. Using Corollary IV.7 we know that $\text{err}_T(H) \leq \text{err}_T(C)$ for all $C \in \mathcal{C}$. Thus for all $C \in \mathcal{C}$

$$\begin{aligned} & \mathcal{P}_{T \sim \mathcal{D}^t} [\text{err}_{\mathcal{D}}(H) - \text{err}_{\mathcal{D}}(C) \leq \varepsilon] \\ &= \mathcal{P}_{T \sim \mathcal{D}^t} [\text{err}_{\mathcal{D}}(H) - \text{err}_T(H) + \text{err}_T(H) \\ & \quad - \text{err}_T(C) + \text{err}_T(C) - \text{err}_{\mathcal{D}}(C) \leq \varepsilon] \\ &\geq \mathcal{P}_{T \sim \mathcal{D}^t} [\text{err}_{\mathcal{D}}(H) - \text{err}_T(H) \leq \varepsilon/2 \text{ and} \\ & \quad \text{err}_T(H) \leq \text{err}_T(C) \text{ and } \text{err}_T(C) - \text{err}_{\mathcal{D}}(C) \leq \varepsilon/2] \\ &= \mathcal{P}_{T \sim \mathcal{D}^t} [\text{err}_{\mathcal{D}}(H) - \text{err}_T(H) \leq \varepsilon/2 \text{ and} \\ & \quad \text{err}_T(C) - \text{err}_{\mathcal{D}}(C) \leq \varepsilon/2] \\ &\geq \mathcal{P}_{T \sim \mathcal{D}^t} [|\text{err}_{\mathcal{D}}(H) - \text{err}_T(H)| \leq \varepsilon/2 \text{ and} \\ & \quad |\text{err}_T(C) - \text{err}_{\mathcal{D}}(C)| \leq \varepsilon/2] \\ &\geq 1 - \delta. \end{aligned}$$

This holds especially for $C = \argmin_{C' \in \mathcal{C}} \text{err}_{\mathcal{D}}(C')$ and hence (1) follows with $s = 4\ell \cdot \log |\Phi^*|$. The statements (2) and (3) follow immediately from Theorem IV.8. \square

V. CONCLUSIONS

We prove that FOCN(P)-definable concepts over structures of polylogarithmic degree can be learned in sublinear time. For structures with no degree bound we show that there is no consistent model-learning algorithm that runs in sublinear time with only local access to the structure. Furthermore, we show how to use a consistent parameter-learning algorithm to solve q -CLIQUE.

It remains open whether one can obtain a similar result for model learning and improve the lower bound on the running time for consistent model-learning algorithms. Our results imply PAC-learnability on structures of bounded degree. It would be interesting to investigate structures of polylogarithmic degree in this context.

In addition to the COUNT operator that we analyze within the logic FOCN(P), another direction for future research is the analysis of learning algorithms for stronger logics that implement other aggregating operators from SQL.

REFERENCES

- [1] M. Grohe and G. Turán, “Learnability and definability in trees and similar structures,” *Theory Comput. Syst.*, vol. 37, no. 1, pp. 193–220, 2004.
- [2] M. Grohe and M. Ritzert, “Learning first-order definable concepts over structures of small degree,” in *32nd Annual ACM/IEEE Symposium on Logic in Computer Science, LICS 2017, Reykjavik, Iceland, June 20-23, 2017*, pp. 1–12, IEEE Computer Society, 2017.
- [3] D. Kuske and N. Schweikardt, “First-order logic with counting: At least, weak hanf normal forms always exist and can be computed!,” *CoRR*, vol. abs/1703.01122, 2017.
- [4] M. Grohe and N. Schweikardt, “First-order query evaluation with cardinality conditions,” in *Proceedings of the 37th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, Houston, TX, USA, June 10-15, 2018* (J. V. den Bussche and M. Arenas, eds.), pp. 253–266, ACM, 2018.
- [5] M. Grohe, D. Neuen, and P. Schweitzer, “A faster isomorphism test for graphs of small degree,” in *59th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2018, Paris, France, October 7-9, 2018* (M. Thorup, ed.), pp. 89–100, IEEE Computer Society, 2018.
- [6] M. Grohe, C. Löding, and M. Ritzert, “Learning MSO-definable hypotheses on strings,” in *International Conference on Algorithmic Learning Theory, ALT 2017, 15-17 October 2017, Kyoto University, Kyoto, Japan* (S. Hanneke and L. Reyzin, eds.), vol. 76 of *Proceedings of Machine Learning Research*, pp. 434–451, PMLR, 2017.
- [7] W. W. Cohen and C. D. Page Jr., “Polynomial learnability and inductive logic programming: Methods and results,” *New Generation Comput.*, vol. 13, no. 3&4, pp. 369–409, 1995.
- [8] J. Kietz and S. Dzeroski, “Inductive logic programming and learnability,” *SIGART Bulletin*, vol. 5, no. 1, pp. 22–32, 1994.
- [9] S. Muggleton, “Inductive logic programming,” *New Generation Comput.*, vol. 8, no. 4, pp. 295–318, 1991.
- [10] S. Muggleton, “Inductive logic programming,” in *Inductive Logic Programming* (S. Muggleton, ed.), vol. 38 of *The APIC Series*, pp. 1–27, Academic Press, 1992.
- [11] S. Muggleton and L. D. Raedt, “Inductive logic programming: Theory and methods,” *J. Log. Program.*, vol. 19/20, pp. 629–679, 1994.
- [12] A. Abouzied, D. Angluin, C. H. Papadimitriou, J. M. Hellerstein, and A. Silberschatz, “Learning and verifying quantified boolean queries by example,” in *Proceedings of the 32nd ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS 2013, New York, NY, USA - June 22 - 27, 2013* (R. Hull and W. Fan, eds.), pp. 49–60, ACM, 2013.
- [13] A. Bonifati, R. Ciucanu, and S. Staworko, “Learning join queries from user examples,” *ACM Trans. Database Syst.*, vol. 40, no. 4, pp. 24:1–24:38, 2016.
- [14] L. Libkin, *Elements of Finite Model Theory*. Texts in Theoretical Computer Science. An EATCS Series, Springer, 2004.
- [15] M. Grohe, *Descriptive Complexity, Canonisation, and Definable Graph Structure Theory*. Lecture Notes in Logic, Cambridge University Press, 2017.
- [16] L. Hella, L. Libkin, J. Nurmonen, and L. Wong, “Logics with aggregate operators,” *J. ACM*, vol. 48, no. 4, pp. 880–907, 2001.
- [17] S. Feferman and R. L. Vaught, “The first order properties of products of algebraic systems,” *Fundamenta Mathematicae*, vol. 47, no. 1, pp. 57–103, 1959.
- [18] D. Lokshtanov, D. Marx, and S. Saurabh, “Lower bounds based on the exponential time hypothesis,” *Bulletin of the EATCS*, vol. 105, pp. 41–72, 2011.
- [19] M. Grohe, “Generalized model-checking problems for first-order logic,” in *STACS 2001, 18th Annual Symposium on Theoretical Aspects of Computer Science, Dresden, Germany, February 15-17, 2001, Proceedings* (A. Ferreira and H. Reichel, eds.), vol. 2010 of *Lecture Notes in Computer Science*, pp. 12–26, Springer, 2001.
- [20] S. Shalev-Shwartz and S. Ben-David, *Understanding Machine Learning: From Theory to Algorithms*. New York, NY, USA: Cambridge University Press, 2014.