

Distribution Free Learning with Local Queries*

Galit Bary-Weisberg[†]

Amit Daniely[‡]

Shai Shalev-Shwartz[§]

March 8, 2016

Abstract

The model of learning with *local membership queries* interpolates between the PAC model and the membership queries model by allowing the learner to query the label of any example that is similar to an example in the training set. This model, recently proposed and studied by Awasthi et al. [5], aims to facilitate practical use of membership queries.

We continue this line of work, proving both positive and negative results in the *distribution free* setting. We restrict to the boolean cube $\{-1, 1\}^n$, and say that a query is q -local if it is of a hamming distance $\leq q$ from some training example. On the positive side, we show that 1-local queries already give an additional strength, and allow to learn a certain type of DNF formulas. On the negative side, we show that even $(n^{0.99})$ -local queries cannot help to learn various classes including Automata, DNFs and more. Likewise, q -local queries for any constant q cannot help to learn Juntas, Decision Trees, Sparse Polynomials and more. Moreover, for these classes, an algorithm that uses $(\log^{0.99}(n))$ -local queries would lead to a breakthrough in the best known running times.

*This paper is based on the M.Sc. thesis [6] of the first author. The thesis offers a more elaborated discussion, as well as experiments.

[†]Matific inc. Most work was done while the author was an M.Sc. student at the Hebrew University, Jerusalem, Israel

[‡]Google inc. Most work was done while the author was a Ph.D. student at the Hebrew University, Jerusalem, Israel

[§]School of Computer Science and Engineering, The Hebrew University, Jerusalem, Israel

1 Introduction

A child typically learns to recognize a cat based on two types of input. The first is given by her parents, pointing at a cat and saying “Look, a cat!”. The second is given as a response to the child’s frequent question “What is that?”. These two types of input were the basis for the learning model originally suggested by Valiant [21]. Indeed, in Valiant’s model, the learner can randomly sample labelled examples from “nature”, but it can also make a *membership query (MQ)* for the label of *any* unseen example. Today, the acronym PAC stands for the restricted model in which MQ are forbidden, while the full model is called PAC+MQ. Much work has been done investigating the limits and the strengths of MQ. In particular, membership queries were proven stronger than the vanilla PAC model [1, 10, 12, 16]. Yet, MQ are rarely used in practice. This is commonly attributed to the fact that MQ algorithms query very artificial examples, that are uninterpretable by humans (e.g., [8]).

Awasthi et al. [5] suggested a solution to the problem of unnatural examples. They considered a mid-way model that allows algorithms to make only *local* queries, i.e., query examples that are close to examples from the sample set. Hopefully, examples which are similar to natural examples will appear natural to humans. Awasthi et al. considered the case where the instance space is $\{-1, 1\}^n$ and the distance between examples is the Hamming distance. They proved positive results on learning sparse polynomials with $O(\log(n))$ -local queries under what they defined as locally smooth distributions¹. They also proposed an algorithm that learns DNF formulas under the uniform distribution in quasi-polynomial time using $O(\log(n))$ -local queries.

Our work follows Awasthi et al. and investigates local queries in the *distribution free* setting, in which no explicit assumptions are made on the underlying distribution, but only on the learned hypothesis. We prove both positive and negative results in this context:

- One of the strongest and most beautiful results in the MQ model shows that automata are learnable with membership queries [1]. We show that unfortunately, this is probably not the case with local queries. Concretely, we show that even $(n^{0.99})$ -local queries cannot help to learn automata. Namely, such an algorithm will imply a standard PAC algorithm for automata. As learning automata is hard under several assumptions [17, 13], our result suggests that it is hard to learn automata even with $(n^{0.99})$ -local queries.
- We prove a similar result for several additional classes. Namely, we show that $(n^{0.99})$ -local queries cannot help to learn DNFs, intersection of halfspaces, decision lists, depth- d circuits for any $d \geq 2$, and depth- d threshold circuits for any $d \geq 2$. Likewise, for any constant q , q -local queries cannot help to learn Juntas, Decision Trees, Sparse polynomials, and Sparse polynomial threshold functions. In fact, we show that even $(\log^{0.99}(n))$ -local queries are unlikely to lead to polynomial time algorithms. Namely, any algorithm that uses $(\log^{0.99}(n))$ -local queries will result with a PAC algorithm whose running time significantly improves the state of the art in these well studied problems.
- On the positive side we show that already 1-local queries are probably stronger than the vanilla PAC model. Concretely, we show that a certain simplistic but natural

¹A distribution is locally α -smooth for $\alpha \geq 1$ if its density function is $\log(\alpha)$ -Lipschitz.

learning problem, which we term learning DNFs with evident examples, is learnable with 1-local queries. Furthermore, we show that without queries, this problem is at least as hard as learning decision trees, that are conjectured to be hard to learn.

2 Previous Work

Membership Queries Several concept classes are known to be learnable only if membership queries are allowed: Deterministic Finite Automata [1], k -term DNF for $k = \frac{\log(n)}{\log(\log(n))}$ [10], decision trees and k -almost monotone-DNF formulas [12], intersections of k -halfspaces [7] and DNF formulas under the uniform distribution [16]. The last result builds on Freund’s boosting algorithm [15] and the Fourier-based technique for learning using membership queries due to [18]. We note that there are cases in which MQ do not help. E.g., in the case of learning DNF and CNF formulas [2], assuming that one way functions exist, and in the case of distribution free agnostic learning [14].

Local Membership Queries Awasthi et al. focused on learning with $O(\log(n))$ -local queries. They showed that t -sparse polynomials are learnable under locally smooth distributions using $O(\log(n) + \log(t))$ -local queries, and that DNF formulas are learnable under the uniform distribution in quasi-polynomial time ($n^{O(\log \log n)}$) using $O(\log(n))$ -local queries. They also presented some results regarding the strength of local MQ. They proved that under standard cryptographic assumptions, $(r + 1)$ -local queries are more powerful than r -local queries (for every $1 \leq r \leq n - 1$). They also showed that local queries do not always help. They showed that if a concept class is agnostically learnable under the uniform distribution using k -local queries (for constant k) then it is also agnostically learnable (under the uniform distribution) in the PAC model.

We remark that besides local queries, there were additional suggestions to solve the problem of unnatural examples. For example, to allow “I don’t know” answers, or to be tolerant to some incorrect answers [3, 4, 11, 20, 9].

3 Setting

We consider binary classification where the instance space is $\mathcal{X} = \mathcal{X}_n = \{-1, 1\}^n$ and the label space is $\mathcal{Y} = \{0, 1\}$. A learning problem is defined by a hypothesis class $\mathcal{H} \subset \{0, 1\}^{\mathcal{X}}$. The learner receives a *training set*

$$S = \{(\mathbf{x}_1, h^*(\mathbf{x}_1)), (\mathbf{x}_2, h^*(\mathbf{x}_2)), \dots, (\mathbf{x}_m, h^*(\mathbf{x}_m))\} \in (\mathcal{X} \times \mathcal{Y})^m$$

where the \mathbf{x}_i ’s are sampled i.i.d. from some *unknown* distribution \mathcal{D} on \mathcal{X} and $h^* : \mathcal{X} \rightarrow \mathcal{Y}$ is some *unknown* hypothesis. The learner is also allowed to make membership queries. Namely, to call an oracle, which receives as input some $\mathbf{x} \in \mathcal{X}$ and returns $h^*(\mathbf{x})$. We say that a membership query for $\mathbf{x} \in \mathcal{X}$ is q -local if there exists a training example \mathbf{x}_i whose Hamming distance from \mathbf{x} is at most q .

The learner returns (a description of) a hypothesis $\hat{h} : \mathcal{X} \rightarrow \mathcal{Y}$. The goal is to approximate h^* , namely to find $\hat{h} : \mathcal{X} \rightarrow \mathcal{Y}$ with *loss* as small as possible, where the loss is defined as

$L_{\mathcal{D}, h^*}(\hat{h}) = \Pr_{\mathbf{x} \sim \mathcal{D}} \left(\hat{h}(\mathbf{x}) \neq h^*(\mathbf{x}) \right)$. We will focus on the so-called realizable case where h^* is assumed to be in \mathcal{H} , and will require algorithms to return a hypothesis with loss $< \epsilon$ in time that is polynomial in n and $\frac{1}{\epsilon}$. Concretely,

Definition 3.1 (Membership-Query Learning Algorithm) *We say that a learning algorithm \mathcal{A} learns \mathcal{H} with q -local membership queries if*

- *There exists a function $m_{\mathcal{A}}(n, \epsilon) \leq \text{poly}\left(n, \frac{1}{\epsilon}\right)$, such that for every distribution \mathcal{D} over \mathcal{X} , every $h^* \in \mathcal{H}$ and every $\epsilon > 0$, if \mathcal{A} is given access to q -local membership queries, and a training sequence*

$$S = \{(\mathbf{x}_1, h^*(\mathbf{x}_1)), (\mathbf{x}_2, h^*(\mathbf{x}_2)), \dots, (\mathbf{x}_m, h^*(\mathbf{x}_m))\}$$

where the \mathbf{x}_i 's are sampled i.i.d. from \mathcal{D} and $m \geq m_{\mathcal{A}}(n, \epsilon)$, then with probability of at least² $\frac{3}{4}$ (over the choice of S), the output \hat{h} of \mathcal{A} satisfies $L_{\mathcal{D}, h^}(\hat{h}) < \epsilon$.*

- *Given a training set of size m*
 - *\mathcal{A} asks at most $\text{poly}(m, n)$ membership queries.*
 - *\mathcal{A} runs in time $\text{poly}(m, n)$.*
 - *The hypothesis returned by \mathcal{A} can be evaluated in time $\text{poly}(m, n)$.*

We remark that learning with 0-local queries is equivalent to PAC learning, while learning with n -local queries is equivalent to PAC+MQ learning.

4 Learning DNFs with Evident Examples

Intuitively, when evaluating a DNF formula on a given example, we check a few conditions corresponding to the formula's terms, and deem the example positive if one of them holds. We will consider the case that for each of these conditions, there is a chance to see a “prototype example”, that satisfies it in a strong or evident way. In the sequel, we denote by $h_F : \{-1, 1\}^n \rightarrow \{0, 1\}$ the function induced by a DNF formula F over n variables.

Definition 4.1 *Let $F = T_1 \vee T_2 \vee \dots \vee T_d$ be a DNF formula. We say that an example $\mathbf{x} \in \{-1, 1\}^n$ satisfies a term T_i (with respect to the formula F) **evidently** and denote $T_i(\mathbf{x}) \equiv 1$ if :*

- *It satisfies T_i . (In particular, $h_F(\mathbf{x}) = 1$)*
- *It does **not** satisfy any other term T_k (for $k \neq i$) from F .*
- *No coordinate change will turn T_i False and another term T_k True. Concretely, if for $j \in [n]$ we denote $\mathbf{x}^{\oplus j} = (x_1, \dots, x_{j-1}, -x_j, x_{j+1}, \dots, x_n)$, then for every coordinate $j \in [n]$, if $\mathbf{x}^{\oplus j}$ satisfies F (i.e. if $h_F(\mathbf{x}^{\oplus j}) = 1$) then $\mathbf{x}^{\oplus j}$ satisfies T_i and only T_i .*

²The success probability can be amplified to $1 - \delta$ by repetition.

Definition 4.2 Let $F = T_1 \vee T_2 \vee \dots \vee T_d$ be a DNF formula. We say that $h^* : \{-1, 1\}^n \rightarrow \{0, 1\}$ is realized by F with evident examples w.r.t. a distribution \mathcal{D} if $h^* = h_F$ and for every term³ T_i , $\Pr_{\mathbf{x} \sim \mathcal{D}}(T_i(\mathbf{x}) \equiv 1 | T_i(\mathbf{x}) = 1) \geq \frac{1}{n}$.

For example, the definition holds whenever h^* is realized by a DNF in which any pair of different terms contains two opposite literals. Also, functions that are realized by a decision tree can be also realized by a DNF in which every positive example satisfies a single term, corresponding to root-to-leaf path that ends with a leaf labelled by 1. Hence, the assumption holds for decision trees provided that for every such path, there is a chance to see an example satisfying it evidently, in the sense that flipping each variable in the path will turn the example negative.

4.1 An algorithm

Algorithm 1

Input: $S_1, S_2 \in (\{-1, 1\}^n \times \{0, 1\})^m$

Output: A DNF formula H

Start with an empty DNF formula H

for all positive examples $(\mathbf{x}, y) \in S_1$ **do**

Define $T = x_1 \wedge \bar{x}_1 \wedge x_2 \wedge \bar{x}_2 \wedge \dots \wedge x_n \wedge \bar{x}_n$

for $1 \leq j \leq n$ **do**

Query $\mathbf{x}^{\oplus j}$ to get $h^*(\mathbf{x}^{\oplus j})$

if $h^*(\mathbf{x}^{\oplus j}) = 1$ **then**

Remove x_j and \bar{x}_j from T

if $h^*(\mathbf{x}^{\oplus j}) = 0$ **then**

if $x_j = 1$ **then**

Remove \bar{x}_j from T

if $x_j = 0$ **then**

Remove x_j from T

$H = H \vee T$

for all T in H **do**

if $T(\mathbf{x}) = 1$ but $y = 0$ for some $(\mathbf{x}, y) \in S_2$ **then**

Remove T from H

Return H

Theorem 4.3 Algorithm 1 learns with 1-local-queries poly-sized DNFs with evident examples.

Idea The algorithm is based on the following claim that follows easily from definition 4.1.

³We note that the quantity $\frac{1}{n}$ in the following equation is arbitrary, and can be replaced by $\frac{1}{n^c}$ for any $c > 0$.

Claim 1 Let $F = T_1 \vee T_2 \vee \dots \vee T_d$ be a DNF formula over $\{-1, 1\}^n$. For every $\mathbf{x} \in \{-1, 1\}^n$ that satisfies a term T_i evidently (with respect to F), and for every $j \in [n]$ it holds that:

$$h_F(\mathbf{x}^{\oplus j}) = 1 \iff \text{the term } T_i \text{ does not contain the variable } x_j$$

By this claim, if \mathbf{x} is an evident example for a certain term T , one can easily reconstruct T . Indeed, by flipping the value of each variable and checking if the label changes, one can infer which variables appear in T . Furthermore, the sign of these variables can be inferred from their sign in \mathbf{x} . Hence, after seeing an evident example for all terms, one can have a list of terms containing all terms in the DNF. This list might have terms that are not part of the DNF. Yet, such terms can be thrown away later by testing if they make wrong predictions.

Proof (of Theorem 4.3) We will show that algorithm 1 learns with 1-local-queries any function that is realized by a DNF of size $\leq n^2$ with evident examples. Adapting the proof to hold with n^c instead of n^2 , for any $c > 0$, is straight-forward. First, it is easy to see that this algorithm is efficient. Now, fix a distribution \mathcal{D} and let $h^* : \{-1, 1\}^n \rightarrow \{0, 1\}$ be a hypothesis that is realized, w.r.t. \mathcal{D} , by a DNF formula $F = T_1 \vee T_2 \vee \dots \vee T_d$, $d \leq n^2$ with evident examples. Let $\epsilon > 0$, and suppose we run the algorithm on two independent samples from \mathcal{D} , denoted $S_1 = \{(\mathbf{x}_i, h^*(\mathbf{x}_i))\}_{i=1}^{m_1}$ and $S_2 = \{(\mathbf{x}'_i, h^*(\mathbf{x}'_i))\}_{i=1}^{m_2}$. We will show that if $m_1 \geq \frac{32n^3}{\epsilon} \log \frac{32n^2}{\epsilon} \geq \frac{32nd}{\epsilon} \log \frac{32d}{\epsilon}$ and $m_2 \geq \frac{32m_1}{\epsilon} \log \frac{32m_1}{\epsilon}$ then with probability of at least $\frac{3}{4}$, the algorithm will return a function \hat{h} with $L_{\mathcal{D}, h^*}(\hat{h}) < \epsilon$. Let H be the DNF formula returned by the algorithm, and let \hat{h} be the function induced by H . We have that

$$\begin{aligned} L_{\mathcal{D}, h^*}(\hat{h}) &= \Pr_{\mathbf{x} \sim \mathcal{D}} \left(h^*(\mathbf{x}) \neq \hat{h}(\mathbf{x}) \right) \\ &= \Pr_{\mathbf{x} \sim \mathcal{D}} \left(h^*(\mathbf{x}) = 1 \text{ and } \hat{h}(\mathbf{x}) = 0 \right) + \Pr_{\mathbf{x} \sim \mathcal{D}} \left(h^*(\mathbf{x}) = 0 \text{ and } \hat{h}(\mathbf{x}) = 1 \right) \end{aligned}$$

The proof will now follow from claims 2 and 3.

Claim 2 With probability at least $\frac{7}{8}$ over the choice of S_1, S_2 we have

$$\Pr_{\mathbf{x} \sim \mathcal{D}} \left(h^*(\mathbf{x}) = 1 \text{ and } \hat{h}(\mathbf{x}) = 0 \right) \leq \frac{\epsilon}{2} \quad (1)$$

Proof We first note that if there is an evident example \mathbf{x} for a term T_i in S_1 , then T_i will be in the output formula. Indeed, in the for-loop that go over the examples in S_1 , when processing the example $(\mathbf{x}, h^*(\mathbf{x}))$, it is not hard to see that T_i will be added. We furthermore claim that the term T_i won't be removed at the for loop that tests the terms collected in the first loop. Indeed, if for some $(\mathbf{x}, h^*(\mathbf{x})) \in S_2$ we have $T_i(\mathbf{x}) = 1$, it must be the case that $h^*(\mathbf{x}) = 1$. Now, we say that the term T_i is *revealed* if we see an evident example for this term in S_1 . We also denote $p_i = \Pr_{\mathbf{x} \sim \mathcal{D}} (T_i(\mathbf{x}) = 1)$. We have

$$\begin{aligned} \Pr_{\mathbf{x} \sim \mathcal{D}} \left(h^*(\mathbf{x}) = 1 \text{ and } \hat{h}(\mathbf{x}) = 0 \right) &\leq \sum_{i=1}^d \Pr_{\mathbf{x} \sim \mathcal{D}} \left(T_i(\mathbf{x}) = 1 \text{ and } \hat{h}(\mathbf{x}) = 0 \right) \\ &\leq \sum_{i: T_i \text{ is not revealed}} p_i \end{aligned}$$

Now, by the assumption that h^* is realized with evident queries, the probability (over the choice of S_1, S_2) that T_i is not revealed is at most $(1 - \frac{p_i}{n})^{m_1}$. Hence, if we denote by A_i the event that T_i is not revealed, we have

$$\begin{aligned}
\mathbb{E}_{S_1 \sim \mathcal{D}^{m_1}} \left[\Pr_{\mathbf{x} \sim \mathcal{D}} \left(h^*(\mathbf{x}) = 1 \text{ and } \hat{h}(\mathbf{x}) = 0 \right) \right] &\leq \mathbb{E}_{S_1 \sim \mathcal{D}^{m_1}} \left[\sum_{i=1}^d p_i \cdot 1_{A_i} \right] \\
&= \sum_{i=1}^d p_i \mathbb{E}_{S_1 \sim \mathcal{D}^{m_1}} [1_{A_i}] \\
&= \sum_{i=1}^d p_i \Pr_{S_1 \sim \mathcal{D}^{m_1}} [A_i] \\
&= \sum_{i=1}^d p_i \left(1 - \frac{p_i}{n} \right)^{m_1} \\
&= \sum_{i|p_i < \frac{\epsilon}{32d}} p_i \left(1 - \frac{p_i}{n} \right)^{m_1} + \sum_{i|p_i \geq \frac{\epsilon}{32d}} p_i \left(1 - \frac{p_i}{n} \right)^{m_1} \\
&\leq \sum_{i|p_i < \frac{\epsilon}{32d}} \frac{\epsilon}{32d} + \sum_{i|p_i \geq \frac{\epsilon}{32d}} \left(1 - \frac{p_i}{n} \right)^{m_1} \\
&\leq d \cdot \frac{\epsilon}{32d} + \sum_{i|p_i \geq \frac{\epsilon}{32d}} e^{-\frac{m_1 p_i}{n}} \\
&\leq \frac{\epsilon}{32} + d \cdot e^{-\frac{m_1 \epsilon}{32dn}}
\end{aligned}$$

Since $m_1 \geq \frac{32dn}{\epsilon} \log \frac{32d}{\epsilon}$ the last expression is bounded by $\frac{\epsilon}{16}$. By Markov's inequality we conclude that the probability over the choice of S_1, S_2 that (1) does not hold is less than $\frac{1}{8}$. \square

Claim 3 *With probability at least $\frac{7}{8}$ over the choice of S_1, S_2 we have*

$$\Pr_{\mathbf{x} \sim \mathcal{D}} \left(h^*(\mathbf{x}) = 0 \text{ and } \hat{h}(\mathbf{x}) = 1 \right) \leq \frac{\epsilon}{2} \tag{2}$$

Proof Let $\hat{T}_1, \dots, \hat{T}_r$ be the terms that were added to H at the first for-loop. Denote $q_i = \Pr_{\mathbf{x} \sim \mathcal{D}} \left(\hat{T}_i(\mathbf{x}) = 1 \text{ and } h^*(\mathbf{x}) = 0 \right)$. We have

$$\Pr_{\mathbf{x} \sim \mathcal{D}} \left(h^*(\mathbf{x}) = 0 \text{ and } \hat{h}(\mathbf{x}) = 1 \right) \leq \sum_{i: \hat{T}_i \text{ is not removed}} q_i$$

Now, the probability that \hat{T}_i is not removed is $(1 - q_i)^{m_2}$. Hence, using an argument similar to the one used in the proof of claim 2, and since $m_2 \geq \frac{32m_1}{\epsilon} \log \frac{32m_1}{\epsilon} \geq \frac{32r}{\epsilon} \log \frac{32r}{\epsilon}$, the claim follows. \square

4.2 A matching Lower Bound

In this section we provide evidence that the use of queries in our algorithm is crucial. We will show that the problem of learning poly-sized decision trees can be reduced to the problem of learning DNFs with evident examples. As learning decision trees is conjectured to be intractable, this reduction serves as an indication that learning DNFs with strongly evident examples is hard without membership queries. In fact, we will show that learning decision trees can be even reduced to the case that *all* positive examples are evident.

Theorem 4.4 *Learning poly-sized DNFs with evident examples is as hard as learning poly-sized decision trees.*

We denote by h_T the function induced by a decision tree T . The proof will use the following claim:

Claim 4 *There exists a mapping (a reduction) $\varphi : \{-1, 1\}^n \rightarrow \{-1, 1\}^{2n}$, that can be evaluated in $\text{poly}(n)$ time so that for every decision tree \mathcal{T} over $\{-1, 1\}^n$ there exists a DNF formula F over $\{-1, 1\}^{2n}$ such that the following holds:*

1. *The number of terms in F is upper bounded by the number of leaves in \mathcal{T}*
2. $h_{\mathcal{T}} = h_F \circ \varphi$
3. $\forall \mathbf{x}$ such that $h_{\mathcal{T}}(\mathbf{x}) = 1$, $\varphi(\mathbf{x})$ satisfies some term in F evidently.

Proof Define φ as follows:

$$\forall \mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathcal{X}_n \quad \varphi(x_1, x_2, \dots, x_n) = (x_1, x_1, x_2, x_2, \dots, x_n, x_n)$$

Now, for every tree \mathcal{T} , we will build the desired DNF formula F as follows: First we build a DNF formula F' over $\{-1, 1\}^n$. Every leaf labeled '1' in \mathcal{T} will define the following term—take the path from the root to that leaf and form the logical AND of the literals describing the path. F' will be a disjunction of these terms. Now, for every term T in F' we will define a term $\phi(T)$ over \mathcal{X}_{2n} in the following way: Let $P_T = \{i \in [n] : x_i \text{ appear in } T\}$ and $N_T = \{i \in [n] : \bar{x}_i \text{ appear in } T\}$. So

$$T = \bigwedge_{j \in P_T} x_j \bigwedge_{j \in N_T} \bar{x}_j$$

Define

$$\phi(T) = \bigwedge_{j \in P_T} x_{2j-1} \bigwedge_{j \in P_T} x_{2j} \bigwedge_{j \in N_T} \overline{x_{2j-1}} \bigwedge_{j \in N_T} \overline{x_{2j}}$$

Finally, define F to be the DNF formula over \mathcal{X}_{2n} given by

$$F = \bigvee_{T \in F'} \phi(T)$$

We will now prove that φ and F satisfy the required conditions. First, φ can be evaluated in linear time in n . Second, it is easy to see that $h_{\mathcal{T}} = h_F \circ \varphi$, and as every term in F matches one of \mathcal{T} 's leaves, the number of terms in F cannot exceed the number of leaves in \mathcal{T} . It is left to show that the third requirement holds. Let there be an \mathbf{x} such that $h_{\mathcal{T}}(\mathbf{x}) = 1$, then \mathbf{x} is matched to one and only one path from \mathcal{T} 's root to a leaf labeled '1'. From the construction of F , \mathbf{x} satisfies one and only one term in F' . Regarding the last requirement, that no coordinate change will make one term from F False and another one True, we made sure this will not happen by ‘‘doubling’’ each variable. By this construction, in order to change a term from False to True at least two coordinate must change their value. \square

We are now ready to prove theorem 4.4.

Proof [of theorem 4.4] Suppose that \mathcal{A} learns size- n DNFs with evident examples. Using the reduction from claim 4 we will build an efficient algorithm \mathcal{B} that learns size- n decision trees. For every training set

$$S = \{(\mathbf{x}_1, h^*(\mathbf{x}_1)), (\mathbf{x}_2, h^*(\mathbf{x}_2)), \dots, (\mathbf{x}_m, h^*(\mathbf{x}_m))\} \in (\mathcal{X}_n \times \{0, 1\})^m$$

we define

$$\varphi(S) := \{(\varphi(\mathbf{x}_1), h^*(\mathbf{x}_1)), (\varphi(\mathbf{x}_2), h^*(\mathbf{x}_2)), \dots, (\varphi(\mathbf{x}_m), h^*(\mathbf{x}_m))\} \in (\mathcal{X}_{2n} \times \{0, 1\})^m$$

The algorithm \mathcal{B} will work as follows: Given a training set S , \mathcal{B} will run \mathcal{A} on $\varphi(S)$, and will return $\hat{h} \circ \varphi$, where \hat{h} is the hypothesis returned by \mathcal{A} . Since φ can be evaluated in $\text{poly}(n)$ time and \mathcal{A} is efficient, \mathcal{B} is also efficient. We will prove that \mathcal{B} learns size- n trees. Since \mathcal{A} learns size- n DNFs with evident examples, there exists a function $m_{\mathcal{A}}(n, \epsilon) \leq \text{poly}(n, \frac{1}{\epsilon})$, such that if \mathcal{A} is given a training sequence

$$S = \{(\mathbf{x}_1, h^*(\mathbf{x}_1)), (\mathbf{x}_2, h^*(\mathbf{x}_2)), \dots, (\mathbf{x}_m, h^*(\mathbf{x}_m))\} \in (\mathcal{X}_n \times \{0, 1\})^m$$

where the \mathbf{x}_i 's are sampled i.i.d. from a distribution \mathcal{D} , h^* is realized by a poly-sized DNF with evident examples, and $m \geq m_{\mathcal{A}}(n, \epsilon)$, then with probability of at least $\frac{3}{4}$ (over the choice of S), the output \hat{h} of \mathcal{A} satisfies $L_{\mathcal{D}, h^*}(\hat{h}) \leq \epsilon$. Let \mathcal{D} be a distribution on \mathcal{X}_n and let $h_{\mathcal{T}}$ be a hypothesis that can be realized by a tree with $\leq n$ leafs. Define a distribution $\tilde{\mathcal{D}}$ on \mathcal{X}_{2n} by,

$$(\tilde{\mathcal{D}})(\mathbf{z}) = \begin{cases} \mathcal{D}(\mathbf{x}) & \text{if } \exists \mathbf{x} \in \mathcal{X}_n \text{ such that } \mathbf{z} = \varphi(\mathbf{x}) \\ 0 & \text{otherwise} \end{cases}$$

Now, since $h_{\mathcal{T}}$ is realized by \mathcal{T} , from the conditions that φ satisfies, we get that $h_{\mathcal{T}} = h \circ \varphi$, where h is realized by a DNF of size $\leq n$ with evident examples w.r.t. $\tilde{\mathcal{D}}$. Now if $S \in (\mathcal{X}_n \times \{0, 1\})^m$ is an i.i.d. sample with $m \geq m_{\mathcal{A}}(2n, \epsilon)$ we have that with probability of at least $\frac{3}{4}$ it holds that

$$\begin{aligned}
L_{\mathcal{D}, h_{\mathcal{T}}}(\mathcal{B}(S)) &= L_{\mathcal{D}, h_{\mathcal{T}}}(\hat{h} \circ \varphi) \\
&= \Pr_{\mathbf{x} \sim \mathcal{D}}[h_{\mathcal{T}}(\mathbf{x}) \neq \hat{h} \circ \varphi(\mathbf{x})] \\
&= \Pr_{\mathbf{x} \sim \mathcal{D}}[h \circ \varphi(\mathbf{x}) \neq \hat{h} \circ \varphi(\mathbf{x})] \\
&= \Pr_{\mathbf{z} \sim \tilde{\mathcal{D}}}[h(\mathbf{z}) \neq \hat{h}(\mathbf{z})] \\
&= L_{\tilde{\mathcal{D}}, h}(\hat{h}) \\
&= L_{\tilde{\mathcal{D}}, h}(\mathcal{A}(\varphi(S))) < \epsilon
\end{aligned}$$

□

5 Lower Bounds

We first present a general technique to prove lower bounds on learning with local queries. For $A \subset \mathcal{X}_n$ and $q > 0$ denote

$$B(A, q) = \{\mathbf{x} \in \mathcal{X}_n \mid \exists \mathbf{a} \in A, d(\mathbf{x}, \mathbf{a}) \leq q\}$$

We say that a mapping $\varphi : \{-1, 1\}^n \rightarrow \{-1, 1\}^{n'}$ is a q -reduction of type A from a class \mathcal{H} to a class \mathcal{H}' if the following holds:

1. φ is efficiently computable.
2. For every $h \in \mathcal{H}$ there is $h' \in \mathcal{H}'$ such that

(a) $h = h' \circ \varphi$

(b) The restriction of h' to $B(\varphi(\mathcal{X}_n), q) \setminus \varphi(\mathcal{X}_n)$ is the constant function 1.

We say that a mapping $\varphi : \{-1, 1\}^n \rightarrow \{-1, 1\}^{n'}$ is a q -reduction of type B if the same requirements hold, except that (2b) is replaced by the following condition: For every $\mathbf{z} \in B(\varphi(\mathcal{X}_n), q)$ there is a unique $\mathbf{x} \in \mathcal{X}_n$ satisfying $d(\mathbf{z}, \varphi(\mathbf{x})) \leq q$, and furthermore, $h'(\mathbf{z}) = h(\mathbf{x})$.

Lemma 5.1 *Suppose that there is a q -reduction φ from \mathcal{H} to \mathcal{H}' . Then, learning \mathcal{H}' with q -local queries is as hard as PAC learning \mathcal{H} .*

Proof (sketch) Suppose that \mathcal{A}' learns \mathcal{H}' with q -local queries. We need to show that there is an algorithm that learns \mathcal{H} . Indeed, by an argument similar to the one in theorem 4.4, it is not hard to verify that the following algorithm learns \mathcal{H} . Given a sample $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\} \subset \{-1, 1\}^n \times \{0, 1\}$, run \mathcal{A}' on the sample $\{(\varphi(\mathbf{x}_1), y_1), \dots, (\varphi(\mathbf{x}_m), y_m)\} \subset \{-1, 1\}^{n'} \times \{0, 1\}$. Whenever \mathcal{A}' makes a q -local query for $\mathbf{z} \in \{-1, 1\}^{n'}$, respond 1 if φ is of type A. If φ is of type B, respond y_i , where \mathbf{x}_i is the unique training sample satisfying $d(\mathbf{z}, \varphi(\mathbf{x}_i)) \leq q$. Finally, if \mathcal{A}' returned the hypothesis h' , return $h' \circ \varphi$. □

We next use Lemma 5.1 to prove that for several classes, local queries are not useful. Namely, if the class is learnable with local queries then it is also learnable without queries. We will use the following terminology. We say that q -local queries cannot help to learn a class \mathcal{H} if \mathcal{H} is learnable if and only if it is learnable with q -local queries.

Corollary 5.2 *For every $\epsilon_0 > 0$, $(n^{1-\epsilon_0})$ -local queries cannot help to learn poly-sized DNFs, intersection of halfspaces, decision lists, depth- d circuits for any $d = d(n) \geq 2$, and depth- d threshold circuits for any $d = d(n) \geq 2$.*

Proof We will only prove the corollary for DNFs. The proofs for the remaining classes are similar. Also, for simplicity, we will assume that $\epsilon_0 = \frac{1}{3}$. Consider the mapping $\varphi : \{-1, 1\}^n \rightarrow \{-1, 1\}^{n^3}$ that replicates each coordinate n^2 times. To establish the corollary, we show that φ is an $\left((n')^{\frac{2}{3}} - 1\right)$ -reduction of type A from poly-sized DNF to poly-sized DNF.

Indeed, let $F = T_1 \vee \dots \vee T_d$ be a DNF formula on n variables, consider the following formula on the n^3 variables $\{x_{i,j}\}_{1 \leq i \leq n, 1 \leq j \leq n^2}$. Let,

$$\begin{aligned} T'_t(\{x_{i,j}\}_{i,j}) &= T_t(x_{1,1}, \dots, x_{n,1}) \\ G'(\{x_{i,j}\}_{i,j}) &= \bigvee_{i=1}^n \bigvee_{j=1}^{n^2-1} (x_{i,j} \wedge \neg x_{i,j+1}) \vee (x_{i,j+1} \wedge \neg x_{i,j}) \\ F' &= (T'_1 \vee \dots \vee T'_d) \vee G' \end{aligned}$$

It is not hard to verify that $h_F = h_{F'} \circ \varphi$. Moreover, if $\mathbf{x} = \{x_{i,j}\}_{i,j} \in B(\varphi(\mathcal{X}_n), n^2 - 1) \setminus \varphi(\mathcal{X}_n)$ then $x_{i,j} \neq x_{i,j+1}$ for some i, j , meaning that $h_G(x) = 1$ and therefore also $h_F(x) = 1$. \square

Corollary 5.3 *For every $\epsilon_0 > 0$, $(n^{1-\epsilon_0})$ -local queries cannot help to learn poly-sized automata.*

Proof Again, for simplicity, we assume that $\epsilon_0 = \frac{1}{3}$, and consider the mapping $\varphi : \{-1, 1\}^n \rightarrow \{-1, 1\}^{n^3}$ that replicates each coordinate n^2 times. To establish the corollary, we show that φ is an $\left((n')^{\frac{2}{3}} - 1\right)$ -reduction of type A from poly-sized automata to poly-sized automata.

Indeed, let A be an automaton on n variables. It is enough to show that there is an automaton A' on the n^3 variables $\{x_{i,j}\}_{1 \leq i \leq n, 1 \leq j \leq n^2}$ such that (i) the size of A' is polynomial, (ii) A' accepts any string in $B(\varphi(\mathcal{X}_n), n^2 - 1) \setminus \varphi(\mathcal{X}_n)$, and (iii) A' accepts $\varphi(x)$ if and only if A accepts x . Now, by the product construction of automata [19], if A'_1, A'_2 are automata that induce the functions $h_{A'_1}, h_{A'_2} : \{-1, 1\}^{n^3} \rightarrow \{0, 1\}$, then the function $h_{A'_1} \vee h_{A'_2}$ can be induced by an automaton of size $|A'_1| \cdot |A'_2|$. Hence, it is enough to show that there are poly-sized automata A'_1, A'_2 that satisfies (ii) and (iii) respectively.

A construction of a size $O(n^2)$ automaton that satisfies (ii) is a simple exercise. We next explain how to construct a poly-sized automaton A'_2 satisfying (iii). The states of A'_2 will be the $S(A) \times [n^2]$ (here, $S(A)$ denotes the set of states of A). The start state of A'_2 will be $(\alpha_0, 1)$, where α_0 is the start state of A , and the accept states of A'_2 will be the cartesian product of the accept states of A with $[n^2]$. Finally if $\delta : S(A) \times \{-1, 1\} \rightarrow S(A)$ is the transition function of A , then the transition function of A'_2 is defined by

$$\delta'((\alpha, i), b) = \begin{cases} (\alpha, i + 1) & 1 \leq i < n^2 \\ (\delta(\alpha, b), 1) & i = n^2 \end{cases}$$

It is not hard to verify that A'_2 satisfies (iii). \square

In the next corollary, a *Junta* of size t is a function $h : \{-1, 1\}^n \rightarrow \{0, 1\}$ that depends on $\leq \log(t)$ variables. The rationale behind this definition of size, is the fact that in order to describe a general function that depends on K variables, at least 2^K bits are required. Likewise, the sample complexity of learning Juntas is proportional to t rather than $\log(t)$.

Corollary 5.4 *For every constant $q_0 > 0$, q_0 -local queries cannot help to learn poly-sized Juntas.*

Proof Consider the mapping $\varphi : \{-1, 1\}^n \rightarrow \{-1, 1\}^{(2q_0+1)n}$ that replicates each coordinate $2q_0 + 1$ times. To establish the corollary, we show that φ is a q_0 -reduction of type B from poly-sized Juntas to poly-sized Juntas. Indeed, let $h : \{-1, 1\}^n \rightarrow \{0, 1\}$ be a function that depends on K variables. It is enough to show that there is a function h' on the variables $\{z_{i,j}\}_{i \in [n], j \in [2q_0+1]}$ that satisfies (i) $h = h' \circ \varphi$, (ii) for every $\mathbf{x} \in \mathcal{X}$, if $\mathbf{z} \in \{-1, 1\}^{(2q_0+1)n}$ is obtained from $\varphi(\mathbf{x})$ by modifying $\leq q_0$ coordinates, then $h'(\mathbf{z}) = h'(\varphi(\mathbf{x}))$ and (iii) h' depends on $(2q_0 + 1)K$ variables. It is not hard to check that the following function satisfies these requirements:

$$h'(\mathbf{z}) = h(\text{MAJ}(z_{1,1}, \dots, z_{1,2q_0+1}), \dots, \text{MAJ}(z_{n,1}, \dots, z_{n,2q_0+1}))$$

\square

We remark that by taking $q_0 = \log^{1-\epsilon_0}(n)$ for some $\epsilon_0 > 0$, and using a similar argument, it can be shown that an efficient algorithm for learning poly-sized Juntas with $(\log^{1-\epsilon_0}(n))$ -local queries would imply a PAC algorithm for poly-sized Juntas that runs in time $n^{O(\log^{1-\epsilon_0}(n))}$.

Corollary 5.5 *For every constant $q_0 > 0$, q_0 -local queries cannot help to learn poly-sized decision trees.*

Proof As with Juntas, consider the mapping $\varphi : \{-1, 1\}^n \rightarrow \{-1, 1\}^{(2q_0+1)n}$ that replicates each coordinate $2q_0 + 1$ times. To establish the corollary, we show that φ is a q_0 -reduction of type B from poly-sized decision trees to poly-sized decision trees. Indeed, let $h : \{-1, 1\}^n \rightarrow \{0, 1\}$ be a function that is realized by a decision tree T . It is enough to show that there is a function h' on the variables $\{z_{i,j}\}_{i \in [n], j \in [2q_0+1]}$ that satisfies (i) $h = h' \circ \varphi$, (ii) for every $\mathbf{x} \in \mathcal{X}$, if $\mathbf{z} \in \{-1, 1\}^{(2q_0+1)n}$ is obtained from $\varphi(\mathbf{x})$ by modifying $\leq q_0$ coordinates, then $h'(\mathbf{z}) = h'(\varphi(\mathbf{x}))$ and (iii) h' can be realized by a tree of size $|T|^{2q_0+1}$. As we explain, this will hold for the following function:

$$h'(\mathbf{z}) = \text{MAJ}(h(z_{1,1}, \dots, z_{n,1}), \dots, h(z_{1,2q_0+1}, \dots, z_{n,2q_0+1}))$$

It is not hard to check that h' satisfies (i) and (ii). As for (iii), consider the following tree T' . First replicate T on the variables $z_{1,1}, \dots, z_{n,1}$. Then, on the obtained tree, replace each leaf by a replica of T on the variables $z_{1,2}, \dots, z_{n,2}$. Then, again, replace each leaf by a replica of T on the variables $z_{1,3}, \dots, z_{n,3}$. Continue doing so, until the leaves are replaced by a replica on the variables $z_{1,2q_0+1}, \dots, z_{n,2q_0+1}$. Now, each leaf in the resulted tree corresponds to $(2q_0 + 1)$ root-to-leaf paths in the original tree T . The label of such leaf will be the majority of the labels of these paths. \square

As with Juntas, we remark that by taking $q_0 = \log^{1-\epsilon_0}(n)$ for some $\epsilon_0 > 0$, and using a similar argument, it can be shown that an efficient algorithm for learning poly-sized trees with $(\log^{1-\epsilon_0}(n))$ -local queries would imply a PAC algorithm for poly-sized trees that runs in time $n^{O(\log^{1-\epsilon_0}(n))}$.

In the sequel, we say that a function $h : \{-1, 1\}^n \rightarrow \{0, 1\}$ is realized by a poly-sized polynomial, if it is the function induced by a polynomial with polynomially many non-zero coefficient and degree $O(\log(n))$. A similar convention applies to the term poly-sized polynomial threshold function. We remark that the following corollary and its proof remain correct also if in our definition, we replace the number of coefficients with the ℓ^1 norm of the coefficients vector.

Corollary 5.6 *For every constant $q_0 > 0$, q_0 -local queries cannot help to learn poly-sized polynomials, as well as poly-sized polynomial threshold functions.*

Proof We will prove the corollary for polynomials. The proof for polynomial threshold functions is similar. Consider the mapping $\varphi : \{-1, 1\}^n \rightarrow \{-1, 1\}^{(2q_0+1)n}$ that replicates each coordinate $2q_0 + 1$ times. To establish the corollary, we show that φ is a q_0 -reduction of type B from poly-sized polynomials to poly-sized polynomials. Indeed, let $h : \{-1, 1\}^n \rightarrow \{0, 1\}$ be a poly-sized polynomial. It is enough to show that there is a poly-sized polynomial h' on the variables $\{z_{i,j}\}_{i \in [n], j \in [2q_0+1]}$ that satisfies (i) $h = h' \circ \varphi$, and (ii) for every $\mathbf{x} \in \mathcal{X}$, if $\mathbf{z} \in \{-1, 1\}^{(2q_0+1)n}$ is obtained from $\varphi(\mathbf{x})$ by modifying $\leq q_0$ coordinates, then $h'(\mathbf{z}) = h'(\varphi(\mathbf{x}))$. As we explain next, this will hold for the following polynomial:

$$h'(\mathbf{z}) = h(\text{MAJ}(z_{1,1}, \dots, z_{1,2q_0+1}), \dots, \text{MAJ}(z_{n,1}, \dots, z_{n,2q_0+1}))$$

It is not hard to check that h' satisfies (i) and (ii). It remains to show that h' is poly-sized. Indeed, the majority function on $2q_0 + 1$ coordinates is a polynomial of degree $\leq 2q_0 + 1$ with at most 2^{2q_0+1} non zero coefficients (this is true for any function on $2q_0 + 1$ coordinates). Hence, if we replace each variable in h by a polynomial of the form $\text{MAJ}(z_{i,1}, \dots, z_{i,2q_0+1})$, the degree is multiplied by at most $(2q_0 + 1)$, while the number of non zero coefficients is multiplied by at most 2^{2q_0+1} . \square

As with Juntas and decision trees, by taking $q_0 = \log^{1-\epsilon_0}(n)$ for some $\epsilon_0 > 0$, and using a similar argument, it can be shown that an efficient algorithm for learning poly-sized polynomials or polynomial threshold functions with $(\log^{1-\epsilon_0}(n))$ -local queries would imply a PAC algorithm for the same problem that runs in time $n^{O(\log^{1-\epsilon_0}(n))}$.

6 Conclusion and Future Work

We have shown that in the distribution free setting, for many hypothesis classes, local queries are not useful. As our proofs show, this stems from the fact that learning these classes without queries can be reduced to a case where local queries are pointless, in the sense that the answer to them is either always 1, or the label of the closest training example. On the other hand, the learning problem of DNFs with evident examples circumvents this property. Indeed, the underlying assumption enforces that local changes can change the label in a

non-trivial manner. While this assumption might be intuitive in some cases, it is certainly very restrictive. Therefore, a natural future research direction is to seek less restrictive assumptions, that still possess this property.

More concrete directions arising from our work concern classes for which we have shown that $(\log^{0.99}(n))$ -local queries are unlikely to lead to efficient algorithms. We conjecture that for some $a > 0$ even (n^a) -local queries won't lead to efficient distribution free algorithms for these classes.

References

- [1] Dana Angluin. Learning regular sets from queries and counterexamples. *Information and computation*, 75(2):87–106, 1987.
- [2] Dana Angluin and Michael Kharitonov. When won't membership queries help? *Journal of Computer and System Sciences*, 50(2):336–355, 1995.
- [3] Dana Angluin and Donna K Slonim. Randomly fallible teachers: Learning monotone dnf with an incomplete membership oracle. *Machine Learning*, 14(1):7–26, 1994.
- [4] Dana Angluin, Mārtiņš Kriķis, Robert H Sloan, and György Turán. Malicious omissions and errors in answers to membership queries. *Machine Learning*, 28(2-3):211–255, 1997.
- [5] Pranjal Awasthi, Vitaly Feldman, and Varun Kanade. Learning using local membership queries. *arXiv preprint arXiv:1211.0996*, 2012.
- [6] Galit Bary. Learning using 1-local membership queries. *arXiv preprint arXiv:1512.00165*, 2015.
- [7] Eric B Baum. Neural net algorithms that learn in polynomial time from examples and queries. *Neural Networks, IEEE Transactions on*, 2(1):5–19, 1991.
- [8] Eric B Baum and Kenneth Lang. Query learning can work poorly when a human oracle is used. In *International Joint Conference on Neural Networks*, volume 8, 1992.
- [9] Laurence Bisht, Nader H Bshouty, and Lawrence Khoury. Learning with errors in answers to membership queries. *Journal of Computer and System Sciences*, 74(1):2–15, 2008.
- [10] Avrim Blum and Steven Rudich. Fast learning of k-term dnf formulas with queries. In *Proceedings of the twenty-fourth annual ACM symposium on Theory of computing*, pages 382–389. ACM, 1992.
- [11] Avrim Blum, Prasad Chalasanani, Sally A Goldman, and Donna K Slonim. Learning with unreliable boundary queries. In *Proceedings of the eighth annual conference on Computational learning theory*, pages 98–107. ACM, 1995.
- [12] Nader H Bshouty. Exact learning boolean functions via the monotone theory. *Information and Computation*, 123(1):146–153, 1995.

- [13] Amit Daniely and Shai Shalev-Shwartz. Complexity theoretic limitations on learning dnf's. *arXiv preprint arXiv:1404.3378*, 2014.
- [14] Vitaly Feldman. On the power of membership queries in agnostic learning. *The Journal of Machine Learning Research*, 10:163–182, 2009.
- [15] Yoav Freund. Boosting a weak learning algorithm by majority. *Information and computation*, 121(2):256–285, 1995.
- [16] Jeffrey Jackson. An efficient membership-query algorithm for learning dnf with respect to the uniform distribution. In *Foundations of Computer Science, 1994 Proceedings., 35th Annual Symposium on*, pages 42–53. IEEE, 1994.
- [17] Michael Kearns and Leslie Valiant. Cryptographic limitations on learning boolean formulae and finite automata. *Journal of the ACM (JACM)*, 41(1):67–95, 1994.
- [18] Eyal Kushilevitz and Yishay Mansour. Learning decision trees using the fourier spectrum. *SIAM Journal on Computing*, 22(6):1331–1348, 1993.
- [19] Michael Sipser. *Introduction to the Theory of Computation*, volume 2.
- [20] Robert H Sloan and György Turán. Learning with queries but incomplete information. In *Proceedings of the seventh annual conference on Computational learning theory*, pages 237–245. ACM, 1994.
- [21] Leslie G Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.