

We consider only the Probably Approximately Correct (PAC) model under the uniform distribution. A learning problem is given by a concept class C of functions $f: \{0,1\}^n \rightarrow \{-1,1\}$.

A learning algorithm is (typically) randomized algorithm with:

- Input: accuracy parameter $\epsilon > 0$.
- Either: ① Random examples $(x, f(x))$ for x chosen uniformly, or
② Membership queries to f . (where $f \in C$)
- Output: a function $h: \{0,1\}^n \rightarrow \{-1,1\}$, known as the hypothesis, described by a circuit.
- Goal: h should be ϵ -close to f with prob. 99% .

- Remarks:
- General-PAC learning requires the algorithm to work under any dist. This affects both choice of random samples and the closeness of h and f . This typically very hard.
 - Often one requires h to be of the same "form" as C . This is called proper learning and is also typically much harder.
 - One can always check if the hypothesis h is ϵ -close to f : simply compare them on enough random examples. This also means that one can always boost the 99% correctness to $1-\delta$ by paying an extra $O(\log 1/\delta)$ factor.

Learning Using Spectral Concentration

Def: For a family \mathcal{S} of subsets of $[n]$, we say that f is ϵ -concentrated on \mathcal{S} if $\sum_{S \in \mathcal{S}} \hat{f}(S)^2 \leq \epsilon$.

Claim: If $f: \{0,1\}^n \rightarrow \{-1,1\}$ is ϵ -concentrated on \mathcal{S} then $g: \{0,1\}^n \rightarrow \mathbb{R}$ given by $g = \sum_{S \in \mathcal{S}} \hat{f}(S) \chi_S$ satisfies $\|f - g\|_2 \leq \epsilon$.

Proof: $f - g = \sum_{S \in \mathcal{S}} \hat{f}(S) \chi_S$ and hence, $\|f - g\|_2^2 = \sum_{S \in \mathcal{S}} \hat{f}(S)^2 \leq \epsilon$. \square

Claim: Let $f: \{0,1\}^n \rightarrow \{-1,1\}$ and $g: \{0,1\}^n \rightarrow \mathbb{R}$ satisfy $\|f - g\|_2^2 \leq \epsilon$. Then, $h: \{0,1\}^n \rightarrow \{-1,1\}$ given by $h(x) = \text{sign}(g(x))$ is ϵ -close to f .

Proof: For each x such that $f(x) \neq h(x)$ we must have $(f(x) - g(x))^2 \geq 1$. Now use $\|f - g\|_2^2 = \mathbb{E}_x [(f(x) - g(x))^2]$. \square

Thm: [Linial, Mansour, Nisan, 89] Suppose we know a set S on which f is $\epsilon/2$ -concentrated. Then, in time $\text{poly}(|S|, 1/\epsilon, n)$ using random examples, we can output a function h that is ϵ -close to f w.p. $\geq 99\%$.

Proof: The algorithm - for each $s \in S$ estimate $\hat{f}(s)$ to within $\pm \sqrt{\frac{\epsilon}{16|S|}}$ and let \tilde{f}_s be the estimate. Let $\tilde{g} = \sum_{s \in S} \tilde{f}_s \cdot \chi_s$. Now output $h = \text{sign}(\tilde{g})$.

This takes $\text{poly}(|S|, n, 1/\epsilon)$ time. By the previous claim, it suffices to show that $\|f - \tilde{g}\|_2 \leq \epsilon$. Define $g = \sum_{s \in S} \hat{f}(s) \chi_s$. Then, $\|f - \tilde{g}\|_2 \leq \|f - g\|_2 + \|g - \tilde{g}\|_2 \leq \sqrt{\epsilon/2} + \sqrt{\sum_{s \in S} (\hat{f}(s) - \tilde{f}_s)^2} \leq \sqrt{\epsilon/2} + \sqrt{|S| \cdot \frac{\epsilon}{16|S|}} \leq \sqrt{\epsilon}$. \square

Cor: (The low degree algorithm [LMN]) If all $f \in C$ is $\epsilon/2$ -concentrated on $S = \{s : |s| \leq d\}$ then C can be learned in time $n^{O(d)} \cdot \text{poly}(1/\epsilon)$ using random examples.

Cor [Kushilevitz & Mansour '91]: If any $f \in C$ is $\epsilon/4$ -concentrated on some set of size $\leq M$ then C can be learned in time $\text{poly}(M, 1/\epsilon, n)$ using membership queries.

Proof: Assume we are given some $f \in C$ and let S' be a set of size $\leq M$ on which it is $\epsilon/4$ -concentrated. Let $S' = \{s \in S \mid \hat{f}(s)^2 \geq \frac{\epsilon}{4M}\}$. Then f is $\epsilon/2$ -concentrated on S' . Using the G-L algorithm we can find a set L containing all S s.t. $\hat{f}(s)^2 \geq \frac{\epsilon}{4M}$ in time $\text{poly}(M, 1/\epsilon, n)$. In particular, $S' \subseteq L$ so f is $\frac{\epsilon}{2}$ -concentrated on L .

Now apply [LMN] to find h that is ϵ -close to f . \square

Thm: The class $C = \{f : I(f) \leq d\}$ is learnable in time $n^{O(d/\epsilon)}$ using random examples.

Proof: Recall that $I(f) = \sum_s |s| \cdot \hat{f}(s)^2$. Notice that $I(f) \leq d$ implies that $\sum_{|s| \geq d/\epsilon} \hat{f}(s)^2 \leq \epsilon$. Hence, we can use the Low Degree Alg. with $S = \{s : |s| \leq d/\epsilon\}$.

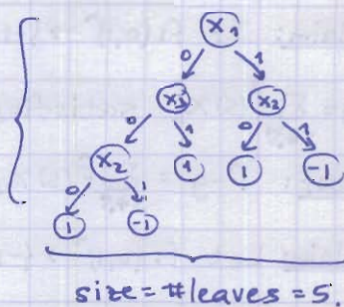
Learning Decision Trees

Remarks: We assume w.l.o.g. that on any path each variable appears at most once.

Notice that the decision trees are universal,

i.e., any Boolean function can be written as a decision tree (of depth $\leq n$ and size $\leq 2^n$).

A decision tree of depth d has size $\leq 2^d$.



Prop: Assume $f: \{0,1\}^n \rightarrow \{-1,1\}$ is computable by a depth- d decision tree. Then,

1. f is of degree $\leq d$ (i.e., $\sum_{|S|>d} \hat{f}(s)^2 = 0$).
2. All Fourier coeff. of f are integer multiples of 2^{-d} .
3. The number of nonzero coeff. is at most 4^d .

Proof: For each path P in the decision tree, let 1_P be the indicator function of P .

Then, we can write $f = \sum_P \hat{f}(P) \cdot 1_P$ because the paths define a partition of $\{0,1\}^n$.

Because 1_P is the AND of at most d literals, its Fourier coeff. are multiples of 2^{-d} ,

and is of degree at most d . (Since, e.g., $\text{AND}(x_1, \bar{x}_3, x_2) = \left(\frac{1-x_{11}}{2}\right) \left(\frac{1+x_{31}}{2}\right) \left(\frac{1-x_{12}}{2}\right)$).

This proves 1 & 2. 3 follows from 2. \square

Cor: 1. Depth- d DTs are exactly learnable (i.e., with $\epsilon=0$) in time $n^{O(d)}$ using random examples.

2. Depth- d DTs are exactly learnable in time $\text{poly}(n, 2^d)$ using membership queries.

3. Depth- $O(\log n)$ DTs are exactly learnable in time $\text{poly}(n)$ using membership queries.

Proof: Estimate all nonzero Fourier coeff. to within $\pm 2^{-d}/4$ and round to nearest multiple of 2^{-d} . \square

Observation: Any DT of size L is ϵ -close to a DT of depth $\log(4/\epsilon)$.

Proof: Cut everything deeper than $\log(4/\epsilon)$. The resulting DT differs from the original on a random input x w.p. $\leq L \cdot 2^{-\log(4/\epsilon)} = \epsilon$. \square

Cor: DTs of size L are 4ϵ -concentrated on a set of size $4^{\log(4/\epsilon)} = (4/\epsilon)^2$ of Fourier coeff. of level $\leq \log(4/\epsilon)$.

Hence, DTs of size L can be learned in time $\text{poly}(L, n, 1/\epsilon)$ using queries and time $n^{O(\log(4/\epsilon))}$ using random examples.

Remark: These are the best known results. It is an open question to do poly-size DT in time $\text{poly}(n)$ using random examples.

Learning DNFs

Def: A DNF is a disjunction $\phi = T_1 \vee T_2 \vee \dots \vee T_L$ where each term is a conjunction of literals (e.g., $(x_{11} \bar{x}_5) \vee (x_{31} x_{41} x_5)$). The size of a DNF is the number of terms L , and the width of a DNF is the maximal # of literals in a term.

Remark: Any DT of size L can be written as a DNF of size L . In other words,

DNF-size(f) \leq DT-size(f). Similarly, DNF-width(f) \leq DT-depth(f).

As we saw in homework, a width w DNF f has $I(f) \leq 2w$ and hence is ϵ -concentrated on levels $\leq \frac{2w}{\epsilon}$. Our goal is to improve this to $O(w \cdot \log 1/\epsilon)$.

Thm: If f is computable by a width- w DNF then for any $d \geq 5$,

$$\sum_{|S| \geq 20dw} \hat{f}(S)^2 \leq 2^{-d+1}$$

Def: A random p -restriction is a pair (I, x) where $I \subseteq [n]$ is chosen by including each coordinate w.p. p and $x \in \{0,1\}^{\bar{I}}$ is chosen uniformly.

Thm [Hastad's switching lemma '86]: If f is computable by a width- w DNF and (I, x) is a random p -restriction, then $\Pr[\text{DT-depth}(f_{x \rightarrow \bar{I}}) > d] \leq (5pw)^d$.

Example: $p = \frac{1}{10w}$, then we get that w.p. $\geq 3/4$ the restriction has DT of depth ≤ 2 .

Claim: For $I \subseteq [n]$, and $x \in \{0,1\}^{\bar{I}}$, $\widehat{f_{x \rightarrow \bar{I}}}(S) = \sum_{T \subseteq \bar{I}} \hat{f}(S \cup T) \chi_T(x)$

Example: $n=4$

If we choose $x \in \{0,1\}^{\bar{I}}$ uniformly, then $\mathbb{E}_x[\widehat{f_{x \rightarrow \bar{I}}}(S)] \geq |\mathbb{E}_x[\widehat{f_{x \rightarrow \bar{I}}}(S)]| = |\hat{f}(S)|$.

$$\begin{aligned} \text{Also, } \mathbb{E}_x[\widehat{f_{x \rightarrow \bar{I}}}(S)^2] &= \sum_{T, T' \subseteq \bar{I}} \hat{f}(S \cup T) \cdot \hat{f}(S \cup T') \cdot \mathbb{E}_x[\chi_T(x) \cdot \chi_{T'}(x)] \\ &= \sum_{T \subseteq \bar{I}} \hat{f}(S \cup T)^2 \end{aligned}$$

by chin

Proof: Let (I, x) be a random restriction with $p = \frac{1}{10w}$. By Hastad's

switching lemma, $f_{x \rightarrow \bar{I}}$ has a decision tree of depth $\leq d$ w.p. $\geq 1 - 2^{-d}$.

In such a case, $\sum_{\substack{S \subseteq I \\ |S| > d}} \widehat{f_{x \rightarrow \bar{I}}}(S)^2 = 0$. Therefore, $2^{-d} \geq \mathbb{E}_{x, I} \left[\sum_{\substack{S \subseteq I \\ |S| > d}} \widehat{f_{x \rightarrow \bar{I}}}(S)^2 \right] =$

$$= \mathbb{E}_I \left[\sum_{\substack{S \subseteq I \\ |S| > d}} \mathbb{E}_x[\widehat{f_{x \rightarrow \bar{I}}}(S)^2] \right] = \mathbb{E}_I \left[\sum_{\substack{S \subseteq I \\ |S| > d}} \sum_{T \subseteq \bar{I}} \hat{f}(S \cup T)^2 \right] = \sum_{U \subseteq [n]} \Pr[|U \cap I| > d] \cdot \hat{f}(U)^2$$

For $|U| \geq 20dw$, we get by Chernoff that $\Pr_I[|U \cap I| > d] \geq 1/2$.

(because $|U \cap I|$ is distributed like Binom($\geq 20dw, \frac{1}{10w}$)).

Hence, $\sum_{|U| \geq 20dw} \hat{f}(U)^2 \leq 2^{-d+1}$ ▣

Thm: If f has width- w DNF, then $\sum_U \left(\frac{1}{20w}\right)^{|U|} \cdot |\hat{f}(U)| \leq 2$.

Solutions - Homework 2

5. (a). $f(x) = \sum_S \hat{f}(s) \chi_S(x) = E_{S \sim D} [\chi_S(x)]$, where D is the distribution on S given by $\Pr[S] = \hat{f}(S)$. By Chernoff, for each fixed x , if we choose $S_1, S_2, \dots, S_{c/n}$ for some large enough c , then $|f(x) - \frac{1}{c/n} \sum_{i=1}^{c/n} \chi_{S_i}(x)| < 0.01$ w.p. $\geq 1 - 2^{-c/n}$.

Therefore, by union bound, this sample is good for all x simultaneously w.p. $\geq 1/2$.

In particular, such sample exists.

(b). Write $f = f^+ - f^-$ where $f^+ = \sum_{\hat{f}(s) \geq 0} \hat{f}(s) \chi_S$ and similarly for f^- . Now apply (a) separately to $\frac{f^+}{\|f^+\|_1}$ and $\frac{f^-}{\|f^-\|_1}$ with accuracies $\frac{0.01}{2\|f^+\|_1}$ and $\frac{0.01}{2\|f^-\|_1}$.

1. Consider $f = \chi_{\{1,2,3\}} = \chi_{\{1\}} \cdot \chi_{\{2\}} \cdot \chi_{\{3\}}$. f is $1/2$ -far from dictators. We claim that the test must accept f w.p. 1 : Assume the test checks that $f(x) \cdot f(y) \cdot f(z) = -1$; $f(x)f(y)f(z) = \chi_{\{1\}}(x) \chi_{\{2\}}(x) \chi_{\{3\}}(x) \cdot \dots \cdot \chi_{\{3\}}(z) = (-1)(-1)(-1) = -1$.

So test must accept f . Similarly, for $f(x) \cdot f(y) \cdot f(z) = 1$.

7.



(a). If $\text{val}(G) \geq 1 - \lambda$ then consider the assignment $f_v = \chi_{\{L(v)\}}$, where L is the assignment to the unique label cover. w.p. $\geq 1 - \lambda$, the constraint $\{u, v\} \in E$ chosen by the tester is s.t. $\sigma_{u \rightarrow v}(L(u)) = L(v)$. In such a case the tester accepts w.p. $1 - \delta$. So overall acceptance prob. is $(1 - \lambda)(1 - \delta) \geq 1 - \lambda\delta$.

(b). Consider the assignment giving some dictators to each f_v . Then the test applies Hastad's to the average of two dictators, $\chi_{\{i\}}$ and $\chi_{\{j\}}$.

If $i=j$ then the test accepts w.p. $1 - \delta$.

If $i \neq j$ the test accepts w.p. $\frac{5}{8} - \frac{\delta}{4}$ (since the acceptance prob. is $\frac{1}{2} + \frac{1}{2} \sum_S (1 - 2\delta)^{|S|} \hat{f}(S)^2$)

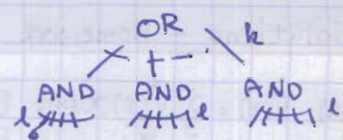
2. (a). A function is 1 -resilient if it is balanced and remains balanced after fixing any one coordinate.

(b). Consider the test that takes x uniformly and $w \sim \mu_\delta$ for $\delta = \frac{1-p}{2}$ and checks that $f(x) \cdot f(x \oplus w) = 1$. Its acceptance prob. is, $\frac{1}{2} + \frac{1}{2} \sum_S p^{|S|} \cdot \hat{f}(S)^2$

If f is 1 -resilient then this is $\leq \frac{1}{2} + \frac{1}{2} p^2$. If $\exists S, |S| \leq 1$ and $|\hat{f}(S)| \geq \epsilon$ then this is $\geq \frac{1}{2} + \frac{1}{2} p \epsilon^2$. Choosing $p = \frac{\epsilon^2}{2}$ makes the first case $\leq \frac{1}{2} + \frac{1}{8} \epsilon^4$

and the second $\geq \frac{1}{2} + \frac{1}{4} \epsilon^4$. We can distinguish the two with 90% confidence using $\Theta(\frac{1}{\epsilon^4})$ repetitions.

3. (a). $2^{-(k-1)} \cdot (1-2^{-k})^{k-1}$



(b). $(1-2^{-k})^k = \Pr[\text{Traces}=0]$. So for any l we can choose $k = \lfloor 2^l \cdot \ln 2 \rfloor$.

(c). All influences are $\Theta(2^{-k}) = \Theta(\frac{\log n}{n})$ where $n = k \cdot l$. For majority this is $\Theta(\frac{1}{\sqrt{n}})$. [KKL] showed that this is best possible.

4. (a). $E[f(x) \cdot g(x)] - E[f] \cdot E[g] = \langle f, g \rangle - \hat{f}(\emptyset) \hat{g}(\emptyset) = \sum_{S \neq \emptyset} \hat{f}(S) \hat{g}(S)$.

(b). Assume f depends on $\{1, \dots, r\}$. $\text{Cov}[h, f] = \sum_{S \neq \emptyset} \hat{f}(S) \hat{h}(S) = \sum_{[r] \supseteq S \neq \emptyset} \hat{f}(S) \hat{h}(S) \leq \sqrt{\sum_{S \neq \emptyset} \hat{f}(S)^2} \cdot \sqrt{\sum_{S \neq \emptyset} \hat{h}(S)^2} \leq 1 \cdot \sqrt{\sum_{S \neq \emptyset} \hat{h}(S)^2}$. By assumption, $\forall i, \sum_{S \ni i} (1-\delta)^{|S|} \hat{h}(S)^2 \leq \epsilon$.

Therefore, $(1-\delta) \cdot \sum_{S \neq \emptyset} \hat{h}(S)^2 \leq \sum_{i=1}^r \sum_{S \ni i} (1-\delta)^{|S|} \cdot \hat{h}(S)^2 \leq \epsilon \cdot r$.

Possible: $\min(1, \frac{\epsilon}{(1-\delta)^{|S|}}) \leq |S| (1-\delta)^{|S|} \quad (1 \leq |S| \leq r)$

6. (a). $\|f^{\text{odd}}\|_2^2 = \sum_{|S| \text{ odd}} \hat{f}(S)^2 \leq \sum_{S \neq \emptyset} \hat{f}(S)^2 \leq \sum_S |S| \hat{f}(S)^2 = I(f)$.

(b). Write $F = (f_1, f_2, \dots, f_m)$ and apply (a) to each f_i and sum the ineq.

(c). Assume F is an embedding with distortion D .

$n^2 \leq E_x [\|F(x) - F(x \oplus (1, \dots, 1))\|_2^2] \leq \sum_{i=1}^m E_x [\|F(x) - F(x \oplus e_i)\|_2^2] \leq n \cdot D^2 \Rightarrow D \geq \sqrt{n}$.

13.3.2008

Observations: 1. $\hat{f}_{x \rightarrow \bar{I}}(s) = \sum_{T \subseteq \bar{I}} \hat{f}(s \cup T) \cdot \chi_T(x)$

2. $E_x [|\hat{f}_{x \rightarrow \bar{I}}(s)|] \geq |\hat{f}(s)|$

3. $E_x [\hat{f}_{x \rightarrow \bar{I}}(s)^2] = \sum_{T \subseteq \bar{I}} \hat{f}(s \cup T)^2$.

Thm: If f is computable by a width- w DNF, then $\forall d \geq 5, \sum_{|S| \geq 20dw} \hat{f}(S)^2 \leq 2^{-d+1}$.

Remark: This gives $n^{O(w \log 1/\epsilon)}$ time alg. to learn DNFs from random examples.

Also $n^{O(\log n \cdot \log 1/\epsilon)}$ for poly-size DNFs.

Thm: If f has width- w DNF then $\sum_U \left(\frac{1}{20w}\right)^{|U|} \cdot |\hat{f}(U)| \leq 8$.

Proof: Let (I, x) be a random restriction with $p = \frac{1}{20w}$. Then,

$E_{I, x} [\|\hat{f}_{x \rightarrow \bar{I}}\|_1] \leq \sum_{d=0}^{\infty} \Pr[\text{DT-depth}(f_{x \rightarrow \bar{I}}) = d] \cdot 2^d \stackrel{\text{Hastad}}{\leq} \sum_{d=0}^{\infty} 4^{-(d-1)} \cdot 2^d = 8$.

$\|f\|_1$ of depth- d DT is $\leq 2^d$

On the other hand, $E_{I, x} [\|\hat{f}_{x \rightarrow \bar{I}}\|_1] = E_I \left[\sum_{S \subseteq \bar{I}} E_x [|\hat{f}_{x \rightarrow \bar{I}}(S)|] \right] \geq E_I \left[\sum_{S \subseteq \bar{I}} |\hat{f}(S)| \right] =$

$= \sum_{U \subseteq [n]} \Pr_{I, x} [U \subseteq \bar{I}] \cdot |\hat{f}(U)| = \sum_{U \subseteq [n]} \left(\frac{1}{20w}\right)^{|U|} \cdot |\hat{f}(U)|$.

This implies that $\sum_{|U| \leq O(w \log 1/\epsilon)} |\hat{f}(U)| \leq w^{O(w \log 1/\epsilon)}$.

Cor: If f is computable by a width w DNF, then it is ϵ -concentrated on a set of size $w^{O(w \log 1/\epsilon)}$.

Proof: Define $S = \{U : |U| \leq O(w \log 1/\epsilon) \text{ and } |\hat{f}(U)| \geq \frac{\epsilon}{w^{O(w \log 1/\epsilon)}}\}$. Then $|S| \leq w^{O(w \log 1/\epsilon)}$.

Moreover, $\sum_{U \notin S} \hat{f}(U)^2 \leq \sum_{|U| > O(w \log 1/\epsilon)} \hat{f}(U)^2 + \sum_{|U| \leq O(w \log 1/\epsilon)} \hat{f}(U)^2 \leq 2\epsilon$.
 $\leq \epsilon$ since $\sum \hat{f}(U)^2 \leq \max |\hat{f}(U)| \cdot \sum |\hat{f}(U)|$. \square

This gives: $w^{O(w \log 1/\epsilon)}$ learning time algorithm for width- w DNF and $n^{O(\log \log n \cdot \log 1/\epsilon)}$ time alg. for poly-size DNF [Mansour].

Summary: (uniform PAC, constant ϵ)

	random examples	queries
poly-size depth circuits	$n^{O(\log^d n)}$ [LMN93]	$\rightarrow \dots$
poly-size DNF	$n^{O(\log n)}$ --	poly(n) [Jackson94]
poly-size DT	$n^{O(\log n)}$ --	poly(n)
$\log n$ -junta	$n^{0.704 \log n}$	poly(n)

- Open problems:
- Are poly(n) size DNFs concentrated on poly(n) coefficients? (this would imply Jackson). Tribes may be counterexample?!
 - Big open question: can monotone polysize DNFs be learned from random examples in polytime? (CCC'06: O'Donnell & Servedio did this for monotone DTs).

Learning Juntas (Mossell, O'Donnell, Servedio, STOC'03)

We want to learn k -juntas. Think of k as being very small, say, $k = \log n$. Because a k -junta has $\leq 2^k$ nonzero Fourier coeff all in levels $\leq k$ and multiple of 2^k , we can learn them exactly in time $\text{poly}(2^k, n)$ using queries, or in time $n^k \cdot \text{poly}(2^k, n)$ using random examples. We will improve this slightly to $n^{0.704k} \cdot \text{poly}(2^k, n)$.

Step 1: Finding a relevant coordinate is enough.

Prop: If there is an algorithm for finding a relevant coordinate in a given (nonconstant) k -junta in time $n^{\alpha(k)} \cdot \text{poly}(2^k, n)$ using random examples then there is an algorithm for learning k -juntas in the same time.

Proof sketch: Find a relevant coordinate, say i . Now recur on $f_{0 \rightarrow i}$ and on $f_{i \rightarrow 1}$.

Notice that we can simulate examples from restrictions of f of r variables using examples from f with an overhead of 2^r (on average). After finding all k relevant coordinates we can determine the junta by observing examples for all 2^k possible settings. This takes $\Theta(k \cdot 2^k)$ examples. \square

If a k -junta has a nonzero Fourier coeff. at level $d > 0$ then we can find it in time $n^d \cdot \text{poly}(n, 2^k)$. This gives a relevant coordinate.

Does a junta always have a nonzero Fourier coeff. in level $0 < d < k$?

No! For instance, parity $\chi_{\{i_1, i_2, \dots, i_k\}}$. But parities are easy to learn!

Given examples $(x^1, f(x^1)), (x^2, f(x^2)), \dots$ write a sequence of linear equalities

over \mathbb{GF}_2 in variables a_1, a_2, \dots, a_n :
 $a_1 x_1^1 + a_2 x_2^1 + \dots + a_n x_n^1 = f(x^1)$
 $a_1 x_1^2 + a_2 x_2^2 + \dots + a_n x_n^2 = f(x^2)$

We only need $\sim n$ equations. \vdots

Step 2: Multilinear polynomials over \mathbb{GF}_2 .

Examples: PARITY $(x_1, \dots, x_n) = x_1 + x_2 + \dots + x_n$ of degree 1.

AND $(x_1, \dots, x_n) = x_1 \cdot x_2 \cdot \dots \cdot x_n$ of degree n .

Prop: Any $f: \mathbb{GF}_2^n \rightarrow \mathbb{GF}_2$ can be uniquely expressed as a multilinear polynomial.

Proof: We can write $f(x) = \sum_{a \in \mathbb{GF}_2^n} f(a) \cdot \prod_{i=1}^n (1 - (x_i - a_i))$. Uniqueness follows because the only multilin. poly. that is constantly 0 is the zero poly. \square

Thm: The class $\{f: \mathbb{GF}_2^n \rightarrow \mathbb{GF}_2 \mid \deg_{\mathbb{GF}_2}(f) \leq e\}$ is learnable from random examples in time $n^{w \cdot e} \cdot \text{poly}(n, 2^e)$ where w is the matrix multiplication constant (currently $w \leq 2.376$ [Coppersmith Winograd 90], Strassen $\log_2 7$).

Proof: Take $m \approx n^e$ random examples (requires proof) $(x^1, f(x^1)), (x^2, f(x^2)), \dots$ and solve equations in variables $\{a_T\}_{|T| \leq e}$ over \mathbb{GF}_2 :

$\left\{ \sum_{|T| \leq e} a_T \cdot \prod_{j \in T} x_j^i = f(x^i) \right\}_{i=1}^m$. This takes time $n^{w \cdot e} \cdot \text{poly}(n, 2^e)$.
const!

Thm [Siegenthaler 84]: If $f: \{0,1\}^k \rightarrow \{-1,1\}$ is s.t. $\hat{f}(S) = 0$ for all S with $1 \leq |S| \leq d$ then over \mathbb{GF}_2 , we can express f as a degree $\leq k-d$ multilinear poly.

Proof: Let $g = f \cdot \chi_{\{1, \dots, k\}}$. If we can represent g as a degree $\leq k-d$ multilin. poly over \mathbb{GF}_2 then we're done because we obtain f by adding $x_1 + x_2 + \dots + x_n$ which is of degree 1. Notice that $\forall S, \hat{g}(S) = \hat{f}(S \Delta \{1, 2, \dots, k\})$, hence $\hat{g}(S) = 0$ for all S with

$k-d \leq |S| \leq k-1$. Consider the multilinear poly. over \mathbb{R} , $h(x_1, \dots, x_n) = \frac{1}{2} - \frac{1}{2} \sum_S \hat{g}(S) \prod_{i \in S} (1-2x_i)$.

Clearly, for all $x_1, \dots, x_n \in \{0,1\}$, $h(x_1, \dots, x_n) = \frac{1}{2} - \frac{1}{2} g(x_1, \dots, x_n) \in \{0,1\}$

This implies that if we expand h , then the coeff of every monomial is integer

(this can be seen by induction on the degree of the monomial: the free coeff. is $h(0, \dots, 0)$ and hence integer; the x_i coeff is $h(0, \dots, 1, 0, \dots, 0) - h(0, \dots, 0) \in \mathbb{Z}$, and is therefore integer). If we reduce each coeff of h modulo 2, we get a representation of g as a multilin. poly over \mathbb{GF}_2 . It remains to prove that the degree is $\leq k-d$.

Notice that $\prod_{i \in S} (1-2x_i)$ is of degree $= |S|$. This means that the coeff of $\prod_{i \in T} x_i$ for some T of size $\geq k-d$ is only affected by $\hat{g}(\{1, \dots, k\})$ and is therefore $-\frac{1}{2} \hat{g}(\{1, \dots, k\}) \cdot (-2)^{|T|}$.

Since all coeff of h are integers, $-\frac{1}{2} \hat{g}(\{1, \dots, k\}) \cdot (-2)^{|T|}$ must be integer. This implies that for T of size $\geq k-d+1$ the coeff is even and therefore disappears when we take modulo 2. \blacksquare

Remark: We could avoid the last complication by using the homework (unbalanced functions have "low" Fourier coeff).

Thm: k -juntas can be learned in time $n^{\frac{nk}{w+1}} \cdot \text{poly}(n, 2^k)$, from random examples.

Proof: Let $d = \frac{nk}{w+1}$ ($\approx 0.7k$). Look for nonzero Fourier coeff up to level d in (in levels $1, \dots, d$) time $n^d \cdot \text{poly}(n, 2^k)$. If found, we have a relevant coordinate, and we're done.

Otherwise, our function is of degree $\leq k-d = \frac{k}{w+1}$ over \mathbb{GF}_2 , so we can learn it in time $n^{\frac{nk}{w+1}} \cdot \text{poly}(n, 2^k)$. \blacksquare

Open Questions: • Learn k -juntas in time $\text{poly}(n, 2^k)$ or even $n^{k/2} \cdot \text{poly}(n, 2^k)$.

• Do something similar over $\{0,1,2\}^n$, or over $\{0,1\}^n$ with μ_p .