

Fourier Learning

Monday, November 1, 2021 6:38 AM

Yuh

Motivated by learning DNF.

Still an open problem to PAC learn DNF or decision trees (lists we can learn).

So we consider special distributions

- uniform
- product
- Gaussian etc

and allow for membership queries, i.e., "what is the label of x ?"

Assume $x \in \{-1, 1\}^n$ $f: \{-1, 1\}^n \rightarrow \{-1, 1\}$

Note that any such $f \in \{-1, 1\}^{2^n}$ (as a table)

Defines inner product of f and g with respect to distribution D as

$$\langle f, g \rangle = \sum D(x) f(x) g(x)$$

$$\langle f, g \rangle_D = \sum_x D(x) f(x) g(x)$$

$$\langle f, f \rangle_D = \|f\|_D^2 = 1 \quad (\text{since } f(x)^2 = 1)$$

Viewing f as a vector, the standard basis is e_1, e_2, \dots, e_{2^n} .

But we can use any basis and write $f(x) = \sum_v \langle f, v \rangle v$ where $\{v\}$ is an orthonormal basis.

What's an interesting basis?

The set of parity functions!

$\forall S \subseteq [n], \chi_S(x) = \prod_{i \in S} x_i$ 2^n functions.

$$\langle \chi_S, \chi_S \rangle_D = 1$$

$$\langle \chi_S, \chi_T \rangle_D = \mathbb{E}_D \left(\prod_{i \in S} x_i \prod_{j \in T} x_j \right) = 0$$

for a product distribution D .

for a "product" distribution \mathcal{D} .

Hence $\{\chi_S\}$ is an orthonormal basis!

So any f can be written as

$$f(x) = \sum_S \hat{f}_S \chi_S(x) \quad \text{where } \hat{f}_S = \langle f, \chi_S \rangle_{\mathcal{D}}.$$

Thm 1. (Parseval) $\langle f, f \rangle_{\mathcal{D}} = \langle \hat{f}, \hat{f} \rangle$

Thm 2. (Plancherel) $\langle f, g \rangle_{\mathcal{D}} = \langle \hat{f}, \hat{g} \rangle.$

Pf.

$$\begin{aligned} & \sum_x \mathcal{D}(x) \sum_S \hat{f}_S \chi_S(x) \sum_T \hat{g}_T \chi_T(x) \\ &= \sum_{S, T} \hat{f}_S \hat{g}_T \mathbb{E}_{\mathcal{D}}(\chi_S(x) \chi_T(x)) \\ &= \sum_S \hat{f}_S \hat{g}_S = \langle \hat{f}, \hat{g} \rangle. \end{aligned}$$

A decision tree is a Boolean function f .

We want to learn f by approximating all

We want to learn f by approximating
of its significant Fourier coefficients \hat{f}_s .

Our approx is g .

$$\Pr_{\mathcal{D}}(g(x) \neq f(x)) \leq \mathbb{E}_{\mathcal{D}}((f(x) - g(x))^2) = \sum_S (\hat{f}_s - \hat{g}_s)^2.$$

We will learn all \hat{f}_s for which $|\hat{f}_s| \geq \tau$.

Note. $\sum_S \hat{f}_s^2 = 1 \Rightarrow |\hat{f}_s| \leq 1$

Lemma ^[DNF] If a decision tree has m leaves

then $\|f\|_1 = \sum |\hat{f}_s| \leq 2m + 1$.

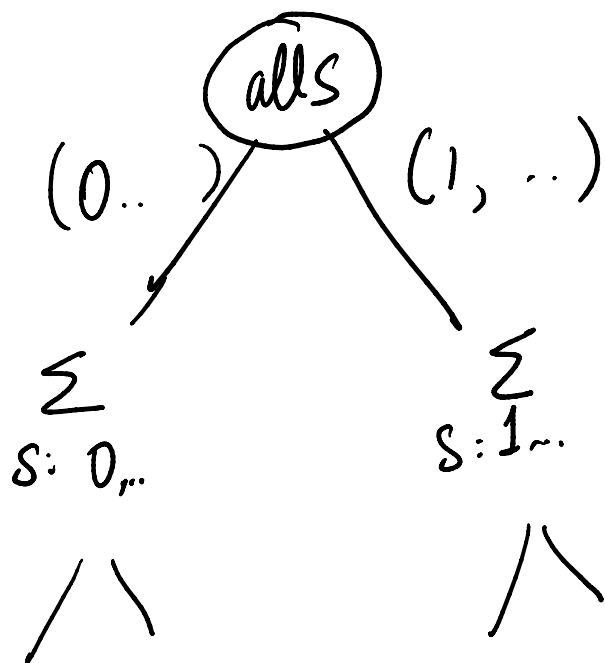
Thm. If we learn all $\hat{f}_s \geq \frac{\epsilon}{\|f\|_1}$

then $\|\hat{f}_s - \hat{g}_s\|^2 \leq \epsilon$.

Pf. $\|\hat{f}_s - \hat{g}_s\|^2 = \sum_{s: |\hat{f}_s| \leq \frac{\epsilon}{\|f\|_1}} \hat{f}_s^2 \leq \sum_S |\hat{f}_s| \cdot \frac{\epsilon}{\|f\|_1}$

$$\|f_s\| \leq \frac{\epsilon}{\|f\|}, \quad = \epsilon.$$

How to learn all large Fourier coefficients.



At each node estimate whether

$$\sum f_s^2 \geq \tau.$$

S: prefix α

width $\leq \frac{1}{\tau}$, depth $\leq n$

$$\# \text{ nodes} \leq \frac{n}{\tau}.$$

How to estimate

$$\sum_{S \in S_\alpha} \hat{f}_S^2 ?$$

Suppose $\alpha = (0, \dots, 0)$

$$\text{Claim } \sum \hat{f}_S^2 \in (f(yx) f(zx))$$

Claim. $\sum_{S \subseteq \Omega} \hat{f}_S^2 = \mathbb{E} (f(yx) f(zx))$
 $x \sim \{0,1\}^{n+k}, y, z \sim \{0,1\}^k$

Suppose f is a parity function

if f agrees with α , then $f(yx) = f(zx)$ so we get 1.

else $\Pr(f(yx) = f(zx)) = \frac{1}{2} \Rightarrow$ we get 0.

Any f can be written as a weighted sum of parities. So,

$$f = \sum_U \hat{f}_U \chi_U$$

$$\mathbb{E}(f(yx) f(zx)) = \mathbb{E} \left(\sum_U \hat{f}_U \chi_U(yx) \sum_V \hat{f}_V \chi_V(zx) \right)$$

$$= \sum_{U, V} \hat{f}_U \hat{f}_V \underbrace{\mathbb{E}(\chi_U(yx) \chi_V(zx))}_{=0 \text{ if } U \neq V}$$

$$= \sum_U \hat{f}_U^2 \underbrace{\mathbb{E}(\chi_U(yx) \chi_U(zx))}_{=0 \text{ if } U \text{ does not}}$$

.. ..

= 0 if U does not agree with $\alpha = (0, \dots, 0)$

$$= \sum_{U \in U_\alpha} \hat{f}_U^2$$

What about general α ?

Lemma.
$$\sum_{S \in S_\alpha} \hat{f}_S^2 = \mathbb{E}_{x \sim \{0,1\}^{n-k}} (f(yx) f(zx) \chi_\alpha(y) \chi_\alpha(z))$$

 $y, z \sim \{0,1\}^k$

Proof [of Lemma [DNF]].

Consider a single conjunction T . Let $T(x) = 1$ if x satisfies it and $T(x) = 0$ o.w.

Then
$$\langle T, T \rangle_D = \mathbb{E}_D (T(x)^2) = \frac{1}{2^{\pi}}$$

$$\hat{f}_S = \langle T, \chi_S \rangle_D$$

$$= \mathbb{E}_D (T(x) \chi_S(x))$$

$$= \frac{1}{2^{\pi}} \mathbb{E} (\chi_S(x) / (T(x) - 1)) = \begin{cases} 0 & \text{if } S \text{ contains } x_i \notin T \end{cases}$$

$$= P_{\mathcal{D}}(\pi(x)=1) \mathbb{E}_{\mathcal{D}}(\chi_S(x) / T(x)=1) = \begin{cases} 1 & x_i \in T \\ \frac{1}{2^{|\pi|}} & \text{o.w.} \end{cases}$$

$$\text{So } \|\hat{T}\|_1 = 1 \quad \|\hat{T}\|_2^2 = 2^{|\pi|} \cdot \frac{1}{2^{2|\pi|}} = \frac{1}{2^{|\pi|}}.$$

For a decision tree with m leaves,
we can write

$$f(x) = 2 \left(T_1(x) + \dots + T_m(x) \right) - 1$$

$$\text{So } \|\hat{f}\|_1 \leq 2 \sum_{i=1}^m \|\hat{T}_i\|_1 + 1 \leq 2m + 1.$$