

Learning Under Distributional Assumptions

As we've seen, many natural learning problems are hard if we make no distributional assumptions

Today: Study learning under simple distributions

There is no known subexponential time alg. for learning intersections of two halfspaces, but...

Theorem 1 [Klivans, O'Donnell, Servedio]

For any constant k , there is a poly. time algorithm for learning intersections of k halfspaces under uniform distribution on hypercube

The main idea is to show:

"intersections of k halfspaces can be approximated by low degree polys."

We'll employ Fourier analysis over the hypercube

Proposition: For any function

$$f: \{\pm 1\}^d \rightarrow \mathbb{R}$$

we can write it as

$$f(x) = \sum_{S \subseteq [d]} \hat{f}(S) \chi_S(x)$$

where $\chi_S(x) = \prod_{i \in S} x_i$ and

$$\hat{f}(S) = \mathbb{E}_{x \sim_{\text{u}} \{\pm 1\}^d} [f(x) \chi_S(x)] \triangleq \langle f, \chi_S \rangle$$

Intuition: The functions χ_S are orthonormal and they span the space of functions on the hypercube

$$\mathbb{E}_{x \sim_{\text{u}} \{\pm 1\}^d} [\chi_T(x) \chi_S(x)] = \mathbb{E}_{x \sim_{\text{u}} \{\pm 1\}^d} [x_3] = 0$$

$T = \{1, 2, 3\} \quad S = \{1, 2\}$

Thus going from one representation of a function

(1) its evaluation at each point of the hypercube

to another

(2) its spectrum, i.e. the collection of its Fourier coefficients $\hat{f}(s)$

is just an orthogonal transformation

Thus we have:

Fact [Parseval's Identity]:

$$\mathbb{E} [f(x)^2] = \sum_{s \subseteq [d]} \hat{f}(s)^2$$

$x \sim_{\text{unif}} \mathbb{S}^d$

Let's prove this to make sure we are comfortable with Fourier analysis over the hypercube

Proof: Substituting in the Fourier representation, we have

$$\mathbb{E} \left[\left(\sum_{s \in [d]} \hat{f}(s) X_s(x) \right) \left(\sum_{t \in [d]} \hat{f}(t) X_t(x) \right) \right] \quad (*)$$

$X \sim \mathcal{U}(\{\pm 1\}^d)$

Now by orthogonality, all the cross terms that do not match are zero,

$$(*) = \mathbb{E} \left[\sum_{s \in [d]} \hat{f}(s)^2 X_s^2(x) \right]$$

$X \sim \mathcal{U}(\{\pm 1\}^d)$ $\uparrow = 1$

$$= \sum_{s \in [d]} \hat{f}(s)^2$$



The starting point for many learning algorithms is truncation

$$g_\ell(x) = \sum_{|s| < \ell} \hat{f}(s) X_s(x)$$

which will give us a good low-degree approximation

Let's make this precise:

def: We say that f is $\alpha(\epsilon, d)$ concentrated

$$\sum_{|s| \geq \alpha(\epsilon, d)} \hat{f}(s)^2 \leq \epsilon$$

Or to put it another way, if we take $g_{\alpha(\epsilon, d)}$ and let r be the residual — i.e.

$$r(x) = f(x) - g_{\alpha(\epsilon, d)}(x)$$

we have

$$\mathbb{E}_{x \sim \mathcal{U}_{\pm 1}^d} [r(x)^2] = \sum_{|s| \geq \alpha(\epsilon, d)} \hat{f}(s)^2 \leq \epsilon$$

The seminal work of Linial, Mansour and Nisan introduced Fourier analysis into computational learning theory

They showed the following meta-theorem

Theorem 2 If a concept class \mathcal{H} has Fourier concentration $\alpha(\epsilon, d)$ then there is a PAC learning algorithm that runs in time $d^{O(\alpha(\epsilon, d))}$ on uniform distribution

The Low Degree Algorithm

For each $S \subseteq [d]$ with $|S| < \alpha(\epsilon, d)$

Take m samples to estimate

$$\hat{g}(s) = \frac{\sum_{i=1}^m f(x) X_S(x)}{m} \approx \hat{f}(s)$$

Output the estimate

$$g(x) = \sum_{|S| < \alpha(\epsilon, d)} \hat{g}(s) X_S(x)$$

Proof: If each estimate satisfies

$$|\hat{g}(s) - \hat{f}(s)|^2 \leq \frac{\epsilon}{d^{O(\alpha(\epsilon, d))}} \quad (\Delta)$$

then we have

$$\mathbb{E}[|g(x) - f(x)|^2] \leq 2\varepsilon$$

And now, given x , if you output the label $y = \text{sgn}(g(x))$ we can check

$$(f(x) - g(x))^2 \geq \mathbb{1}_{\text{sgn}(g(x)) \neq f(x)}$$

Thus we have

$$\mathbb{P}_{x \sim \mathcal{U}(\mathbb{S}^d)} [\text{sgn}(g(x)) \neq f(x)] \leq 2\varepsilon$$

Finally, standard tail bounds imply we can choose

$$m = \frac{d^{O(\alpha(\varepsilon, d))}}{\varepsilon^2} \log \frac{d^{O(\alpha(\varepsilon, d))}}{\delta}$$

s.t. (Δ) holds for all S with $|S| < \alpha(\varepsilon, d)$ with probability $\geq 1 - \delta$. \square

Main Question: what kinds of functions have Fourier Concentration?

We will see that stable functions have Fourier concentration, and are thus learnable

def: Let $T_\eta(x)$ denote the noise operator which flips each bit independently with prob. η

def: Let $NS_\eta(x)$ denote the noise sensitivity of f , defined as

$$NS_\eta(x) = \mathbb{P}_{x \sim \mathcal{U}^{\pm\{0,1\}^d}} [f(x) \neq f(T_\eta(x))]$$

The noise sensitivity has a closed-form expression:

Proposition 2: For $\eta < \frac{1}{2}$

$$NS_\eta(x) = \frac{1}{2} - \frac{1}{2} \sum_s (1 - 2\eta)^{|s|} \hat{f}(s)^2$$

Proof: First we claim (\square)

$$NS_n(f) = \frac{1}{2} - \frac{1}{2} \mathbb{E} [f(x) f(T_n(x))]$$

Again, we can substitute in the Fourier representation:

$$(\square) = \mathbb{E} \left[\left(\sum_s \hat{f}(s) \chi_s(x) \right) \left(\sum_u \hat{f}(u) \chi_u(T_n(x)) \right) \right]$$

And again, only the terms where $s = u$ contribute. Thus we have:

$$(\square) = \sum_s \hat{f}(s)^2 \underbrace{\mathbb{E} [\chi_s(x) \chi_s(T_n(x))]}_{(c)}$$

Now we can think about the noise operator equivalently as:

(1) Select bits independently with probability 2^{-n}

② set them u.a.r. to ± 1

Thus we have

$$(\circ) = \mathbb{P} \left[\begin{array}{l} \text{no bit in} \\ s \text{ is selected} \end{array} \right] = (1 - 2n)^{|s|}$$

Putting it all together

$$NS_n(f) = \frac{1}{2} - \frac{1}{2} \sum_s \hat{f}(s)^2 (1 - 2n)^{|s|}$$



Note: there are many manifestations of these same principles in other learning settings

gaussian
noise sensitivity



learnability
under gaussian
distributions

which applies for functions on \mathbb{R}^d

With these tools in hand, the key ideas are

- ① halfspaces have low noise sensitivity
- ② combining functions with low noise sensitivity, does not increase noise sensitivity too much
- ③ low noise sensitivity \Rightarrow Fourier concentration

For ①, we have:

Theorem 3 [Peres]: Let f be any halfspace
Then $NS_{\pi}(f) \leq O(\sqrt{\pi})$

We will not prove this. But the intuition is related to the central limit theorem:

In particular, consider

$$\text{MAJ} \triangleq \text{sgn} \left(\sum_{i=1}^d x_i \right)$$

Now if we consider the quantity

$$X \triangleq \frac{\sum_{i=1}^d x_i}{\sqrt{d}}$$

it behaves like $N(0, 1)$ - a mean zero and variance one Gaussian

Now let $y = T_n(x)$ and $Y = \frac{\sum_{i=1}^d y_i}{\sqrt{d}}$

Again by the central limit theorem

$$Y - X \sim N(0, 4n)$$

because we flip a n fraction of bits in expectation

And the dominant contribution to the event $\text{sgn}(Y) \neq \text{sgn}(X)$ comes from:

$$\textcircled{1} \quad -c\sqrt{n} \leq X \leq c\sqrt{n}$$

$\textcircled{2}$ $|Y - x| > c\sqrt{n}$ and the sign goes the right way

You can check that the probability all these things happen is about \sqrt{n}

Next we will translate this into a bound on the Fourier concentration

Lemma 1: For any $0 < \eta < \frac{1}{2}$, we have

$$\sum_{|s| \geq \frac{1}{\eta}} \hat{f}(s)^2 \leq \left(\frac{2}{1 - \frac{1}{e^2}} \right) NS_{\eta}(f)$$

Proof: From Proposition 2, we have

$$2NS_{\eta}(f) = 1 - \sum_s (1 - 2\eta)^{|s|} \hat{f}(s)^2$$

Applying Parseval's identity

$$= \sum_s \hat{f}(s)^2 - \sum_s (1-2\pi)^{|s|} \hat{f}(s)^2$$

$$= \sum_s \hat{f}(s)^2 (1 - (1-2\pi)^{|s|})$$

$$\geq \sum_{|s| \geq \frac{1}{\pi}} \hat{f}(s)^2 (1 - (1-2\pi)^{|s|})$$

$$\geq \sum_{|s| \geq \frac{1}{\pi}} \hat{f}(s)^2 (1 - (1-2\pi)^{\frac{1}{\pi}})$$

$$\geq \sum_{|s| \geq \frac{1}{\pi}} \hat{f}(s)^2 (1 - \frac{1}{e^2})$$

Rearranging completes the proof. \square

Finally, we need an elementary fact

Fact 2: Suppose we consider composite functions $h(x) = g(f_1(x), \dots, f_k(x))$

$$\text{Then } NS_n(h) \leq \sum_{i=1}^k NS_n(f_i)$$

Proof: This follows from the definition of noise sensitivity + union bound



Now we are ready to prove Theorem 1

Proof: Let f be an intersection of k halfspaces. Using Lemma 1:

$$\sum_{|S| \geq \frac{1}{m}} \hat{f}(S)^2 \leq 2.32 NS_n(f)$$

And applying Fact 2 and Theorem 3:

$$\leq C R \sqrt{n}$$

Now set $n = \frac{\epsilon^2}{ck^2}$, which implies

$$\sum_{|S| \geq \frac{ck^2}{\epsilon^2}} \hat{f}(S)^2 \leq \epsilon$$

Finally, using the low degree algorithm we get a $d^{O(\frac{k^2}{\epsilon^2})}$ time algorithm, as desired. \square

Is this bound tight?

For intersections of halfspaces, can beat the union bound substantially

Theorem [Kane]: For an intersection of k halfspaces

$$NS_n(f) \leq O(\sqrt{n \log k})$$

Corollary: There is a $d^{O(\frac{\log k}{\epsilon^2})}$ time alg.
for learning intersections of k halfspaces
over the uniform distribution on $\{\pm 1\}^d$

Note: All of these algorithms even work
in the agnostic setting and get error $\text{opt} + \epsilon$

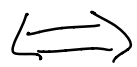
In Gaussian space, can do even better

Theorem [Vempala]: There is a
 $\text{poly}(d, k, \frac{1}{\epsilon}) + k^{O(\frac{\log k}{\epsilon^4})}$

time algorithm for learning intersections of
 k halfspaces over a Gaussian

Intuition: If you restrict to just the
positive / negative examples, you get

span of
normals



directions of
smallest variance

Statistical Query Lower Bounds

Recall, we discussed cryptographic lower bounds for learning

Main Question: But what about lower bounds for "nice" distributions?

Today, we'll introduce a powerful framework

def [Kearns]: In the statistical query model, in each step

(i) The algorithm specifies a query

$$q: \mathbb{R}^d \times Y \rightarrow [-1, 1]$$

and a tolerance τ

② It gets a response r s.t.

$$|r - \mathbb{E}_{(x,y) \sim D} [a(x,y)]| \leq \tau$$

This is a restriction on what you are allowed to do algorithmically

And you can turn SQ algorithms into actually learning algorithms in the following sense

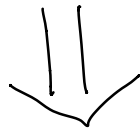
SQ Algorithm

Assumptions

① at most polynomially many queries

② each tolerance is at least inverse polynomially big

③ can compute the next query,
and evaluate it on samples
efficiently



There is a polynomial time learning
algorithm

The important point is SQ algorithms
do not look directly at the data, but
rather use inexact statistics / summaries

Meta Claim: By dropping ③, can prove
many tight upper / lower bounds

Let's start with a key example:

Learning Parity Functions

① There is an unknown $T \subseteq [d]$

② $X \sim_{\mu} \{\pm 1\}^d$, and $y = X_T(x)$

Observe that parities are not at all Fourier concentrated / stable

Theorem 3 [Blum, Furst, Jackson, Kearns, Mansour, Rudich] Any SQ algorithm for learning parities must make at least $2^{\Omega(d)}$ queries or have $\tau = 2^{-\Omega(d)}$

Proof (sketch): It turns out that it suffices to allow only correlational queries:

$$q(x, y) \triangleq y q(x)$$

Now, once again, we use the Fourier representation

$$q(x) = \sum_s \hat{q}(s) \chi_s(x)$$

Thus we have

$$\begin{aligned} \mathbb{E}[y q(x)] &= \mathbb{E}\left[\chi_T(x) \sum_s \hat{q}(s) \chi_s(x)\right] \\ &= \hat{q}(T) \end{aligned}$$

Now by Parseval's identity

$$\sum_s \hat{q}(s)^2 = \mathbb{E}[q(x)^2] \leq 1$$

Thus the intuition is

- (i) we can only have $|\hat{q}(s)| \geq \tau$
for at most $\frac{1}{\tau^2}$ parities

② If we have $|\hat{q}(T)| < \tau$ the oracle can just reply "zero"

Thus it will take us about $\tau^2 2^d$ queries to find T . \square

Is there a non-SQ algorithm for learning parities?

Fact 3: There is a polynomial time algorithm for learning parities

Main Idea: solve for T by setting up a linear system over \mathbb{F}_2

$$\begin{matrix} m \times d & & 1 \times 1 \\ \nearrow & \downarrow & \\ A & C & = y \\ \left. \begin{matrix} -x_1- \\ -x_2- \\ \vdots \end{matrix} \right\} \end{matrix}$$

But this algorithm is highly non-robust

In fact:

i.e. flip label with prob. $\frac{1}{3}$

Theorem [Blum, Kalai, Wasserman]

There is a $2^{d/\log d}$ time algorithm for learning noisy parities

This is the best known algorithm

Are there other natural learning problems that separate SQ and general learning algorithms?

Next time: SGD as an SQ algorithm, and SQ lower and upper bounds for learning deep nets on Gaussian inputs