

## Lecture 18

### Fourier-based learning algorithms

- the low degree algorithm
- Fourier Concentration
- Noise sensitivity

Recall Fourier Transform:

$$\chi_s(x) = \prod_{i \in S} x_i$$

$$\langle f, g \rangle = \frac{1}{2^n} \sum_x f(x) g(x)$$

$$\hat{f}(s) = \langle f, \chi_s \rangle = 1 - 2 \cdot \Pr[f(x) \neq \chi_s(x)]$$

← lemma

$$\forall f, f(x) = \sum_s \hat{f}(s) \chi_s(x)$$

$$\text{Plancherel } \langle f, g \rangle = \sum_s \hat{f}(s) \hat{g}(s)$$

# Learning via Fourier Representation

will look at learning algorithms that are based on estimating Fourier representation of fctn  $f$   
(similar to polynomial interpolation)

Approximating one Fourier coefficient:

lemma for any  $S \subseteq [n]$ , can approx  $\hat{f}(s)$  to within additive  $\delta$   
(i.e.  $|\text{output} - \hat{f}(s)| \leq \delta$ )  
with prob  $\geq 1 - \delta$  in  $O\left(\frac{1}{\delta^2} \log \frac{1}{\delta}\right)$  samples.  
no queries needed!

(Proved last time)

# The low degree algorithm

definition of fctns for which low degree

Fourier coeffs pretty much suffice to describe fctn:

def  $f: \pm 1^{\mathbb{S}^n} \rightarrow \mathbb{R}$  has  $\alpha(\epsilon, n)$ -Fourier concentration

$$\text{if } \sum_{\substack{S \subseteq [n] \\ \text{s.t.} \\ |S| > \alpha(\epsilon, n)}} \hat{f}(S)^2 \leq \epsilon \quad \forall 0 < \epsilon < 1$$

for Boolean  $f$ , this implies

$$\sum_{\substack{S \subseteq [n] \\ \text{s.t.} \\ |S| \leq \alpha(\epsilon, n)}} \hat{f}(S)^2 \geq 1 - \epsilon$$

## examples

1) fctn  $f$  which depends on  $\leq k$  vars } if  $f$  doesn't depend on  $x_i$  then all  $\hat{f}(S)$  for which  $i \in S$  satisfy  $\hat{f}(S) = 0$

has  $\sum_{\substack{S \text{ s.t.} \\ |S| > k}} \hat{f}(S)^2 = 0$

# Low degree algorithm

approximates fctns with  $d \equiv \Omega(\frac{1}{\epsilon}, n)$  Fourier concentrations

Given:  $d$  degree  
 $\gamma$  accuracy  
 $\delta$  confidence

Algorithm:

• Take  $m = O\left(\frac{n^d}{\gamma} \ln \frac{n^d}{\delta}\right)$  samples

• For each  $S$  s.t.  $|S| \leq d$ :

$C_S \leftarrow$  estimate of  $\hat{f}(S)$

• let  $h(x) \equiv \sum_{|S| \leq d} C_S \cdot X_S(x)$

• output  $\text{sign}(h)$  as hypothesis

$\left(\binom{n}{d}\right)$  of these  
reuse samples

Why does this work?

Two stages:

1) Show that  $f$  has low F.C.

$$\Rightarrow E_x [(f(x) - h(x))^2] \text{ small}$$

2) Show that  $\Pr [f(x) \neq \text{sign}(h(x))] \leq E_x [(f(x) - h(x))^2]$

↑  
Hamming dist

put together:  
 $f$  has low F.C.  
 $\Rightarrow \text{sign}(h(x))$   
is good approximation  
of  $f$

First "stage":

Thm 1 if  $f$  has  $d = \alpha(\epsilon, n)$ -F.C. then

$h$  satisfies  $E_x [(f(x) - h(x))^2] \leq \epsilon + \gamma$

with prob  $\geq 1 - \delta$  (Proved last time)

2nd "stage":

Thm 2  $f: \{\pm 1\}^n \rightarrow \{\pm 1\}$

$h: \{\pm 1\}^n \rightarrow \mathbb{R}$

then  $\Pr [f(x) \neq \text{sign}(h(x))] \leq E_{x \in \mathcal{X}} [(f(x) - h(x))^2]$

Proof.

$$E[(f(x) - h(x))^2] = \frac{1}{2^n} \sum_x (f(x) - h(x))^2 \quad \text{defn.}$$

$$P_r[f(x) \neq \text{sign}(h(x))] = \frac{1}{2^n} \sum_x \mathbb{1}_{\{f(x) \neq \text{sign}(h(x))\}}$$

compare these terms by term to get Thm.

Consider " $(f(x) - h(x))^2$ " vs. " $\mathbb{1}_{\{f(x) \neq \text{sign}(h(x))\}}$ " :

Case 1 if  $f(x) = \text{sign}(h(x))$ :

$$\mathbb{1}_{f(x) \neq \text{sign}(h(x))} = 0$$

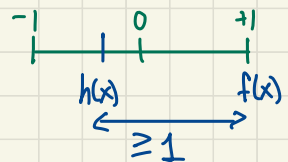
$$(f(x) - h(x))^2 \geq 0$$

Case 2 if  $f(x) \neq \text{sign}(h(x))$ :

$$\mathbb{1}_{f(x) \neq \text{sign}(h(x))} = 1$$

$$(f(x) - h(x))^2 \geq 1$$

Why? e.g.  
if  $f(x) = +1$  then in this case  $h(x) < 0$ :



So,  $\forall x$

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

(other case is analogous)

# Correctness of learning algorithm

Thm if  $\mathcal{C}$  has Fourier concentration  $d = \alpha(\epsilon, \eta)$

then there is a  $q = O\left(\frac{n^d}{\epsilon} \log \frac{n^d}{\delta}\right)$  sample  
uniform distribution learning algorithm for  $\mathcal{C}$

ie. algorithm gets  $q$  samples + with prob  $\geq 1 - \delta$   
outputs  $h'$  st.  $\Pr[f \neq h'] \leq 2\epsilon$

Pf.

run low degree alg with  $\gamma = \epsilon$

thm 1  $\Rightarrow$  get  $h$  st.  $E[(f-h)^2] \leq \epsilon + \epsilon = 2\epsilon$

output  $h' = \text{sign}(h)$

$\uparrow$   
thm 2  $\Rightarrow h'$  has error  $\leq 2\epsilon$





# Applications

1) Bounded depth decision trees

$$f(x) = \sum_{l \in \text{leaves of } T} \underbrace{f_l(x)}_{\substack{\text{fctn} \\ \text{which} \\ \text{depends on} \\ \leq \text{depth many} \\ \text{vars}}} \cdot \underbrace{\text{val}(l)}_{\text{const}}$$

$$\hat{f}(s) = \sum \text{val}(l) \underbrace{\hat{f}_l(s)}_{\substack{0 \text{ for} \\ |s| > \text{depth}}} \quad \text{linearity}$$

$$\Rightarrow \forall s \text{ st. } |s| > \text{depth}, \hat{f}(s) = 0$$

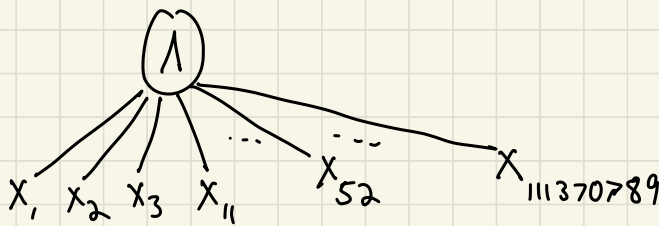
$$\text{so } O\left(\frac{n}{\varepsilon} \log \frac{n}{\delta}\right)^{\text{depth}} \text{ suffices}$$

## 2) Constant depth ckt

def. "Boolean Ckt  $C$ " is DAG

gates:  $\wedge, \vee, \neg, \perp, 0, X_1, \dots, X_n$   
operations    consts    vars

how many inputs? const, poly, unbounded?



can we compute parity of  $n$  bits  
(xor)  
in const depth?

yes! can compute any fctn on  $n$  bits  
in const depth "Karnaugh maps"

parity in const depth, poly size?

no! [Furst Saxe Sipser]  $\xi$  lemon  
Switching lemma

lemons  $\Rightarrow$  lemonade:

Thm [Hastad, Linial Mansour Nisan]

$\forall f$  computable via size  $s$  depth  $d$  ckt

$$\sum_{|S| > t} \hat{f}^2(S) \leq \alpha \quad \text{for } t = O\left(\log \frac{s}{\alpha}\right)^{d-1}$$

$$\left. \begin{array}{l} \text{take } s = \text{poly}(n) \\ d = \text{const} \\ \alpha = O(\epsilon) \end{array} \right\} \Rightarrow t = O\left(\log^d\left(\frac{n}{\epsilon}\right)\right)$$

yields  $n^{O(\log^d(\frac{n}{\epsilon}))}$  sample algorithm

(can improve to  $n^{O(\log \log n)}$  [Jackson])

(recall parity of  $s$  will have 1 large Fourier coeff of degree  $|S|$ )

### 3) Learning halfspaces

def.  $h(x) = \text{sign}(w \cdot x - \theta)$  is "halfspace function"

$$\text{sign}(y) = \begin{cases} +1 & \text{if } y \geq 0 \\ -1 & \text{o.w.} \end{cases}$$

Thm Let  $h$  be halfspace over  $\{\pm 1\}^n$   
then  $h$  has f.c.  $\alpha(\epsilon) = \frac{C}{\epsilon^2}$

$$\left(\text{i.e. } \sum_{|S| \geq \frac{C}{\epsilon^2}} \hat{h}(S)^2 \leq \epsilon\right)$$

(will prove soon)

Corr low degree alg learns halfspaces  
under unif dist with  $n^{O(1/\epsilon^2)}$   
unif. samples.

(actually  $O(n^5)$  sample algorithms exist,  
but this approach will have  
"big win" soon)

key idea:

## Noise Sensitivity

← use to bound  
Fourier  
Concentration

def. "Noise operator"  $0 < \varepsilon < 1/2$

$N_\varepsilon(x)$  = randomly flip each bit of  $x$   
with prob  $\varepsilon$

def "Noise sensitivity"

$$NS_\varepsilon(f) = \Pr_{\substack{x \in \{\pm 1\}^n \\ \uparrow \\ \text{noise}}} [f(x) \neq f(N_\varepsilon(x))]$$

## Examples

1.  $f(x) = X_1$        $NS_\varepsilon(f) = \varepsilon$

2.  $f(x) = X_1 X_2 \dots X_k$        $NS_\varepsilon(f) = \Pr [f(x) = F \wedge f(N_\varepsilon(x)) = T]$   
 $\quad \quad \quad + \Pr [f(x) = T \wedge f(N_\varepsilon(x)) = F]$

by symmetry  $\rightarrow = 2 \cdot \Pr [f(x) = T \wedge f(N_\varepsilon(x)) = F]$

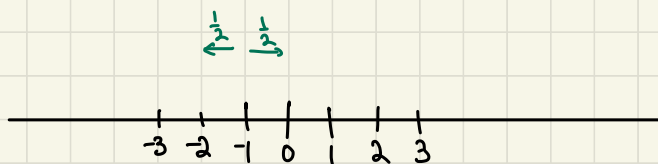
if  $\varepsilon \ll \frac{1}{k}$ ,  $\approx \frac{1}{2^{k-1}} (\varepsilon k)$   
if  $\varepsilon \gg \frac{1}{k}$ ,  $\approx \frac{1 - e^{-k\varepsilon}}{2^{k-1}}$   $\rightarrow = \frac{2}{2^k} (1 - (1-\varepsilon)^k)$

$$3. f(x) = \text{Maj}(x_1, \dots, x_n)$$

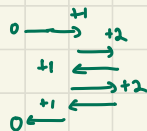
$$ns_\varepsilon(f) = O(\sqrt{\varepsilon})$$

Sketch:

$\text{Maj}(x) \sim$  random walk on line starting at 0



e.g.  $x = \begin{pmatrix} 1 & 1 & -1 & 1 & -1 \\ 1 & 2 & 1 & 2 & 1 & 0 \end{pmatrix}$



well known fact:

$$E[X_1 + X_2 + \dots + X_n]$$

$$= \sqrt{n}$$

† likely to be close to  $\sqrt{n}$

$N_\varepsilon(x) \sim$  random walk on  $\varepsilon n$  bits

each flip displaces by  $\pm 2$

(-1  $\rightarrow$  +1 or +1  $\rightarrow$  -1)

$$E[\text{displacement}] = 2\sqrt{\varepsilon n}$$

Process: take walk specified by  $X$  + continue walk according to  $N_\varepsilon(x) \cdot 2$

heuristic argument:

pretend first walk leaves us at  $\sqrt{n}$

$\Pr$  [2<sup>nd</sup> walk takes us across 0]

$$= \frac{1}{2} \Pr [2^{\text{nd}} \text{ displacement} > \sqrt{n}]$$

$$\underbrace{\frac{1}{2\sqrt{\varepsilon}} \cdot 2\sqrt{\varepsilon n}}$$

$$\leq 2\sqrt{\varepsilon} \text{ by Markov's } \neq$$

4. any LTF ( $\frac{1}{2}$  space)

$$\underline{\text{Thm}} \text{ (Peres)} \quad NS_\varepsilon(\text{LTF}) < 8.8\sqrt{\varepsilon}$$

best possible since  $NS_\varepsilon(\text{Maj}) = \Theta(\sqrt{\varepsilon})$

5. Parity fctns  $\chi_S(x)$  for  $|S|=k$

$$NS_\varepsilon(F) = \Pr [\text{odd \# bits in } S \text{ flipped by } N_\varepsilon]$$

$$= \frac{1 - (1 - 2\varepsilon)^k}{2}$$

for  $|S|=1: \varepsilon$