

## Lecture 17

### Fourier-based learning algorithms

- Fourier Concentration via Noise sensitivity
- Learning heavy Fourier coeffs (with queries)

Recall Fourier Transform:

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

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

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

=  $2 \cdot \Pr[f(x) = \chi_s(x)] - 1$

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

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

Parseval's  $1 = \langle f, f \rangle = \sum_s \hat{f}(s)^2$

for Boolean  
 $f: \{\pm 1\}^n \rightarrow \{\pm 1\}$

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

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

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

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.  $\Pr[f \neq h'] \leq 2\epsilon$

# Applications

- 1) Bounded depth decision trees
- 2) Const depth ckts
- 3) halfspaces (linear threshold fctns)

key idea:

## Noise Sensitivity

← use to bound  
Fourier  
Concentration

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

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

def "Noise sensitivity"

$$NS_\epsilon(f) = \Pr_{\substack{x \in \{0,1\}^n \\ \uparrow \\ \text{noise}}} [f(x) \neq f(N_\epsilon(x))]$$

## Examples

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

2.  $f(x) = X_1 X_2 \dots X_k$        $NS_\epsilon(f) = \frac{2}{2^k} (1 - (1 - \epsilon)^k)$

3.  $f(x) = \text{Maj}(X_1, \dots, X_n)$        $NS_\epsilon(f) = O(\sqrt{\epsilon})$

4.  $f(x)$  is any LTF       $NS_\epsilon(f) \leq 8.8\sqrt{\epsilon}$

5. Parity fctns       $NS_\epsilon(f) = \Pr[\text{odd \# bits flipped by noise}] = \frac{1 - (1 - 2\epsilon)^k}{2}$

(End review)

6. Any  $f$

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

$$NS_{\varepsilon}(f) = \frac{1}{2} - \frac{1}{2} \sum_S (1-2\varepsilon)^{|S|} \hat{f}(S)^2$$

for parity  $\chi_S$  fctns:  $\frac{1}{2} - \frac{1}{2}(1-2\varepsilon)^{|S|}$

pf. homework?

# Noise Sensitivity vs. Fourier Concentration

Thm  $\forall f: \{\pm 1\}^n \rightarrow \{\pm 1\} \quad 0 < \gamma < \frac{1}{2}$

$$\sum_{|s| \geq \frac{1}{\gamma}} \hat{f}(s)^2 < 2.32 \, ns_{\gamma}(f)$$

Pf  $2 \cdot ns_{\gamma}(f) = 1 - \sum_s (1-2\gamma)^{|s|} \hat{f}(s)^2$  *previous thm*

$$= \sum_s \hat{f}(s)^2 - \sum_s (1-2\gamma)^{|s|} \hat{f}(s)^2$$
 *Booleen Parseval*

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

$$\geq \sum_{\substack{s \text{ st.} \\ |s| \geq \gamma\gamma}} [1 - (1-2\gamma)^{\gamma\gamma}] \hat{f}(s)^2$$

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

$$\text{So } \sum_{|s| \geq \gamma\gamma} \hat{f}(s)^2 < \underbrace{\left(\frac{2}{1-e^{-2}}\right)}_{2.32} \cdot ns_{\gamma}(f)$$

▣

Corr for halfspace  $h: \{\pm 1\}^n \rightarrow \{\pm 1\}$

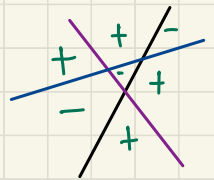
$$\sum_{|s| \geq O(\frac{1}{\epsilon^2})} \hat{f}(s) \leq \epsilon$$

(pf omitted - some calculations + bound on NS)

$\Rightarrow$  can learn any halfspace from  $n^{O(1/\epsilon^2)}$   
random examples

(actually can do a lot better)

Corr any function of  $k$  halfspaces  
can be learned with  $n^{O(k/\epsilon^2)}$  samples



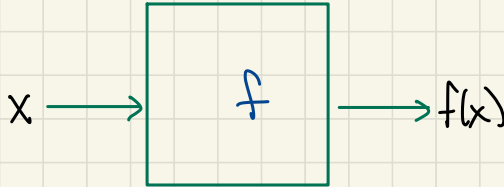
e.g. parity  
of  $k$  vars,  
 $\wedge$  of  $\frac{k}{2}$  spaces

Pf idea noise sensitivity  $\leq 8.8 k \epsilon$  by union bound.

# Learning Heavy Fourier Coeffs

[Goldreich Levin]

[Kushilevitz Mansour]



not just low degree  $S$



all close linear fctns

Given  $f, \theta$

• Output all coeffs  $S$  st.  $|\hat{f}(S)| \geq \theta$

• Only output  $S$  st.  $|\hat{f}(S)| \geq \frac{\theta}{2}$  ← no junk

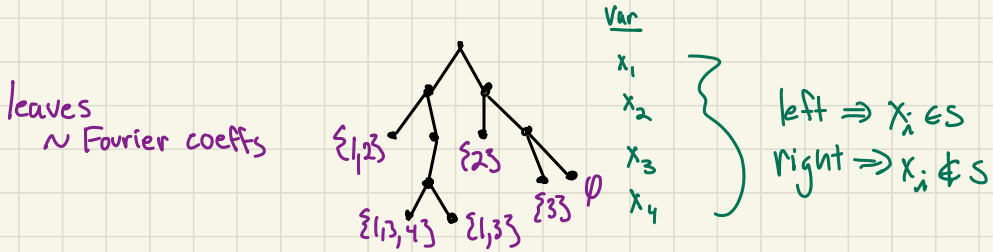
Probably can't do it with only random examples

What if can query  $f$  at any input?

Observation if there is a really big  $\hat{f}(S)$ , can use self-correction to find it

*say  $> 3/4$*

# Main Idea: "exhaustive search with good pruning"



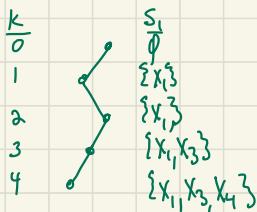
ONLY OUTPUT THOSE THAT REACH BOTTOM LEVEL

recursive algorithm:

- each node  $\sim$  setting of  $x_1, \dots, x_i$
- estimate "total energy" of subtrees  $x_1 \dots x_i^o(x_{i+1} = +1)$   
 $\& x_1 \dots x_i^o(x_{i+1} = -1)$
- only go down paths with high enough energy

How to prune?

Define quantity:



Fix  $0 \leq k \leq n$   
 $S_1 \subseteq [k]$

current "level" of search  
 current "node" of search

↓ doesn't depend on 1<sup>st</sup> k bits

2<sup>k</sup> such fctns (for each S<sub>1</sub>)

$$f_{k, S_1}: \{\pm 1\}^{n-k} \rightarrow \mathbb{R}$$

all Fourier coeffs which agree on first k elements

$$\text{s.t. } f_{k, S_1}(x) = \sum_{T_2 \subseteq \{k+1, \dots, n\}} \hat{f}(S_1, T_2) \chi_{T_2}(x)$$

all extensions of S<sub>1</sub> to indices in {k+1...n}

could be S<sub>1</sub>, T<sub>2</sub> but no need since  $\chi_{S_1, T_2} = \chi_{S_1} \cdot \chi_{T_2}$  same for all

notation: index 1 → prefix 2 → suffix

Sanity Checks:

1) k=0

$$f_{0, \emptyset}(x) = \sum_{T_2 \subseteq [n]} \hat{f}(T_2) \chi_{T_2}(x) = f(x)$$

↑ since k=0      ↑ since S<sub>1</sub>=∅

2) k=n

$$f_{n, S_1}(x) = \hat{f}(S_1) \cdot \chi_{\emptyset}(x)$$

↑ sum over T<sub>2</sub>=∅      since T<sub>2</sub>=∅

Plan Only go down paths with  $E[f^2(x)] \geq \theta^2$   
 $K_S$

1. can we compute it?

2. does it bring us to right leaves?

- do we get to all heavy leaves?

- do we get junk? (light leaves)

3. how many paths do we take?

lots of dead ends?

is runtime good?

# Not too many paths! (answer to 3)

Lemma "not too many" ← at any stage in algorithm

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

$$(1) \leq \frac{1}{\theta^2} \text{ s's satisfy } |\hat{f}(s)| \geq \theta$$

$$(2) \forall 0 \leq k \leq n, \leq \frac{1}{\theta^2} \text{ fctns } f_{k, s_1}$$

$$\text{have } E_x [f_{k, s_1}^2(x)] \geq \theta^2$$

Pf

$$(1) \text{ Boolean Parseval's } 1 = \sum_s \hat{f}(s)^2$$

$$\text{so if } > \frac{1}{\theta^2} \text{ s's satisfy } |\hat{f}(s)| \geq \theta$$

$$\text{then } \sum_s \hat{f}(s)^2 > \frac{1}{\theta^2} \cdot \theta^2 > 1$$

→ ←

(2) For given  $k$ :

Claim:  $\forall k, s_1 \leq k$

$$E_x [f_{k, s_1}^2(x)] = \sum_{T_2 \in \{k+1, \dots, n\}} \hat{f}(s_1, \nu_{T_2})^2$$

pf of claim:

$$E_x [f_{k, s_1}^2(x)] = E_x \left[ \left( \sum_{T_2} \hat{f}(s_1, \nu_{T_2}) \chi_{T_2}(x) \right)^2 \right] \quad \text{def.}$$

$$= E_x \left[ \sum_{\substack{T_2, T_2' \\ \in \{k+1, \dots, n\}}} \hat{f}(s_1, \nu_{T_2}) \cdot \hat{f}(s_1, \nu_{T_2'}) \chi_{T_2}(x) \chi_{T_2'}(x) \right]$$

$$= \sum_{T_2, T_2'} \hat{f}(s_1, \nu_{T_2}) \hat{f}(s_1, \nu_{T_2'}) E \left[ \chi_{T_2}(x) \cdot \chi_{T_2'}(x) \right]$$

$$= \sum_{T_2} \hat{f}(s_1, \nu_{T_2})^2$$

$$\begin{aligned} &= 1 \text{ if } T_2 = T_2' \\ &= 0 \text{ o.w.} \end{aligned}$$

Using Claim:

$$1 = \sum_S \hat{f}(s)^2 \stackrel{\text{Booleam Parseval's}}{=} \sum_{S_1 \subseteq K} \sum_{T_2 \subseteq [K \setminus S_1]} \hat{f}(S_1 \cup T_2)^2$$

$$= \sum_{S_1} E_x [f_{K, S_1}^2(x)] \quad \text{claim}$$

So  $\leq \frac{1}{\theta^2}$   $S_1$ 's can have  $E_x [f_{K, S_1}^2(x)] > \theta^2$

~~□~~

Does algorithm bring us to good leaves?

(answer to 2)

Fact: "not missing out"  $\Rightarrow$  find all big  
Fourier coeffs

For any  $S_1$ , if  $\exists T_2$  st.

$$|\hat{f}(S_1 \vee T_2)| > \theta$$

then  $E_x [f_{k_{S_1}}^2(x)] = \sum_{T_2} \hat{f}(S_1 \vee T_2)^2$  via claim

$$\geq \theta^2$$

$\Rightarrow E_x [f_{k_{S_1}}^2(x)]$  is a good measure  
to use when deciding whether  
to investigate subtree!

So we find all good leaves  
we don't spend too much time

but do we output junk?

No junk (answer to 2)

Simple fix:

For each "candidate" to  $S$ ,  
estimate its Fourier coeff & make sure it is  
big enough before outputting it

Can we estimate  $f_{k, s_1}(x)$ ?

recall:

$$f: \{\pm 1\}^n \rightarrow \{\pm 1\}$$
$$0 \leq k \leq n$$
$$S_1 \subseteq [k]$$

(answer to 1)

$$f_{k, s_1}(x) = \sum_{T_2 \subseteq [kH, \dots, n]} \hat{f}(S_1, VT_2) \chi_{T_2}(x)$$

$$E_x [f_{k, s_1}(x)^2] = \sum_{T_2 \subseteq [kH, \dots, n]} \hat{f}(S_1, VT_2)^2$$

Bad idea: estimate each

$$\hat{f}(S_1, VT_2) \quad \forall T_2 \leftarrow \text{too much time}$$

Another challenge:

heavy  $\hat{f}(S_1, VT_2)$  can differ  
greatly depending on how  
you fix vars in  $S_1$

e.g.  $S_1 = \{x_i\}$

$$f(x) = \begin{cases} 1 & \text{if } x_i = 1 \\ \prod_i x_i & \text{if } x_i = -1 \end{cases}$$

Can we estimate  $f_{k,s_1}(x)$ ?

recall:

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

$$0 \leq k \leq n$$

$$s_1 \in [k]$$

(answer to 1)

$$f_{k,s_1}(x) = \sum_{T_2 \subseteq \{k+1..n\}} \hat{f}(s_1, T_2) \chi_{T_2}(x)$$

$$E_x [f_{k,s_1}(x)^2] = \sum_{T_2 \subseteq \{k+1..n\}} \hat{f}(s_1, T_2)^2$$

Lemma " $f_{k,s_1}(x)$  Estimation lemma"

for  $x \in \{\pm 1\}^{n-k}$

$$f_{k,s_1}(x) = E_{y \in \{\pm 1\}^k} [f(yx) \chi_{s_1}(y)]$$

↑ concatenation
} Use this to estimate  $f_{k,s_1}(x)$

"agreement"

think of as  
fctn of  $y$  since  
 $x$  is fixed  
throughout

pf  
Fourier representation  $\Rightarrow$

$$f(yx) = \sum_T \hat{f}(T) \chi_T(yx)$$

$$T = T_1 \cup T_2 \quad T_1 \subseteq [k] \quad T_2 \subseteq \{k+1..n\}$$

$$\text{so } \chi_T(yx) = \chi_{T_1}(y) \cdot \chi_{T_2}(x)$$

$$E_y[f(y|x)\chi_{S_1}(y)]$$

$$= E_y\left[\sum_{T_1} \sum_{T_2} \hat{f}(T_1, T_2) \chi_{T_1}(y) \chi_{T_2}(x) \cdot \chi_{S_1}(y)\right] \quad \text{above}$$

$$= \sum_{T_1} \sum_{T_2} \hat{f}(T_1, T_2) \chi_{T_2}(x) E_y[\chi_{T_1}(y) \chi_{S_1}(y)]$$

$$= \begin{cases} 0 & \text{if } T_1 \neq S_1 \\ 1 & \text{if } T_1 = S_1 \end{cases}$$

$$= \sum_{T_2} \hat{f}(S_1, T_2) \chi_{T_2}(x)$$

$$= f_{k, S_1}(x)$$



## Overall Algorithm:

Pick random  $x$ 's  $\in \{\pm 1\}^{n-k}$

For each  $x$ , pick random  $y$ 's  $\in \{\pm 1\}^k$

Estimate  $E_y [f(yx) \chi_{S_i}(y)]$

which gives estimate of  
 $f_{k, S_i}(x)$

Estimate  $E_x [f_{k, S_i}(x)^2]$

Chernoff + samples

$\Rightarrow$  Can get  $\gamma$ -additive estimate with  
prob  $\geq 1 - \delta$  in  $O\left(\frac{1}{\gamma^2} \log \frac{1}{\delta}\right)$  queries

$\Rightarrow$  Can get  $\frac{\gamma^2}{2}$ -additive est of  $f_{k, S_i}^2(x)$

Thm  $\forall \theta > 0$ , KM-alg outputs

$S = \{s_1, \dots, s_k\}$  st.  $k = O(1/\theta^2)$  + with prob  $\geq 1 - \delta$

$$\forall s_i \in S \quad |\hat{f}(s_i)| \geq \frac{\theta}{2} \quad \text{no junk}$$

$$\forall s \notin S \quad |\hat{f}(s)| \leq \theta \quad \text{no misses}$$

+ query time is poly  $(n, \frac{1}{\theta}, \log \frac{1}{\delta})$

Pf. if  $\hat{f}(s) < \frac{\theta}{2}$ ,  $\hat{f}(s)^2 \leq \frac{\theta^2}{4}$

+ test  $\times$  prevents it from being output

if  $\hat{f}(s) > \theta$ ,  $\forall k$  if  $s_1, \dots, s_k$  agree on  $[1..k]$

"not missing out" fact  $\Rightarrow$

$$E_x [f_{k, s}^2(x)] \geq \theta^2$$

total # nodes explored at level  $k \leq \frac{1}{\theta^2}$

$$\Rightarrow \text{total nodes explored} \leq \frac{n}{\theta^2}$$