

## Lecture 17

Fourier-based learning algorithms

- learning one Fourier coeff
- the low degree algorithm

Review from last time:

def uniform distribution learning algorithm for concept class  $\mathcal{C}$  is algorithm  $A$  st.

- $A$  given  $\epsilon, \delta > 0$   
access to  $E_X(f)$  for  $f \in \mathcal{C}$
- $A$  outputs  $h$  st. with prob  $\geq 1 - \delta$   
error( $h$ ) w.r.t.  $f$  is  $\leq \epsilon$   
 $h$  is  $\epsilon$ -close to  $f$

Parameters of interest:

- $m$  # samples used by  $A$  "Sample Complexity"
- $\epsilon$  accuracy parameter
- $\delta$  confidence parameter
- runtime? hope for poly( $\log(\text{domainsize}), \frac{1}{\epsilon}, \frac{1}{\delta}$ )

• description of  $h$ ?

• should it be similar to description of fctns in  $\mathcal{C}$ ? "proper learning"

• at least should be relatively  
compact + efficient to evaluate  
 $O(\log |\mathcal{C}|)$

# 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!

Pf. Chernoff +  $\hat{f}(s) = 2 \Pr_x [f(x) = \chi_s(x)] - 1$   
estimate this



Can we find **any** or **all** heavy coefficients?

there are exponentially many coeffs

Can use same samples to estimate each coeff, but must union bound prob of error (error = bad approx) on any of them.

Need  $\delta \ll \frac{1}{2^n}$ , which needs

$O(\frac{1}{\delta^2} \cdot n)$  samples, but

exponential runtime. ← turns out that queries help a lot

What if we "know where to look" for heavy coeffs?

e.g. all heavy coeffs are in "low degree" coeffs? If so, can search!

# Fourier Representations of Important Examples

1)  $\overline{\text{AND}}$  on  $T \subseteq N$  st.  $|T| = k$

$$\overline{\text{AND}}(x) = \begin{cases} 1 & \text{if } \forall i_j \in T = \{i_1, \dots, i_k\} \\ & x_{i_j} = -1 \\ -1 & \text{o.w.} \end{cases}$$

$$f(x) = \begin{cases} 1 & \text{if } \forall i \in T \quad x_i = -1 \\ 0 & \text{o.w.} \end{cases} \quad \left. \vphantom{f(x)} \right\} \text{AND over range } \{0,1\}$$

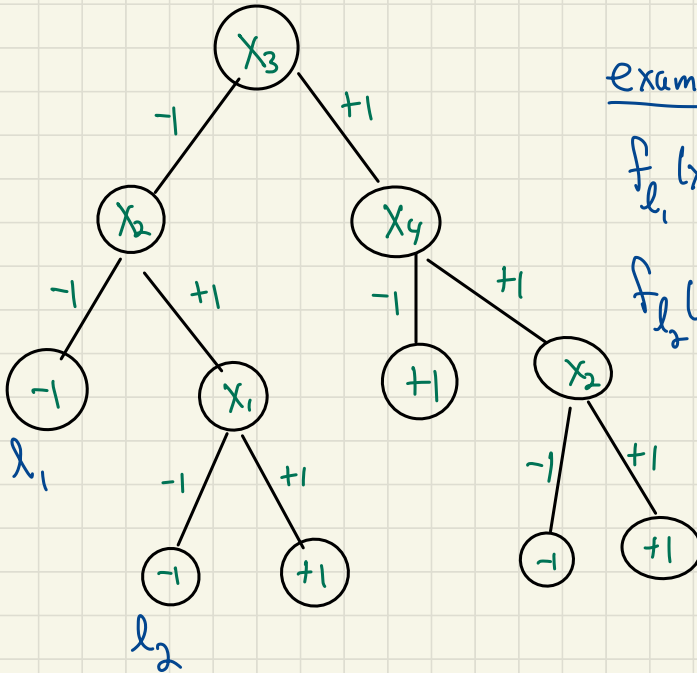
$$= \frac{(1-x_{i_1})}{2} \cdot \frac{(1-x_{i_2})}{2} \cdot \dots \cdot \frac{(1-x_{i_k})}{2}$$

$$= \sum_{S \subseteq T} \frac{(-1)^{|S|}}{2^k} \chi_S$$

$$\overline{\text{AND}} = 2f(x) - 1 = -1 + \frac{2}{2^{|T|}} + \sum_{\substack{S \subseteq T \\ S \neq \emptyset}} \frac{(-1)^{|S|}}{2^{|T|-1}} \chi_S$$

Note: all Fourier coeffs containing vars not in  $T$  are 0

## 2) Decision trees



examples

$$f_{l_1}(x) = \frac{(1-x_3)}{2} \cdot \frac{(1-x_2)}{2}$$

$$f_{l_2}(x) = \frac{(1-x_3)}{2} \cdot \frac{(1+x_2)}{2} \cdot \frac{(1-x_1)}{2}$$

First consider path functions:

$$f_l(x) = \prod_{i \in V_l} \frac{(1 \pm x_i)}{2}$$

← left or right

← vars visited on path to leaf l

$$= \frac{1}{2^{|V_l|}} \sum_{S \subseteq V_l} (\pm 1)^{|S|} x_S$$

←  $(-1)^{\# \text{ left turns taken in } S}$

$$= \begin{cases} 1 & \text{if } x \text{ takes path to } l \\ 0 & \text{o.w.} \end{cases}$$

So  $f(x) = \sum_{l \in \text{leaves of } T} f_l(x) \cdot \text{val}(l)$

exactly one of these is 1.  
all others are 0.

Comment only coeffs corresponding to  $S$  st.  
 $|S| \leq \text{max path length}$  have a hope of  
being non-zero.



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

2)  $f = \text{AND}$  on  $T \subseteq \{1..n\}$  has  $\log(\frac{4}{\epsilon})$ -F.C.

- all  $\hat{f}(s)^2 = 0$  for  $|s| > |T|$
- if  $|T| \leq \log \frac{4}{\epsilon}$  then ✓
- if  $|T| \geq \log \frac{4}{\epsilon}$  then:

$$\begin{aligned}\hat{f}(\varphi)^2 &= (1 - 2 \Pr(f(x) \neq \chi_{\varphi}(x)))^2 \\ &= \left(1 - \frac{2}{2^{|\pi|}}\right)^2 \\ &> 1 - \epsilon\end{aligned}$$

$$\text{so } \sum_{s \neq \emptyset} \hat{f}(s)^2 \leq \epsilon \quad \& \quad f \text{ has } 0\text{-F.C.}$$

Idea: can we approximate  $f$   
by only considering  
low degree Fourier  
coeffs?

# Low degree algorithm

approximates fctns with  $d \equiv \Omega(\varepsilon, n)$  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  $\geq 1 - \delta$

Pf (1) each low degree Fourier coeff is well approximated:

Claim with prob  $\geq 1 - \delta$ ,  $\forall s$  st  $|s| \leq d$

$$|C_s - \hat{f}(s)| \leq \gamma \quad \text{for } \gamma \leftarrow \sqrt{\frac{\epsilon}{nd}}$$

Pf of claim (Chernoff + union bnd)

note,  $\frac{1}{\gamma^2} = \frac{n^d}{\tau}$

Chernoff bnd  $\Rightarrow$

$$O\left(\frac{n^d}{\tau} \ln \frac{n^d}{\delta}\right) = O\left(\frac{1}{\gamma^2} \ln \frac{n^d}{\delta}\right) \text{ samples}$$

yields  $\Pr [ |C_s - \hat{f}(s)| > \gamma ] < \frac{\delta}{n^d}$

union bnd over all  $\binom{[d]}{s}$   $s$ 's  $\Rightarrow$

$$\Pr [ \exists s \text{ s.t. } |C_s - \hat{f}(s)| > \gamma ] < \delta$$

(2) all low degree Fourier coeffs well approx  $\Rightarrow$  low  $l_2$  error:

Assume  $\forall s$  s.t.  $|s| \leq d$ ,  $|C_s - \hat{f}(s)| \leq \gamma$ :

define  $g(x) \equiv f(x) - h(x)$

Fourier transform linear  $\Rightarrow \forall s \quad \hat{g}(s) = \hat{f}(s) - \hat{h}(s)$

by defn,  $\forall s$  s.t.  $|s| > d$ ,  $\hat{h}(s) = 0 \Rightarrow \hat{g}(s) = \hat{f}(s)$

$$|s| \leq d, \hat{h}(s) = C_s$$

$$\Rightarrow \hat{g}(s) = \hat{f}(s) - C_s$$

$$\text{so } \sum \hat{g}(s)^2 \leq \gamma^2$$

$$\text{So } E[(f(x) - h(x))^2] = E[g(x)^2]$$

$$= \sum_s \hat{g}(s)^2 \quad \text{Parseval}$$

$$= \underbrace{\sum_{|s| \leq d} \hat{g}(s)^2}_{\leq \gamma^2} + \underbrace{\sum_{|s| > d} \hat{g}(s)^2}_{\leq \varepsilon \text{ by F.C.}}$$

$$\leq \gamma + \varepsilon$$



2nd "stage":

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

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

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