# Lecture 13:

## Learning via Fourier Coeffs

- Some fctns & their Fourier representation

- the low degree algorithm

- applications

Fourier    Representations  of    Important Examples

Two   examples

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

$$\overline{\text{AND}} \, (X_{i_j} \cdots X_{i_k}) = 1 \qquad \text{if} \qquad \forall \, i_j \in T = \{ i_1 \cdots i_k \}$$

$$X_{i_j} = -1$$

$$-1 \qquad\qquad\qquad o.w.$$

define  $\overset{\text{AND}}{f}(x) = \begin{cases} 1 & \text{if} \quad \forall \, i \in T \quad X_i = -1 \\ 0 & o.w. \end{cases}$    corresponds to AND fctn over $\{0,1\}$

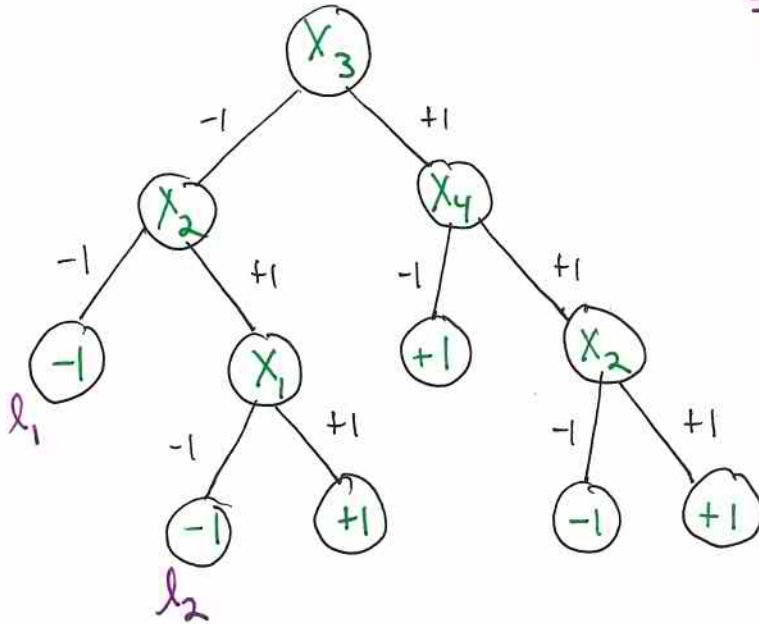$$= \frac{(1 - X_{i_1})}{2} \cdot \frac{(1 - X_{i_2})}{2} \cdots \frac{(1 - X_{i_k})}{2}$$

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

+ so    $\overline{\text{AND}}(x) = 2 f(x) - 1$

$$= -1 + \frac{2}{2^k} \cdot 1 + \sum_{\substack{S \subseteq T \\ |S| > 0}} \frac{(-1)^{|S|}}{2^{k-1}} \chi_S$$

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

2) Decision trees



examples

$$f_{\ell_1}(x) = \frac{(1-X_3)}{2} \cdot \frac{(1-X_2)}{2}$$

$$f_{\ell_2}(x) = \frac{(1-X_3)}{2} \frac{(1+X_2)}{2} \frac{(1-X_1)}{2}$$

First, consider path fctns:

$$f_\ell(x) = \prod_{i \in V_\ell} \frac{(1 \pm X_i)}{2}$$

← left or right

vars visited on path to leaf ℓ

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

$$f_\ell(x) = \begin{cases} 1 & \text{if } X \text{ takes path } \ell \\ 0 & \text{o.w.} \end{cases}$$

$$= \frac{1}{2^{|V_\ell|}} \sum_{S \subseteq V_\ell} (\pm 1) X_S$$

So   $f(x) = \sum_{\ell \in \text{leaves of } T} f_\ell(x) \, val(\ell)$

exactly one of these is 1 others are 0

Comment only coeffs corresponding to S s.t. $|S| \leq$ max path length can be non zero.

# Low degree algorithm

__def__  $f: \{\pm 1\}^n \to \mathbb{R}$  has  $\alpha(\varepsilon, n)$- Fourier concentration

if  $\displaystyle\sum_{\substack{S \subseteq [n] \\ st. \\ |s| > \alpha(\varepsilon, n)}} \hat{f}(s)^2 \leq \varepsilon$   $\forall \; 0 < \varepsilon < 1$

for   Boolean $f$,   this implies  $\displaystyle\sum_{\substack{S \subseteq [n] \\ st. \\ |s| \leq \alpha(\varepsilon, n)}} \hat{f}(s)^2 \geq 1 - \varepsilon$

# examples

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

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

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

   • all  $\hat{f}(s)^2 = 0$  for  $|s| > |T|$

   • if  $|T| \leq \log\frac{4}{\varepsilon}$  then  ✓

   • if  $|T| \geq \log\frac{4}{\varepsilon}$  then  :

   $\hat{f}(\phi)^2 = (1 - 2 Pr(f(x) \neq \chi_\varphi(x)))^2 = \left(1 - \frac{2}{2^{|T|}}\right)^2 > 1 - \varepsilon$

   so  $\displaystyle\sum_{S \neq \varphi} \hat{f}(s)^2 \leq \varepsilon$  + $f$ has  0-F.C.

Now, lets approximate fctns with $d \equiv \alpha(\varepsilon, n)$ F.c.:

### Low Degree Algorithm

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

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

- $C_s \leftarrow$ estimate of $\hat{f}(s)$ (for each $s$ s.t. $|s| \leq d$)

- output $h(x) = \sum_{|s| \leq d} C_s \chi_s(x)$

$\leq \binom{n}{d}$ of these
Can reuse same samples for each!

$\not\!\!+$ use $\operatorname{sign}(h(x))$ as hypothesis!

Why does this work?

Two stages:
1) show that if $f$ has low F.C. $\swarrow^{\frac{L_2 \text{ dist}}{2^n}}$
   then $E_x[(f(x) - h(x))^2]$ small

2) show that $\Pr[f(x) \neq \operatorname{sign}[h(x)]] \leq E_x[(f(x) - h(x))^2]$
   $\uparrow$
   Hamming dist

__Thm__ if   f   has   $d = \alpha(\varepsilon, n) - F.c,$   then

h   satisfies   $E_x\left[ (f(x) - h(x))^2 \right] \leq \varepsilon + \Upsilon$

with   prob $\geq 1 - \delta$

__Pf__

__Claim__   with   prob $\geq 1 - \delta$,   $\forall$ s   s.t.   $|s| \leq d$,   $|c_s - \hat{f}(s)| \leq \gamma$

for   $\gamma \leftarrow \sqrt{\dfrac{\Upsilon}{nd}}$

Pf of claim

note,   $\dfrac{1}{\gamma^2} = \dfrac{nd}{\Upsilon}$

Chernoff bnd $\Rightarrow$ $O\left( \dfrac{n^d}{\Upsilon} \ln \dfrac{n^d}{\delta} \right) = O\left( \dfrac{1}{\gamma^2} \ln \dfrac{n^d}{\delta} \right)$ samples

yields $Pr\left[ |c_s - \hat{f}(s)| > \gamma \right] < \dfrac{\delta}{n^d}$

union bnd $\Rightarrow$ $Pr\left[ \exists\ s\ \text{s.t.}\ |c_s - \hat{f}(s)| > \gamma \right] < \delta$

only $\binom{n}{d} < n^d$ such

s's of size $\leq d$

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

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

$|s| \leq d,\ \hat{h}(s) = c_s$   $\Rightarrow$   $\hat{g}(s) = \hat{f}(s) - c_s$

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

so $\quad E\left[ (f(x) - h(x))^2 \right] = E\left[ g(x)^2 \right]$

$$= \sum_S \hat{g}(s)^2 \qquad \text{Parseval}$$

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

$$\leq \Upsilon + \varepsilon \qquad \blacksquare$$

__Thm__ $\quad f : \{\pm 1\}^n \to \{\pm 1\}$

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

then $\quad Pr\left[ f(x) \neq \text{sign}(h(x)) \right] \leq E\left[ (f(x) - h(x))^2 \right]$

__Pf.__ $\quad E\left[ (f(x) - h(x))^2 \right] = \frac{1}{2^n} \sum_x (f(x) - h(x))^2 \qquad \text{defn} \quad \Big\}$ show term by term

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

But $\quad$ if $f(x) = \text{sign}(h(x))$

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

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

if $\quad f(x) \neq \text{sign}(h(x))$

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

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

So $\forall x, \quad \left( f(x) - h(x) \right)^2 \geq \mathbb{1}_{f(x) \neq \text{sign}(h(x))} \qquad \blacksquare$

Correctness of learning algorithm :

<u>Thm.</u> if $C$ has Fourier concentration $d = \alpha(\varepsilon, n)$
then there is a $q = O(\frac{n^d}{\varepsilon} \log \frac{n^d}{\delta})$ sample
uniform distribution learning algorithm for $C$
ie. algorithm gets $q$ samples & with prob $\geq 1 - \delta$
outputs $h'$ s.t. $\Pr[f \neq h'] \leq 2\varepsilon$

<u>Pf.</u> run low degree alg with $\gamma = \varepsilon$
get $h$ st. $E[(f-h)^2] \leq \varepsilon + \varepsilon = 2\varepsilon$
output $\text{sign}(h)$
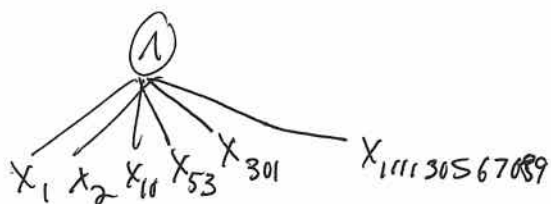
# Applications

1) Bounded depth decision trees

$$f(x) = \sum_{\ell \in \text{leaves of } T} f_\ell(x) \cdot \text{val}(\ell)$$

$\underbrace{}_{\text{const}}$

$\underbrace{}_{\substack{\text{fctn which} \\ \text{depends on} \leq \text{depth} \\ \text{many vars}}}$

by linearity, $\hat{f}(s) = \sum \text{val}(\ell) \cdot \hat{f}_\ell(s)$ which is $0$ if $|s| > \text{depth}$
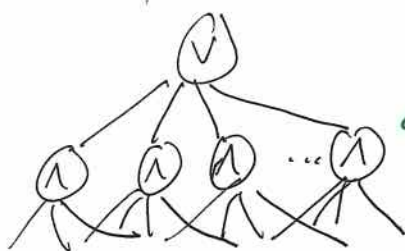
2) Constant depth ckts:

<u>Def</u>. "Boolean Ckt C" is a DAG

gates: $\wedge, \vee, \neg, 1, 0, X_1 \cdots X_n$
and, or, not, or $\pm 1$, vars
how many inputs?
const? poly? unbounded?



$X_1 \quad X_2 \quad X_{10} \quad X_{53} \quad X_{301} \quad X_{1111305567089}$

Can we compute parity of n bits in const depth?

yes! can compute any n-bit fctn in constdepth



each "$\wedge$" picks an arbitrary sat setting & checks if input matches

Can we compute parity of n bits in const

depth, poly size?

No! [Furst Saxe Sipser] } lemma

Switching lemma

## Lemons ⟹ Lemonade

Thm [Hastad, Linial Mansour Nisan]

prove via
random
restrictions
as in
[FSS] parity
result

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

$$\sum_{|s| > t} \hat{f}^2(s) \leq \alpha \qquad \text{for} \qquad t = O\left(14 \log \frac{2s}{\alpha}\right)^{d-1}$$

Take Advanced Complexity!

take $s = \text{poly}(n)$
$d = \text{const}$
$\alpha = O(\varepsilon)$ $\Big\} \Rightarrow t = O\left(\log^d\left(\frac{n}{\varepsilon}\right)\right)$

Gives $n^{O\left(\log^d\left(\frac{n}{\varepsilon}\right)\right)}$ sample query algorithm

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