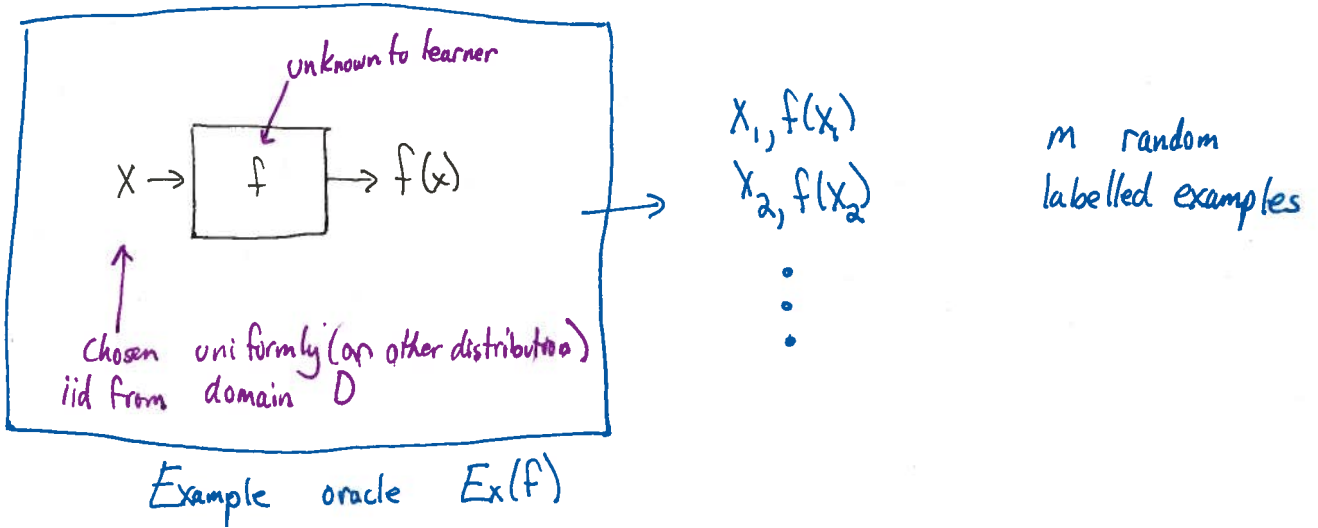


## Learning via Fourier Coeffs

- Some fctns & their Fourier representation
- the low degree algorithm
- applications

# Learning

Learn from random, uniform examples: How do we formalize?



After seeing several examples, learner should output hypothesis  $h$ .

- hopefully  $h=f$
- is that asking too much?
- how about  $\text{dist}(h, f) < \epsilon$ ?

what is distance?

e.g.  $\Pr_{x \in D} [h(x) \neq f(x)]$ ?

but then what distribution on  $D$ ?

Today: uniform

In general: match distribution of examples

Valiant's PAC model

"Probably Approximately Correct"

def. given hypothesis  $h$ , error of  $h$  wrt  $f$  is

$$\text{error}(h) = \Pr_{x \in D} [f(x) \neq h(x)]$$

Note: this is defn wrt uniform. In general, this is

often will use:  
 $f$  is  $\epsilon$ -close to  $h$  wrt  $D$  if  $\Pr_{x \in D} [f(x) \neq h(x)] \leq \epsilon$

the same as distribution on  $D$  from example one

Note if  $f$  is arbitrary, there is nothing you can do! (ie. can't learn a random fctn)  
 However, if you know something about  $f$ , there may be hope!

What if you know that  $f$  is from a family of functions

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

- $A$  is given  $\epsilon, \delta > 0$  access to  $E_x(f)$  for  $f \in \mathcal{C}$
- $A$  outputs  $h$  st. with prob  $\geq 1 - \delta$  error( $h$ ) wrt.  $f$  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  $(\log(\text{domain size}), \frac{1}{\epsilon}, \frac{1}{\delta})$
- description of  $h$ ?
  - should it be similar to description of  $f$ ?  
(proper learning)
  - at least should be relatively compact  
 $O(\log |C|)$  + efficient to evaluate

## Remarks

- as before, dependence on  $\delta$  needn't be more than  $O(\log(1/\delta))$ .  
why?
- Uniform case is special case of PAC-model:
  - Given  $EX_{\mathcal{D}}(f)$  for unknown  $\mathcal{D}$
  - output  $h$  with small error according to same  $\mathcal{D}$  (some  $\mathcal{D}$  can be harder than others)

Ignoring RuntimeOccam's Razor

learning is easy!

i.e. can easily achieve small sample complexity

Brute Force Algorithm

- Draw  $M = \frac{1}{\epsilon} (\ln |\mathcal{C}| + \ln \frac{1}{\delta})$  uniform examples
- Search over all  $h \in \mathcal{C}$  until find one that labels all examples correctly & output it.  
(choose arbitrarily if  $\geq 1$  such  $h$  works)

Behavior:

What should behavior be?

- $f$  is a good thing to output ✓
- what is a bad thing to output?

$h$  is "bad" if  $\text{error}(h) \text{ wrt } f \geq \epsilon$

$\Pr$  [bad  $h$  consistent with examples]

$$\leq (1-\epsilon)^M$$

$\Pr$  [any bad  $h$  consistent with examples]

$$\leq |\mathcal{C}| (1-\epsilon)^M \quad \leftarrow \text{union bound}$$

$$\leq |\mathcal{C}| (1-\epsilon)^{\frac{1}{\epsilon} (\ln |\mathcal{C}| + \ln \frac{1}{\delta})}$$

$$\leq \delta$$

$\therefore$  unlikely to output any bad  $h$

[Does the Bible really predict JFK's assassination?]

Comments

• proof didn't use anything special about uniform distribution

works for any  $\mathcal{D}$ ,  
as long as error defined w.r.t. same  $\mathcal{D}$  as  
sample generator

• once we have a good  $h$

1) can predict values of  $f$  on new

random inputs since  $\Pr_{x \in \mathcal{D}} [f(x) = h(x)] \geq 1 - \delta$   
according to  $\mathcal{D}$

2) can compress description of samples

$(x_1, f(x_1)) (x_2, f(x_2)) \dots (x_m, f(x_m))$   $m(\log |D| + \log |R|)$   
range of  $f$

↓

$x_1 \dots x_m$ , description of  $h$   $m \log |D| + \log |C|/k$

so learning, prediction & compression are related.

learning  $\Rightarrow$  prediction & compression

formal relations in other direction too

Occam's Razor: simplest explanation is best

# An efficient learning algorithm

$C^0$  = conjunctions over  $\{0,1\}^n$

ie.  $f(x) = x_i x_j \bar{x}_k$

• can't hope for 0-error from subexponential # of random examples

eg. how to distinguish  $f(x) = x_i \dots x_n$   
from  $f(x) = 0$ ?

• Brute force:  $M = \frac{1}{\epsilon} (\ln(2^n) + \ln \frac{1}{\delta})$  examples takes much time

• Poly time algorithm:

• draw  $\text{poly}(1/\epsilon)$  random examples to estimate

$\Pr[f(x)=1]$  to additive error  $\pm \frac{\epsilon}{4}$

if estimate  $< \epsilon/2$ , output " $h(x)=0$ "

• since estimate  $\geq \epsilon/2$  + error  $\leq \epsilon/4$

$\Pr[f(x)=1] \geq \epsilon/4$

so, every  $O(1/\epsilon)$  examples see new random "positive" example (expected)

} just look at these

• in set of positive examples

let  $V = \{ \text{vars set same way in each example} \}$

output  $h(x) = \bigwedge_{i \in V} x_i^{b_i}$

$\leftarrow b_i$  tells us if  $i$  complemented or not

behavior of algorithm:

for  $i$  in conjunction:

must be set same way in each  
positive example  $\Rightarrow$  in  $V$

for  $i$  not in conjunction:

$\Pr [i \in V] \leq \Pr [i \text{ set same in each  
of } k \text{ positive examples}]$

$$\leq \frac{1}{2^{k-1}}$$

$\Pr [\text{any } i \text{ that not in conjunction manages to survive}]$

$$\leq \frac{n}{2^{k-1}}$$

$$\leq \delta \quad \text{if pick } k = \log \frac{n}{\delta}$$

So  $\Omega(\log \frac{n}{\delta})$  positive examples

+  $\Omega(\frac{1}{\epsilon} \log \frac{n}{\delta})$  total examples suffice!



## Learning via Fourier Representation

learning algorithms based on estimating Fourier representation of fctn  $f$  (similar to poly interpolation)

### Approximating one Fourier coefficient:

lemma can approx any specific Fourier coeff  $s$  to w/in additive  $\gamma$  (i.e.  $|\text{output} - \hat{f}(s)| \leq \gamma$ ) with prob  $\geq 1 - \delta$  in  $O(\frac{1}{\gamma^2} \log \frac{1}{\delta})$  samples

Note no queries needed!!

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

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

there are exponentially many coefficients.

Can use same samples for all coeffs, but must Union bnd prob of error on any of them

Using  $\delta = \frac{1}{2^n}$ , gives  $O(\frac{1}{\gamma^2} \cdot n)$  samples, but exp runtime.

queries can help a lot!

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

e.g. all heavy coeffs are in "low degree" coeffs?

If so, can search!

# Fourier Representations of Important Examples

Two examples

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

$$\overline{\text{AND}}(x_{i_1} \dots x_{i_k}) = 1 \quad \text{if } \forall i_j \in T = \{i_1, \dots, i_k\} \\ x_{i_j} = -1$$

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

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

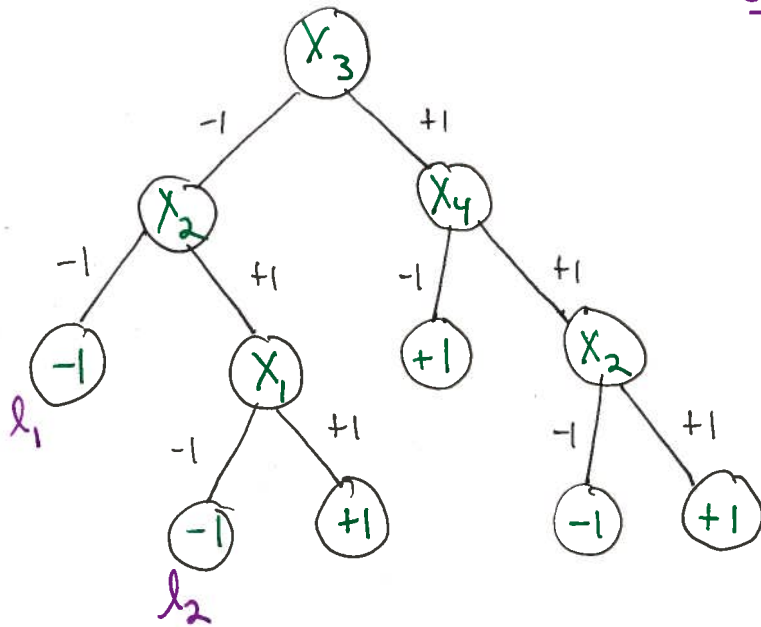
$$= \sum_{S \subseteq T} \frac{(-1)^{|S|}}{2^k} \chi_S$$

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

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

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

## 2) Decision trees



examples

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

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

First, consider path fctns:

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

vars visited on path to leaf l

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

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

exactly one of these is  $\pm 1$   
others are 0

$$f_l(x) = \begin{cases} 1 & \text{if } x \text{ takes } p \\ 0 & \text{o.w.} \end{cases}$$

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

## Low degree algorithm

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

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

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

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


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 :

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

so  $\sum_{S \neq \emptyset} \hat{f}(S)^2 \leq \epsilon$  +  $f$  has 0-F.C. 

Now, let's approximate  $f_{\text{true}}$  with  $d = \alpha(\epsilon, n)$  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!

Use  $\text{sign}(h(x))$  as hypothesis!

Why does this work?

Two stages:

1) show that if  $f$  has low F.C. then  $E_x [(f(x) - h(x))^2]$  small

$L_2 \text{ dist} / 2^n$

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

↑  
 Hamming dist

Thm if  $f$  has  $d = d(\epsilon, n)$  - F.c. then  
 $h$  satisfies  $E_x [(f(x) - h(x))^2] \leq \epsilon + \gamma$   
 with  $\text{prob} \geq 1 - \delta$

PF

Claim with  $\text{prob} \geq 1 - \delta$ ,  $\forall s$  st.  $|s| \leq d$ ,  $|C_s - \hat{f}(s)| \leq \gamma$   
 for  $\gamma \leftarrow \sqrt{\frac{\gamma}{nd}}$

Pf of claim

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

Chernoff bnd  $\Rightarrow O\left(\frac{nd}{\gamma} \ln \frac{nd}{\gamma}\right) = O\left(\frac{1}{\gamma^2} \ln \frac{nd}{\gamma}\right)$  samples  
 yields  $\Pr [ |C_s - \hat{f}(s)| > \gamma ] < \frac{\delta}{nd}$

union bnd  $\Rightarrow \Pr [ \exists s \text{ st. } |C_s - \hat{f}(s)| > \gamma ] < \delta$   
 $\uparrow$   
 only  $\binom{n}{d} < nd$  such  
 sets of size  $\leq d$

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

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,  $\forall s$  st.  $|s| > d$ ,  $\hat{h}(s) = 0 \Rightarrow \hat{g}(s) = \hat{f}(s)$   
 $|s| \leq d$ ,  $\hat{h}(s) = C_s \Rightarrow \hat{g}(s) = \hat{f}(s) - C_s$   
 so  $\hat{g}(s)^2 \leq \gamma^2$

$$\begin{aligned}
 \text{so } E[(f(x) - h(x))^2] &= E[g(x)^2] \\
 &= \sum_s \hat{g}(s)^2 \quad \text{Parseval} \\
 &= \underbrace{\sum_{|s| \leq d} \hat{g}(s)^2}_{\leq \gamma^2} + \underbrace{\sum_{|s| > d} \hat{g}(s)^2}_{\leq \epsilon \text{ by F.C.}} \\
 &\leq n^d \cdot \gamma^2 + \epsilon \\
 &\leq \tau + \epsilon \quad \blacksquare
 \end{aligned}$$

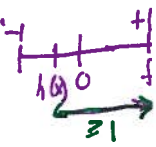
Thm  $f: \{\pm 1\}^n \rightarrow \{\pm 1\}$   
 $h: \{\pm 1\}^n \rightarrow \mathbb{R}$   
then  $\Pr[f(x) \neq \text{sign}(h(x))] \leq E[(f(x) - h(x))^2]$

Pf.  $E[(f(x) - h(x))^2] = \frac{1}{2^n} \sum_x (f(x) - h(x))^2$  defn } show term by term

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

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

eg: 

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

## Correctness of learning algorithm:

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

Pf. run low degree alg with  $\gamma = \epsilon$   
 get  $h$  s.t.  $E[(f-h)^2] \leq \epsilon + \epsilon = 2\epsilon$   
 output  $\text{sign}(h)$  ■

## Applications

1) Bounded depth decision trees

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

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

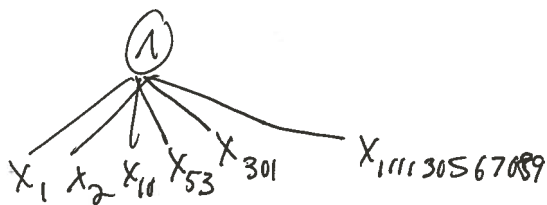
by linearity,  $\hat{f}(s) = \sum \text{val}(\ell) \cdot \hat{f}_{\ell}(s)$  which is 0 if  $|s| > \text{depth}$



2) Constant depth ckt:

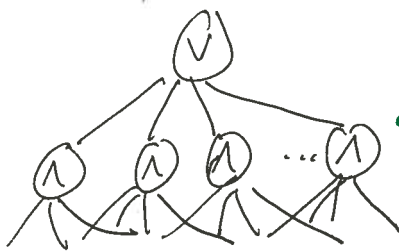
Def. "Boolean Ckt C" is a DAG

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



Can we compute parity of  $n$  bits in const depth?

yes! can compute any  $n$ -bit fctn in const depth



each "1" picks an arbitrary sat setting & checks if input matches

Can we compute parity of  $n$  bits in const depth, poly size?

No! [Fürst Saxe Sipser] } lemma  
 Switching lemma

Lemons  $\Rightarrow$  Lemonade

Thm [Hastad, Linial Mansour Nisan]

prove via  
random  
restrictions  
as in  
[FSS] parity  
result

Take Advanced Complexity!

$\forall f$  computable via size  $s$  depth  $d$  ckt  
 $\sum_{|s| \geq t} \hat{f}^2(s) \leq \alpha$  for  $t = O(14 \log \frac{2s}{\alpha})^{d-1}$

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

Gives  $n^{O(\log^d(\frac{n}{\epsilon}))}$  sample query algorithm  
 (note: can improve to  $n^{O(\log \log n)}$  [Jackson])

Application 2

## Learning halfspaces (linear threshold fctns)

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

$$\begin{array}{c} \uparrow \\ \text{sign}(x) \begin{cases} +1 & \text{if } x \geq 0 \\ -1 & \text{o.w} \end{cases} \end{array}$$

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

$$\text{(ie. } \sum_{|S| \geq \frac{c}{\epsilon^2}} \hat{h}(S)^2 \leq \epsilon \text{)}$$

(Will prove soon)

Corr low degree algorithm learns halfspaces under  
uniform distribution with  $n^{O(1/\epsilon^2)}$  uniform samples.

(Actually can learn in  $O(n^5)$  but we'll get a "big win"  
from this approach soon...)

Key idea: Noise sensitivity

def.

(1)  $0 < \epsilon < 1/2$  "Noise operator"

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

(2) "Noise Sensitivity"

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

Examples

1)  $f(x) = x_1$        $NS_\epsilon(f) = \epsilon$

2)  $f(x) = x_1 x_2 \dots x_k$        $NS_\epsilon(f) = \Pr [f(x) = F \text{ \& } f(N_\epsilon(x)) = T]$   
 $+ \Pr [f(x) = T \text{ \& } f(N_\epsilon(x)) = F]$

$= 2 \cdot \Pr [f(x) = T \text{ \& } f(N_\epsilon(x)) = F]$

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

$\leftarrow$  symmetric since  $N_\epsilon(x)$

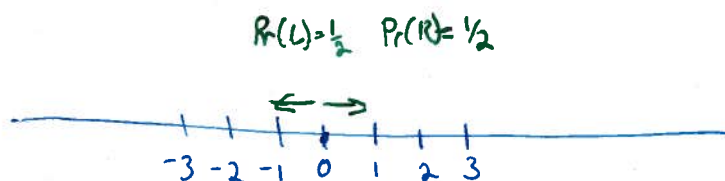
$= \begin{cases} \frac{1}{2^{k-1}} (\epsilon k) & \text{if } \epsilon \ll 1/k \\ 1 - \frac{e^{-k\epsilon}}{2^{k-1}} & \text{if } \epsilon \gg 1/k \end{cases}$  } small!  
but fcn is rarely 1

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

$nS_\epsilon(f) = O(\sqrt{\epsilon})$

Sketch

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



$x_i = 1 \Rightarrow R \times 2$   
 $x_i = -1 \Rightarrow L \times 2$   
 $\sum x_i = \text{end pt of walk}$   
 (if start at 0)

eg.  $x = (11-1-1-1)$



ends at 0

Equivalent process:

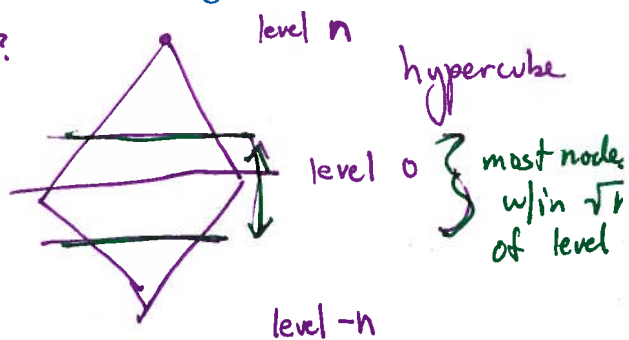
pick random pt on hypercube & output the level

Well known fact:

expected distance from startpt ( $|x_1 + x_2 + \dots + x_n|$ ) after  $n$  steps is  $\sqrt{n}$

& likely to be close to expectation

why?



•  $N_\epsilon(x) \sim$  random walk on  $\epsilon n$  bits each has twice value taken

expected displacement  $2\sqrt{\epsilon n}$

↑ since chosen  
 $+1 \rightarrow +1$   
 $-1 \rightarrow +1$

Another view:  
 • take walk according to  $x$   
 • then take walk according to  $N_\epsilon(x)$

• heuristic argument:

pretend 1st walk leaves us at  $\sqrt{n}$

$$\Pr[2^{\text{nd}} \text{ walk takes us across 0}] = \frac{1}{2} \Pr[2^{\text{nd}} \text{ displacement} > \sqrt{n}]$$

$$= \frac{1}{2\sqrt{\epsilon}} 2 \cdot \sqrt{\epsilon n}$$

$$< 2\sqrt{\epsilon} \quad \text{by Markov's } \neq$$

4) any LTF (1/2 space)

Thm [Peres]  $NS_{\epsilon}(\text{LTF}) < 8.8\sqrt{\epsilon}$  ← best possible since  
 Maj is LTF & has  
 $NS_{\epsilon}(\text{MAJ}) = \Theta(\sqrt{\epsilon})$

5) parity fctns  $\chi_S(x)$  for  $|S|=k$  ← we've done  
 $k=1$  (dictators)

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

$$= \binom{k}{1} (1-\epsilon)^{k-1} \epsilon + \binom{k}{3} (1-\epsilon)^{k-3} \epsilon^3 + \dots$$

but  $1^k = ((1-\epsilon) + \epsilon)^k = (1-\epsilon)^k + \binom{k}{1} (1-\epsilon)^{k-1} \epsilon + \binom{k}{2} (1-\epsilon)^{k-2} \epsilon^2 + \dots$

$-(1-2\epsilon)^k = ((1-\epsilon) - \epsilon)^k = (1-\epsilon)^k - \binom{k}{1} (1-\epsilon)^{k-1} \epsilon + \binom{k}{2} (1-\epsilon)^{k-2} \epsilon^2 - \dots$

÷ 2

$$\frac{1 - (1-2\epsilon)^k}{2} = NS_{\epsilon}(f)$$

↑ high for large  $k$  (if  $S$  likely to get "hit")

6) any f:

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

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

note for parity fctns, this gives

$$NS_\epsilon(\chi_s) = \frac{1}{2} - \frac{1}{2} \cdot (1-2\epsilon)^{|s|} \quad \checkmark$$

Pf.

$$NS_\epsilon(f) = \Pr_{x, y \in N_\epsilon(x)} [f(x) \neq f(y)]$$

$$= E_{x, y \in N_\epsilon(x)} [1_{f(x) \neq f(y)}]$$

$$= E_{x, y} \left[ \frac{(f(x) - f(y))^2}{4} \right] = \frac{1}{4} E [ \underbrace{f(x)^2 + f(y)^2}_{\text{always 1}} - 2f(x)f(y) ]$$

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

$$= \frac{1}{2} - \frac{1}{2} E \left[ \sum_s \hat{f}(s) \chi_s(x) \sum_T \hat{f}(T) \chi_T(y) \right]$$

$$= \frac{1}{2} - \frac{1}{2} \sum_{s, T} \hat{f}(s) \hat{f}(T) E[\chi_s(x) \chi_T(y)]$$

$= 0$  if  $s \neq T$  } show using standard techniques using pairs on  $i \in SA$ .

but, if  $s = T$ ?

$$E_{x, y} [\chi_s(x) \chi_T(y)] = 1 \cdot \Pr[\chi_s(x) = \chi_T(y)] + (-1) \Pr[\chi_s(x) \neq \chi_T(y)]$$

$$= 1 - 2 \Pr[\chi_s(x) \neq \chi_T(y)]$$

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

$$NS_\epsilon(\chi_s) = \frac{1}{2} - \frac{1}{2} (1-2\epsilon)^{|s|}$$

$$= (1-2\epsilon)^{|s|}$$

□

# Noise Sensitivity w. Fourier Concentration :

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

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

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

previous thm

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

Booleam Parseval's

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

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

drop lower terms  
raise to lower power

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

so  $\sum_{|s| > \frac{1}{\gamma}} \hat{f}(s)^2 < \frac{2}{1-e^{-2}} ns_\gamma(f) < 2.32$

Which gives:

Corr 1

$f : \{\pm 1\}^n \rightarrow \{\pm 1\}$ .  $\beta : [0, 1/2] \rightarrow [0, 1/2]$  is fctn st.  $ns_\gamma(f) \leq \beta(\gamma)$

Then  $\sum_{|s| \geq m} \hat{f}(s)^2 \leq \epsilon$  for  $m = \frac{1}{\beta^{-1}(\epsilon/2.32)}$  }  $f$  has Fourier concentra  $m$

pf

$$\sum_{|s| \geq m} \hat{f}(s)^2 \leq \sum_{|s| \geq \frac{1}{\beta^{-1}(\epsilon/2.32)}} \hat{f}(s)^2 \leq 2.32 ns_{\frac{1}{\beta^{-1}(\epsilon/2.32)}}(f) \leq 2.32 \beta\left(\beta^{-1}\left(\frac{\epsilon}{2.32}\right)\right) = \epsilon$$

↑ thm
↑ set to  $\gamma$ 
↑ def of  $\beta$



Back to  $\frac{1}{2}$  spaces:

Corr 2 for  $\frac{1}{2}$  space  $h: \{\pm 1\}^n \rightarrow \{\pm 1\}$

$$\sum_{|S| \geq n \left(\frac{\epsilon}{8.8}\right)} \hat{f}(S)^2 \leq \epsilon$$

PF

for  $\frac{1}{2}$  space  $ns_\epsilon(h) \leq 8.8 \sqrt{\epsilon}$  so use

$$\beta(\epsilon) = 8.8 \sqrt{\epsilon}$$

$$\beta^{-1}(x) = \left(\frac{x}{8.8}\right)^2$$

$$\beta^{-1}\left(\frac{\epsilon}{2.32}\right) \leq \left(\frac{\epsilon}{2.32 \times 8.8}\right)^2$$

$$\leq \left(\frac{\epsilon}{20}\right)^2 \quad (20.416)$$

so for  $n \geq \frac{1}{\left(\frac{\epsilon}{20}\right)^2}$

□

⇒ Can learn any  $\frac{1}{2}$  space over uniform with  $n = O\left(\frac{1}{\epsilon^2}\right)$  random examples but, can do a lot better! still, technique extends ...

Learn any fctn of  $k$   $\frac{1}{2}$  spaces:

$h_1, \dots, h_k$  are  $k$  half spaces

$g: \{\pm 1\}^k \rightarrow \{\pm 1\}$  Boolean fctn

$f: \{\pm 1\}^n \rightarrow \{\pm 1\}$  st  $f(x) = g(h_1(x), \dots, h_k(x))$

Thm  $ns_\epsilon(f) = 8.8k\sqrt{\epsilon}$

PF

$$\begin{aligned}
 \Pr_{S_\varepsilon}(f) &= \Pr[f(x) \neq f(N_\varepsilon(x))] \\
 &= \Pr[g(h_1(x), \dots, h_k(x)) \neq g(h_1(N_\varepsilon(x)), \dots, h_k(N_\varepsilon(x)))] \\
 &\leq \Pr[h_1(x) \neq h_1(N_\varepsilon(x)) \text{ or } h_2(x) \neq h_2(N_\varepsilon(x)) \text{ or } \dots] \\
 &\leq k \cdot 8.8 \cdot \sqrt{\varepsilon} \quad \text{union bnd}
 \end{aligned}$$

Take  $\beta(\varepsilon) = 8.8 k \sqrt{\varepsilon}$

so  $\beta^{-1}(x) = \left(\frac{x}{8.8 k}\right)^2$

+ corr 1  $\Rightarrow$

Thm  $\sum_{|S| \geq \Omega\left(\frac{k^2}{\varepsilon}\right)} \hat{f}(S)^2 \leq \varepsilon$

Thm low degree alg learns any Boolean fcn  
of  $k$   $\frac{1}{2}$  spaces with  $n(O(k^2/\varepsilon^2))$  samples

e.g.  $\wedge$  of 2  $\frac{1}{2}$  spaces  
parity of  $k$  vars