# Lecture 21

- weak learning of monotone fctns

- begin: distribution-free weak learning
$$\Rightarrow \text{strong learning}$$

# Boolean Cube



$111\cdots1$

$000\cdots0$

hypercube

level $k$:
   nodes labeled by
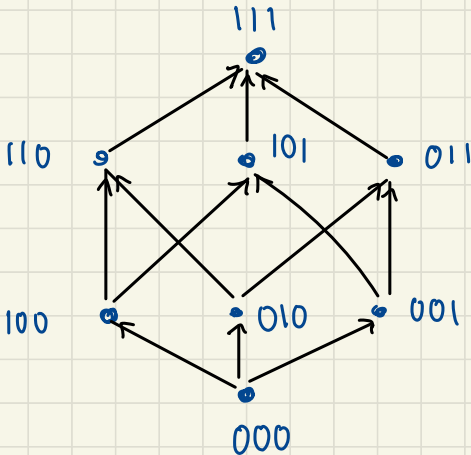      $k$ 1's & $n-k$ 0's

#nodes on level $k$:

$$\binom{n}{k}$$

edges:

   $x \rightarrow y$
   if flip one 0 in $x$ to a 1
        to get $y$

#nodes: $2^n$
#edges: $\dfrac{n \cdot 2^n}{2}$



111

110    101    011

100    010   001

000

example for $n=3$

# Monotone Functions

def. partial order $\leq$ : $x \leq y$ iff $\forall i \; x_i \leq y_i$

monotone fctn $f$ : $x \leq y \implies f(x) \leq f(y)$

Are there fast learning algorithms for the class
of monotone functions?
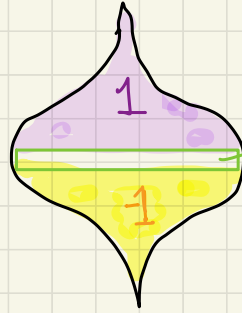
Occam's razor:

poly $(\log |C|)$ samples suffice

↳ Class of monotone fctns

$\geq 2^{2/\sqrt{n}}$ monotone fctns

So only gives exponential bound

Why so many monotone fctns?



Consider "slice" fctns:

*set middle row in all possible ways w/o violating monotonicity*

*Note: on uniform dist, easy to learn slice fctns. ie. Output "Majority" ⟹ Occam is "weak" on this class/distribution*

$2^{\binom{n}{\frac{n}{2}}}$ *options*

*all are monotone!*

H.W.: $2^{O(\sqrt{n})}$ random samples suffice for unif dist

T̲oday: what if you compromise on error?

Can get very slight "win"

All monotone fctns have weak agreement with S̲ome dictator fctn.

**Thm** $\forall f$ monotone, $\exists g \in \{\pm 1, x_1, x_2, \ldots x_n\} \equiv S$

s.t. $\Pr_x [f(x) = g(x)] \geq \frac{1}{2} + \Omega(\frac{1}{n})$

↖ uniform
distribution

↗ slightly
better than
random
guessing

(can get $\frac{1}{2} + \Omega(\frac{1}{\sqrt{n}})$
if add majority)

<u>note</u> Slice fctns have weak
agreement with all dictators
on uniform dist

$\Rightarrow$ learning algorithm:
estimate agreement of $f$ with all members of $S$
output best

<u>Pf.</u>

Case 1: $f(x)$ has weak agreement with $+1$ or $-1$ ✓

Case 2: otherwise $\Pr[f(x)=1] \in [\frac{1}{4}, \frac{3}{4}]$

Let's first look at monotone $\Big\}$ excuse for
fctns in a different way: $\Big\}$ a
detour

# Monotone Functions on Boolean Cube :

## A "graph" view



111···1

red  +1

blue  −1

000···0

monotone ⟹ no blue above any red

$$x \le y \quad \text{if} \quad \forall i \quad x_i \le y_i$$

$f$ monotone if

$$\forall \, x \le y, \quad f(x) \le f(y)$$

## Influence of f :

$$\text{Inf}_i(f) = \frac{\#\text{ red-blue edges in } i^{th} dir}{2^{n-1}}$$

$$= \Pr_x \left[ f(x) \ne f(x^{\oplus i}) \right]$$

↰ $x$ with $i^{th}$ bit flipped

$$\text{Inf}(f) = \frac{\#\text{ red-blue edges}}{2^n}$$

$$= \sum_{i=1}^{n} \text{Inf}_i(f)$$

**Thm 1** $f$ monotone $\Rightarrow$ $\inf_i(f) = \hat{f}(\{i\})$

**Thm 2** majority fctn $f(x) \equiv sign\left(\sum_{i=1}^{n} x_i\right)$ (odd $n$)

maximizes influence among

monotone fctns


Pfs on h.w.


## Plan:

note: $\inf_i(f) = \hat{f}(\{i\})$     (Thm 1)

<span style="color:green">early fourier lecture:
agreement
vs.
Fourier coeffs</span>

$$= 2 \cdot Pr[f(x) = \chi_{\{i\}}(x)] - 1$$

$$\underbrace{\qquad}_{\chi_i}$$

So showing $\inf_i(f) \geq \Omega(\frac{1}{n})$

<span style="color:purple">weak learner</span>

is equivalent to showing

$$Pr[f(x) = x_i] \geq \frac{1}{2} + \frac{\inf_i(f)}{2} \geq \frac{1}{2} + \Omega(\frac{1}{n})$$

<span style="color:purple">such an $i$ would give us our theorem!</span>

To show that such an $i$ exists, will use a cool tool:

## Canonical Path Argument

**Plan** (1) define canonical path for every red-blue pair of nodes

(such a path _must_ cross at least one red-blue edge)

(2) Show upper bound on # of c.p.'s passing through any edge

(in particular, any red-blue edge)

(3) Conclude lower bound on # of red-blue edges.

<u>Part I</u>: define canonical path for every
red-blue pair of nodes

<u>def</u> $\forall (x,y)$ s.t. x <span style="color:red">red</span> & y blue

"Canonical path from x to y" is:
scan bits left to right
flipping where needed
each flip ⟿ step in path

<u>example</u>:

| dimension | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| x = | $-1$ | $+1$ | $+1$ | $+1$ |
| w = | $+1$ | $+1$ | $+1$ | $+1$ |
| z = | $+1$ | $-1$ | $+1$ | $+1$ |
| y = | $+1$ | $-1$ | $+1$ | $-1$ |

$x \to w \to z \to y$
each step
has Hamming
distance 1

note: c.p.'s can go up & down
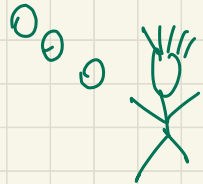e.g. $x \to w$ is up step    $w \to z$ is downstep

# Big question:

How many red-blue $x, y$ pairs have canonical paths?

recall, $\Pr[f(x) = 1] \in [\frac{1}{4}, \frac{3}{4}]$
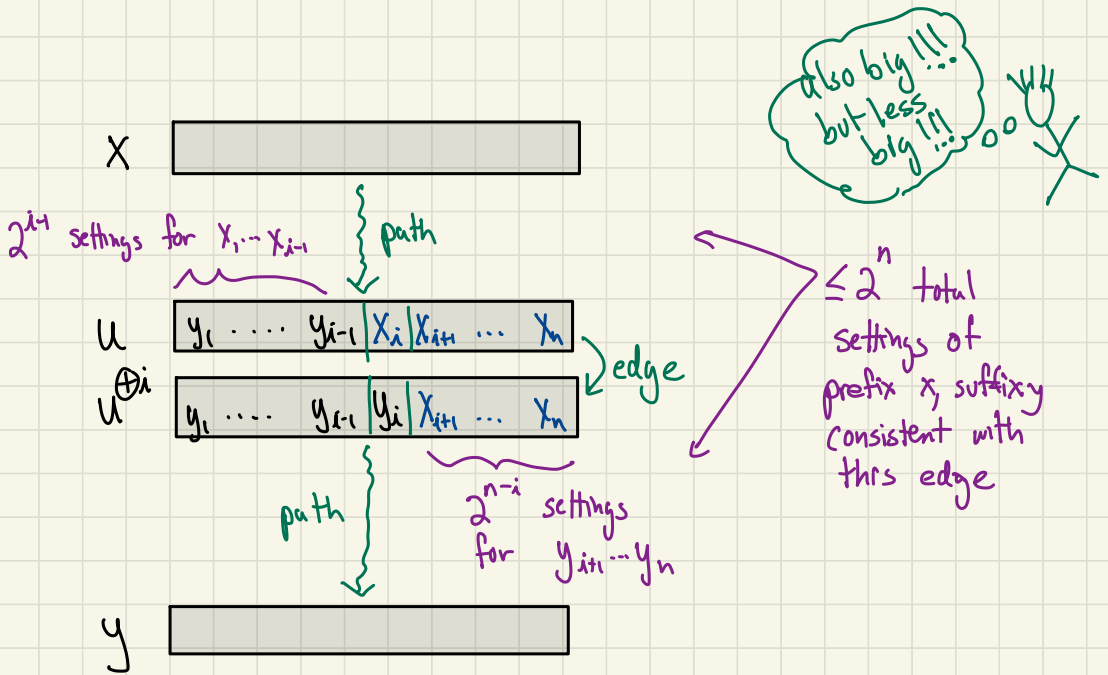
$$\#paths \geq \underbrace{\frac{1}{4} \cdot 2^n}_{\substack{\text{l.b. on} \\ \#red}} \cdot \underbrace{\frac{1}{4} \cdot 2^n}_{\substack{\text{l.b. on} \\ \#blue}} = \frac{1}{16} \cdot 2^{2n}$$

woa, that's a big number!

<u>Part II:</u> Show upper bound on # of c.p.'s

passing through any edge

for any red-blue edge $e$, how many $x$-$y$ pairs

can cross it with canonical $x$-$y$ path?

$X$ [bar]

*also big!!! but less big!!! oo*

$2^{i-1}$ settings for $X_1 \cdots X_{i-1}$  { path

$U$ $\boxed{y_1 \cdots y_{i-1} | X_i | X_{i+1} \cdots X_n}$

$U^{\oplus i}$ $\boxed{y_1 \cdots y_{i-1} | y_i | X_{i+1} \cdots X_n}$  } edge

$\leq 2^n$ total settings of prefix $x$, suffix $y$ consistent with this edge

path {

$2^{n-i}$ settings for $y_{i+1} \cdots y_n$

$y$ [bar]

<u>Main point:</u> all canonical paths crossing $U, U^{\oplus i}$

<u>agree</u> on $y_1 \cdots y_{i-1}$ & $X_{i+1} \cdots X_n$

$\Rightarrow \leq 2^n$ possible paths for each $\begin{smallmatrix} X_1 \cdots X_i \\ y_i \cdots y_n \end{smallmatrix}$

example :  ie.
started at ++++

$(-+++)$  or  $(++++)$  must come from node with   this suffix

two options
for x

$$+\ +\ +\ +$$
$$+\ -\ ++$$
$$e = (+++\text{+}, +\ -++)$$

must go
to node with
this prefix

difference in $i=2$

$(+\ -\ -\ -)$  or  $(+--+)$  or  $(+-+-)$  or  $(+-++)$

4 options for y

<u>Part III</u>: Conclude lower bound on # of red-blue edges.

$(\text{\# red-blue edges}) \times (\max \text{ \# canonical paths that use each edge})$

$\geq \text{\# red-blue canonical paths}$

$\qquad \qquad \qquad \uparrow$ since each crosses $\geq 1$ red-blue edge

$\overset{\text{l.b. on \# r-b pairs}}{\overbrace{\phantom{xxxx}}}$

$\Rightarrow \text{\# red blue edges} \geq \dfrac{\frac{1}{16} \cdot 2^{2n}}{2^n} = \frac{1}{16} \cdot 2^n$

$\qquad\qquad\qquad\qquad\qquad \nwarrow$ u.b. on \# canonical paths crossing any edge

$\Rightarrow \exists\ i \text{ s.t. } \geq \dfrac{2^n}{16} \cdot \dfrac{1}{n}$ red-blue edges in direction $i$

$$\implies \exists i \text{ st. } \text{Inf}_i(f) = \hat{f}(\{i\}) = 2 \cdot \Pr[f(x) = x_i] - 1$$

$$\geq \frac{2^n}{\frac{16n}{2^{n-1}}} = \frac{1}{8n}$$

↑ total # edges in dir $i$

$$\implies \exists i \text{ st. } \Pr[f(x) = x_i] \geq \frac{1}{2} + \frac{1}{16n}$$

✓

▨

Other uses of canonical path arguments:

- routing
- expansion/conductance of hypercube/other
    Markov chains

# What good is weak learning?

unclear

here can only weakly learn on
uniform distribution

ability to weakly learn on
all distributions
$\implies$ ability to strongly learn
[Schapire]
"boosting"

# Weak vs. Strong Learning

**Def.** Algorithm $\mathcal{A}$ "weakly PAC learns" concept

class $\mathcal{C}$ if $\exists \gamma > 0$

    s.t. $\forall c \in \mathcal{C}$ + $\forall$ dists $\mathcal{D}$

    $\forall \delta > 0$     $\longleftarrow$     ($\delta = \frac{1}{4}$ or $\frac{1}{n^2}$ doesn't

                                            affect )

    with prob $\geq 1 - \delta$

    given examples of c

    $\mathcal{A}$ outputs h s.t. $\Pr_{\mathcal{D}}[h(x) = c(x)] \geq \frac{1}{2} + \frac{\gamma}{2}$

                             *not good*

                             *Compared*        advantage

                             *to*                over

                           *$1 - \varepsilon$ or 99%*     guessing

It was first conjectured that weak learning is

      easier than strong (i.e. $\exists$ fctns that can

                         weakly learn but not

                         strongly learn)

## Surprise!!

Can "boost" a weak learner

**Thm** if $C$ can be weakly learned on

_any_ dist $\mathcal{D}$ then $C$ can be

(strongly) learned

ie. $\forall \varepsilon \; \exists A$

dependence on $\gamma$ ?

$\delta$ ?

$\varepsilon$ ?

Will prove for case of $\mathcal{D}_0 = U$

# Applications:

1) "theoretical"

- uniform distribution algorithms for
  poly term DNF
  weight- w poly threshold fctns
  (Boosting + KM)

  } low degree alg doesn't work well

- Ave case vs. worst case complexity

2) practical: "Boosting"
  Freund- Schapire

# <u>Good & Bad Ideas</u>

1) Simulate weak learner several times on same distribution & take
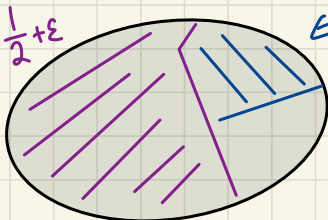
   majority answer
   <u>or</u>
   best answer

- gives better confidence
- but doesn't reduce error — what if always get same answer?

2) filter out examples on which current hypothesis does well & run weak learner on part where you do badly

$\frac{1}{2} + \varepsilon$ &larr; $\frac{1}{2} + \varepsilon$ of non-purple

<u>Problem</u>: given <u>new</u> example, how do you know which section it is in?

3) Keep some samples on which you are
ok in your filtering.
Always use majority vote on previous
hypotheses to predict value of new
samples.

history: Schapire, Freund-Schapire, Impagliazzo-
                                    Servedio-Klivans

# Filtering Procedures:

- decide which samples to keep vs. throw out

- samples on which you guess

  Correctly: needed for checking future
                        hypotheses
  incorrectly: needed for improvement