

Lecture 19

- weak learning of monotone fctns
- begin: distribution-free weak learning
⇒ strong learning

Monotone Functions

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

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

Are there fast learning algorithms for the class of monotone functions?

Occam's razor:

poly $(\log |\mathcal{C}|)$ samples suffice

Well studied partial order:

Boolean Cube (aka hypercube)

will use 0/1 instead of +1/-1 for a little bit

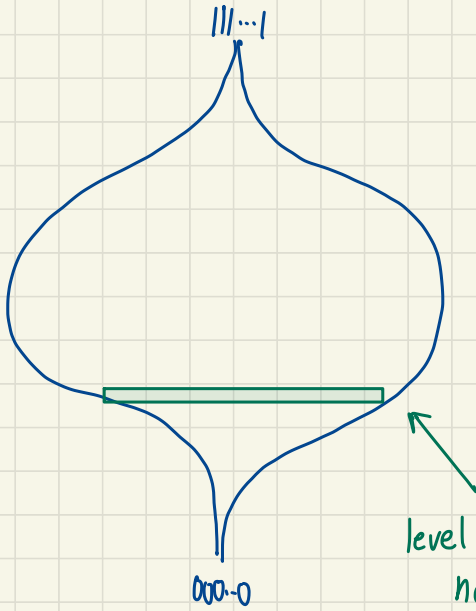
nodes labeled by n -bit strings

edges: (directed)

$x \rightarrow y$

if flip one 0 in x to 1 to get y

eg. $000 \rightarrow 010$



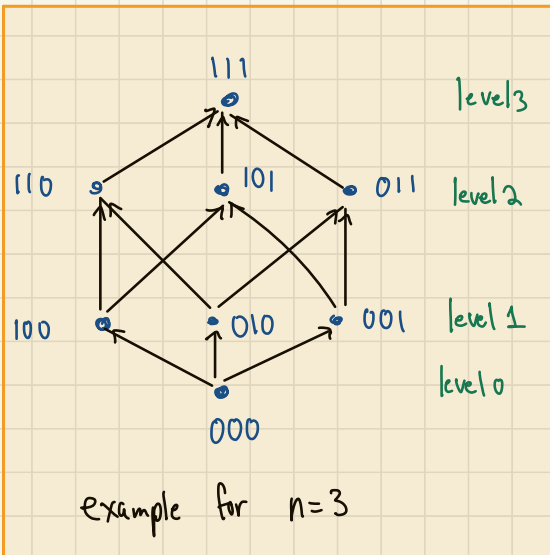
level k :

nodes labeled by k 1's + $(n-k)$ 0's

nodes: 2^n

edges: $\frac{n \cdot 2^n}{2}$

nodes on level k : $\binom{n}{k}$



how many fctns on Boolean cube?

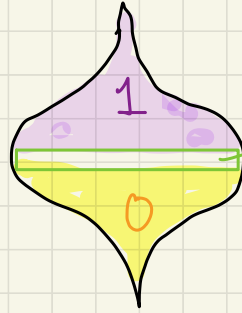
$$\geq 2^{2^{1/n}}$$

so Occam bnd is

only $O(2^{1/n})$

Why so many monotone fctns?

Consider "slice" fctns;



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

$2^{\binom{n}{2}}$ options
all are monotone!

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

Today: what if you have queries?

can get very slight "win"

All monotone fctns have weak
agreement with some dictator
fctn. (or const fctn)

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

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

constant fctns

dictator fctns

uniform distribution

slightly better than random guessing

note slice fctns have weak agreement with all dictators on uniform dist

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

\Rightarrow learning algorithm:

estimate agreement of f with all members of \mathcal{S}
output best

Pf.

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 fctns in a different way: } excuse for a detour

Monotone Functions on Boolean Cube: (also switch

A "graph" view

0 \rightarrow -1
1 \rightarrow 1

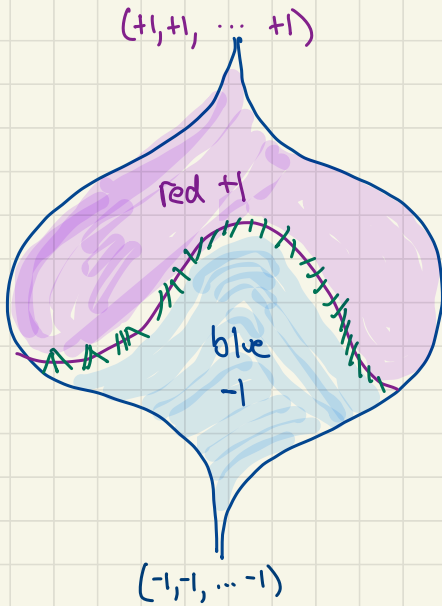
note this is not what we did in past but equivalent in this case

monotone \Rightarrow no blue above any red

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

f monotone if

$$\forall x \leq y, f(x) \leq f(y)$$



Influence of f : [these notions are interesting also for nonmonotone fctns]

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

$$= \Pr_x [f(x) \neq f(x^{\oplus i})]$$

\leftarrow x with i^{th} bit flipped

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

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

} measures size of "border"

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

Thm 2 majority fctn $f(x) \equiv \text{sign}(\sum_{i=1}^n x_i)$ (odd n)
maximizes influence among
monotone fctns

Pfs on h.w.

Plan:

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

$$= 2 \cdot \Pr[f(x) = \underbrace{\chi_{\{\pm 1\}}(x)}_{x_i}] - 1$$

early Fourier lecture:
agreement
vs.
Fourier coeffs

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

is equivalent to showing

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

such an i would give us our theorem!

weak learner

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.

Part I: define canonical path for every red-blue pair of nodes

def $\forall (x, y)$ s.t. x red & y blue

"canonical path from x to y " is:

scan bits left to right

flipping where needed

each flip \rightsquigarrow step in path

Important points
must have
at least one
red-blue edge,
possibly more

example:

	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

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

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

Big question:

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

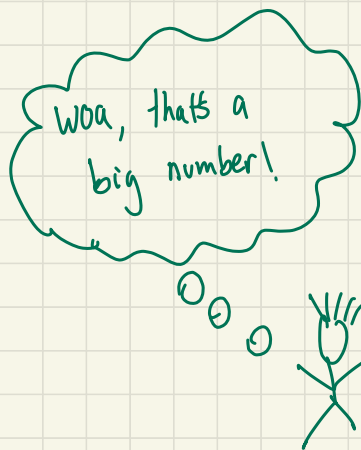
recall, $\Pr[f(x)=1] \in [\frac{1}{4}, \frac{3}{4}]$

$$\# \text{paths} \geq \underbrace{\frac{1}{4} \cdot 2^n}_{\text{l.b. on } \# \text{red}} \cdot \underbrace{\frac{1}{4} \cdot 2^n}_{\text{l.b. on } \# \text{blue}} = \frac{1}{16} \cdot 2^{2n}$$

recall the plan:

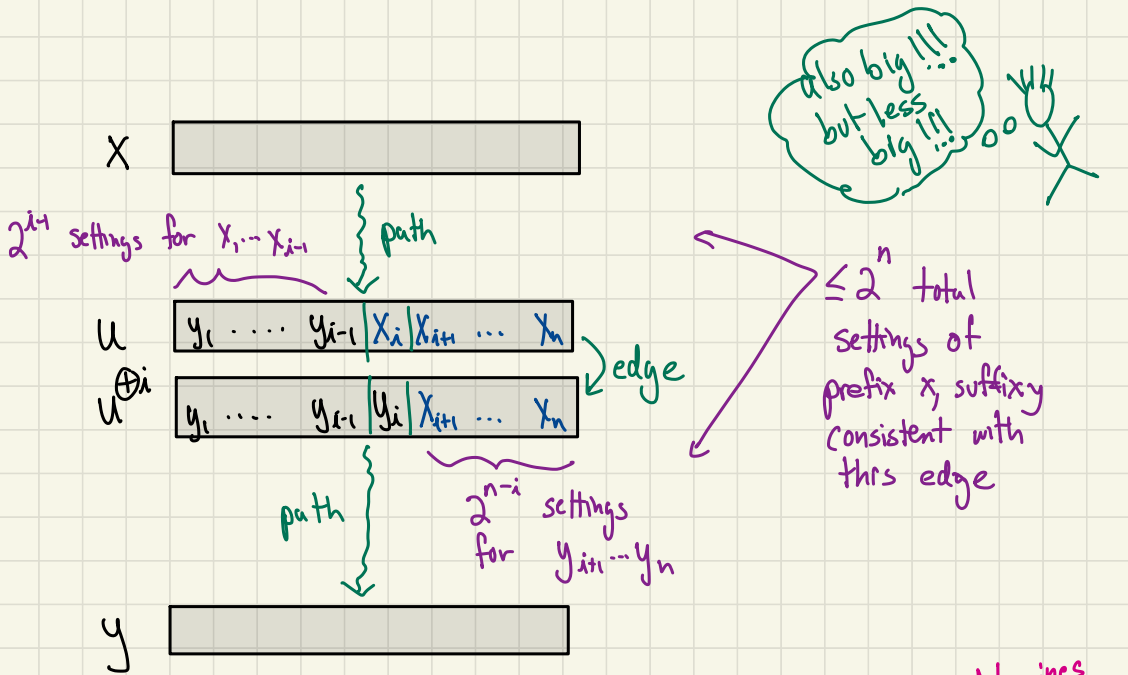
- to show: lots of red-blue edges
- each path has ≥ 1 red-blue edge, but that doesn't mean we have lots of red-blue edges. Maybe they all use the same red-blue edge?

- in the following: bound $\#$ of paths that can share a red-blue edge



Part II: Show upper bound on # of c.p.'s passing through any edge

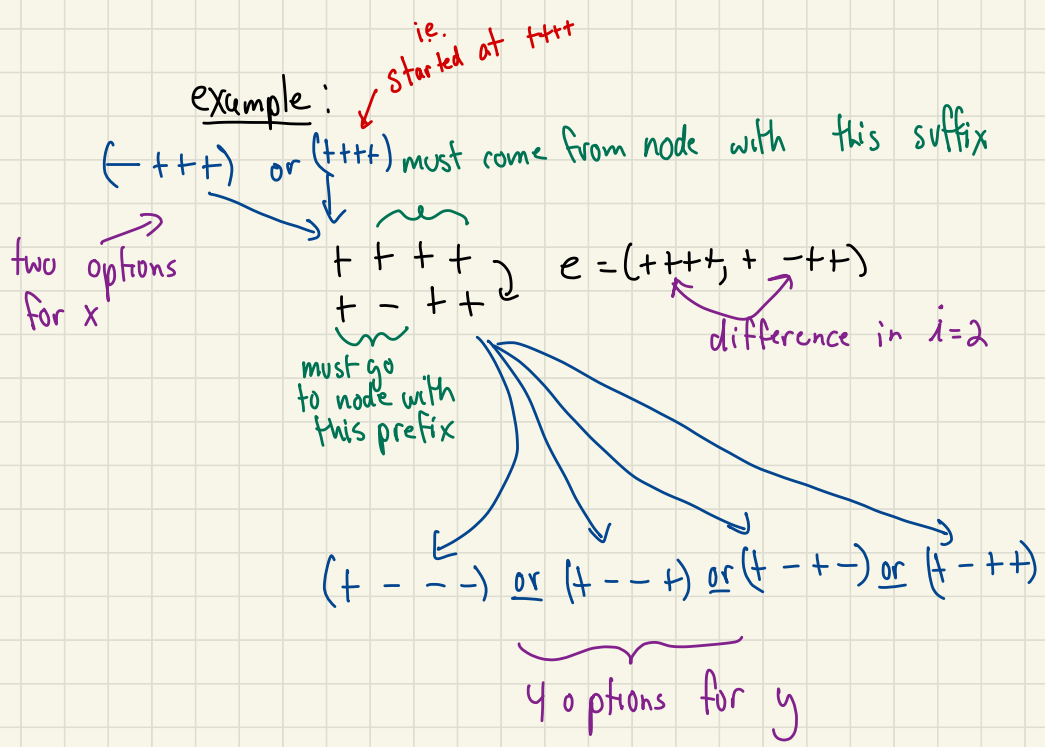
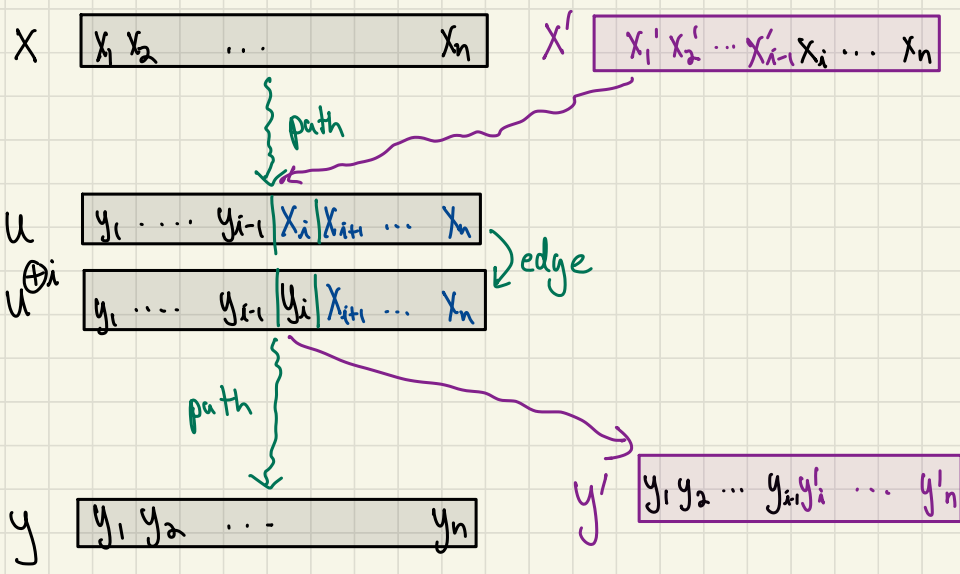
for any red-blue edge e , how many x - y pairs can cross it with canonical x - y path?



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

agree on $y_1 \dots y_{i-1} \oplus x_{i+1} \dots x_n$
 $\Rightarrow \leq 2^n$ possible paths for each $X_1 \dots X_i, y_i \dots y_n$

determines (X_i, y_i)



Part III: Conclude lower bound on # of red-blue edges.

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

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

↑ since each crosses ≥ 1 red-blue edge

$$\Rightarrow \# \text{red blue edges} \geq \frac{\overset{\text{l.b. on \# r-b pairs}}{1}{16} \cdot 2^{2n}}{2^n} = \frac{1}{16} \cdot 2^n$$

↑ u.b. on # canonical paths crossing any edge

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

$$\Rightarrow \exists i \text{ st. } \ln f_i(f) = \hat{f}(\{i\}) = 2 \cdot \Pr[f(x) = x_i] - 1$$

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

↑ total #
edges in dir i

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



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

⇒ ability to strongly learn
[Schapire]

"boosting"

Weak vs. Strong Learning

Def. Algorithm A "weakly PAC learns" concept class \mathcal{C} if $\exists \epsilon > 0$

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

$\forall \delta > 0$

$\leftarrow (\delta = \frac{1}{4}$ or $\frac{1}{n^2}$ doesn't affect)

with prob $\geq 1 - \delta$

given examples of c

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

not good
compared
to
 $1 - \epsilon$ or 99%

↑
advantage
over
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 \mathcal{C} can be weakly learned on

any dist \mathcal{D} then \mathcal{C} can be

(strongly) learned

ie. $\forall \epsilon$

dependence on γ ?

δ ?

ϵ ?

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

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

Good & Bad Ideas

1) simulate weak learner several times
on same distribution & take

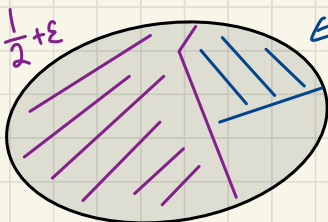
majority answer

or

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



$\leftarrow \frac{1}{2} + \epsilon$ of non-purple

Problem: given new 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