# Lecture 16

fast weak learning of monotone functions

# Monotone Functions

<u>def</u>    partial order $\leq$

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

monotone function $f$

$$x \leq y \implies f(x) \leq f(y)$$

Learning algorithms for the <u>class</u> of monotone functions?

in homework we saw $2^{O(\sqrt{n})}$ random samples suffice for uniform distribution

why is this nontrivial?

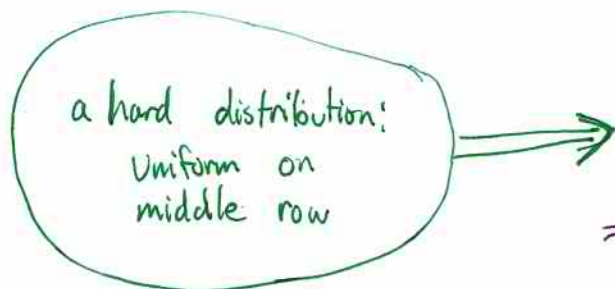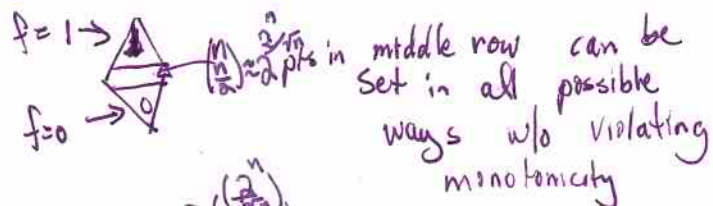we said poly samples is easy, the problem is computation time? poly in what?

but need $\text{poly}(\log |C|)$ samples
↑
all monotone fctns

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

why $\geq 2^{2^{n}}$ monoatone fctns?

consider slice fctns:

$f=1 \rightarrow$

$f=0 \rightarrow$

$\binom{n}{\frac{n}{2}} \approx \frac{2^n}{\sqrt{n}}$ pts in middle row can be set in all possible ways w/o violating monotonicity

a hard distribution:
Uniform on
middle row

$\implies$ learning needs $\Omega\left(2^{\left(\frac{2^n}{\sqrt{n}}\right)}\right)$
even with queries in PAC model

Today's question:

what about learning monotone distributions,

on uniform distribution,

with queries?

here we will get a very slight "win":  slightly better than random guess ↓

All monotone fctns have weak agreement $(\frac{1}{2}+\theta(\frac{1}{n}))$

with some dictator fctn.

(can get $\frac{1}{2}+\theta(\frac{1}{\sqrt{n}})$ with majority + dictators)

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

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

Why does this give weak learning algorithm? estimate agreement of $f$ with all members of $S$ & output member with max agreement
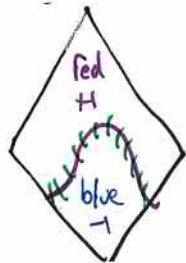
Pf.

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

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

First a break,

before we prove case 2 ...

what is another way to think of influence of monotone fctns?



· # nodes $= 2^n$,   # edges $= \frac{n \cdot 2^n}{2}$

· each level has $\binom{n}{j}$ weight $j$ nodes

· monotone $\Rightarrow$ no blue above any red

· slicing cube in roughly half cuts many edges + many in same direction

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

recall H.W. :

Thm $f$ monotone

$$\mathrm{Inf}_i(f) = \hat{f}(\{i\}) \equiv 2\Pr[f(x) = \chi_{\{i\}}(x)] - 1$$

↑ H.V.   ↑ Known   $\underbrace{\chi_{\{i\}}}_{\chi_i}$

Plan :

Show $\mathrm{Inf}_i(f) \geq \Omega\left(\frac{1}{n}\right)$

$$\Rightarrow \Pr[f(x) = \chi_i] \geq \frac{1}{2} + \frac{\mathrm{Inf}_i(f)}{2}$$

$$\geq \frac{1}{2} + \Omega\left(\frac{1}{n}\right)$$

Will use following tool :

<u>Canonical Path Argument</u>

<u>Plan</u> 1) define canonical path for every red-blue pair
of nodes (note such a path <u>must</u> cross at
least one red-blue edge)

2) show upper bnd on # of c.p.'s
passing through <u>any</u> edge
(in particular, any red-blue edge)

3) conclude lower bnd. on # of red-blue edges

Part 1 of plan:

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

example        direction

|  | $\underline{1}$ | $\underline{2}$ | $\underline{3}$ | $\underline{4}$ |
|---|---|---|---|---|
| $x =$ | $-1$ | $+1$ | $+1$ | $+1$ |
| $w =$ | $\searrow +1$ | $+1$ | $+1$ | $+1$ |
| $z =$ | $+1$ | $\searrow -1$ | $+1$ | $+1$ |
| $y =$ | $+1$ | $-1$ | $+$ | $-1$ |

$x \to w \to z \to y$

each step is

Hamming distance $\underline{1}$

How many $\overset{\text{red-blue}}{\stackrel{\vee}{}} x, y$ pairs have canonical paths?

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

\# paths $\geq \frac{1}{4} \cdot 2^n \cdot \frac{1}{4} \cdot 2^n = \frac{1}{16} \cdot 2^{2n}$

Part 2 of plan:

For any (red-blue) edge $e$, how many $x, y$ pairs can cross it with canonical $xy$-path?



$x$

$e = (u, u^{\oplus i})$   $u$   $y_1 \cdots y_{i-1} | x_i | x_{i+1} \cdots x_n$   edge in $i^{th}$ direction
red blue   $u^{\oplus i}$   $y_1 \cdots y_{i-1} | y_i | x_{i+1} \cdots x_n$

$\leq 2^{i-1}$ settings for prefix of $x$

$y_1 \cdots y_{i-1} \boxed{x_i} \; x_{i+1} \cdots x_n$
$y_1 \cdots y_{i-1} \boxed{y_i} \; x_{i+1} \cdots x_n$

$\underbrace{\qquad}_{\leq 2^{n-i} \text{ settings for suffix of } y}$

$y$

**Main point:**
all canonical paths crossing $u, u^{\oplus i}$ agree on $y_1 \cdots y_{i-1} + x_{i+1} \cdots x_n$

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

Part 3 of plan:

(# red-blue edges)(max # canonical-paths that use it) $\geq$ # red-blue canonical paths

since each uses $\geq 1$ red-blue edge

l.b. on # r-b pairs

So

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

u.b. on # canonical paths crossing any edge

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

so $\qquad \text{Inf}_i(f) \geq \dfrac{\frac{2^n}{16n}}{2^{n-1}} = \dfrac{1}{8n} = \hat{f}(\{i\}) = 2\Pr[f(x) = x_i] - 1$

$\underset{\text{total \# edges}}{\underset{\text{in dir } i}{\uparrow}}$

$\therefore \Pr[f(x) = x_i] \geq \dfrac{1}{2} + \dfrac{1}{16 \cdot n}$

Canonical Path argument also used in

· routing

· expansion / conductance of hypercube / other Markov Chains

What good is weak learning?

unclear

here only uniform distribution

if can learn in all distributions,

can do much more

(next result does not apply to monotone

function learning ...
$\qquad\qquad\qquad \swarrow$ i.e. $\frac{1}{n}$ agreement

in particular, this weak notion of learning $\underset{\text{learning}}{\underset{\leftarrow}{\text{i.e. const}}} \overset{\geq \frac{1}{2}}{\text{agreement}}$

provably doesn't give anything for stronger learning)

# Weak vs. Strong Learning

Def. Algorithm $A$ weakly "PAC learns" concept class $C$

if $\forall$ $c \in C$ & $\forall$ dists $\mathcal{D}$ $\qquad \exists \gamma > 0$

$\forall \epsilon, \delta > 0$ $\quad (\delta = \frac{1}{4}$ or $\frac{1}{n^2}$ doesn't affect$)$

with prob $\geq 1-\delta$

given examples of $c$

$A$ outputs $h$ s.t. $\Pr_{\mathcal{D}}\left[ h(x) \neq c(x) \right] \leq \cancel{\epsilon}$

$$\frac{1}{2} - \frac{\gamma}{2}$$

$\uparrow$
advantage

It was conjectured that distribution free weak learning
was really weaker but surprise!

Can "boost" a weak learner

Thm if $C$ can be weakly learned on
<u>any</u> dist $\mathcal{D}$ then $C$ can be
(strongly) learned.

## Applications

1) "theoretical"
- Unif dist Algorithms for poly term DNF
  weight w - poly threshhold fctns   } low degree
                                        alg doesn't
                                        work well

  ∴   (Boosting + KM)

- Ave  case  vs.  worst case

2) practical - Boosting
          Freund-Schapire

## Good & Bad   Ideas
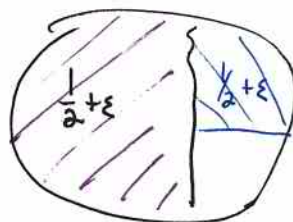
1) simulate  weak  learner  several  times  on
   same  distribution  & take  majority answer
                                    -or-
                              best  answer

   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.



Problem: given a new
example, how do you
know which section it
is in?

3) Keep some samples on which you are ok

always use majority vote on all previous hypotheses
to predict value of new samples

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

## Filtering Procedures

- decide which samples to keep, which to throw out

- samples on which so far you guess correctly ← need for checking future hypotheses

  incorrectly ← need to improve on these