

Lecture 13

- Basics of Fourier Analysis
on Boolean cube (some review,
some new)
- Analysis of linearity test
- learning Boolean functions - a model

Recall from last time:

def. $f: \{\pm 1\}^n \rightarrow \{\pm 1\}$ is "linear" (homomorphism) if

$$\forall x, y \in \{\pm 1\}^n \quad f(x) \cdot f(y) = f(x \odot y)$$

Coordinatewise multiplication
 $x \odot y = (x_1 y_1, x_2 y_2, \dots)$

def f is " ϵ -linear" if \exists linear g
s.t. f & g agree on $\geq 1-\epsilon$ fraction
of inputs.

A useful observation:

$$\forall a, y \in \{\pm 1\}^n \quad \Pr_x [y = a \odot x] = \frac{1}{2^n}$$

Linear fctns are:

$$\begin{aligned} S &\subset \{1, \dots, n\} \\ \chi_S(x) &= \prod_{i \in S} x_i \end{aligned}$$

"Parity fctns"

1 if $\# i \in S$ s.t. $x_i = -1$ even
-1 if " " " " " odd

e.g. $f(x) = x_1 \cdot x_3 \cdot x_5 \cdot x_7$
 $g(x) = x_1 \cdot x_2$

def. "linearity tester"

given: query access to fct f & parameter ϵ

requirements:

• if f linear, $\Pr[\text{tester passes}] = 1$

• if f not ϵ -linear, $\Pr[\text{tester fails}] \geq 3/4$

arbitrary
constant

Proposed test:

Pick random $x, y \in \mathbb{Z}_2^n$

Test $f(x) \cdot f(y) = f(xy)$

repeat
how
many
times?

Express event that test passes as algebraic fctn:

$$f(x) \cdot f(y) \cdot f(x \odot y) = \begin{cases} 1 & \text{if test accepts} \\ -1 & \text{" " rejects} \end{cases}$$

$$f(x) \cdot f(y) = f(x \odot y)$$

↕
test accepts

" " rejects

$$\begin{matrix} \updownarrow \\ f(x) \cdot f(y) \neq f(x \odot y) \end{matrix}$$



0/1 indicator var }
$$\frac{1 - f(x) \cdot f(y) \cdot f(x \odot y)}{2} = \begin{cases} 0 & \text{if accepts} \\ 1 & \text{o.w.} \end{cases}$$

Now we have a new way to express rejection probability:

rejection probability

$$\begin{aligned} \delta_f &\equiv \Pr_{x,y} [f(x) \odot f(y) \neq f(x \odot y)] \\ &= E_{x,y} \left[\frac{1 - f(x) \cdot f(y) \cdot f(x \odot y)}{2} \right] \end{aligned}$$

Fourier Analysis on Boolean Cube

want basis to describe all fctns. $f: \{\pm 1\}^n \rightarrow \{\pm 1\}$

Linear fctns:

$$S \subseteq \{1..n\}$$

for $x \in \{\pm 1\}^n$,

$$\chi_S(x) = \prod_{i \in S} x_i$$

parity fctns

define $\langle f, g \rangle = \frac{1}{2^n} \sum_{x \in \{\pm 1\}^n} f(x) g(x)$ inner product (but normalized)

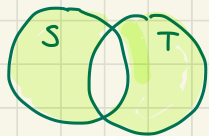
Fact parity (linear) fctns $\{\chi_S\}$ is orthonormal basis w.r.t. inner product!

Proof of fact:

$$\bullet \langle \chi_S, \chi_S \rangle = \frac{1}{2^n} \sum_x \underbrace{(\chi_S(x))^2}_{\substack{+1 \\ +1}} = \frac{2^n}{2^n} = 1 \quad \text{normal}$$

• if $S \neq T$:

$S \Delta T$:



$$\langle \chi_S, \chi_T \rangle = \frac{1}{2^n} \sum_x \chi_S(x) \cdot \chi_T(x)$$

if $i \in S \Delta T$
then $\chi_i \cdot \chi_i = 1$
"drops out"
so can ignore

$$= \frac{1}{2^n} \sum_x \chi_{S \Delta T}(x)$$

nonempty since $S \neq T$
so pick $j \in S \Delta T$

$$= \frac{1}{2^n} \sum_{\text{pairs } x, x^{\oplus j}}$$

$x^{\oplus j} = x$ with j th
bit flipped

$$= \frac{1}{2^n} \sum_{\text{pairs } x, x^{\oplus j}} \underbrace{x_j \cdot \prod_{i \in (S \Delta T) \setminus \{j\}} x_i}_{\text{equal}} + \underbrace{\bar{x}_j \cdot \prod_{i \in (S \Delta T) \setminus \{j\}} x_i}_{\text{equal}}$$

$$= \frac{1}{2^n} \sum_{\text{pairs}} 0$$

$$= 0$$

one is +1
the other is -1
so sum to 0

Orthogonal!

So $\{\chi_S\}$ is an orthonormal basis

Thm f is uniquely expressible as linear comb. of χ_s .

$$\begin{aligned}\underline{\text{Def.}} \quad \hat{f}(s) &\equiv \langle f, \chi_s \rangle \\ &= \frac{1}{2^n} \sum_{x \in \{\pm 1\}^n} f(x) \chi_s(x)\end{aligned}$$

Fourier
Coefficients
of
 f

$$\underline{\text{Thm}} \quad \forall f \quad f(x) = \sum \hat{f}(s) \chi_s(x)$$

Fourier coeffs of linear fctns:

$$\underline{\text{Fact}} \quad f \text{ linear} \Leftrightarrow \exists s \subseteq [n] \text{ st. } \begin{cases} \hat{f}(s) = 1 \\ \hat{f}(T) = 0 \end{cases} \quad \begin{array}{l} \leftarrow \text{one is} \\ \text{really} \\ \text{big} \\ \leftarrow \text{others} \\ \text{are } 0 \end{array}$$

e.g. if $f(x) = x_1 \cdot x_2$

$$f(x) = 0 \cdot \chi_{\emptyset} + 0 \cdot \chi_{\{1\}} + 0 \cdot \chi_{\{2\}} + 1 \cdot \chi_{\{1,2\}}$$

Fourier coeffs characterize distance to linear:

Lemma $\forall S \subseteq [n]$

$$\hat{f}(s) = 1 - 2 \operatorname{dist}(f, \chi_S)$$
$$= 1 - 2 \Pr_{x \in \{\pm 1\}^n} [f(x) \neq \chi_S(x)]$$

Pf $2^n \cdot \hat{f}(s) = \sum_x f(x) \chi_S(x)$ def of Fourier coeff

$$= \sum_{\substack{x \text{ st.} \\ f(x) = \chi_S(x)}} 1 + \sum_{\substack{x \text{ st.} \\ f(x) \neq \chi_S(x)}} -1$$

$$= 2^n [1 - \operatorname{dist}(f, \chi_S)] - 2^n [\operatorname{dist}(f, \chi_S)]$$

$$= 2^n (1 - 2 \cdot \operatorname{dist}(f, \chi_S))$$

~~□~~

example: $f =$ all -1 's

for $s = \emptyset$: $\operatorname{dist}(f, \chi_\emptyset) = 1$ so $\hat{f}(\emptyset) = -1$

$\forall s \neq \emptyset$ $\operatorname{dist}(f, \chi_s) = \frac{1}{2}$ why? so $\hat{f}(s) = 0$

Observation: Any two distinct linear fctns differ on exactly $\frac{1}{2}$ of inputs

pf. let $f = \chi_T$ for $T \neq S$
 $g = \chi_S$

note $\hat{f}(s) = \hat{\chi}_T(s) = \langle \chi_T, \chi_s \rangle$

$0 = \langle \chi_T, \chi_S \rangle = 1 - 2 \text{dist}(\chi_T, \chi_S)$

↑
since orthonormal

↑
lemma

$\Rightarrow \text{dist}(\chi_T, \chi_S) = \frac{1}{2}$

Corollary: if $S \neq \emptyset$, $\chi_S(x) = +1$ on exactly $\frac{1}{2}$ the inputs

Very useful tools:

Plancherel's identity

$\langle f, g \rangle = \left\langle \sum_{S \subseteq [n]} \hat{f}(s) \chi_s, \sum_{T \subseteq [n]} \hat{g}(T) \chi_T \right\rangle$

$= \sum_{S, T} \hat{f}(s) \hat{g}(T) \langle \chi_s, \chi_T \rangle$

bilinearity of $\langle \cdot, \cdot \rangle$

$= \sum_s \hat{f}(s) \hat{g}(s)$

$\underbrace{\quad}_0$ if $S \neq T$
 $\underbrace{\quad}_1$ if $S = T$

Parseval's identity:

$$\langle f, f \rangle = \sum_s \hat{f}(s)^2$$

"Boolean Parseval's"

$$\text{if } f: D \rightarrow \{\pm 1\} \quad \langle f, f \rangle = \frac{1}{2^n} \sum_x \underbrace{f(x) \cdot f(x)}_{=+1 \text{ since } f \text{ is Boolean}} = \frac{1}{2^n} \cdot 2^n = 1$$

$$\text{so } \sum_s \hat{f}(s)^2 = 1$$

Analysis of linearity test

if δ_f is small, can we conclude that f is close to linear?

doesn't this contradict Oppenheim's lower bound? this

is only for $f: \{\pm 1\}^n \rightarrow \{\pm 1\}$

YES! rejection probability gives upper bound on distance.

Thm. $f: \{\pm 1\}^n \rightarrow \{\pm 1\}$ is δ_f -close to some linear fctn

Pf.

$$\begin{aligned} & E_{x,y} [f(x)f(y)f(x \circ y)] \\ &= E_{x,y} \left[\left(\sum_s \hat{f}(s) \chi_s(x) \right) \cdot \left(\sum_T \hat{f}(T) \chi_T(y) \right) \cdot \left(\sum_u \hat{f}(u) \chi_u(x \circ y) \right) \right] \\ &= E_{x,y} \left[\sum_{s,T,u} \hat{f}(s) \hat{f}(T) \hat{f}(u) \underbrace{\chi_s(x) \chi_T(y) \chi_u(x \circ y)}_{\text{what is this?}} \right] \end{aligned}$$

$$\begin{aligned} \text{if } s=T=u : & \chi_s(x) \chi_T(y) \chi_u(x \circ y) \\ &= \prod_{i \in S} x_i \cdot \prod_{i \in S} y_i \cdot \prod_{i \in S} \overbrace{(x_i y_i)}^{x_i \cdot y_i} \\ &= \prod_i x_i \cdot y_i \cdot (x_i \cdot y_i) = \prod_i x_i^2 y_i^2 = \prod_i (1 \cdot 1) = 1 \end{aligned}$$

if $\neg (S=T=U)$:

$$E_{x,y} [\chi_S(x) \chi_T(y) \chi_U(xoy)]$$

$$= E_{x,y} \left[\prod_{i \in S} \chi_i \prod_{j \in T} y_j \prod_{k \in U} (\chi_k \cdot y_k) \right]$$

$$= E_{x,y} \left[\prod_{i \in S \Delta U} \chi_i \cdot \prod_{j \in T \Delta U} y_j \right]$$

$$= E_{x,y} \left[\prod_{i \in S \Delta U} \chi_i \right] \cdot E_{x,y} \left[\prod_{j \in T \Delta U} y_j \right] \quad \text{since } x,y \text{ indep}$$

assumption
 \Rightarrow one
of these
holds
so \rightarrow

$$\text{if } S \neq U \text{ then } S \Delta U \neq \emptyset \Rightarrow E_{x,y} \left[\prod_{i \in S \Delta U} \chi_i \right] = 0$$

$$\text{analogously if } T \neq U, E_{x,y} \left[\prod_{j \in T \Delta U} y_j \right] = 0$$

$$= 0 \quad \text{so all these terms drop out!}$$

$$E_{x,y} [f(x)f(y)f(xoy)]$$

$$= E_{x,y} \left[\left(\sum_s \hat{f}(s) \chi_s(x) \right) \cdot \left(\sum_T \hat{f}(T) \chi_T(y) \right) \cdot \left(\sum_u \hat{f}(u) \chi_u(xoy) \right) \right]$$

$$= E_{x,y} \left[\sum_{s,T,u} \hat{f}(s) \hat{f}(T) \hat{f}(u) \chi_s(x) \chi_T(y) \chi_u(xoy) \right]$$

$$= \sum_{S=T=U} \hat{f}(s)^3$$

$$\leq \max_s \hat{f}(s) \cdot \underbrace{\sum_s \hat{f}(s)^2}_{=1 \text{ by "Parseval's"}} = 1$$

$$= \max_s \hat{f}(s)$$

$$= \max_s (1 - 2 \cdot \text{dist}(f, \chi_s))$$

$$= 1 - 2 \cdot \min_s (\text{dist}(f, \chi_s))$$

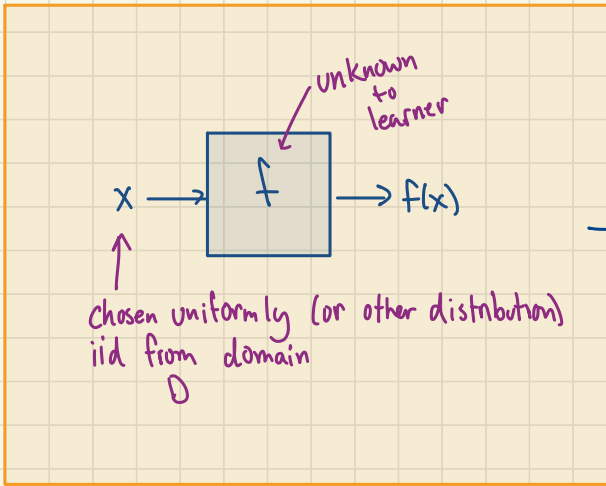
$$\begin{aligned} \text{So } \delta_f &\geq \frac{1 - (1 - 2 \min_s (\text{dist}(f, \chi_s)))}{2} \\ &= \min_s \text{dist}(f, \chi_s) \end{aligned}$$

$\Rightarrow \exists s$ st f is δ_f -close to χ_s

Learning

Learn from random, uniform examples ↙ How to formalize?

lots of other models -
will mention some
others later



$x_1, f(x_1)$
 $x_2, f(x_2)$
 \vdots
 $x_m, f(x_m)$

m random, labelled examples

Example oracle $E_x(f)$

After seeing several examples, learner should output hypothesis h .

what do we hope h satisfies?

- hopefully $h = f$ ← asking too much?

- at least $\text{dist}(h, f) \leq \epsilon$

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

what distribution on inputs do we use?

today uniform

in general, match distribution of example oracle

Valiant's
PAC
model
"probably
approximately
correct"

Common terms for same thing:

- $\text{dist}(h, f)$
 $x \in D$

- $\text{error}(h)$ (w.r.t. f is understood from context)
 $x \in D$

- h is ϵ -close to f (dist D understood from context)

Note in above:

$X \in \mathcal{D}$ can be chosen according to uniform or any other prespecified distribution

Note if f is arbitrary, there is nothing you can do that is "efficient" in terms of sample complexity (e.g. you can't learn a random fctn f without seeing the value of f for most inputs)

However, if you know something about f , there may be hope.

here: what if you know that f is a member of fctn family \mathcal{C} ?

e.g. \mathcal{C} = linear fctns
K-term DNF
⋮
⋮

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

- A given $\epsilon, \delta > 0$
access to $Ex(f)$ for $f \in \mathcal{C}$
- A outputs h st. with prob $\geq 1 - \delta$
error (h) w.r.t. f is $\leq \epsilon$
 h is ϵ -close to f

Parameters of interest:

- m # samples used by A "Sample Complexity"
- ϵ accuracy parameter
- δ confidence parameter
- runtime? hope for poly($\log(\text{domainsize}), \frac{1}{\epsilon}, \frac{1}{\delta}$)

• description of h ?

- should it be similar to description of fctns in \mathcal{C} ? "proper learning"
- at least should be relatively compact & efficient to evaluate
 $O(\log |\mathcal{C}|)$

Remarks

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

Efficient learning algorithm for conjunctions:

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

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

Note:

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

e.g. how to distinguish $f(x) = x_1 x_2 \dots x_n$
from $f'(x) = 0 \quad \forall x$?

Behavior of poly time algorithm:

for i in conjunction:

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

for i not in conjunction:

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

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

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

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

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

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

or $\Omega(\frac{1}{\epsilon} \log \frac{n}{\delta})$ total examples, will suffice
to rule out all $i \notin$ conjunction.