Learning via Fourier Coeffs

- Some fcthns & their Fourier representation

- the low degree algorithm

- applications

# Learning via Fourier Representation

learning algorithms based on estimating fourier representation of fcn $P$ (similar to poly interpolation)

## Approximating one Fourier coefficient:

**lemma** can approx any specific Fourier coeff $S$ to w/in additive $\gamma$

(ie. $|\text{output} - \hat{f}(s)| \leq \gamma$)

with prob $\geq 1 - \delta$ in $O\left(\frac{1}{\gamma^2} \log \frac{1}{\delta}\right)$ samples

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

Note no queries needed!!

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 $\delta = \frac{1}{2^n}$, give $O\left(\frac{1}{\gamma^2} \cdot n\right)$ 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) $\overline{AND}$ on $T \subseteq N$ s.t. $|T| = k$

$$\overline{AND} (x_{i_1} \cdots x_{i_k}) = 1 \qquad \text{if} \quad \forall_{i \in T} = \{\lambda_1 \cdots \lambda_k\}$$
$$X_{i_j} = -1$$

$$\text{o.w.} \qquad \forall_{i \in T} \quad X_i = -1 \Big\} \begin{array}{l} \text{corresponds} \\ \text{to} \\ \text{AND fctn} \\ \text{over } \{0,1\}^S \end{array}$$

define $f(x) = \begin{cases} 1 & \text{if} \quad -1 \\ 0 & \text{o.w.} \end{cases}$

$$= \frac{(1-X_{i_1})}{2} \cdot \frac{(1-X_{i_2})}{2} \cdots \frac{(1-X_{i_k})}{2}$$

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

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

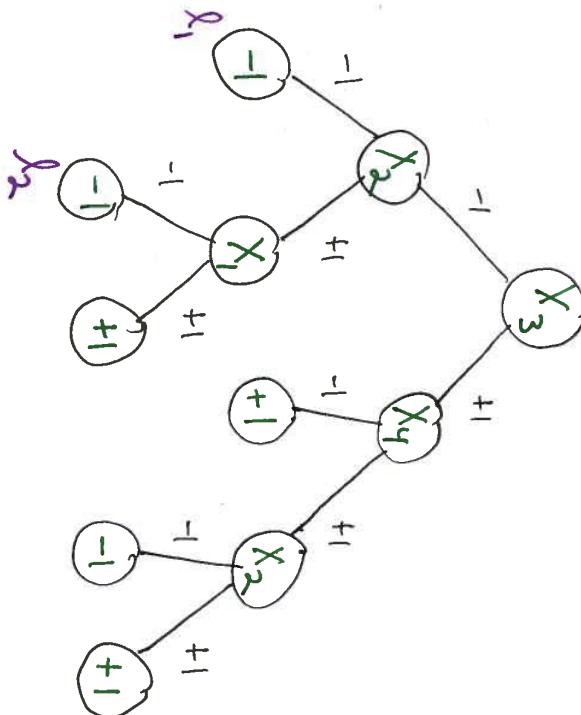$$= -1 + \frac{2}{2^k} + \sum_{\substack{S \subseteq T \\ |S| > 0}} \frac{(-1)^{|S|}}{2^{k-1}} \chi_S$$

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:

$$f_l(x) = \prod_{i \in V_l} \frac{(1 \pm X_i)}{2}$$  ← left or right

vars visited on path to leaf l

$$= \frac{1}{2^{|V_l|}} \sum_{S \subseteq V_l} (\pm 1) X_S$$  ← $(-1)^{\# \text{left turns taken in } S}$

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

So $f(x) = \displaystyle\sum_{l \text{ leaves of } T} f_l(x) \, val(l)$

exactly one of these is 1 / others are 0

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

## Low degree algorithm

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

if $\displaystyle\sum_{\substack{S \subseteq [n] \\ \text{st.} \\ |S| > \alpha(\varepsilon,n)}} \hat{f}(S)^2 \le \varepsilon \qquad \forall\, 0 < \varepsilon < 1$

for Boolean $f$, this implies $\displaystyle\sum_{\substack{S \subseteq [n] \\ \text{st.} \\ |S| \le \alpha(\varepsilon,n)}} \hat{f}(S)^2 \ge 1-\varepsilon$

examples

1) fcn $f$ which depends on $\le k$ vars has $\log(\tfrac{4}{\varepsilon})$-F.C.

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

$\begin{cases} \text{if } f \text{ doesn't depend} \\ \text{on } X_i, \text{ then all} \\ \hat{f}(S) \text{ for which } i \in S \\ \text{satisfy } \hat{f}(S)=0 \end{cases}$

2) $f = $ AND on $T \subseteq \{1..n\}$ has $\log(\tfrac{4}{\varepsilon})$-F.C.

• all $\hat{f}(S)^2 = 0$ for $|S| > |T|$

• if $|T| \le \log\tfrac{4}{\varepsilon}$ then ✓

• if $|T| \ge \log\tfrac{4}{\varepsilon}$ then :

$\hat{f}(\phi)^2 = \left(1 - 2\Pr(f(x) \ne \chi_\phi(x))\right)^2 = \left(1 - \tfrac{2}{2^{|T|}}\right)^2 > 1-\varepsilon$

so $\displaystyle\sum_{S \ne \phi} \hat{f}(S)^2 \le \varepsilon \quad + f$ has 0-F.C.

Now, let's approximate fctns with $d = \alpha(\varepsilon, n)$ F.C.:

## Low Degree Algorithm

Given    $d \rightarrow$ degree

       $\gamma \rightarrow$ accuracy

       $\delta \rightarrow$ confidence

Algorithm

- Take $m = O\left(\frac{n^d}{\gamma} \ln \frac{n^d}{\delta}\right)$ samples

- $\hat{c}_s \leftarrow$ estimate of $\hat{f}(s)$ (for each $S$   s.t. $|S| \leq d$ )

- output $h(x) = \sum\limits_{|S| \leq d} \hat{c}_s \chi_s(x)$   $\underbrace{\leq \binom{n}{d}}_{}$ of these

       $|S| \leq d$       Can reuse same samples for each!

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

Why does this work?

Two stages:

1) show that if $f$ has low F.C. $\swarrow$ $L_2$ dist $\frac{1}{2^n}$

     then    $E_x[(f(x)-h(x))^2]$ small $\swarrow$

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

               $\uparrow$

            Hamming dist

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

**Pf**

**Claim** with prob $\geq 1-\delta$, $\forall S$ s.t. $|S| \leq d$, $|c_S - \hat{f}(S)| \leq \gamma$
for $\gamma \leftarrow \sqrt{\frac{\varepsilon}{n^d}}$

**Pf of claim**

note, $\frac{1}{\gamma^2} = \frac{n^d}{\varepsilon}$

Chernoff bnd $\Rightarrow O\left(\frac{n^d}{\varepsilon} \ln \frac{n^d}{\delta}\right) = O\left(\frac{1}{\gamma^2} \ln \frac{n^d}{\delta}\right)$ samples
yields $\Pr\left[|c_S - \hat{f}(S)| > \gamma\right] < \frac{\delta}{n^d}$

Union bnd $\Rightarrow \Pr\left[\exists S \text{ s.t. } |c_S - \hat{f}(S)| > \gamma\right] < \delta$
only (a) $< n^d$ such [S of size $\leq d$]

Assume $\forall S$ s.t. $|S| \leq d$, $|c_S - \hat{f}(S)| \leq \gamma$

define $g(x) \equiv f(x) - h(x)$
Fourier transform is linear $\Rightarrow \forall S \quad \hat{g}(S) = \hat{f}(S) - \hat{h}(S)$
by defn, $\forall S$ s.t. $|S| > d, \hat{h}(S) = 0$
$|S| \leq d, \hat{h}(S) = c_S$

$\Rightarrow \hat{g}(S) = \hat{f}(S) - \hat{h}(S)$
$\hat{g}(S) = \hat{f}(S)$
$\hat{g}(S) = \hat{f}(S) - c_S \Rightarrow$
so $\hat{g}(S)^2 \leq \gamma^2$

so $E[(f(x)-h(x))^2] = E[g(x)^2]$

$$= \sum_s \hat{g}(s)^2 \qquad \text{Parseval}$$

$$= \sum_{|s|\leq d} \hat{g}(s)^2 + \sum_{|s|>d} \hat{g}(s)^2$$

$\underbrace{\leq n^d \cdot \gamma^2}_{\leq \gamma^2}$ $\qquad \underbrace{\leq \varepsilon}_{\text{by F.C.}}$

$$\leq T + \varepsilon \qquad \blacksquare$$

Thm  $f: \{\pm1\}^n \to \{\pm1\}$
  $h: \{\pm1\}^n \to \mathbb{R}$
  then  $\Pr[f(x) \neq \text{Sign}(h(x))] \leq E[(f(x)-h(x))^2]$

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

$\Pr[f(x) \neq \text{Sign}(h(x))] = \frac{1}{2^n}\sum_x 1_{\{f(x)\neq \text{Sign}(h(x))\}} \leq E[(f(x)-h(x))^2]$

defn $\left\{ \text{show term by term} \right.$

But  if  $f(x) = \text{Sign}(h(x))$
  $\underbrace{(f(x)-h(x))^2}_{\geq 0} \qquad \underbrace{\frac{1}{f(x)\neq\text{Sign}(h(x))}}_{=0}$

  if  $f(x) \neq \text{Sign}(h(x))$
  $\underbrace{(f(x)-h(x))^2}_{\geq 1} \qquad \underbrace{\frac{1}{f(x)\neq\text{Sign}(h(x))}}_{=1}$

So $\forall x$, $\left(f(x)-h(x)\right)^2 \geq \frac{1}{f(x)\neq\text{Sign}(h(x))}$ $\qquad \blacksquare$

Correctness of learning algorithm:

Thm. if $C$ has fourier concentration $d = \alpha(\varepsilon, n)$

then there is a $q = O(\frac{n^d}{\varepsilon} \log \frac{n}{\delta})$ sample

uniform distribution learning algorithm for $C$

i.e. algorithm gets $q$ samples & with prob $\geq 1-\delta$

outputs $h'$ s.t. $\Pr[f \neq h'] \leq 2\varepsilon$

Pf. run low degree alg with $\tau = \varepsilon$

get $h$ s.t. $E[(f-h)^2] \leq \varepsilon + \varepsilon = 2\varepsilon$

output sign $(h)$

∎

Applications

1) Bounded depth decision trees

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

$f_\ell(x)$ : $\underbrace{\text{const}}$ ← fctn which depends on ≤ depth many vars

by linearity, $\hat{f}(S) = \sum_\ell \text{val}(\ell) \cdot \hat{f_\ell}(S)$  which is 0 if $|S| > \text{dept}$