

Today's lecture:

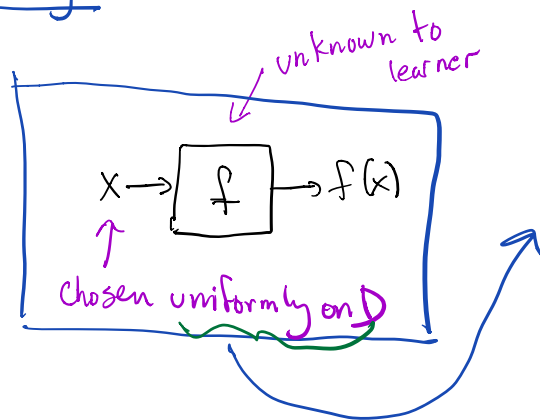
The PAC learning model
motivation
definition

Occam's razor

Learning conjunctions

(if time: begin learning via Fourier representation)

Learning how to formalize?



labelled examples

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

m random
labelled
examples

Example oracle $E_x(f)$

Goal: ~~output f~~
output h

too hard?

is ϵ -close to f
e.g. $\Pr_{x \in D} [f(x) = h(x)] \geq 1 - \epsilon$

which distribution
today assume
uniform!

def given hypothesis h , error of h with respect to f is $\text{error}(h) = \Pr_{x \in D} [f(x) \neq h(x)]$

$\underbrace{x \in D}_u \quad \uparrow$
 f is ϵ -close to h
 wrt. uniform on D

Observe if f arbitrary then nontrivial learning is impossible

What if f is in a class of fctns \mathcal{C}

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

• A is given $\epsilon, \delta > 0$ access to $E_x(f)$ for $f \in \mathcal{C}$ ↙ according to f

• A output h st, with prob $\geq 1 - \delta$ $\text{error}(h)$ wrt. f is $\leq \epsilon$

↘ according to f

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 $\text{poly}(\log(\text{domain size}), \frac{1}{\epsilon}, \frac{1}{\delta})$
- description of h ? $|C|$
 - similar to description of all $f \in C$?
(proper learning)
 - at least should be "compact"
 $O(\log|C|)$ + efficient to evaluate

Remarks

- dependence on δ needn't be more than $O(\log(\frac{1}{\delta}))$
- uniform dist is a special case

Occam's Razor

learning is easy!
wrt sample complexity
not runtime

brute force algorithm

- draw $M = \frac{1}{\epsilon} (\ln |\mathcal{C}| + \ln \frac{1}{\delta})$ samples
- search over all $h \in \mathcal{C}$ until
find one that labels all examples
correctly. Output h .
(choose arbitrarily if > 1)

behavior:

examples come from $f \in \mathcal{C}$
good to output f
bad to output h st.
 $h + f$ not ϵ -close

h is "bad" if $\text{error}(h) \geq \epsilon$

$$\Pr[\text{bad } h \text{ consistent with examples}] \leq (1-\epsilon)^M$$

$\Pr[\text{any bad } h \text{ consistent with examples}]$

$$\leq |\mathcal{C}| \cdot (1-\epsilon)^M \quad \text{union bound}$$

$$\leq \cancel{|\mathcal{C}|} \cdot \underbrace{(1-\epsilon)^{\frac{1}{\epsilon}}}_{e^{-1}} (\ln \cancel{|\mathcal{C}|} + h \frac{1}{\epsilon})$$

$$\leq \delta$$

\Rightarrow unlikely to output any bad h ~~is~~

Proof applies to learning under any distribution

Once we have a good hypothesis h :

1) can predict values of f on
new random inputs $\Pr_{x \in \mathcal{X}} [f(x) = h(x)] \geq 1 - \epsilon$
according to \mathcal{D}

2) can compress description of samples

$$\begin{array}{l} (x_1, f(x_1)) \quad (x_2, f(x_2)) \quad \dots \quad (x_m, f(x_m)) \quad \overset{\# \text{ bits}}{m(\log|D| + \log|R|)} \\ \Downarrow \\ x_1 \dots x_m, \text{ description of } h \quad m \cdot \log|D| + \log|C| \end{array}$$

learning \Rightarrow prediction + "compression"

Occam's Razor: simplest explanation is best

An efficient learning algorithm

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

ie. $f(x) = X_i X_j \bar{X}_k$
(x_i, x_n)

Observe: how to distinguish

$f(x) = X_i, \dots, X_n$ } need $\sim 2^n$
from } samples
 $f(x) = 0$

\Rightarrow can't hope for poly time \pm 0-error

Brute force algorithm: (ie. alg in Occam's razor)

try each $f \in \mathcal{C}$ $|\mathcal{C}| \geq 2^n$

union bound \Rightarrow need $\Omega(\frac{1}{\epsilon} \ln 2^n + \ln \frac{1}{\delta})$
Samples

Poly time algorithm

Simplifying assumption!

assume $\Pr_{x \in \{0,1\}^n} [f(x)=1] > \epsilon$

\Rightarrow in a sample of size m , in expectation many $\geq \epsilon m$ "positive" examples $= 1$

Algorithm:

1	2	3	4		$f(x)$
0	1	1	0	+	1
1	1	1	1	-	0
0	1	0	1	+	1
1	1	0	0	-	0

Take M examples, K of which are "positive"
 $f(x)=1$

let $V = \{ \text{vars set same way} \}$
in each positive example

$$V = \{1, 2\}$$

output $h(x) = \bigwedge_{i \in V} x_i^b$

$$h(x) = \bar{x}_1 x_2$$

Behavior:

$$f(x) = \bar{x}_1$$

for $i \in$ conjunction:
must be set same way in
each positive example \Rightarrow in V

for $i \notin$ conjunction:

$$\begin{aligned} \Pr[i \in V] &\leq \Pr[i \text{ set same in each} \\ &\quad \text{of } K \text{ positive examples}] \\ &\leq \frac{1}{2^K} + \frac{1}{2^K} = \frac{1}{2^{K-1}} \end{aligned}$$

$$\begin{aligned} \Pr[\text{any } i \text{ not in conjunction survives}] \\ &\leq \frac{n}{2^{K-1}} \\ &\leq \delta \quad \text{if pick } K = \log \frac{n}{\delta} \end{aligned}$$

$\Rightarrow \Omega(\log \frac{n}{\delta})$ positive examples
 $+ \Omega(\frac{1}{\epsilon} \log \frac{n}{\delta})$ total examples suffice.

More general algorithm:

using
 $\text{poly}(1/\epsilon)$
 samples

estimate $\Pr[f(x)=1]$ to additive error $\pm \frac{\epsilon}{4}$
 if estimate $< \epsilon/2$, output $h=0$

$$\Rightarrow \Pr[f(x)=1] \leq \frac{\epsilon}{2} + \frac{\epsilon}{4} < \epsilon$$

good answer

O.w. continue

$$\Rightarrow \Pr[f(x)=1] \geq \frac{\epsilon}{2} - \frac{\epsilon}{4} \geq \frac{\epsilon}{4}$$

\Rightarrow see positive example
 every $\frac{4}{\epsilon}$ samples

\Rightarrow above algorithm is
 efficient

QED