


Learning with
Classification Noise

The Model

- Same as PAC, but now we see $\langle x, \hat{y} \rangle$, $x \in P$ and:

$$\hat{y} = \begin{cases} c(x) & \text{with prob. } 1-\eta \\ \tau c(x) & \text{with prob. } \eta \end{cases}$$

$c \in \mathcal{H}$ target fn.

$\eta \in [0, 1/2]$ noise rate

- Goal: still want
 $E(h) = P[h(x) \neq c(x)] \leq \epsilon$

Say C CN learnable by \mathcal{H}
if PAC criteria met
by algo using noisy \hat{y} 's.

- What C are CN-learnable?
- How much harder is CN?
- Are there "general principles" for CN?

Note: Can't use
consistency any more...

An Observation

$$\Pr[h(x) \neq \tilde{y}] \stackrel{\Delta}{=} \hat{\epsilon}(h)$$

$$= (1 - \pi) \epsilon(h) + \pi (1 - \epsilon(h))$$

$$= \underbrace{\pi}_{\text{fixed}} + \underbrace{(1 - 2\pi)}_{\text{increasing}} \epsilon(h)$$

fixed increasing
with $\epsilon(h)$

\therefore ordering by $\tilde{\epsilon}(h) \equiv$
ordering by $\epsilon(h)$

$\therefore \tilde{\epsilon}(h)$ minimization \equiv
 $\epsilon(h)$ minimization

Note: $\tilde{\epsilon}(h_1) - \tilde{\epsilon}(h_2) =$

$$(1-2\eta)(\epsilon(h_1) - \epsilon(h_2))$$

shrinks need resolution
 ϵ

\Rightarrow sample sizes $\geq \frac{1}{(1-2\eta)\epsilon}$

CN learning (monotone) conjunctions

• Original algo:

- start with $h = x_1 x_2 \dots x_n$
- $\langle x_i, + \rangle \Rightarrow$ delete any
 $x_i \leftarrow 0$
↑
Won't work on \tilde{y}

Algo too brittle/sensitive.

Define $\text{fl}(x_i)$:

$$\text{pol}(x_i) = P[x_i = 0]$$

$$(\geq) \text{pol}_1(x_i) = P[x_i = 0, y = 1]$$

↗
true label $c(x)$

Say x_i is

significant: $p_{01}(x_i) \geq \epsilon/8n$

harmful: $p_{01}(x_i) \geq \epsilon/8n$

Let h contain all x_i that
are sig. but not harmful.

$\Pr[h(x)=1, y=0]$: must be

some $x_i \in C, x_i \notin h, x_i = 0$
not \Rightarrow not $\Rightarrow p_{01}(x_i) \leq \epsilon/8n$
harmful sig.

$\therefore \Pr[h(x)=1, y=0] \leq n\left(\frac{\epsilon}{8n}\right) = \epsilon/8$

$\Pr[h(x)=0, y=1]$: must be

some $x_i \notin C, x_i \in h, x_i = 0$

$$p_{0,i}(x_i) \leq \varepsilon/8n$$

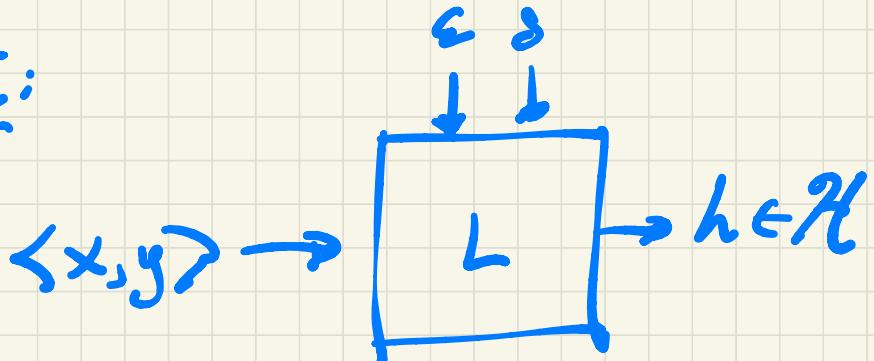
$$\therefore \Pr[h(x)=0, y=1] \leq n \left(\frac{\varepsilon}{8n} \right) = \frac{\varepsilon}{8}$$
$$\not\models \varepsilon(h) \leq \varepsilon/8 + \varepsilon/8 = \varepsilon/4.$$

The $p_{0,i}(x_i)$ are easy to estimate from $x^n P$.

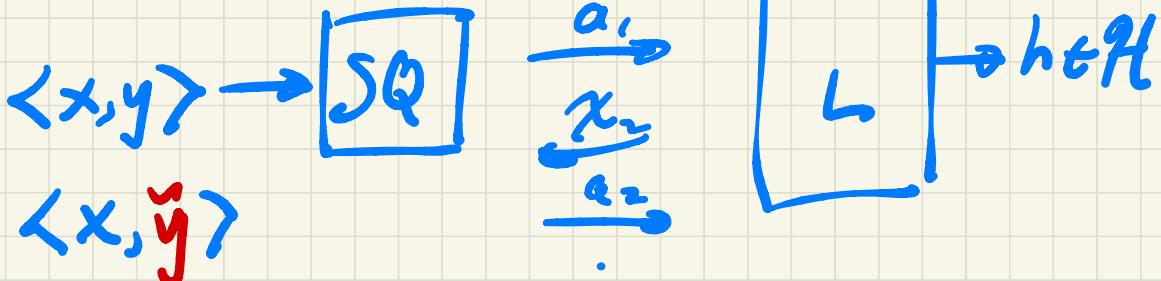
How can we estimate the $p_{0,i}(x_i)$ from $\langle x, \bar{y} \rangle$?

The Statistical Query (SQ) Learning Model

PAC:



SQ:



"statistical
queries"

Queries are predicates χ
over $\langle x, y \rangle$ (no noise):

$$\chi(x, y) \in \{0, 1\}$$

E.g. $\chi(x, y) = 1 \Leftrightarrow \begin{aligned} x_5 &= 1, \\ y &= 0 \\ (\rho_{01}(x_5)) \end{aligned}$

$$\chi(x, y) = 1 \Leftrightarrow \begin{aligned} x_1 &\leq 1.7, \\ x_2 &\geq 2.9, \\ y &= 1 \end{aligned}$$

Assume χ

"reasonable" / computable

Viz w $X(x,y)$ as a request for

$$P_x \stackrel{a}{=} \Pr[X(x,y) = 1]$$

But SQ won't answer exactly, but with tolerance Σ :

$$\hat{P}_x \in P_x + [-\Sigma, \Sigma]$$

Full query: X, Σ

Get back: \hat{P}_x

Say C is SQ-learnable by \mathcal{H}

if \exists algo L

$\forall c \in C, \forall P, \forall \epsilon > 0:$

- L makes $\text{poly}(1/\epsilon, n, \text{size}(c))$ queries, each of tolerance $\epsilon' \geq \frac{1}{\text{poly}(\cdot)}$ (e.g. ϵ^2/n^3)
- L runs in poly time
- L outputs $h \in \mathcal{H}$ s.t.
 $\epsilon(h) \leq \epsilon.$

Theorem

$C \text{ SQ learnable by } \mathcal{H} \Rightarrow$
 $\underline{C \text{ PAC learnable by } \mathcal{H}}.$

Proof. Let L be SQ algo.

On each query x, τ of L ,
sample $m < x, y >$ pairs to
get estimate \hat{P}_x^n . If
 $m = \frac{1}{\epsilon^2} \log\left(\frac{\ell}{\delta}\right)$ then
w.p. $\geq 1 - \frac{\delta}{\ell}$ $\hat{P}_x^n \in P_x + [-\tau, \tau]$.

Now let $\ell = \# \text{ queries of } L$
 $\leq \text{ runtime of } L$.

More Interesting Theorem

C SQ learnable by $\mathcal{H} \Rightarrow$

C CN learnable by \mathcal{H} .



How can we

simulate SQ

from $\langle x, \hat{y} \rangle$?

Decomposing P_x

$$\cdot P_x = \Pr[X(x,y) = 1]$$

$$\cdot \hat{P}_x = \Pr[X(x, \tilde{y}) = 1]$$

Can estimate directly

Now define disjoint X_1, X_2

$$X_1 \cap X_2 = \emptyset, X_1 \cup X_2 = X:$$

$$X_1 \triangleq \{x \in X : X(x, 0) \neq X(x, 1)\}$$

$$X_2 = X - X_1$$

Can test x for membership

Lct $p_1 = P[X_1]$, $p_2 = P[X_2]$

$$p_2 = 1 - p_1$$

Induced distributions:

$$P_1 \text{ over } X_1: P_1[x] = P[x]/p_1$$

$$P_2 \text{ over } X_2: P_2[x] = P[x]/p_2$$

And finally:

$$P_{\chi^1} = P_1[\chi(x,y) = 1]$$

$$P_{\chi^2} = P_2[\chi(x,y) = 1]$$

Note: The following can
be estimated from $\langle \tilde{x}, \tilde{y} \rangle$:

$$\hat{P}_x, P_x^2, P_1, P_2$$

Goal: formula where

$$P_x = \text{some function}$$
$$\text{of } \hat{P}_x, P_x^2, P_1, P_2, \eta$$

$$\tilde{P}_x = (\overset{\check{y}=y}{1} - \tilde{n}) P_x +$$

$$\tilde{n} (p_1 P_1 [x(x, \overset{\check{y}}{y}) = 1] \\ + p_2 P_2 [x(x, \overset{\check{y}}{y}) = 1])$$

$$\tilde{P}_x = \overset{\tilde{y}=y}{(1-\eta)} P_x +$$

$$\eta (p_1 P_1 [Z(x, \tilde{y}) = 1] \\ + p_2 P_2 [Z(x, \tilde{y}) = 1])$$

$$= P_1 P_x^1 + P_2 P_x^2$$

$$\tilde{P}_x = (\overset{\check{y}=y}{1-\eta}) P_x +$$

$$\eta \left(p_1 P_1 [X(x, \overset{\check{y}=y}{y}) = 1] + p_2 P_2 [X(x, \overset{\check{y}=y}{y}) = 1] \right)$$

$$= P_1 [X(x, y) = 0]$$

$$= 1 - P_x'$$

$$\tilde{P}_x = (\overset{\check{y}=y}{1} - \pi) P_x +$$

$$\pi (p_1 P_1 [X(x, \overset{\check{y}}{y}) = 1] + p_2 \overset{\text{circled}}{P_2 [X(x, \overset{\check{y}}{y}) = 1]})$$

$$\overset{\text{red arrow}}{=} P_2 [X(x, \overset{\check{y}}{y}) = 1]$$

$$= P_x^2$$

$$\tilde{P}_x = (\overset{\text{y}=\text{y}}{1-\eta}) P_x +$$

$$\eta \left(p_1 P_1 [x(x, \overset{\text{y}=\text{y}}{y}) = 1] + p_2 P_2 [x(x, \overset{\text{y}=\text{y}}{y}) = 1] \right)$$

$$= (1-\eta) (p_1 P_x' + p_2 P_x'') + \eta (p_1 (1 - P_x') + p_2 P_x'')$$

$$\tilde{P}_x = (\overset{\check{y}=y}{1-\eta}) P_x +$$

$$\eta (p_1 P_1 [Z(x, \overset{\check{y}}{y}) = 1] \\ + p_2 P_2 [Z(x, \overset{\check{y}}{y}) = 1])$$

$$= (1-\eta) (p_1 P_x' + p_2 P_x^2) \\ + \eta (p_1 (1 - P_x') + p_2 P_x^2)$$

$$= (1-2\eta) p_1 P_x' \\ + p_2 P_x^2 + \eta p_1$$

Now solve for P_x'

$$P_x' = \frac{\tilde{P}_x - p_2 P_x^2 - n \tilde{P}_1}{(1-2n)p_1}$$

Plug into $P_x = p_1 P_x' + p_2 P_x^2$:

$$P_x = \frac{\tilde{P}_x - p_2 P_x^2 - n \tilde{P}_1}{1-2n} + p_2 P_x^2$$

$$= \frac{\tilde{P}_x}{1-2n} + \left(1 - \frac{1}{1-2n}\right) p_2 P_x^2$$

$$\frac{-n}{1-2n} P_1$$

Sensitivity Analysis (informal)

If we estimate \hat{P}_x
within $\pm \rho$, error

$$\ln \frac{\hat{P}_x}{1-2\eta} \text{ is } \pm \frac{\rho}{1-2\eta}$$

So better make

$$\frac{\rho}{1-2\eta} < \tau, \rho < \tau(1-2\eta)$$

Sensitivity Analysis (informal)

If error in p_2 is $\pm a$,

P_x^2 is $\pm b$, error in

$$P_2 P_x^2 \leq a + b + ab,$$

make $a, b < \epsilon(1-2n)$

etc.

Guessing \hat{N} (informed)

- Suppose we are only given an

upper bound:

$$\eta \leq \hat{\bar{N}} \leq \frac{1}{2}$$

- Divide $[0, \hat{\bar{N}}]$ into values $0, \Delta, 2\Delta, \dots, \frac{\Delta}{\hat{\bar{N}}}$

- Run simulation for each guess $\hat{\eta} = i \cdot \Delta$
- Verify hypothesis has $E(h) \leq \Sigma$
- One guess will be $\pm \Delta$ of true η
- Make Δ small wrt $\Gamma, \Sigma, \frac{1}{1-2\hat{\eta}}$
- Poly. overhead

Say x_i is

significant: $\text{pol}(x_i) \geq \epsilon/8n$

harmful: $\text{pol}_1(x_i) \geq \epsilon/8n$

Let h contain all x_i that
are sig. but not harmful.

$\Pr[h(x)=1, y=0]$: must be
some $x_i \in C, x_i \notin h, x_i = 0$
not harmful \Rightarrow not sig. $\Rightarrow \text{pol}(x_i) \leq \epsilon/8n$

$\therefore \Pr[h(x)=1, y=0] \leq n\left(\frac{\epsilon}{8n}\right) = \epsilon/8$

estimate all these
(non-adaptive)

Claims

- "Almost" all C in PAC are also in SQ (with a different algo), and thus in CN.
- conjunctions, rectangles, decision lists, linear class. etc.

Claim

- Virtually every "unpractical" ML algo has an SQ variant.
- backprop, dec. tree algos, gradient des., boosting, etc.

So what's

hard in SQ?

Parity Functions

- $X = \{0, 1\}^n$

- If set $S \subseteq \{1, 2, \dots, n\}$ define:

$$f_S(x) \triangleq \sum_{i \in S} x_i \bmod 2$$

$$= \bigoplus_{i \in S} x_i$$

- So $f_S(x) = 1 \Leftrightarrow \# \text{ bits} = 1$
in $x|_S$ is
odd.

- Let PARITY denote the
class of all (2^n) f_S .

PARITY is PAC-learnable

- Note: $f_{\alpha}(x) = \sum_i \alpha_i x_i \bmod 2$

where $\alpha, x \in \{0, 1\}^n$ and

$$\alpha_i = 1 \iff i \in S$$

- Inner product $\alpha \cdot x$ in $GF^n(2)$
- Given $\langle x_1, y_1 \rangle, \dots, \langle x_m, y_m \rangle$:

$$\begin{bmatrix} -x_1- \\ -x_2- \\ \vdots \\ -x_m- \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

Solve sys. of linear eqns!

Why is PARITY hard for SQ?

- Let $P = \text{uniform on } \{0, 1\}^n$
- Look at χ 's of form:
 $\chi(x, y) = 1 \Leftrightarrow x|T = z$
and $y = 1$
- Draw $S \subseteq \{1, \dots, n\}$ randomly
- If T is small ($|T| \leq n/2$)
then whp over draw of S ,
 $S - T \neq \emptyset$
and thus
- $\Pr_x [\chi(x, f_S(x)) = 1] = 1/2$
- "no information" conveyed

Theorem. For $\epsilon = 1/2$, and
any polynomial $p(n)$,
the number of queries
of tolerance $\xi = 1/p(n)$
required to learn
PARITY is exponential
in n .

So $\text{PARITY} \in \text{PAC}$

but $\text{PARITY} \notin \text{SQ}$.

Is $\text{PARITY} \in \text{CN}$?

Open problem,
but believed hard.

- "Nearby" NP-hardness
- Crypto proposals

Query Complexity in SQ

- target class \mathcal{C} , dist P
- view $f \in \mathcal{C}$ as ± 1 -valued
- If $f_i, f_j \in \mathcal{C}$, inner product:

$$\langle f_i, f_j \rangle_p \stackrel{\Delta}{=} \sum_x P(x) f_i(x) f_j(x)$$
$$= P[f_i = f_j] - P[f_i \neq f_j]$$

- For every $f, f' \in \text{PARITY}$,
- $$f \neq f' \Rightarrow \langle f, f' \rangle_p = 0$$

where $P = \text{uniform}$.

• Define

$$SQ(C, P) = \max d \text{ s.t.}$$

$f_1, \dots, f_d \in C$, $\forall 1 \leq i \neq j \leq d$:

$$\langle f_i, f_j \rangle_P \leq 1/d^3$$

• Largest # of "almost orthogonal" fns/vectors
in C w.r.t. P

$$• SQ(PARITY, uniform) = 2^n$$

Theorem Let $SQ(\mathcal{C}, P) = d$.

Then at least $d^{1/3}/2$

queries with $\Sigma \geq 1/d^{1/3}$

are needed to learn

\mathcal{C} wrt P in SQ model.

(even "weakly")

- d superpoly in $n \Rightarrow$ superpoly # queries

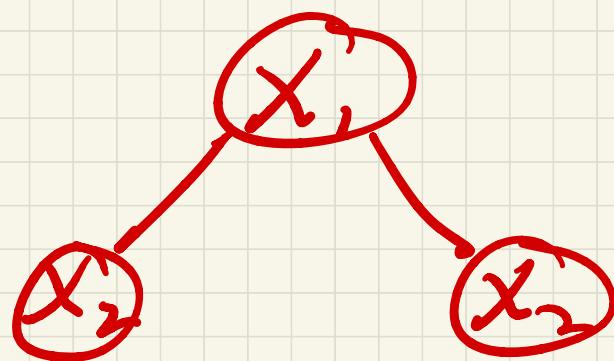
An Application

- Consider "small" parities, i.e. f_S where $|S| \leq \log(n)$
- E.g. parity of first $\log(n)$:

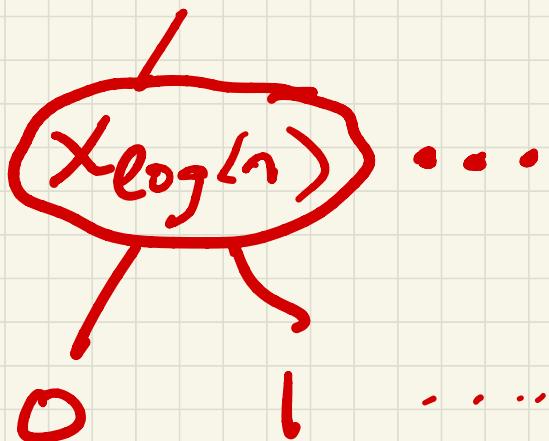
x_1	x_2	\dots	$\times \log(n)$	f_S
0	0	0	...	0
0	0	0	.. 0 1	1
0	0	0	.. 1 0	1
0	0	0	.. 1 1	0
⋮				

→ only $2^{\log(n)} = n$ rows

Encode as decision tree:



⋮



size = n nodes

Encode as DNF:

$$\{ \overline{x_1} \overline{x_2} \dots \overline{x_{\log(n)-1}} x_n \vee \\ \overline{x_1} \overline{x_2} \dots \overline{x_{\log(n)-1}} \overline{x_n} \vee \dots \} \quad \text{all terms where } f_S = 1$$

size $\leq n$ terms

But there are

$$\binom{n}{\log(n)} \approx n^{\log(n)}$$

such small f_S .

So decision trees & DNF
are not learnable
in the SQ model.

For informational, not
computational, reasons.

No assumptions.

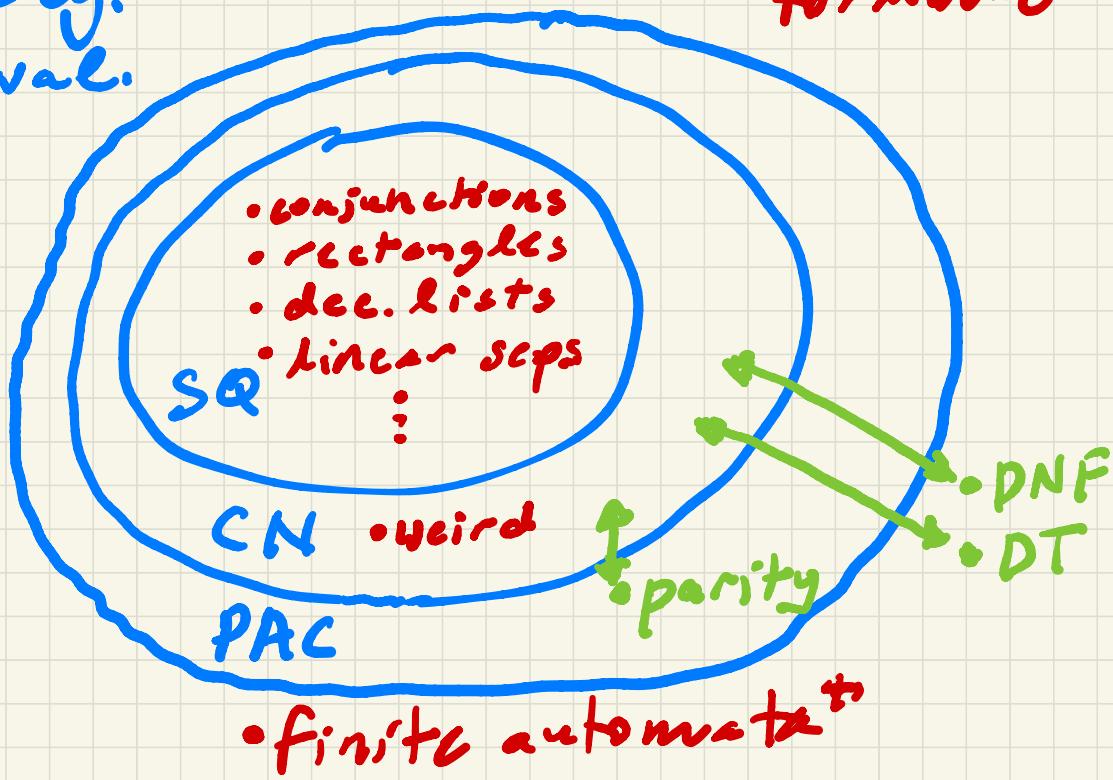
Status in PAC: open.

The PAC Solar System

- status known
- (partially) open

- no eval
- not S*
- boolean formulas*

all
poly.
eval.



$SQ \subsetneq CN \subseteq PAC \subsetneq_{*} \text{poly. eval}$

What about the *?

- There are functions that are easy to compute but hard to invert.
- Example: multiplication.

$$f(p, q) = p \cdot q$$

Over domain p, q prime,

f^{-1} is factorization.

In fact, there are
large families of
such functions.

If primes p, q define

$$f_{p,q}(x) = x^2 \bmod p \cdot q$$

Inverting for unknown
 $p, q \equiv$ factors in $\underline{P \cdot q}$.

Encode as PAC learning
neural nets, finite
automata, boolean
formulas...

Note: Hardness
holds even for weak
PAC learning, where
 $\epsilon = \frac{1}{2} - \frac{1}{\text{poly}(n)}$

Does
weak PAC \Rightarrow PAC?

One final twist...

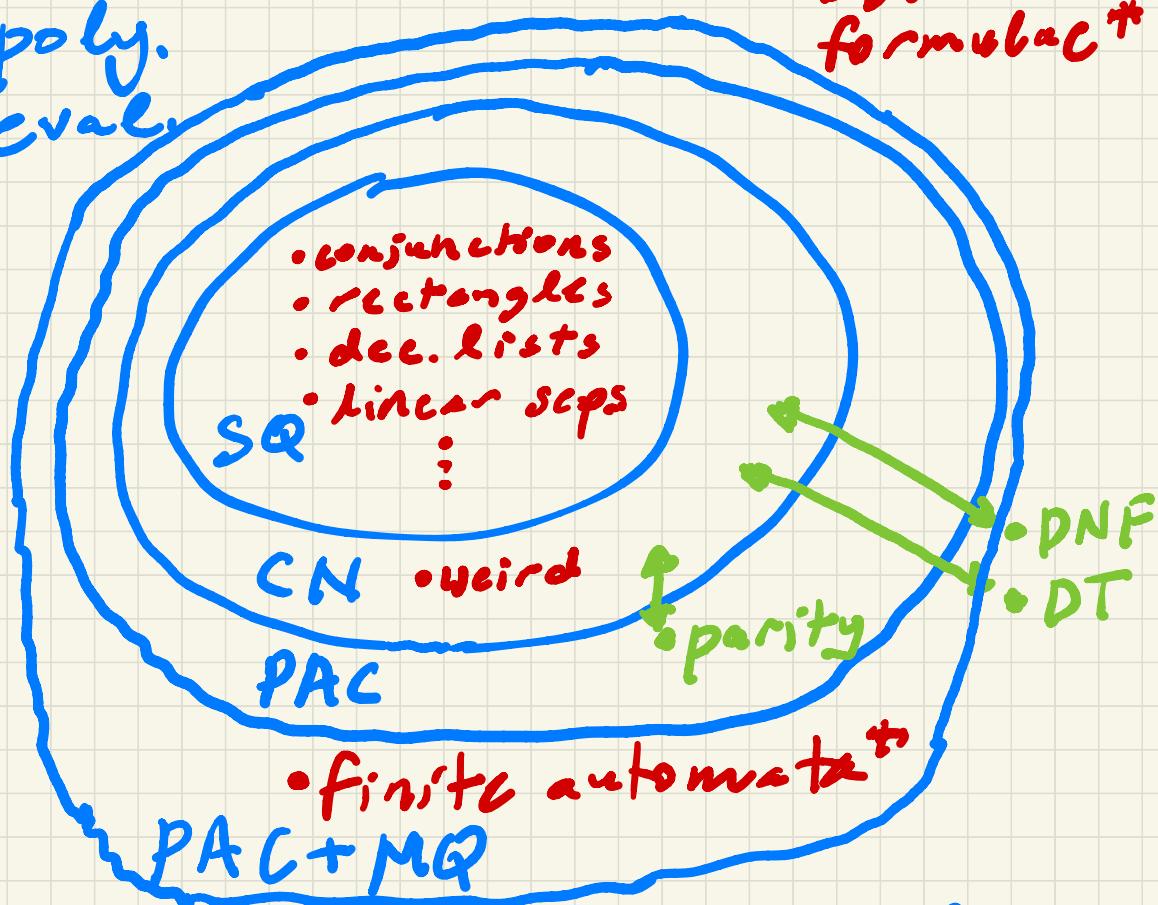
- Let's make PAC more powerful by adding membership queries
- In addition to $\langle x, y \rangle, x \in P, \text{algo}$ can get $c(x)$ for any chosen x .
- Goal remains same

The PAC Solar System

- status Known
- (partially) open

- no eval
nots*
- boolean
formulas*

all
poly.
evals



$$SQ \subsetneq CN \subseteq PAC \subseteq_{*} \text{poly. eval}$$

Next Up:
Boosting
and
Beyond

