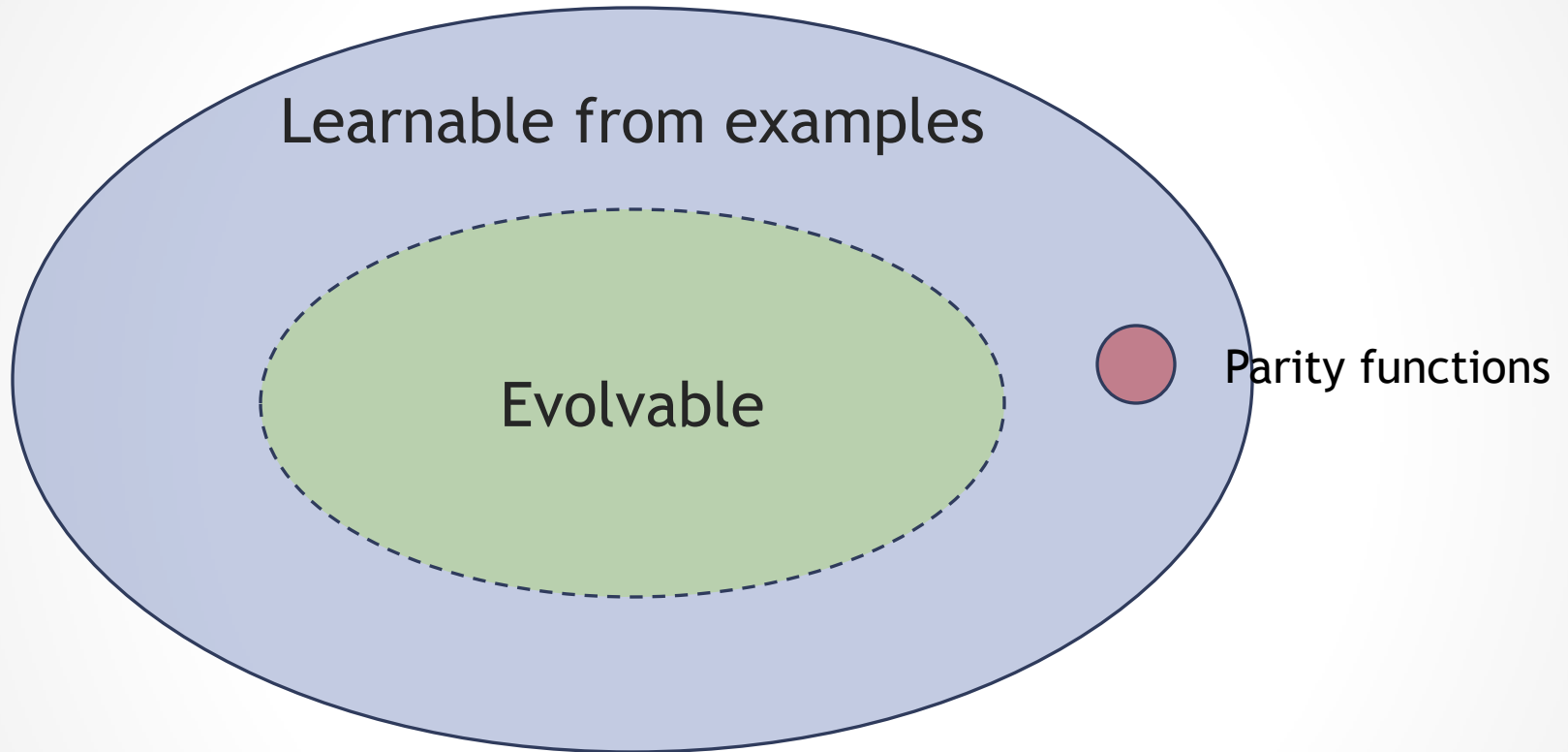


# On the power and the limits of evolvability



Vitaly Feldman  
Almaden Research Center

# Learning from examples vs evolvability



# The core model: PAC [Valiant 84]

- Learner observes random examples:  $(x, f(x))$
- Assumption: unknown Boolean function  $f: X \rightarrow \{-1, 1\}$  labeling the examples comes from a known class of functions  $\mathcal{C}$
- Distribution  $D$  over  $X$  (e.g.  $R^n$  or  $\{-1, 1\}^n$ )

every distribution  $D$

For every  $f \in \mathcal{C}, \epsilon > 0$ , w.h.p. output  $h: X \rightarrow [-1, 1]$   
s.t.  $\mathbb{E}_{x \sim D} [f(x)h(x)] \geq 1 - \epsilon$

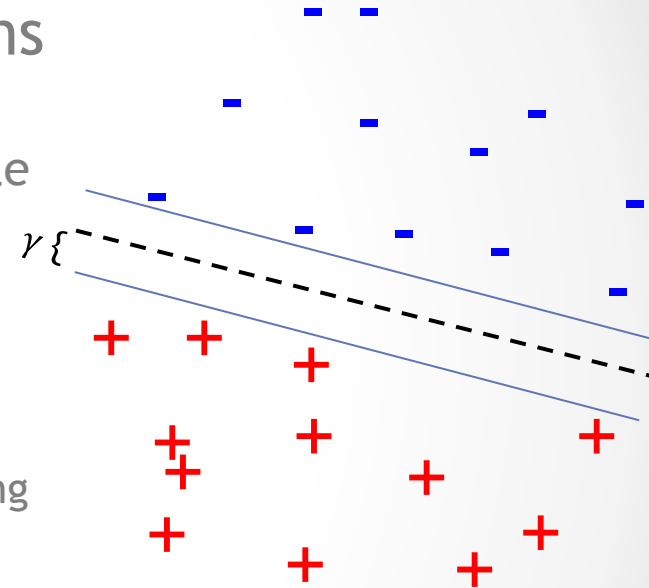
For Boolean  $h$

$\Pr_{x \sim D} [f(x) = h(x)] \geq 1 - \epsilon/2$

Efficient:  $\text{poly}(\frac{1}{\epsilon}, |x|)$  time

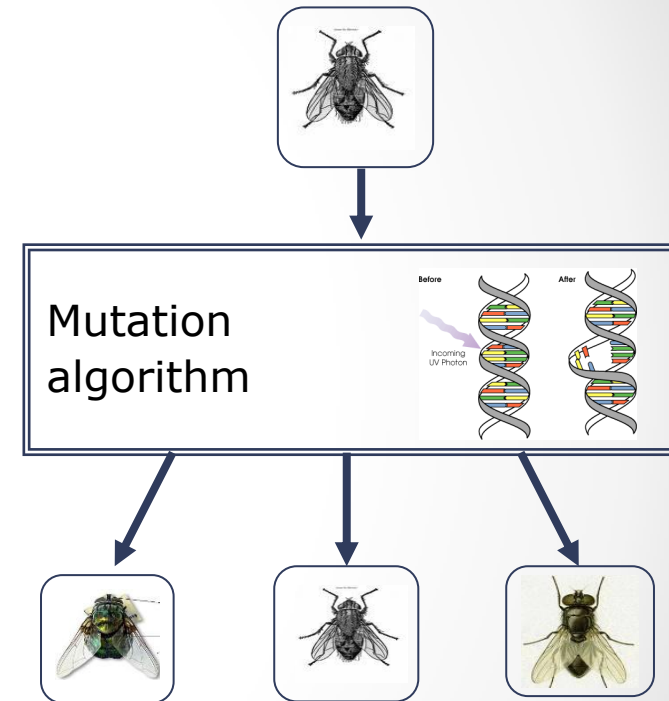
# Classical example

- $C$  halfspaces/linear threshold functions
  - $\text{sign}(\sum_i w_i x_i - \theta)$  for  $w_1, w_2, \dots, w_n, \theta \in \mathbf{R}$
  - equivalent to examples being linearly separable
- Perceptron algorithm
  - [Rosenblatt 57; Block 62; Novikoff 62]
  - Start with LTF  $h^0$  defined by  $w^0 = (0, 0, \dots, 0); \theta^0 = 0$
  - Get a random example  $(x, \ell)$ . If  $h^t(x) = \ell$  do nothing
  - Else let  $h^{t+1}$  be LTF defined by
$$w^{t+1} = w^t + \ell \cdot x; \theta^{t+1} = \theta^t + \ell$$
- Gives PAC learning if  $f$  has significant margin  $\gamma$  on the observed data points



# Evolution algorithm

- $R$  - representation class of functions over  $X$ 
  - E.g. all linear thresholds over  $R^n$
- $M$  - randomized *mutation algorithm* that given  $r \in R$  outputs (a random mutation)  $r' \in R$ 
  - Efficient: poly in  $\frac{1}{\epsilon}, n$
  - E.g. choose a random  $i$  and adjust  $w_i$  by 0, +1 or -1 randomly





# Evolvability

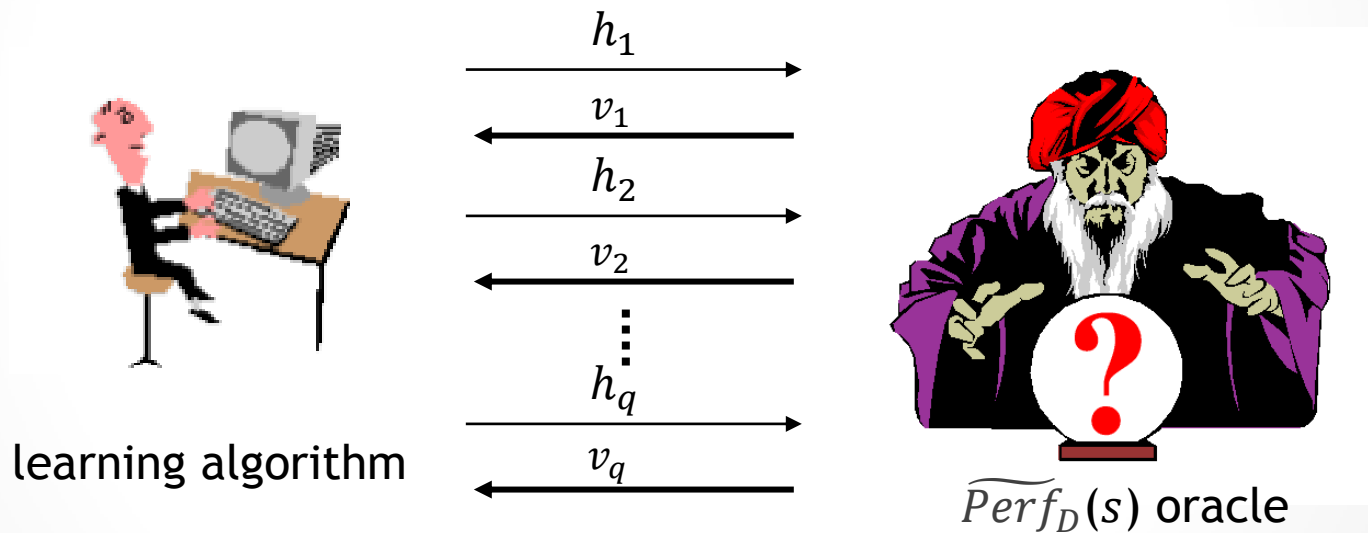
- Class of functions  $C$  is evolvable over  $D$  if exists an evolution algorithm  $(R, M)$  and a polynomial  $g(\cdot, \cdot)$  s.t.

For every  $f \in C, r \in R, \varepsilon > 0$  and a sequence  $r_0 = r, r_1, r_2, \dots$  where  $r_{i+1} \leftarrow \text{Select}(R, M, r_i)$  it holds:  $\text{Perf}_D(f, r_{g(n, \frac{1}{\varepsilon})}) \geq 1 - \varepsilon$  w.h.p.

- Evolvable (*distribution-independently*)
  - Evolvable for all  $D$  by the same mutation algorithm  $(R, M)$

# Limits of evolvability

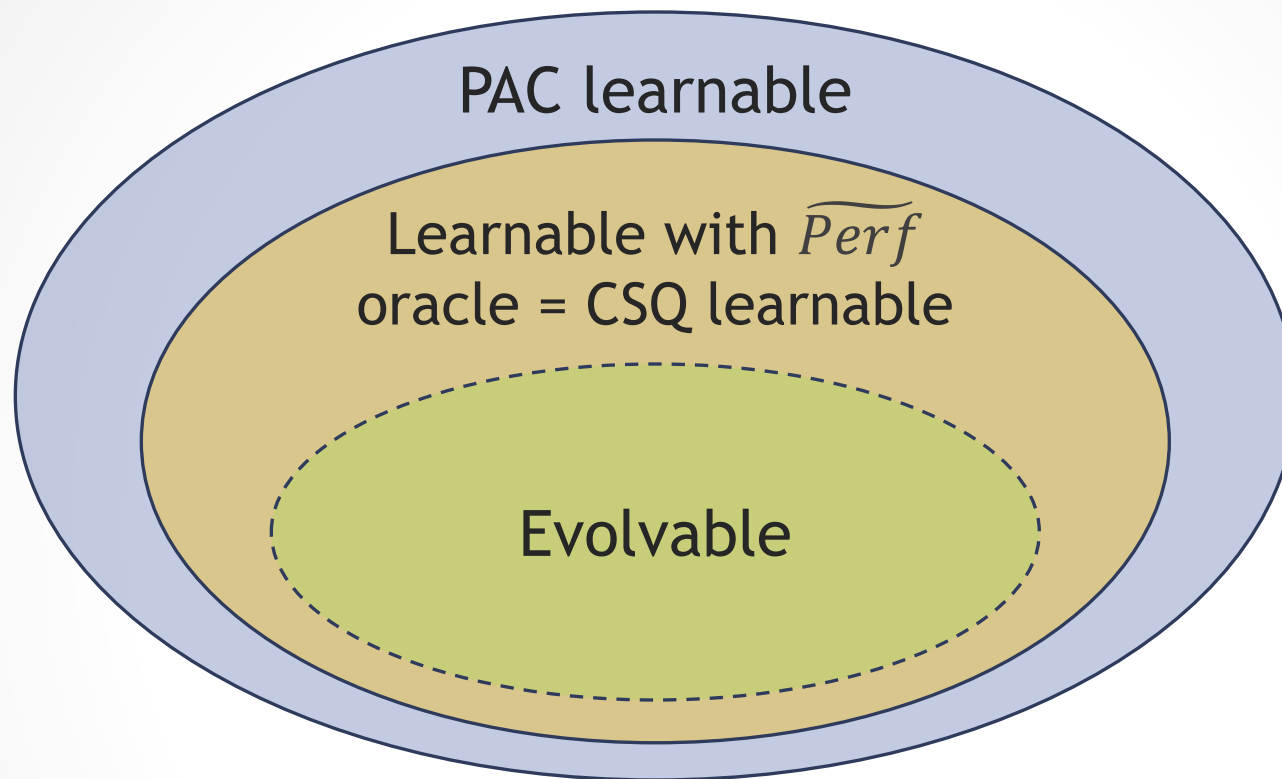
- Feedback is restricted to values of  $\widetilde{Perf}_D(f, r_i)$  for some polynomial number of samples  $s$



$$v_i = \widetilde{Perf}_D(f, h_i) \text{ evaluated on } s \text{ fresh examples}$$



# Evolvable $\subseteq$ CSQ learnable



Correlational Statistical Query:

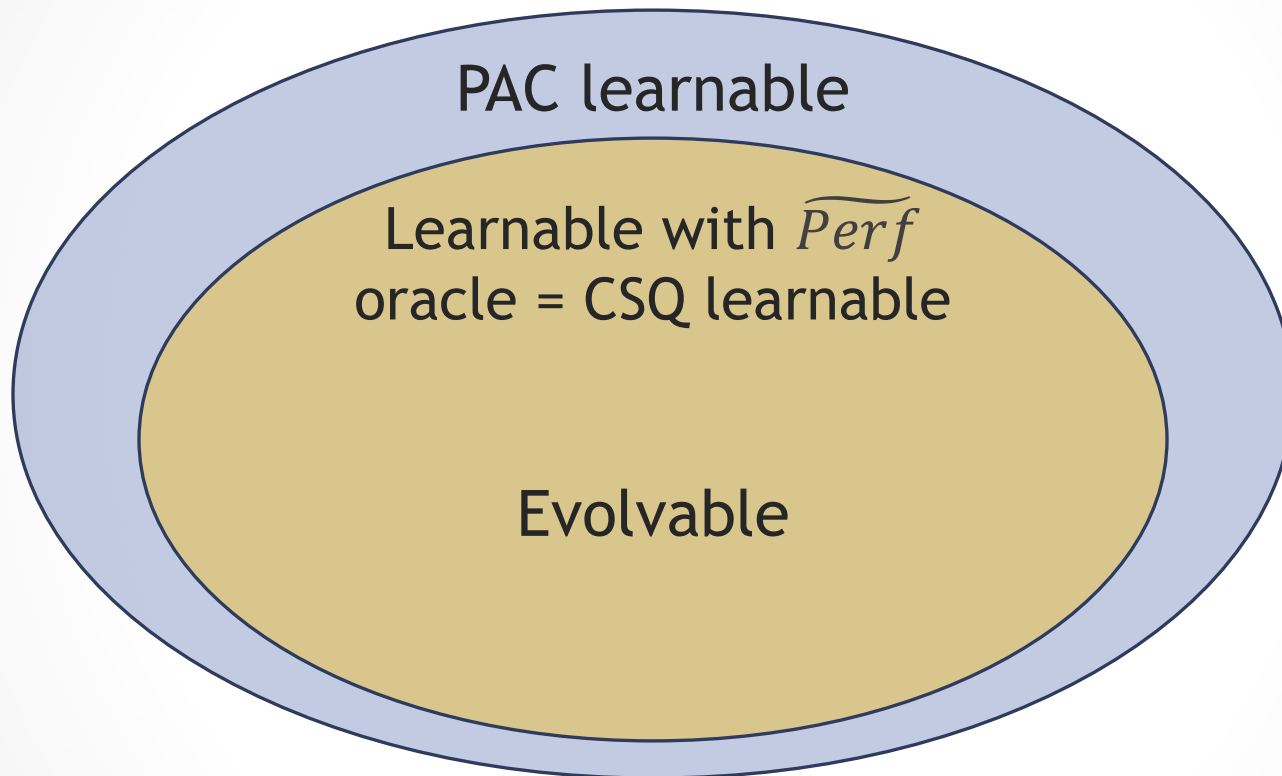
To query  $h$  CSQ oracle responds with any value  $v$

$$|v - \mathbf{E}_D[f(x)h(x)]| \leq \tau \text{ for } \tau \geq \frac{1}{\text{poly}(n, \frac{1}{\epsilon})}$$

Learning by Distances [Ben-David, Itai, Kushilevitz '90]

Restriction of SQ model by Kearns [93]

# CSQ learnable $\subseteq$ Evolvable [F. 08]



# Proof outline

Replace queries for performance values  
with approximate comparisons

For hypothesis  $h: X \rightarrow [-1,1]$ , tolerance  $\tau > 0$  and  
threshold  $t \geq \tau$ , CSQ<sub>></sub> oracle returns:

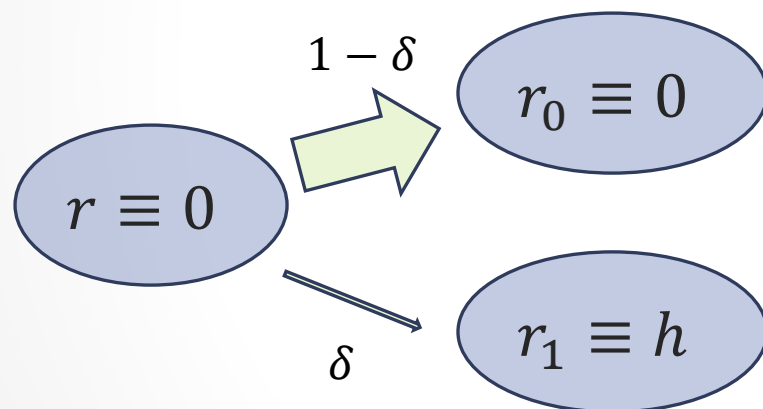
1	if $E_D[f(x)h(x)] \geq t + \tau$
0	if $E_D[f(x)h(x)] \leq t - \tau$
0 or 1	otherwise

Design evolution algorithm with mutations  
that simulate comparison queries

# From comparisons to mutations

For hypothesis  $h: X \rightarrow [-1,1]$ , tolerance  $\tau > 0$  and threshold  $t \geq \tau$ , CSQ<sub>></sub> oracle returns:

1            if  $\mathbf{E}_D[f(x)h(x)] \geq t + \tau$   
0            if  $\mathbf{E}_D[f(x)h(x)] \leq t - \tau$   
0 or 1            otherwise



Beneficial/neutral threshold =  $t$

Mutation pool size  $p = O\left(\frac{\log(\frac{1}{\delta})}{\delta}\right)$

Performance sample size  $s = O\left(\frac{\log(\frac{1}{\delta})}{\tau^2}\right)$

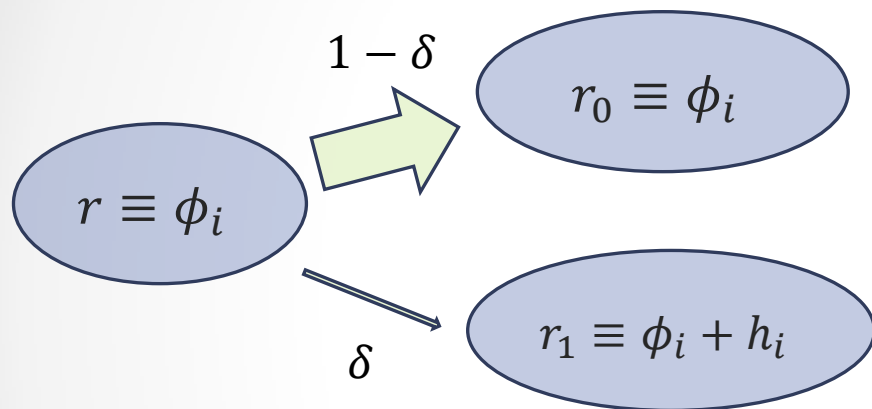
With probability at least  $1 - \delta$ ,  $\text{Select}(r) = r_b$   
where  $b$  is valid response to comparison query

# Simulating $CSQ_{\geq}$ algorithm

Need to answer  $q$  queries

$\phi_i$  is the function obtained after answering  $i$  queries

Need to answer query  $h_i$  with threshold  $t_i$



Beneficial/neutral threshold =  $t_i$

Mutation sample size  $p = O\left(\frac{\log(\frac{1}{\delta})}{\delta}\right)$

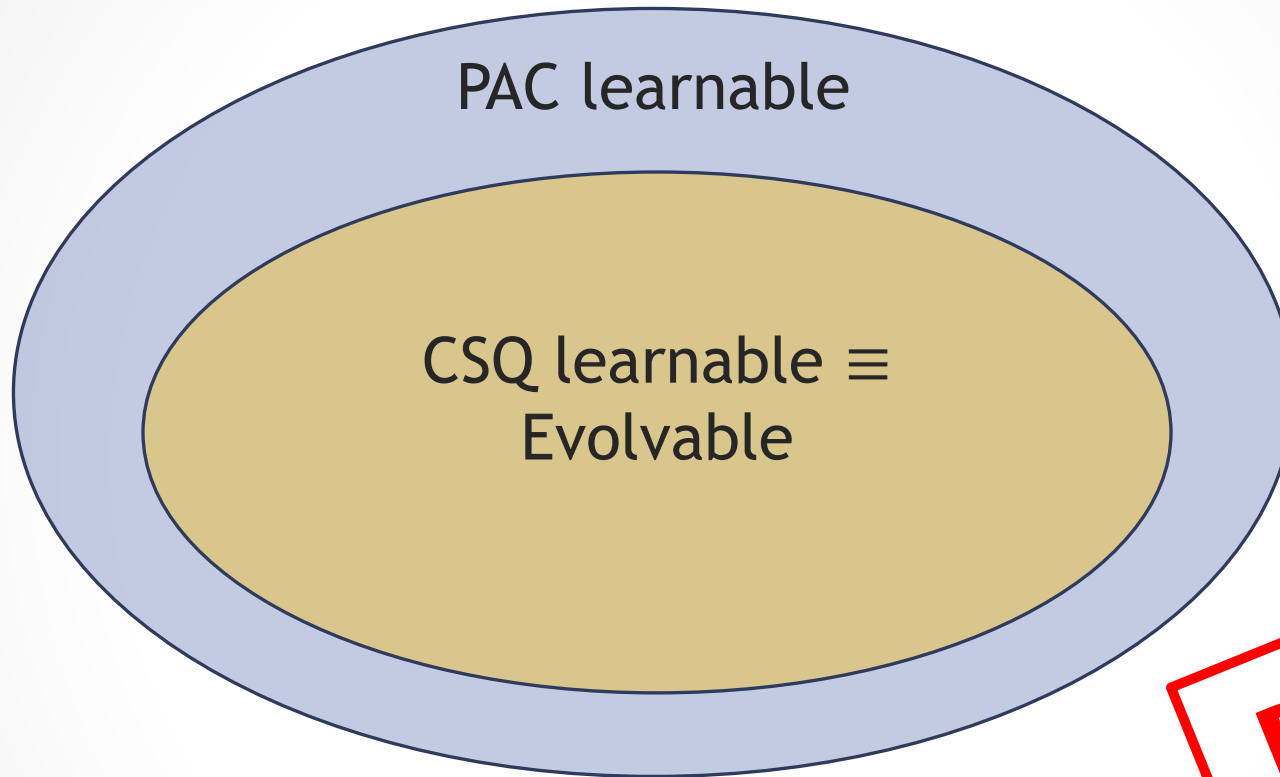
Performance sample size  $s = O\left(\frac{\log(\frac{1}{\delta})}{\tau^2}\right)$

Leads to representations with values in  $[-q, q]$ !

Rescale by  $1/q$  to get functions with values in  $[-1, 1]$

Given answers to queries can compute  $h$  such that  
 $Perf_D(f, h) \geq 1 - \epsilon$  and mutate into it

# CSQ learnable $\equiv$ Evolvable [F.08; F. 09]



**ROBUST**

E.g. optimizing selection; recombination [Kanade 11]; changing thresholds; number of mutations

# How powerful is CSQ learning?

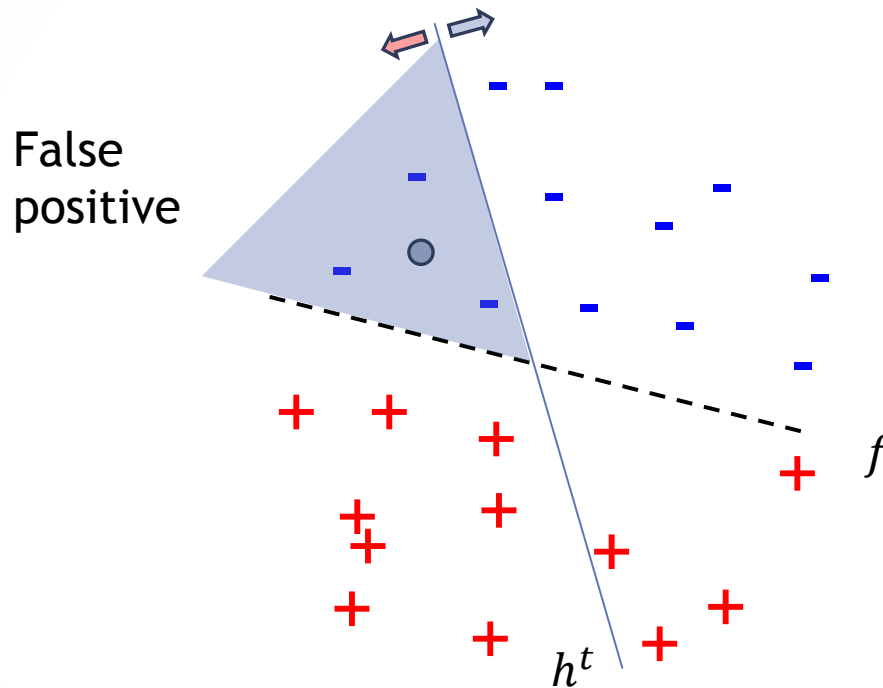
SQ learning [Kearns 93] : learner submits  $\psi(x, \ell)$

SQ oracle returns  $v$  such that  $|v - E_D[\psi(x, f(x))]| \leq \tau$

- Many known algorithms are essentially described in SQ or can be easily translated to SQ
  - Boolean conjunctions, decision lists, simple geometric concepts,  $AC^0$  [Kearns 93]
- Several other were found with more effort
  - Halfspaces with large margin [Bylander 94]
  - General LTFs [BlumFrie.Kann.Vemp. 96; Duna.Vemp. 04]
- General ML techniques
  - Nearest neighbor
  - Gradient descent
  - SVM
  - Boosting

# Perceptron in SQ

- Recall the Perceptron algorithm:
  - Add false negatives, subtract false positive examples




- Use SQs to find the centroid of false positives
- $E_D[x_1 \cdot I(f(x) = -1) \cdot I(h^t(x) = 1)]$  gives the first coordinate of the centroid
- Use the centroid for the Perceptron update
-



# If $D$ is fixed then $SQ \equiv CSQ$

- Decompose SQ into CSQ and a constant

$$\begin{aligned}\mathbf{E}_D[\psi(x, f(x))] &= \mathbf{E}_D \left[ \psi(x, -1) \frac{1 - f(x)}{2} + \psi(x, 1) \frac{1 + f(x)}{2} \right] \\ &= \mathbf{E}_D \left[ \frac{\psi(x, 1) - \psi(x, -1)}{2} f(x) \right] + \mathbf{E}_D \left[ \frac{\psi(x, 1) + \psi(x, -1)}{2} \right]\end{aligned}$$

 CSQ

- Corollary: linear threshold functions are evolvable for any fixed distribution  $D$

# Distribution-independent CSQ

- Single points are learnable [F. 09]
- Characterization of weak-learning [F. 08]

Better than random guessing:  $\mathbb{E}_{x \sim D} [f(x)h(x)] \geq \frac{1}{\text{poly}(n, \frac{1}{\epsilon})}$

$C$  is weakly CSQ learnable if and only if all functions in  $C$  can be represented as linear threshold functions with “significant” margin over a poly-size set of Boolean features

- General linear thresholds are not weakly CSQ learnable [Goldmann, Hastad, Razborov 95] (but are SQ learnable)
- Conjunctions are not CSQ learnable [F. 11]

# Further directions

- Characterize (strong) evolvability (CSQ learning)
  - Strengthen the lower bound for conjunctions
- Are thresholds on a line evolvable distribution independently
- $Perf_D(f, r) = -E_D \left[ (f(x) - r(x))^2 \right]$  then all of SQ is evolvable [F. 09]

