## Two Types of Error

① True Error (aka. expected risk) (aka. Generalization Error)

$$R(h) = P_{x \sim p^*(x)}\left(c^*(x) \neq h(x)\right) \quad \longleftarrow \quad \text{always unknown.}$$

② Train Error (aka. empirical risk)

$$\hat{R}(h) = P_{x \sim S}\left(c^*(x) \neq h(x)\right) \qquad S = \{x^{(1)}, \ldots, x^{(N)}\}$$

$$= \frac{1}{N} \sum_{i=1}^{N} \mathbb{1}\left(c^*(x^{(i)}) \neq h(x^{(i)})\right)$$

$$= \frac{1}{N} \sum_{i=1}^{N} \mathbb{1}\left(y^{(i)} \neq h(x^{(i)})\right)$$

known, computable

## Three Decisions Fns. of Interest

✓ ① True function (oracle), $c^*$

$$y^{(i)} = c^*(x^{(i)}) \qquad \forall i$$

✓ ② Expected Risk Minimizer (lowest true error)

$$h^* = \operatorname*{argmin}_{h \in H} R(h)$$

hypothesis class.

③ Empirical Risk Minimizer (lowest training error)

$$\hat{h} = \operatorname*{argmin}_{h \in H} \hat{R}(h) = \operatorname*{argmin}_{h \in H} \frac{1}{N} \sum_{i} \mathbb{1}\left(y^{(i)} \neq h(x^{(i)})\right)$$

Q: Which of these are unknown?

# PAC Learning

A: Yes!

PAC stands for Probably Approximately Correct

PAC learner yields hypothesis $h$, which is approximately correct $R(h) \approx 0$
with high probability $Pr(R(h) \approx 0) \approx 1$

## Def: PAC Criterion

$$Pr(\forall h, |R(h) - \hat{R}(h)| \leq \epsilon) \geq 1 - \delta$$

↑ small          ↑ small

Q: What is random?
A: $\hat{R}$ is measured on a __random__ sample $S$ of traing examples

## Def: Sample complexity is the minimum value $N$ of traing examples s.t. the PAC criterion holds for a given $\epsilon, \delta$
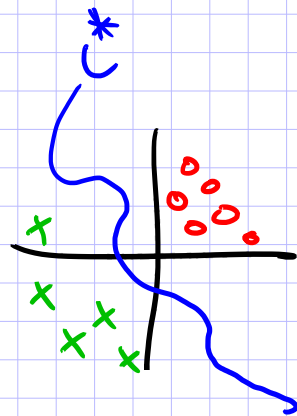
## Def: Hypothesis space $H$ is __PAC learnable__ if sample complexity $N$ is a polynomial function of $\frac{1}{\epsilon}$ and $\frac{1}{\delta}$ for some learning algorithm.

# 4 Bounds

## Two cases for $c^*$

A) Realizable case: $c^* \in H$

B) Agnostic case: $c^*$ is not necessarily in $H$

## Two Cases for $H$:

A) Finite $|H| < +\infty$

B) Infinite $|H| = +\infty$

**Thm 1:** Sample Complexity (Realizable, Finite $|H|$)

$$N \geq \frac{1}{\epsilon} \left[ \ln(|H|) + \ln\left(\frac{1}{\delta}\right) \right] \quad \text{labeled examples}$$

are sufficient to ensure that with prob. $\geq (1-\delta)$

all $h \in H$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$

* Bound is inversely linear in $\epsilon$
  (Halving the error requires double the examples)
* Bound is logarithmic in $|H|$
  (Doubling the hyp. space size requires only logarithmically more training exs. $N$)

**Ex: Conjunctions**

$H$ = class of conjunctions over $\vec{x} \in \{0,1\}^M$

e.g. $h(\vec{x}) = x_1 \bar{x}_3 x_4 = x_1(1-x_3)x_4$

$h(\vec{x}) = x_1 \bar{x}_2 x_4 \bar{x}_6$

Q: Suppose $M = 10$, $\epsilon = 0.1$, $\delta = 0.01$, what is the samp. comp.?

~~$|H| = 2^M$~~     $|H| = 3^M$

Then $N \geq \frac{1}{0.1} \left[ \ln(3^M) + \ln\left(\frac{1}{0.01}\right) \right]$

$= \frac{1}{0.1} \left[ M \ln(3) + \ln\left(\frac{1}{0.01}\right) \right]$

$= \frac{1}{0.1} \left[ 10 \ln(3) + \ln\left(\frac{1}{0.01}\right) \right] \approx 156$ examples

**Corollary 1:**

With prob. $(1-\delta)$ for all $h \in H$ s.t. $\hat{R}(h) = 0$
we have that: $R(h) \leq \frac{1}{N} \left[ \ln(|H|) + \ln\left(\frac{1}{\delta}\right) \right]$

(probably)  (approximately correct) assuming a consistent learner

Thm 2: Sample Complexity (Agnostic, Finite $|H|$)

$$N \geq \frac{1}{2\epsilon^2}\left[\ln(|H|) + \ln\left(\frac{2}{\delta}\right)\right] \text{ labeled examples}$$

are sufficient so that with prob. $\geq (1-\delta)$
all $h \in H$ have $|R(h) - \hat{R}(h)| \leq \epsilon$

★ Bound is inversely quadratic in $\epsilon$ — "quadratically worse"
(halving the error requires $4x$ examples)
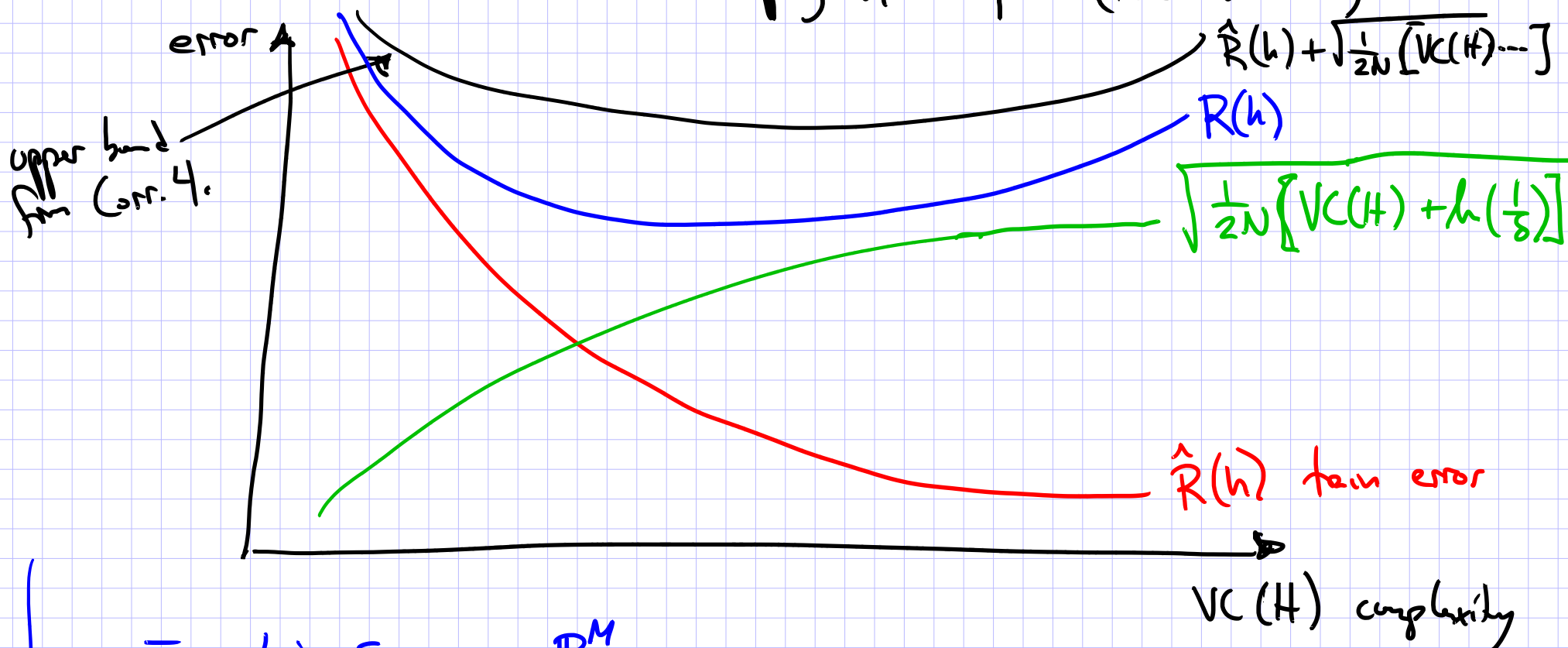
Corollary 4:
with prob. $(1-\delta)$ for all $h \in H$

$$R(h) \leq \hat{R}(h) + \sqrt{\frac{1}{2N}\left[VC(H) + \ln\left(\frac{1}{\delta}\right)\right]}$$

Structural Risk Minimization

Q: Should Corr. 4 inform how do we model selectn?
A: Yes!

Model Selection: tradeoff between low train error
and keeping $H$ simple (low VCDim)



error

upper bound
from Corr. 4.

$\hat{R}(h) + \sqrt{\frac{1}{2N}[VC(H) \cdots]}$

$R(h)$

$\sqrt{\frac{1}{2N}[VC(H) + \ln(\frac{1}{\delta})]}$

$\hat{R}(h)$ train error

$VC(H)$ complexity

Ex: Lin Sep in $\mathbb{R}^M$
$VC(H) = M+1$
How to tradeoff?
Train w/ L1 Regularizer
↳ prefers 0 in $\vec{\theta}$