# Machine Learning Theory

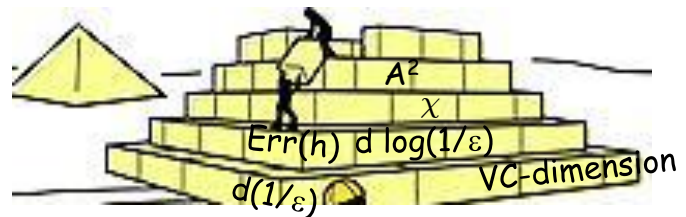## Maria-Florina (Nina) Balcan

February 9th, 2015

M2:
(7.1.  7.2,  7.31,  7.41 − 7.4 3. )

林轩田. ← M2F

# Goals of Machine Learning Theory
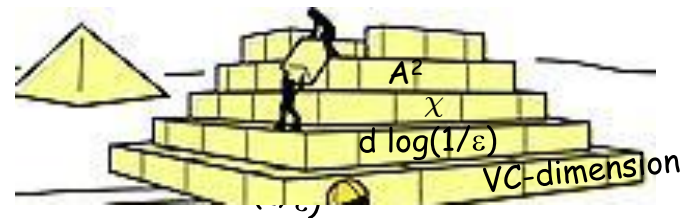
Develop & analyze models to understand:

- what kinds of tasks we can hope to learn, and from what kind of data; what are key resources involved (e.g., data, running time)

- prove guarantees for practically successful algs (when will they succeed, how long will they take?)

- develop new algs that provably meet desired criteria  (within new learning paradigms)

Interesting  tools & connections to other areas:

- Algorithms, Probability & Statistics, Optimization, Complexity Theory, Information Theory, Game Theory.

Very vibrant field:

- Conference on Learning Theory

- NIPS, ICML

(NeurIPS)

# Today's focus: Sample Complexity for Supervised Classification (Function Approximation)

- **Statistical Learning Theory** (Vapnik)  *SLT*

- **PAC** (Valiant)

  *Probably   Approximately   Correct*

- Recommended reading: Mitchell: Ch. 7
  - Suggested exercises: 7.1, 7.2, 7.7

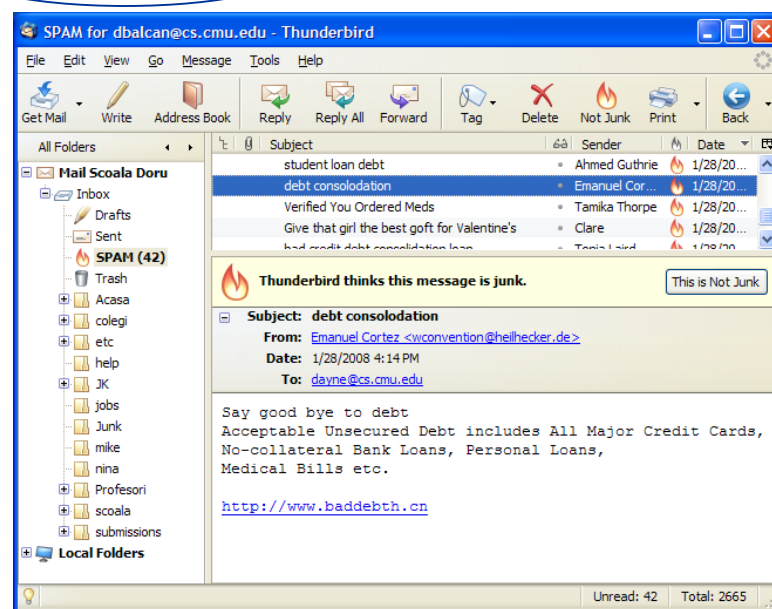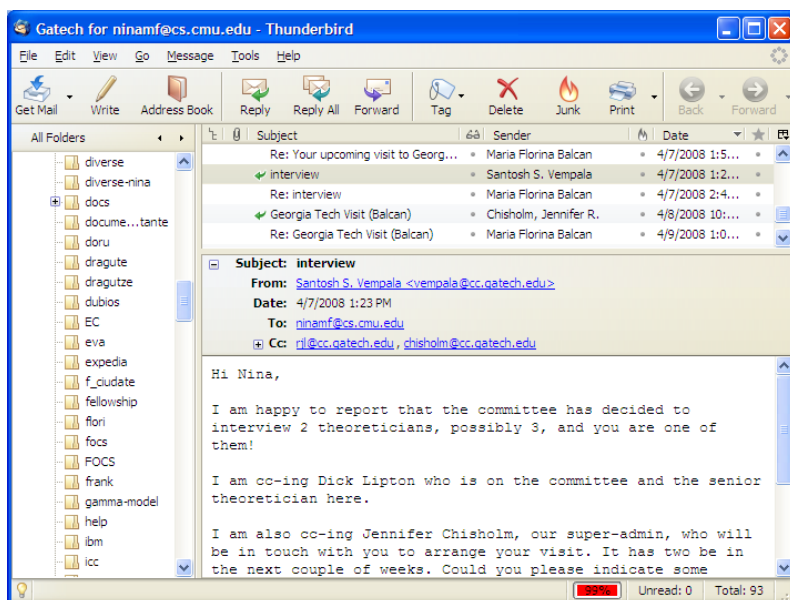- Additional resources: my learning theory course!

# Supervised Classification

Decide which emails are spam and which are important.

Supervised classification

Not spam

spam



**Goal: use emails seen so far to produce good prediction rule for future data.**

$h(\hat{x}) \to \hat{y}$

$ML \begin{cases} Data: \{(x_j, y_j)\}_{j=1}^{M} \\ Alg: \ h: X \to y \end{cases}$

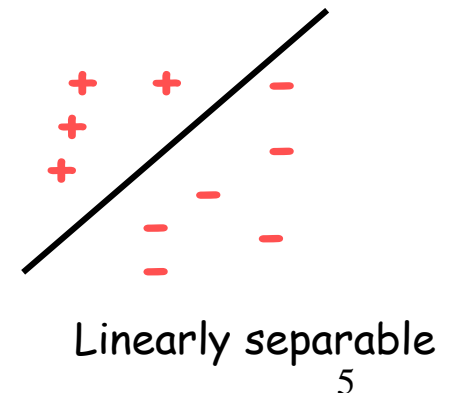Goal: future incoming emails

4

# Example: Supervised Classification

Represent each message by features. (e.g., keywords, spelling, etc.)

| | "money" | "pills" | "Mr." | bad spelling | known-sender | spam? |
|---|---|---|---|---|---|---|
| | Y | N | Y | Y | N | Y |
| | N | N | N | Y | Y | N |
| | N | Y | N | N | N | Y |
| example | Y | N | N | N | Y | N | label |
| | N | N | Y | N | Y | N |
| | Y | N | N | Y | N | Y |
| | N | N | Y | N | N | N |

Reasonable RULES:

Predict SPAM if unknown AND (money OR pills)

Predict SPAM if 2money + 3pills –5 known > 0



Linearly separable

# Two Core Aspects of Machine Learning

**Algorithm Design. How to optimize?**     Computation

Automatically generate rules that do well on observed data.

- E.g.: logistic regression, SVM, Adaboost, etc.
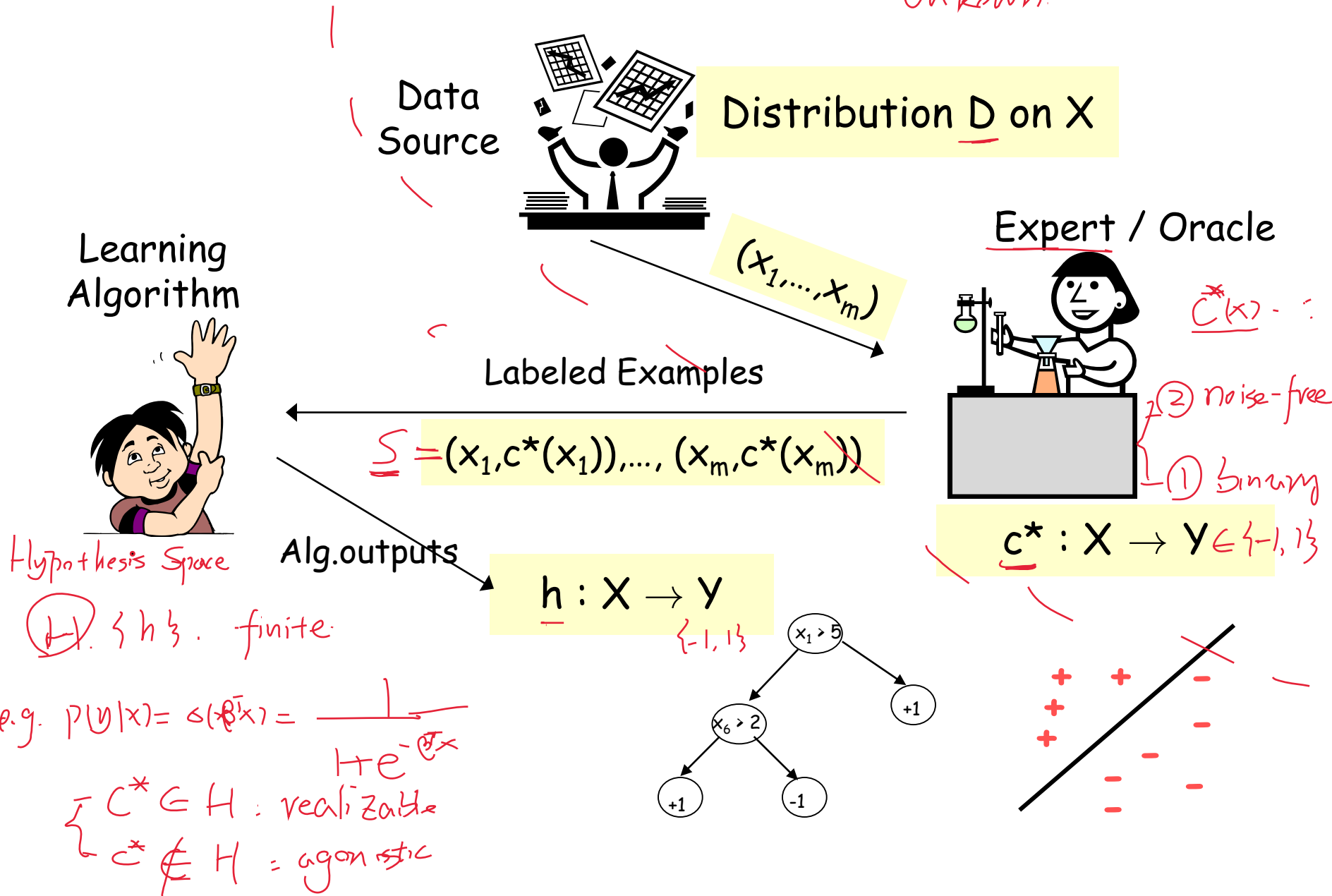
**Confidence Bounds, Generalization**     (Labeled) Data

Confidence for rule effectiveness on future data.

- Very well understood: Occam's bound, VC theory, etc.
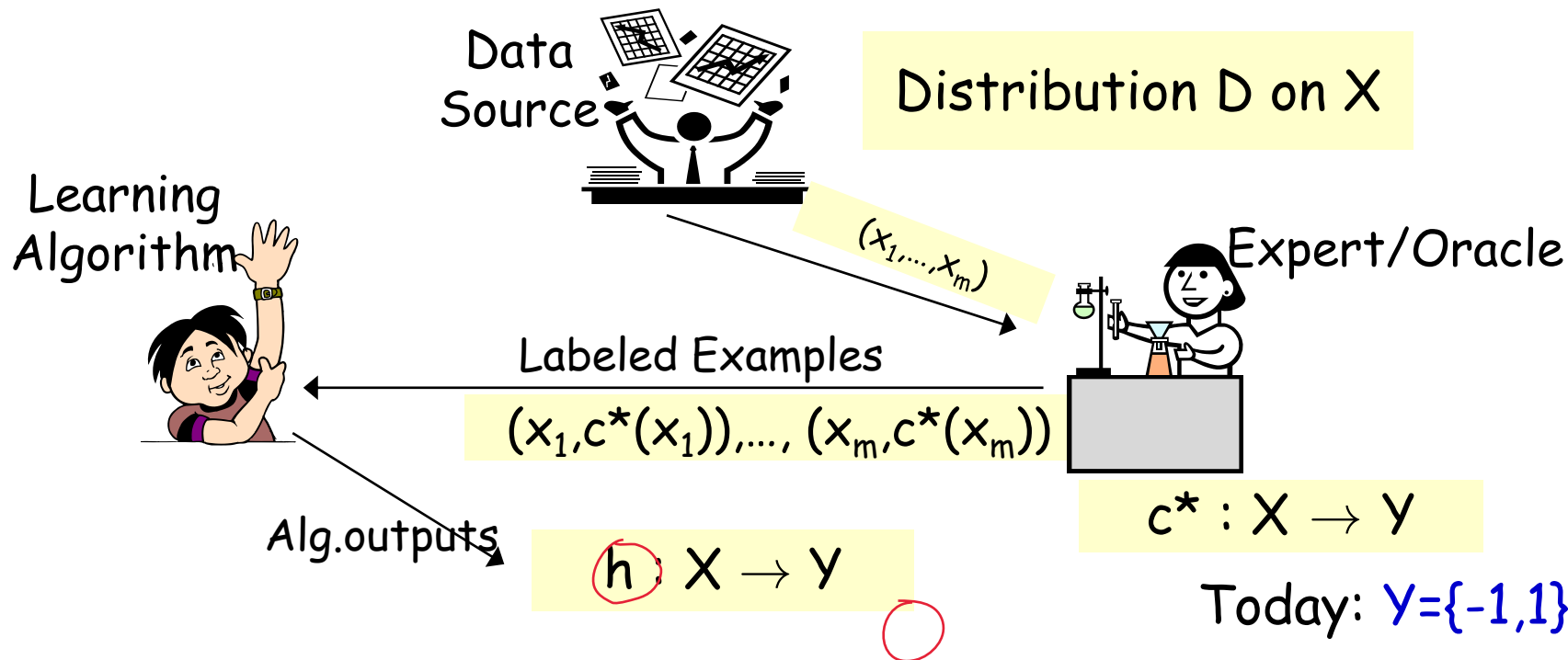- Note: to talk about these we need a precise model.

# PAC/SLT models for Supervised Learning

Unkown.

## Data Source

Distribution D on X

$(x_1,...,x_m)$

## Expert / Oracle

$C^*(x) = $

② noise-free

① binary

Labeled Examples

$S = (x_1, c^*(x_1)),..., (x_m, c^*(x_m))$

$c^* : X \to Y \in \{-1, 1\}$

## Learning Algorithm

Hypothesis Space

$(H) \cdot \{h\} \cdot$ finite.

e.g. $P(y|x) = \sigma(\beta^T x) = \dfrac{1}{1+e^{-\beta^T x}}$

$\begin{cases} c^* \in H : \text{realizable} \\ c^* \notin H : \text{agonistic} \end{cases}$

Alg. outputs

$h : X \to Y$

$\{-1, 1\}$

$x_1 > 5$

$x_6 > 2$

+1

+1

-1

+ + −
+ + −
+ − −
− −

$error_D(h) \leq error_S(h) + \epsilon$

# PAC/SLT models for Supervised Learning

Data Source

Distribution D on X

Learning Algorithm

$(x_1,...,x_m)$

Expert/Oracle

Labeled Examples

$(x_1,c^*(x_1)),..., (x_m,c^*(x_m))$

$c^* : X \rightarrow Y$

Alg.outputs

$h : X \rightarrow Y$

Today: Y={-1,1}

- Algo sees training sample S: $(x_1,c^*(x_1)),..., (x_m,c^*(x_m))$, $x_i$ independently and identically distributed (i.i.d.) from D; labeled by $c^*$   $h \in H$

- Does optimization over S, finds hypothesis h (e.g., a decision tree) $\min_{h \in H} error_S(h)$

$error_D(h) \sim error_S(h) \sim 0$

- Goal:  h has small error over D.

# PAC/SLT models for Supervised Learning

- $X$ – feature or instance space; distribution $D$ over $X$

    e.g., $X = R^d$ or $X = \{0,1\}^d$

- Algo sees training sample S: $(x_1, c^*(x_1)), \ldots, (x_m, c^*(x_m))$, $x_i$ i.i.d. from $D$

    - labeled examples - assumed to be drawn i.i.d. from some distr. $D$ over $X$ and labeled by some target concept $c^*$
    - labels $\in \{-1,1\}$ - binary classification

- Algo does optimization over S, find hypothesis $h$.

$$error_D(h) \neq 0$$
$$error_S(h) = 0$$

- Goal: h has small error over $D$.    *pointwise*

$$err_D(h) = \Pr_{x \sim D}(h(x) \neq c^*(x))$$

$$= \mathbb{E}_D\left[h(x) \neq c^*(x)\right]$$



Instance space $X$

## Need a bias: no free lunch.

# Function Approximation: The Big Picture

$$(C^* \in H)$$

$$H: \quad h \quad X \to \{-1, 1\}$$

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

$$S$$

$|H|$ finite



Hypothesis Space

$$|X| = 2^n$$

$$c^*$$

Feature Space

$$h(S) = \{ (h(x_1), \ldots, h(x_m)) \}$$

dichotomy

$$|H| = 2^{2^n} = 2^{|X|}$$

Additional assumption:

(Complexity)

Q: How many labeled examples are needed to determine which
  of $2^{2^n}$ hypos is $c^*$?

A: $2^n$ labeled examples

$$2^n - 1 : \begin{cases} +1 \to h_i \\ -1 \to h_j \end{cases}$$

$$2^n - 2 : \begin{cases} +1 +1 \to h_a \\ +1 -1 \to h_b \\ -1 +1 \to h_c \\ -1 -1 \to h_d \end{cases}$$

# PAC/SLT models for Supervised Learning
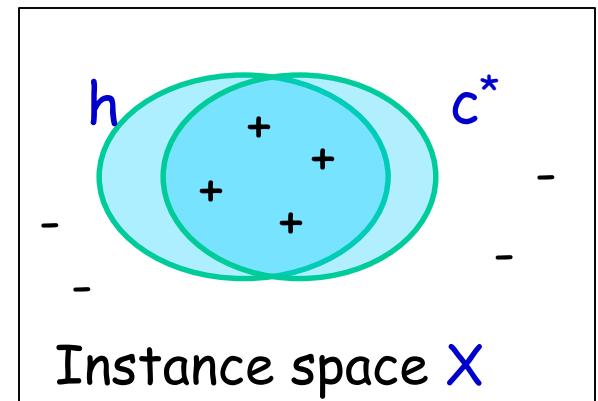
- X – feature or instance space; distribution D over X

    e.g., $X = R^d$ or $X = \{0,1\}^d$

- Algo sees training sample S: $(x_1, c^*(x_1)), \ldots, (x_m, c^*(x_m))$, $x_i$ i.i.d. from D

    - labeled examples - assumed to be drawn i.i.d. from some distr. D over X and labeled by some target concept $c^*$
    - labels $\in \{-1, 1\}$ - binary classification

- Algo does optimization over S, find hypothesis $h$.

- Goal: h has small error over D.

$$err_D(h) = \Pr_{x \sim D}(h(x) \neq c^*(x))$$

Bias: Fix hypotheses space H .

(whose complexity is not too large).

Realizable: $c^* \in H$.

Agnostic: $c^*$ "close to" H. $(c^* \notin H)$



h          $c^*$

+
+
+
-
+
-
-

Instance space X

$err_D(h) \leq err_S(h) + \varepsilon$

# PAC/SLT models for Supervised Learning

- Algo sees training sample $S$: $(x_1, c^*(x_1)), ..., (x_m, c^*(x_m))$, $x_i$ i.i.d. from $D$

- Does optimization over $S$, find hypothesis $h \in H$.

- Goal: $h$ has small error over $D$.   $E_D\left[ h(x) \neq c^*(x) \right]$

  1. True error: $err_D(h) = \Pr_{x \sim D}(h(x) \neq c^*(x))$

  2. Expected error

  3. Generalization error

  How often $h(x) \neq c^*(x)$ over future instances drawn at random from D

- But, can only measure:

  1. Training error: $err_S(h) = \frac{1}{m} \sum_i I\left( h(x_i) \neq c^*(x_i) \right) = \begin{cases} 1, & h(x) \neq c^*(x) \\ 0, & \text{otherwise.} \end{cases}$

  2. Empirical error:

  • Empirical Error Minimization.

  $\min_{h \in H} err_S(h)$

  How often $h(x) \neq c^*(x)$ over training instances

**Sample complexity: bound $err_D(h)$ in terms of $err_S(h)$**

$err_D(h) \leq err_S(h) + \underbrace{\sum}_{(\text{computable})}$

# Sample Complexity for Supervised Learning

- Consistent Learner
    - outputs hypothesis $h$ that perfectly fits the training data $S$,
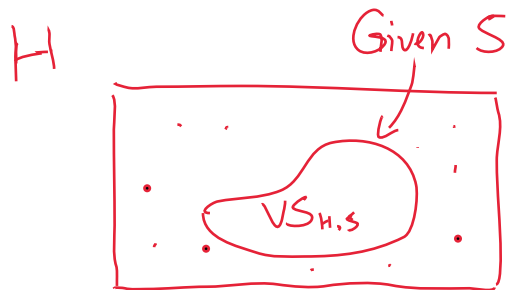
$$h(x) = c^*(x), \qquad \forall x \in S.$$

- Version Space (VS)
    - set of all hypotheses $h \in H$ that correctly classify the training data $S$,

$$VS_{H,S} = \{h \in H | \forall x \in S, h(x) = c^*(x)\}.$$

$err_S(h) = 0$

$H$

Given $S$

$VS_{H,S}$

$\forall$ : /for all

$\exists$ : /exist

# Sample Complexity for Supervised Learning

$(\overset{*}{C} \subseteq H)$

**Definition:** Consider a hypothesis space $H$, target concept $c$, instance distribution $\mathcal{D}$, and set of training examples $S$ of $c$. The version space $VS_{H,D}$ is said to be $\epsilon$-**exhausted** with respect to $c$ and $\mathcal{D}$, if every hypothesis $h$ in $VS_{H,D}$ has error less than $\epsilon$ with respect to $c$ and $\mathcal{D}$.

$(0 < \epsilon < \frac{1}{2})$

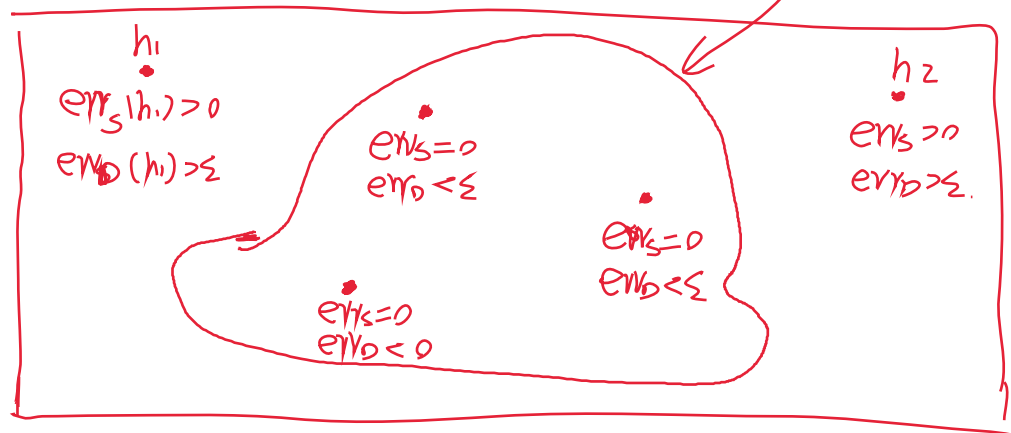$$(\forall h \in VS_{H,D})\,(error_{\mathcal{D}}(h) < \epsilon)$$

$$VS_{H.S} = \{ h \mid h \in H,\ err_S(h) = 0 \}$$

$$\epsilon\text{-exhausted } VS_{H.S} = \{ h \mid h \in H,\ err_S(h) = 0,\ err_D(h) < \epsilon \}$$

$$= \{ h \mid h \in VS_{H.S},\ err_D(h) < \epsilon \}$$

$VS_{H.S} \leftarrow \epsilon\text{-exhausted}$

# Sample Complexity for Supervised Learning

① $c^* \in H$

② $H$ finite

**Theorem 7.1.** $\epsilon$-**exhausting the version space.** If the hypothesis space $H$ is finite, and $D$ is a sequence of $m \geq 1$ independent randomly drawn examples of some target concept $c$, then for any $0 \leq \epsilon \leq 1$, the probability that the version space $VS_{H,D}$ is not $\epsilon$-exhausted (with respect to $c$) is less than or equal to

bad

$|H|e^{-\epsilon m}$

$$P\left(\exists h \in VS_{H,S}, \; err_D(h) \geq \epsilon\right) \leq |H|e^{-\epsilon m}$$

$$1 - P\left(\forall h \in VS_{H,S}, \; err_D(h) < \epsilon\right) \leq \boxed{|H|e^{-\epsilon m} \leq \delta} \qquad (0 < \delta < \tfrac{1}{2})$$

$$P\left(\forall h \in VS_{H,S}, \; err_D(h) < \epsilon\right) \geq 1 - \delta \qquad \text{with high prob.}$$

good ( $\epsilon$-exhausted )

( w hp )

|

$$\ln |H| - \epsilon m \ln e \leq \ln \delta$$

Sample Complexity

$$\Rightarrow \quad m \geq \frac{1}{\epsilon}\left(\ln |H| + \ln \frac{1}{\delta}\right)$$

# Sample Complexity for Supervised Learning

$$P(\exists h \in VS_{H,S}, \; err_D(h) \geq \varepsilon) \leq |H| e^{-\varepsilon m}$$

Proof: if $\exists h \in H$, $err_S(h) = 0$, $err_D(h) \geq \varepsilon$, then VS is not $\varepsilon$-exhausted.

$$(h \in VS_{H,S})$$

Calculate the prob.

$$P(A \cup B) \leq P(A) + P(B)$$

$$err_D(h) \geq \varepsilon \implies P(h(x) \neq C^*(x)) \geq \varepsilon$$

$$\implies 1 - P(h(x) = C^*(x)) \geq \varepsilon$$

$$\implies P(h(x) = C^*(x)) \leq 1 - \varepsilon \qquad \swarrow P(h(s) = C^*(s))$$

$$\implies P(h(x_1) = C^*(x_1), \; \cdots, \; h(x_m) = C^*(x_m)).$$

$$\underline{\underline{i.i.d.}} \quad \prod_{i=1}^{m} P(h(x_i) = C^*(x_i)) \leq (1 - \varepsilon)^m$$

$$1 - x \leq e^{-x}$$

$$P\big(h_1(s) = C^*(s) \cup h_2(s) = C^*(s) \cup \cdots \cup h_{|H|}(s) = C^*(s)\big)$$

$$\leq \sum_{\ell=1}^{|H|} P(h_\ell(s) = C^*(s)) \leq \sum_{\ell=1}^{|H|} (1-\varepsilon)^m = |H|(1-\varepsilon)^m \leq \cdot |H| e^{-\varepsilon m}$$

# Sample Complexity for Supervised Learning

**Consistent Learner**

$(c^* \in H!, \quad err_S(h) = 0)$

- Input: S: $(x_1, c^*(x_1)), ..., (x_m, c^*(x_m))$

- Output: Find h in H consistent with the sample (if one exits).

$err_D(h) \leq err_S(h) + \epsilon$

**Theorem**

$$m \geq \frac{1}{\varepsilon}\left[\ln(|H|) + \ln\left(\frac{1}{\delta}\right)\right]$$

$\Rightarrow \epsilon \geq \frac{1}{m}\left[\ln|H| + \ln\frac{1}{\delta}\right]$

labeled examples are sufficient so that with prob. $1 - \delta$, all $h \in H$ with $err_D(h) \geq \varepsilon$ have $err_S(h) > 0$.

$err_S(h) = 0 \Rightarrow err_D(h) < \epsilon$

A ⟹ B

¬B ⟹ ¬A

¬B        ¬A

Contrapositive: if the target is in H, and we have an algo that can find consistent fns, then we only need this many examples to get generalization error $\leq \epsilon$ with prob. $\geq 1 - \delta$

$err_D(h) \geq \epsilon \geq \frac{1}{m}\left[\ln|H| + \ln\frac{1}{\delta}\right]$

A

¬A: $\quad err_D(h) \leq \epsilon \leq \frac{1}{m}\left[\ln|H| + \ln\frac{1}{\delta}\right]$

(S.L.T)

# Sample Complexity for Supervised Learning

**Consistent Learner**

- Input: S: $(x_1, c^*(x_1)), ..., (x_m, c^*(x_m))$

- Output: Find h in H consistent with the sample (if one exits).

**Theorem**

Bound inversely linear in $\epsilon$

$$m \geq \frac{1}{\varepsilon} \left[ \ln(|H|) + \ln\left(\frac{1}{\delta}\right) \right]$$

labeled examples are sufficient so that with prob. $1 - \delta$, all $h \in H$ with $err_D(h) \geq \varepsilon$ have $err_S(h) > 0$.

Bound only logarithmic in |H|

- $\epsilon$ is called error parameter
  - D might place low weight on certain parts of the space

- $\delta$ is called confidence parameter
  - there is a small chance the examples we get are not representative of the distribution

# Sample Complexity for Supervised Learning

**Consistent Learner**

- Input: S: $(x_1, c^*(x_1)), ...., (x_m, c^*(x_m))$

- Output: Find h in H consistent with the sample (if one exits).

**Theorem**

$$m \geq \frac{1}{\varepsilon}\left[\ln(|H|) + \ln\left(\frac{1}{\delta}\right)\right]$$

labeled examples are sufficient so that with prob. $1 - \delta$, all $h \in H$ with $err_D(h) \geq \varepsilon$ have $err_S(h) > 0$.

**Example:** H is the class of conjunctions over $X = \{0,1\}^n$.   $|H| = 3^n$

E.g., $h = x_1 \overline{x_3} x_5$ or $h = x_1 \overline{x_2} x_4 x_9$

Then $m \geq \frac{1}{\epsilon}\left[n \ln 3 + \ln\left(\frac{1}{\delta}\right)\right]$ suffice

$n = 10, \epsilon = 0.1, \delta = 0.01$  then $m \geq 156$ suffice

# Sample Complexity for Supervised Learning

**Theorem**

$$m \geq \frac{1}{\varepsilon}\left[\ln(|H|) + \ln\left(\frac{1}{\delta}\right)\right]$$

labeled examples are sufficient so that with prob. $1 - \delta$, all $h \in H$ with $err_D(h) \geq \varepsilon$ have $err_S(h) > 0$.

**Proof** Assume k bad hypotheses $h_1, h_2, \ldots, h_k$ with $err_D(h_i) \geq \epsilon$

1) Fix $h_i$. Prob. $h_i$ consistent with first training example is $\leq 1 - \epsilon$.

Prob. $h_i$ consistent with first m training examples is $\leq (1 - \epsilon)^m$.

2) Prob. that at least one $h_i$ consistent with first m training examples is $\leq k\,(1 - \epsilon)^m \leq |H|(1 - \epsilon)^m$.

3) Calculate value of m so that $|H|(1 - \epsilon)^m \leq \delta$

3) Use the fact that $1 - x \leq e^{-x}$, sufficient to set $|H|\,e^{-\epsilon m} \leq \delta$

# Sample Complexity: Finite Hypothesis Spaces

Realizable Case

**Theorem**

$$m \geq \frac{1}{\varepsilon}\left[\ln(|H|) + \ln\left(\frac{1}{\delta}\right)\right]$$

labeled examples are sufficient so that with prob. $1 - \delta$, all $h \in H$ with $err_D(h) \geq \varepsilon$ have $err_S(h) > 0$.

Probability over different samples
of m training examples

# Sample Complexity: Finite Hypothesis Spaces Realizable Case

1) PAC: How many examples suffice to guarantee small error whp.

**Theorem**

$$m \geq \frac{1}{\varepsilon}\left[\ln(|H|) + \ln\left(\frac{1}{\delta}\right)\right]$$

labeled examples are sufficient so that with prob. $1 - \delta$, all $h \in H$ with $err_D(h) \geq \varepsilon$ have $err_S(h) > 0$.

2) Statistical Learning Way:

With probability at least $1 - \delta$, for all $h \in H$ s.t. $err_S(h) = 0$ we have

$$err_D(h) \leq \frac{1}{m}\left(\ln|H| + \ln\left(\frac{1}{\delta}\right)\right).$$

# Supervised Learning: PAC model (Valiant)

- $X$ - instance space, e.g., $X = \{0,1\}^n$ or $X = R^n$
- $S_l = \{(x_i, y_i)\}$ - labeled examples drawn i.i.d. from some distr. $D$ over $X$ and labeled by some target concept $c^*$
  - labels $\in \{-1,1\}$ - binary classification

- Algorithm $A$ PAC-learns concept class $H$ if for any target $c^*$ in $H$, any distrib. $D$ over $X$, any $\varepsilon, \delta > 0$:
  - $A$ uses at most poly($n$,$1/\varepsilon$,$1/\delta$,size($c^*$)) examples and running time.
  - With probab. $1-\delta$, $A$ produces $h$ in $H$ of error at $\leq \varepsilon$.

What if $c^* \notin H$?

# Uniform Convergence

**Theorem**

$$m \geq \frac{1}{\varepsilon}\left[\ln(|H|) + \ln\left(\frac{1}{\delta}\right)\right]$$

labeled examples are sufficient so that with prob. $1 - \delta$, all $h \in H$ with $err_D(h) \geq \varepsilon$ have $err_S(h) > 0$.

- This basic result only bounds the chance that a bad hypothesis looks perfect on the data. What if there is no perfect h∈H (agnostic case)?

- What can we say if $c^* \notin H$?

- Can we say that whp all h∈H satisfy |err$_D$(h) – err$_S$(h)| ≤ ε?

  - Called "uniform convergence".
  - Motivates optimizing over S, even if we can't find a perfect function.

$$err_S(h) - \varepsilon \leq err_D(h) \leq err_S(h) + \varepsilon$$

# Sample Complexity: Finite Hypothesis Spaces

Realizable Case  $c^* \in H$

**Theorem**

$$m \geq \frac{1}{\varepsilon}\left[\ln(|H|) + \ln\left(\frac{1}{\delta}\right)\right]$$

labeled examples are sufficient so that with prob. $1 - \delta$, all $h \in H$ with $err_D(h) \geq \varepsilon$ have $err_S(h) > 0$.  $\implies$  $err_S(h) = 0 \implies err_D(h) \leq \varepsilon$

Agnostic Case  $c^* \notin H$

What if there is no perfect h?

**Theorem** After $m$ examples, with probab. $\geq 1 - \delta$, all $h \in H$ have $|err_D(h) - err_S(h)| < \varepsilon$, for

$$m \geq \frac{1}{2\varepsilon^2}\left[\ln(|H|) + \ln\left(\frac{2}{\delta}\right)\right]$$

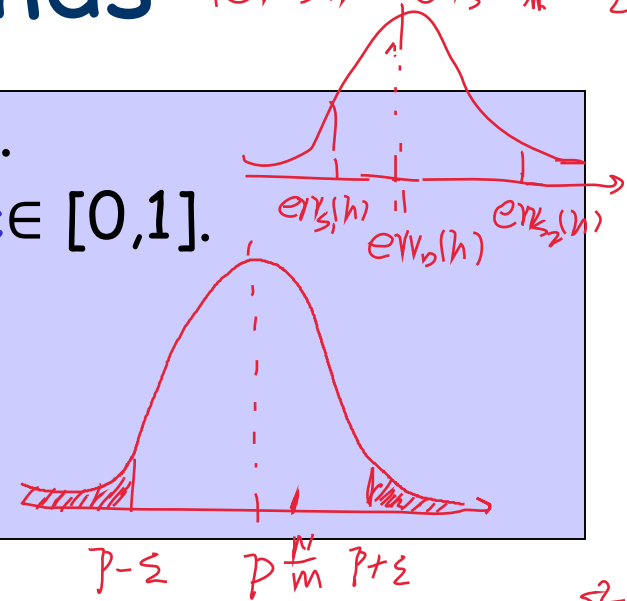To prove bounds like this, need some good tail inequalities.

# Hoeffding bounds

$\lvert err_D(h) - err_S(h)\rvert < \varepsilon$

Consider coin of bias p flipped m times.
 Let N be the observed # heads.  Let $\varepsilon \in [0,1]$.
Hoeffding bounds:

- $\Pr[N/m > p + \varepsilon] \leq e^{-2m\varepsilon^2}$, and
- $\Pr[N/m < p - \varepsilon] \leq e^{-2m\varepsilon^2}$.

$err_S(h)$   $err_S(h)$
$err_D(h)$   $err_D(h)$

Exponentially decreasing tails

$p-\varepsilon$   $p \dfrac{N}{m}$ $p+\varepsilon$

$\Rightarrow P\left(\lvert \tfrac{N}{m} - p\rvert > \varepsilon\right) \leq 2e^{-2m\varepsilon^2}$

- **Tail inequality:** bound probability mass in tail of distribution (how concentrated is a random variable around its expectation).

$P\left(\lvert err_S(h) - err_D(h)\rvert > \varepsilon\right) \leq 2e^{-2m\varepsilon^2}$

$err_D(h) = \Pr_D(h(x) \neq C^*(x))$

$= E_D\{h(x) \neq C^*(x)\}$

$err_S(h) = \dfrac{1}{m}\sum_{i=1}^{m} \mathbb{1}\{h(x_i) \neq C^*(x_i)\}$

$P\left(err_D(h) > err_S(h) + \varepsilon\right) \leq e^{-2m\varepsilon^2}$

- if $\exists h \in H,\ err_D(h) > err_S(h) + \varepsilon$

$P(\exists h \in H,\ err_D(h) > err_S(h) + \varepsilon) \leq \lvert H\rvert \cdot P(err_D(h) > err_S(h) + \varepsilon)$
$\leq \lvert H\rvert e^{-2m\varepsilon^2}$  Calculate the prob.

# Sample Complexity: Finite Hypothesis Spaces
## Agnostic Case

**Theorem** After $m$ examples, with probab. $\geq 1 - \delta$, all $h \in H$ have $|err_D(h) - err_S(h)| < \varepsilon$, for

$$m \geq \frac{1}{2\varepsilon^2}\left[\ln(|H|) + \ln\left(\frac{2}{\delta}\right)\right]$$

- **Proof:** Just apply Hoeffding.
  - Chance of failure at most $2|H|e^{-2|S|\varepsilon^2}$.
  - Set to $\delta$. Solve.
- So, whp, best on sample is $\varepsilon$-best over D.
  - Note: this is worse than previous bound ($1/\varepsilon$ has become $1/\varepsilon^2$), because we are asking for something stronger.
  - Can also get bounds "between" these two.

# What you should know

- Notion of sample complexity.

- Understand reasoning behind the simple sample complexity bound for finite H.