学习问题: learning problem.

unknown target function
$f : X \to Y$

↓

training data set
$D = \{(x_1, y_1), \cdots, (x_N, y_N)\}$

↓

hypothesis set H

→ learning algorithm A → final hypothesis
$g \approx f$.
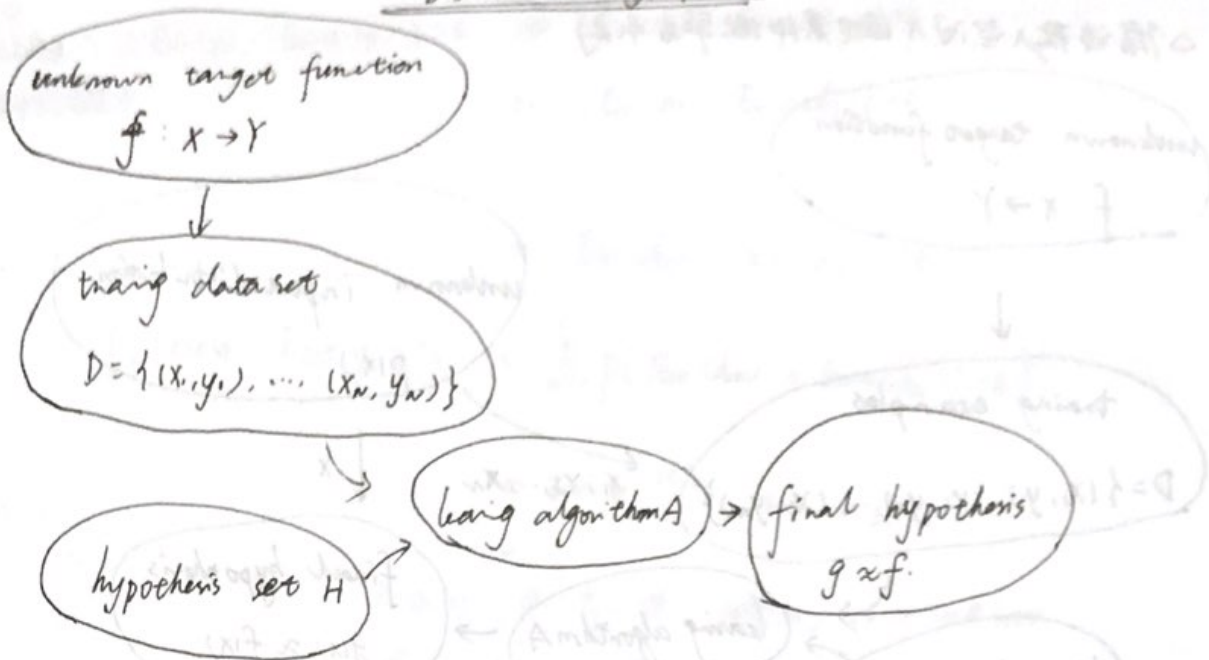
} hypothesis set: 假设目标函数是线性函数.

learning algorithm: 梯度下降. 牛顿法等.

学习是否可信: Is learning feasible?
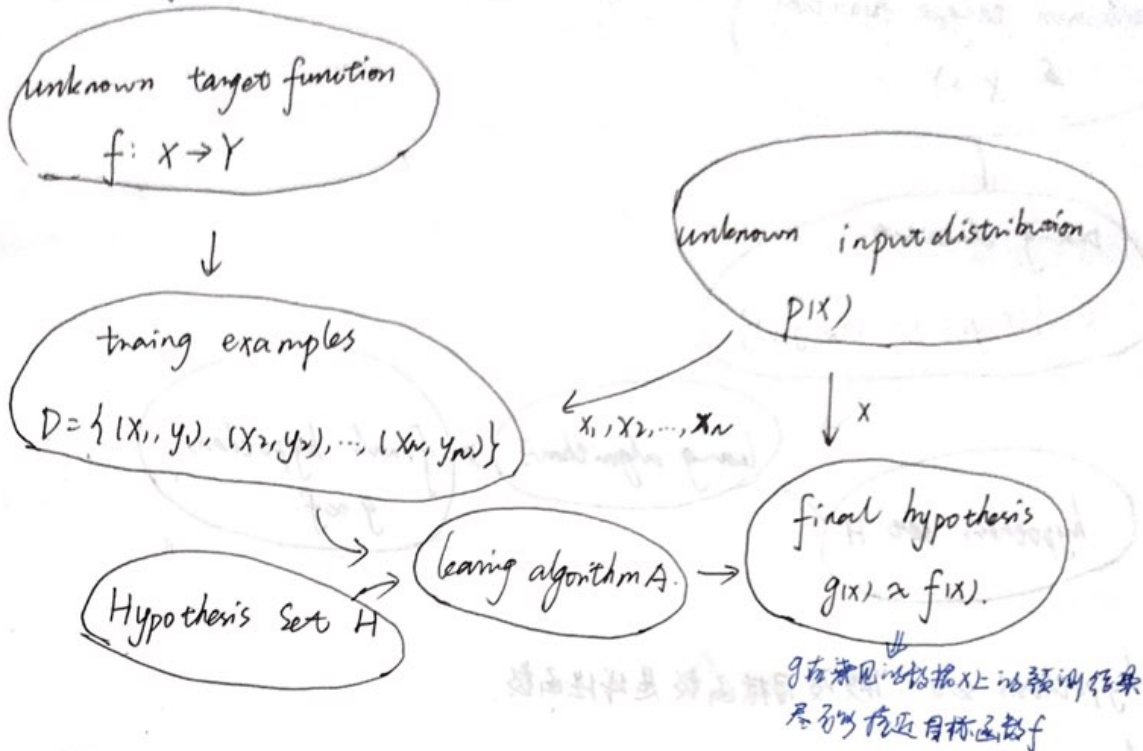
* How could a limited data set reveal enough information to pin down the entire target function?

* learning is feasible ⟹ we can learn something from seen dataset, which we didn't known before.

↓

reveal in probability way

probabilistic

△假设输入空间 X 满足某种概率分布则.

unknown target function
$f: X \to Y$

unknown input distribution
$p(x)$

$\downarrow$

traing examples
$D = \{(x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N)\}$

$x_1, x_2, \ldots, x_N$

$\downarrow X$

Hypothesis Set H

leaning algorithm A.

final hypothesis
$g(x) \approx f(x)$.

g 在常见的数据 X 上的预测结果
尽可能接近目标函数 f

## Hoeffding Inequality

$$p\left[|\nu - \mu| < \varepsilon\right] \leq 2e^{-2\varepsilon^2 N}. \quad \to \text{# samples.}$$

tolerance

$\varepsilon > 0$.

hypothesis   target function

in-sample error      out-sample error.

使用某种方式来
衡量 g 而 f 的接近程度.

$$\Rightarrow p\left[|E_{in}(g) - E_{out}(g)| > \varepsilon\right] \leq 2e^{-2\varepsilon^2 N}, \quad \varepsilon > 0. \quad (1)$$

$$E_{in}(g) = \frac{1}{N} \sum_{n=1}^{N} I\left[g(x_n) \neq f(x_n)\right] \qquad \text{in-sample error}$$

$$E_{out}(g) = P\left[g(x) \neq f(x)\right] \qquad \text{out-sample error.}$$

upper bound



$N$

$\Rightarrow$ 为了压低上界，使得 $E_{in}$ 和 $E_{out}$ 尽可能接近，所以增大 N.

$\Rightarrow$ 从 $e^{-x}$ 的曲线来看，当 N 达到一定数量，增加 N 带来的好处很少.

我们不知道g, $\because |E_{in}(g) - E_{out}(g)| > \varepsilon$ $\Rightarrow |E_{in}(h_1) - E_{out}(h_1)| > \varepsilon$

但是我们知道h.

$$or \quad |E_{in}(h_2) - E_{out}(h_2)| > \varepsilon$$

$$\cdots$$

$$or \quad |E_{in}(h_m) - E_{out}(h_m)| > \varepsilon$$

$$\therefore P[|E_{in}(g) - E_{out}(g)| > \varepsilon] \leq \sum_{m=1}^{M} P[|E_{in}(h_m) - E_{out}(h_m)| > \varepsilon]$$

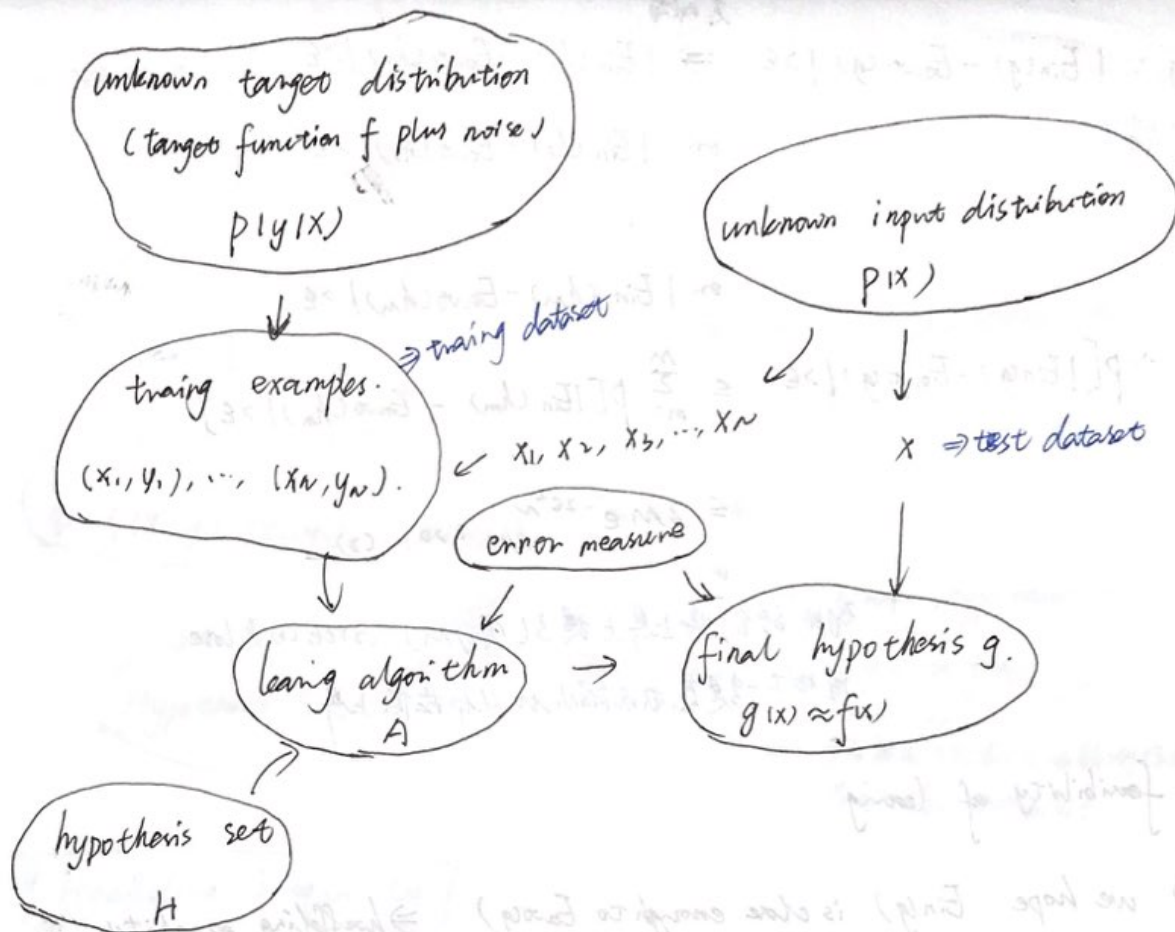$$\leq 2M e^{-2\varepsilon^2 N}, \quad \varepsilon > 0. \quad (2)$$

到此,我们将上界往上提了(因为M). (2)比(1)更 lose.

用比下一步是想.办法缩小M,从而在低上界

$\Downarrow$ feasibility of learning

$\begin{cases} 1) \text{ we hope } E_{in}(g) \text{ is close enough to } E_{out}(g) & \Rightarrow \text{hoeffding equality} \\ 2) \text{ we hope } E_{in}(g) \text{ small enough} & \Rightarrow \text{test } g \text{ with training set} \end{cases}$

$\Rrightarrow E_{out}(g)$ is small enough.

unknown target distribution
( target function f plus noise )
$p(y|x)$

unknown input distribution
$p(x)$

training examples.
$(x_1, y_1), \cdots, (x_N, y_N).$

⇒ training dataset

$x_1, x_2, x_3, \cdots, x_N$

$x$ ⇒ test dataset

error measure

learning algorithm
A

final hypothesis g.
$g(x) \approx f(x)$

hypothesis set
H

由 Hoeffding Inequality 有:

$$P\left[\,|E_{in}(g) - E_{out}(g)| > \varepsilon\,\right] \le 2M e^{-2\varepsilon^2 n}, \quad \varepsilon > 0.$$

∴ 至少有 $1 - 2Me^{-2\varepsilon^2 n}$ 的概率, $E_{out}(g) \le E_{in}(g) + \varepsilon$.

解令 $\delta = 1 - 2Me^{-2\varepsilon^2 n}$. $\Rightarrow \varepsilon = \sqrt{\dfrac{1}{2n} \ln \dfrac{2M}{\delta}}$.

∴ 至少有 $1 - \delta$ 的概率.

$$E_{out}(g) \le E_{in}(g) + \sqrt{\dfrac{1}{2n} \ln \dfrac{2M}{\delta}} \qquad (3)$$

↳ generalization bound.

⇒ 目标: 使用一个更小的值来替换 $M$, 从而缩小 $E_{out}(g)$ 的上界.

1) 定义 growth function. ⇒ 表征假设空间中的有效假设个数.

2) 找到 growth function 的上界.

3) 使用 growth function 替换 $M$ 得到更 tight 的 generalization bound.

1) Define Growth Function.

假设 H 将输入空间 X 映射为 {-1, +1}, 那么有:

$$(x_1, x_2, \ldots, x_n) \xrightarrow{h_1} (h_1(x_1), h_1(x_2), \ldots, h_1(x_n)) \Rightarrow \text{1个 dichotomy}$$

$$(x_1, x_2, \ldots, x_n) \xrightarrow{h_2} (h_2(x_1), h_2(x_2), \ldots, h_2(x_n)) \Rightarrow \text{1个 dichotomy}$$

1° 由一个假设 h 生成的划分结果 (1个 N-tuple) 是一个 dichotomy.
因为它将 n 个点二分为二 (-1 或 ±1 类).

2° $h_1$ 和 $h_2$ 不是 dichotomy, 它们的输出结果才是. 2个假设可能生成相同的 dichotomy.
因此 dichotomy 的数量小于 M.

> **Definition** 由假设空间 H 在 $x_1, \ldots, x_n \in X$ n个点生成的 dichotomies 定义为:
>
> $$H(x_1, x_2, \ldots, x_n) = \{ (h(x_1), h(x_2), \ldots, h(x_n)) \mid h \in H \}$$

> **Definition** [growth function]
>
> $$m_H(n) = \max_{x_1, \ldots, x_n \in X} |H(x_1, \ldots, x_n)|$$

growth function 定义为假设集 H 在给定的 n 个样本上所生成的最多个
dichotomy 的数量.

因为 H 将 X 映射为 2 个类别, 因此: $m_H(n) \leq 2^n < M$

↳ 所以使用 $m_H(n)$ 替换 M 得到更 tight 的上界.

我们说 H 能 shatter $x_1, \ldots x_n$ 意味着 H 能生成所有的 dichotomies, 即且 $m_H(n) = 2^n$.

**Definition** 如果 H 不能够 shatter 任何 由 包含 k 个点的 任意 抽取集，称此 k 就是 H 的一个 break point。于是 $m_H(k) < 2^k$.

**Definition** 给定 N 个点，定义 $B(N, k)$ 为这 N 个 点上 的最大 dichotomies 个数，且 这些 dichotomies 无法 shatter 任意 包含 k 个点的子集。由此 有 $m_H(N) \leq B(N, k)$, 如果 k 是 H 的 break point.

## Bounding Growth Function

定义

$$\begin{cases} B(N, 1) = 1 \\ B(1, k) = 2 \end{cases}, \text{ for } k > 1$$

$\because \quad B(N, k) = \alpha + 2\beta$.

$\alpha + \beta \leq B(N-1, k)$ &

$\beta \leq B(N-1, k-1)$.

$\therefore B(N, k) = \alpha + 2\beta \leq B(N-1, k) + B(N-1, k-1)$

**Lemma**

$$B(N, k) \leq \sum_{i=0}^{k-1} \binom{N}{i}$$

Proof. 已知 $B(N_0, k)$ 设 $N = N_0 + 1$, 有:

$$B(N_0+1, k) \leq B(N_0, k) + B(N_0, k-1)$$

$$\leq \sum_{i=0}^{k-1} \binom{N_0}{i} + \sum_{i=0}^{k-2} \binom{N_0}{i} = 1 + \sum_{i=1}^{k-1} \binom{N_0}{i} + \sum_{i=1}^{k-1} \binom{N_0}{i-1}$$

$$= 1 + \sum_{i=1}^{k-1} \left[ \binom{N_0}{i} + \binom{N_0}{i-1} \right] = 1 + \sum_{i=1}^{k-1} \binom{N_0+1}{i} = \sum_{i=0}^{k-1} \binom{N}{i}$$

$$\boxed{\text{Theorem}}$$

$$m_H(N) \le \sum_{i=0}^{k-1} \binom{N}{i}$$

$$\boxed{\text{Definition}}\ \text{VC dimension. 一个假设集的 VC dimension 定义为 H所能够 shatter 的}$$
最大数据点数，即 $m_H(d_{vc}) = 2^{d_{vc}}$.

因此 $k = d_{vc}+1$ 是一个 break point，有: $\boxed{m_H(N) \le \sum_{i=0}^{d_{vc}} \binom{N}{i}}$.

$$\boxed{m_H(N) \le N^{d_{vc}} + 1}$$

$$\boxed{\text{Definition}}\ (VC\ generalization\ bound).\ \text{For any tolerance } \delta > 0,$$

$$E_{out}(g) \le E_{in}(g) + \sqrt{\frac{8}{N} \ln \frac{4 m_H(2N)}{\delta}}$$

with probability $\ge 1 - \delta$.

$$\boxed{\text{Sample Complexity}}\ \text{假设给定 } \delta > 0,\text{ 代价泛化误差最多只有 } \varepsilon, \text{则}.$$

$$N \ge \frac{8}{\varepsilon^2} \ln \left( \frac{4((2N)^{d_{vc}} + 1)}{\delta} \right).$$