## Recap:

↳ <u>Growth function:</u>

$$m_H(N) = \max_{x_1, \cdots x_N \in \mathbb{R}^d} |H(x_1, \cdots x_N)|$$

worst case No. of dich. vectors $\equiv$ effective #. of hypothesis

↳ <u>VC - dimension:</u>

to bound $m_H(N)$:

$$d_{vc}(H) = N \quad s.t. \begin{cases} m_H(N) = 2^N \\ \\ m_H(N+1) < 2^{N+1} \end{cases}$$

i.e. $N+1$ is the first breakpoint

e.g. Linear classifier $d = 2$

$$m_H(3) = 8 = 2^3, \quad m_H(4) = 14 < 2^4$$

$$\therefore d_{vc}(H) = 3$$

---

↳ <u>Theorem:</u>

For any hypo. set $H$ with $d_{vc}(H) < \infty$,

① $m_H(N) \leq \sum_{i=0}^{d_{vc}(H)} \binom{n}{i} \leq N^{d_{vc}(H)} + 1$

② $Pr\{\triangle(g) > \varepsilon\} \leq 4 \, m_H(2N) \, e^{-\frac{1}{8}N\varepsilon^2}$

Recall: $\triangle(g) = |E_{in}(g) - E_{out}(g)|$

↑
generalization error

// 用处: previously we do $m_H(N) \leq 2^N$, exponential in $N$
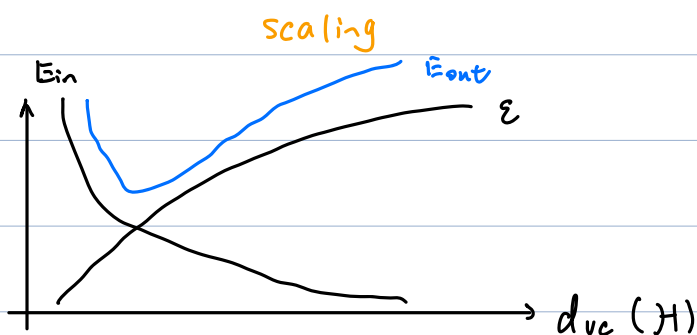
Not the best way to bound

But above is poly. bound

↳ _Note_ :

set $\delta = 4 \, m_H(2N) \, e^{-\frac{1}{8} N \varepsilon^2}$ , with prob. $1 - \delta$

$E_{out}(g) \leq E_{in}(g) + \varepsilon$

$\therefore \varepsilon = \sqrt{\dfrac{8}{N} \log \dfrac{4 \, m_H(2N)}{\delta}} \leq \sqrt{\dfrac{8}{N} \log \dfrac{4 \, ((2N)^{d_{vc}(H)} + 1)}{\delta}}$

$$= \theta \left( \sqrt{\dfrac{d_{vc}(H)}{N} \log \dfrac{N}{\delta}} \right)$$

↑
scaling



e.g.1. Linear classifier in $\mathbb{R}^d$

$\qquad d_{vc}(H) = d + 1$

e.g.2. $2^{nd}$ order model

$\qquad \underline{x} = [1, x_1, x_2] \qquad d = 2$

$\qquad$ Transform to new space:

$\qquad \underline{z} = [1, x_1, x_2, x_1 x_2, x_1^2, x_2^2] \qquad d = 5$

$\qquad$ lin. classification in $\underline{z}$ space $\qquad d_{vc}(H) = 6$

↳ In general, $\underline{x} \in \mathbb{R}^d$, $k^{th}$ order polynomial

$\qquad d_{vc}(H) = \dbinom{k+d}{d} = O(k^d)$

↳ In practice, a rule of thumb is to choose $H$ s.t.

$\qquad d_{vc}(H) \approx \dfrac{N}{10}$

# Generalization Bound in Regression

## 1. Squared error:

↳ $E_{in}(g) = \frac{1}{N} \sum_{n=1}^{N} (y_n - g(\underline{x}_n))^2$

$E_{out}(g) = \mathbb{E}\left[ (y - g(\underline{x}_n)^2 \right]$

↳ We use the bias-variance trade off.

## 2. Bias-variance trade off   (ch. 2.3)

Today : setup for learning model          例 : trade off.

↳ **training set:**

$\mathcal{D} = \{ (\underline{x}_1, y_1) \cdots (\underline{x}_N, y_N) \}$

$\sim P(\mathcal{D}) = P(\underline{x}_1) P(\underline{x}_2) \cdots P(\underline{x}_N)$

↳ **Output hypothesis:**    $g^D$

$\hat{y} = g^D(\underline{x})$
↑
output label

→ 把 —↑ D to tell that we have train the hypo. w/ data.

↳ **Average hypothesis:**

$\bar{g}(\underline{x}) = \mathbb{E}_D\left[ g^D(\underline{x}) \right]$

e.g.  $d = 1$,  $x \in [-1, 1]$

$y = \sin(\pi x)$  // unknown target func.

$\mathcal{D} = \{(u, v)\}$,   $u \sim U(-1, 1)$  ← $P(x)$

∴ $V = \sin(\pi u)$

Assume $H \equiv$ constant hypothesis

( meaning that whatever the model is, the output is the same)

∴ $g^D(x) = V$,  for any $x$.

$$\therefore \bar{g} = \mathbb{E}_u \left[ g^D(x) \right]$$

Avg. Hypo. $\nearrow$

$$= \mathbb{E}_u \left[ v \right]$$

$$= \mathbb{E}_u \left[ \sin(\pi u) \right]$$

$$= 0$$

↳ "best" approx. to $f(x)$ given infinite amount of data

But cannot find $\bar{g}$ in practice.

↳ **Def:**

**Bias of learning model** for any $x$:

$$\text{bias} (\underline{x}) = \left( \bar{g}(x) - f(x) \right)^2$$

**Variance** of `` `` `` ``

$$\text{var} (\underline{x}) = \mathbb{E}_D \left[ \left( g^D(\underline{x}) - \bar{g}(x) \right)^2 \right]$$

e.g. In the prev example.

$$\text{bias} (\underline{x}) = \left( 0 - \sin(\pi x) \right)^2 = \sin^2(\pi x)$$

$$\text{var} (\underline{x}) = \mathbb{E}_u \left[ (v - 0)^2 \right]$$

$$= \mathbb{E}_u \left[ \sin^2(\pi u) \right]$$

$$= \int_{-1}^{1} \frac{1}{2} \sin^2(\pi u) \, du$$

$$= \frac{1}{2}$$