1. Intro:

↳ For supervised learning:

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

$$\mathcal{D} \longrightarrow \boxed{\text{Learning Algo}} \longrightarrow g$$

· Training error:

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

· Test error:

$$E_{out}(g) = \mathbb{E}_{\underline{x}}[e(g(\underline{x}), y)].$$

·



due to overfitting

$E_{out}$

$E_{in}$

time

stop here is good. But how to know this pt?

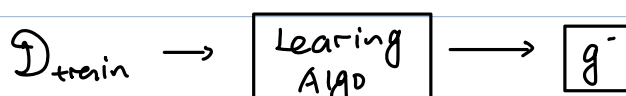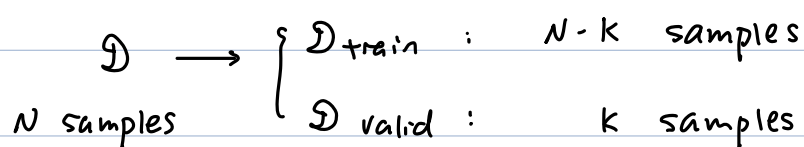↳ Want to estimate $E_{out}(g)$ using training data.

__Problem:__ $E_{in}(h)$ is a good estimation for $E_{out}(h)$   (Hoeffding)

when h is given __independent__ of $\mathcal{D}$

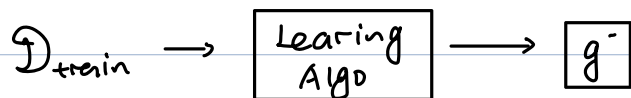But not necessarily good for final hypo. g : VC dim

bias-var trade off

__Idea:__ use validation dataset

$$\mathcal{D} \longrightarrow \begin{cases} \mathcal{D}_{train} : & N-k \text{ samples} \\ \mathcal{D}_{valid} : & k \text{ samples} \end{cases}$$

N samples

$$\mathcal{D}_{train} \longrightarrow \boxed{\begin{array}{c}\text{Learing} \\ \text{Algo}\end{array}} \longrightarrow \boxed{g^-}$$

## 2. Validation

$\mathcal{D} \longrightarrow \begin{cases} \mathcal{D}_{train} : & N-k \text{ samples} \\ \mathcal{D}_{valid} : & k \text{ samples} \end{cases}$

$\mathcal{D}_{train} \longrightarrow \boxed{\begin{array}{c} \text{Learing} \\ \text{Algo} \end{array}} \longrightarrow \boxed{g^-}$

↳ $E_{val}(g^-) = \frac{1}{k} \sum_{(x_n, y_n) \in \mathcal{D}} e(g^-(x_n), y_n)$

We want: $E_{val}(g^-) \approx E_{out}(g^-) \approx E_{out}(g)$

↳ Prf $E_{val}(g^-) \approx E_{out}(g^-)$:

- Properties:

  ① $\mathbb{E}_{\mathcal{D}_{val}}[E_{val}(g^-)] = E_{out}(g^-)$

    ⇒ $E_{val}$ is an unbiased estimate for $E_{out}(g^-)$

  ② $Var[E_{val}(g^-)] = \frac{1}{k}(\sigma^2)$, where

    $\sigma^2 = var[e(g^-(x), y)]$

    As $k \to \infty$, $Var \to 0$

    ∴ $E_{val}(g^-)$ is a consistent estimate for $E_{out}(g^-)$

  ③ With prob. $1-\delta$,

    $E_{out}(g^-) \leq E_{val}(g^-) + \sqrt{\frac{1}{2k} \log \frac{2}{\delta}}$

    $\text{// } O(\frac{1}{\sqrt{k}}) \text{ for binary classification}$

- Fact:

  Relation b/t $E_{val}(g^-)$ and $E_{out}(g^-)$ is nearly identical to

  ˵  ˵  $E_{in}(h)$ and $E_{out}(h)$

- **Prf:**

  Denote $\mathcal{D}_{val} = \{(x_1, y_1) \cdots (x_k, y_k)\}$

  ∵ Property ①:

  $$\mathbb{E}_{\mathcal{D}val}[E_{val}(g^-)] = \mathbb{E}_{val}\left[\underbrace{\frac{1}{k}\sum_{n=1}^{k} e(g^-(x_n), y_n)}_{\color{blue}E_{val}(g^-)}\right]$$

  $$= \frac{1}{k}\sum_{n=1}^{k} \underbrace{\mathbb{E}_{x_n}[e(g^-(x_n), y_n)]}_{\color{blue}def. \; of \; E_{out}(g^-)}$$

  $$= E_{out}(g^-)$$

  ∵ Property ②:

  $$Var[E_{val}(g^-)] = Var\left[\frac{1}{k}\sum_{n=1}^{k} e(g^-(x_n), y_n)\right]$$

  $$= \frac{1}{k^2}\sum_{n=1}^{k} \sigma^2$$

  $$= \frac{1}{k}\sigma^2$$

  ∵ Property ③:

  $g^-$ is independent of $\mathcal{D}_{val}$

  ⇒ Hoeffding bound is valid with $k$ samples

↳ **How to select $k$?**

  $E_{val}(g^-) \approx E_{out}(g^-)$  for large $k$

  $E_{out}(g^+) \approx E_{out}(g)$  for small $k$

  In practise  $k \approx \frac{N}{5}$  is reasonable  ⇒ 20%

## 3. Model selection by Validation:

↳ **Problem:** Given a dataset $\mathcal{D}$,

select the best model from $H_1 \cdots H_M$

e.g. $H_i$ : the class of ith-order polynomial

**step 1:**

train each model $H_m$ on the training set $\mathcal{D}_{train}$ to output the

final hypo. $g_m^-$

**step 2:**

use $\mathcal{D}_{val}$ to compute $E_{val}(g_m^-)$ for $1 \le m \le M$

**step 3:**

select the best model $H_m{}^*$

$$m^* = \underset{1 \le m \le M}{\arg\min} \; E_{val}(g_m^-)$$

( can do better if have time)

**step 4:**

use the complete dataset to train $H_m{}^*$ and output the

final hypo. $g_{m^*}$.