

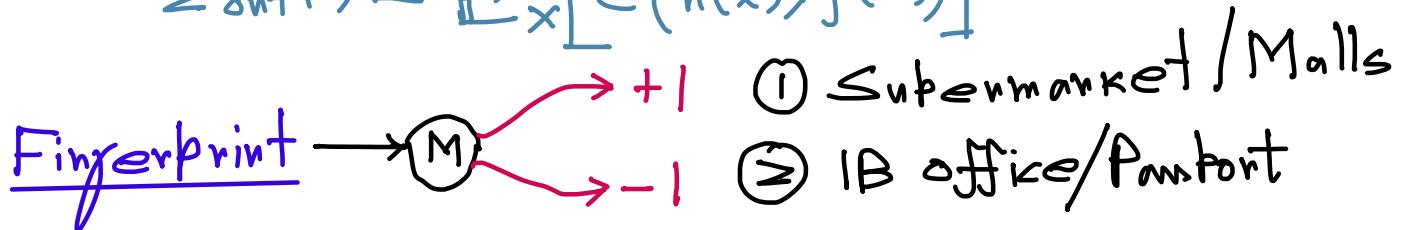
Error!

Regression :- $e(h(x_n), f(x_n)) = (h(x_n) - f(x_n))^2$

Classification :- Pointwise $= [h(x_n) \neq f(x_n)]$

$$\text{Err}_h(h) = \frac{1}{N} \sum_{n=1}^N e(h(x_n), f(x_n))$$

$$\text{Err}_{\text{out}}(h) = \mathbb{E}_x [e(h(x), f(x))]$$



$\textcircled{1}$	Predicted Class		$\textcircled{2}$	Predicted Class	
Actual Class	Yes	No	Actual Class	Yes	No
Yes		-100	Yes		-1
No	-10		No	-10	

Given domain Knowledge

⇒ if not Known Domain Knowledge

① Plausible

② Freindly

Noise :- → Unknown Distribution
 $\Rightarrow P(y|x) \leftarrow$ Points being Generated

Age
Salary
:
:

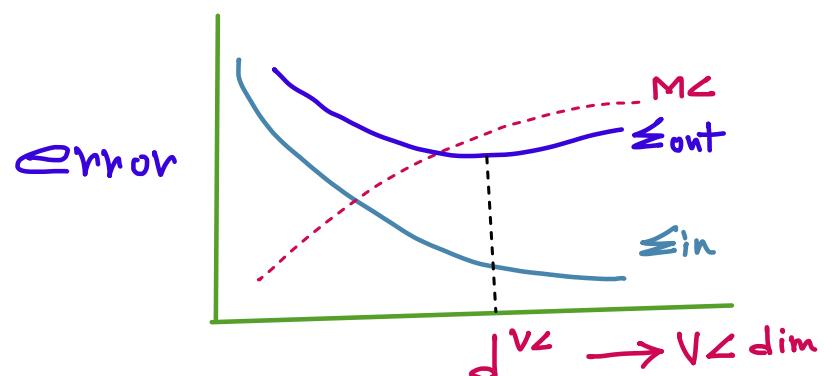
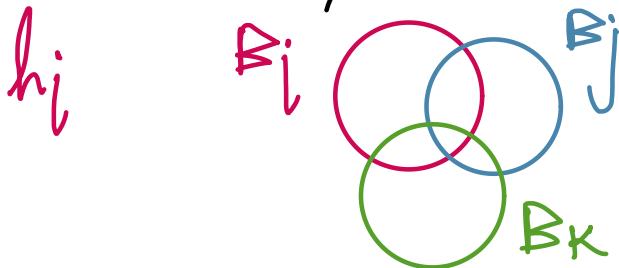
→ +1
-1

$$\text{Noisy Targets} = \frac{\text{Def}}{\text{Target}} + \text{Noise}$$

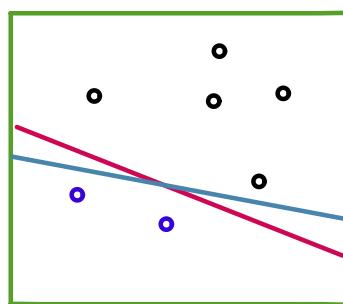
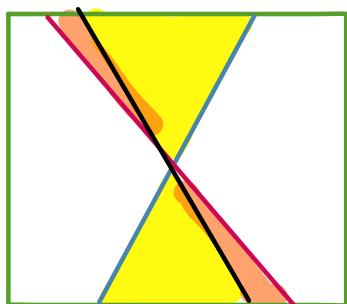
$$f(x) = \mathbb{E}(y|x)$$

$$\frac{\text{Def}}{\text{Target}} P(y|x) = 0 \text{ except } y = f(x)$$

$$(x, y) \sim P(y|x) P(x)$$



$$P(\text{Bad}) \leq \underbrace{P(B_1) + \dots + P(B_M)}_{M \text{ factor}}$$



$$h: X \rightarrow \{+1, -1\}$$

$$d: \{x_1, \dots, x_N\} \rightarrow \{+1, -1\}$$

Dichotomy

$$m_H(n) = \max_{x_1, \dots, x_n \in X} |H(x_1, \dots, x_n)|$$

SHATTER

$$m_p(1) = \geq$$

$$m_p(2) = 4$$

$$m_p(3) = 8$$

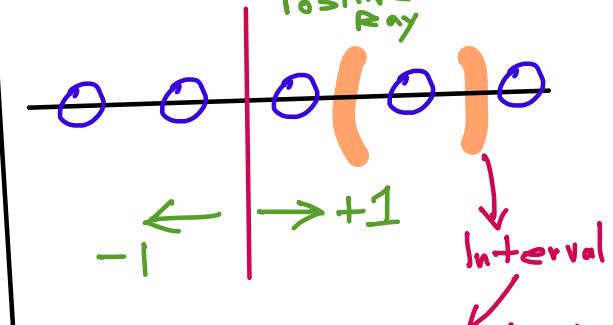
$$m_p(4) = 14$$

Breakpoint
 $K \geq 4$

$$\begin{matrix} -1 \\ \circ \end{matrix} \quad \begin{matrix} +1 \\ \circ \end{matrix}$$

$$\begin{matrix} -1 \\ \circ \end{matrix} \quad \begin{matrix} +1 \\ \circ \end{matrix}$$

$$m_H(N) = N + 1 \rightarrow \text{B.P.} = \geq$$

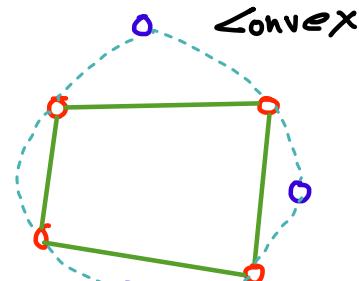


$$m_H(N) = \binom{N+1}{2} + 1$$

$$= \frac{1}{2}(N^2 + N) + 1$$

$$\begin{matrix} +1 \\ \circ \end{matrix} \quad \begin{matrix} -1 \\ \circ \end{matrix} \quad \begin{matrix} +1 \\ \circ \end{matrix}$$

Convex functions:-



$$m_H(N) = \geq^N$$

$$K = \infty (N_0)$$

► If BP = finite, then $m_H(N)$ is poly

$$m_H(N) \leq \dots \leq \dots \leq \text{Poly}$$

$B(N, K)$ = Max number of dichotomy given N points
with B.P.K.

$$m_H(N)$$

With B.P.K.

$$\leq \alpha + \beta$$

$$= (\alpha + \beta) + \beta$$

$$\leq B(N-1, K-1)$$

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

$$N^{d_{VC}}$$

$$\text{Perception} \leq O(N^3)$$

$$d_{VC} = K-1$$

V < -1 inequality

$$\Pr [|E_{in}(h) - E_{out}(h)| > \epsilon] \leq 4 m_H(\geq N) e^{-\frac{\epsilon^2 N}{8}}$$

$$S = 4 m_H(\geq N) e^{-\frac{1}{8} \epsilon^2 N} \quad PAC$$

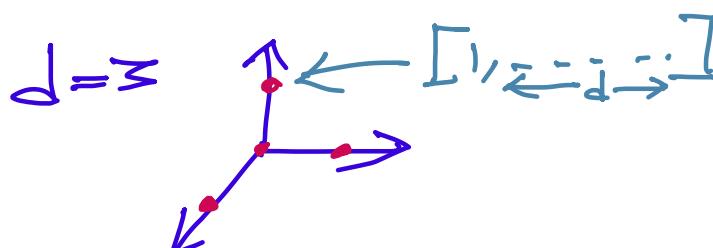
$$\Pr [|E_{in} - E_{out}| < \epsilon] \geq 1 - S$$

$$\Rightarrow E_{out} \leq E_{in} + \epsilon \leq \sqrt{\frac{8}{N} \log \frac{4 m_H(\geq N)}{S}} \quad \text{Generalization Bound}$$

$$\underline{d\text{-dim}} \quad d_{VC} = d+1$$

$d_{VC} \leq d+1 \rightarrow d+1 \text{ can be shattered}$

$d_{VC} \geq d+1 \rightarrow d+2 \text{ cannot be } "$



$$X = \begin{bmatrix} 1 & 0 & 0 & 0 & \cdots & 0 \\ 1 & 1 & 0 & 0 & \cdots & 0 \\ 1 & 0 & 1 & 0 & \cdots & 0 \\ \hline 1 & 0 & 0 & 0 & \cdots & 1 \end{bmatrix}_{(d+1) \times (d+1)}$$

$$y = \begin{bmatrix} \pm & | & \\ \pm & | & \\ \vdots & | & \\ \pm & | & \end{bmatrix} \xrightarrow{\text{XW=Y}} \Rightarrow W = X^{-1}Y$$

$$\Rightarrow d+1 \text{ points } x_j = \sum_{i \neq j} a_{ij} x_i \quad i \neq 0 \text{ for some } i$$

$$\Rightarrow a_{ij} \neq 0 \xrightarrow{>0 \Rightarrow +1} \xrightarrow{<0 \Rightarrow -1} \frac{x_j \xrightarrow{+1 > 0} x_j \xrightarrow{-1 < 0}}{j \xrightarrow{-1}}$$

$$w^T x_j \Rightarrow >0 \quad (\text{Contradiction})$$

Bias vs Variance Tradeoff

$$H = \{h_1, \dots, h_n\} \rightarrow g \approx f$$

① H complex $\Rightarrow \|f\| \downarrow$

How to get g ?



$$\begin{aligned} E_{out}(g) &= E_D \left[E_x \left[(y^D(x) - f(x))^2 \right] \right] \\ &= E_x \left[E_D \left[(y^D(x) - f(x))^2 \right] \right] \end{aligned}$$

any hypothesis

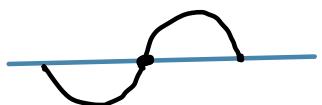
$$\bar{y}(x) = E_D(y^D(x))$$

$$= E_D \left[(y^D(x) - \bar{y}(x))^2 \right] + (\bar{y}(x) - f(x))^2$$

$\uparrow \text{Var}(x)$ $\uparrow \text{Bias}(x)$

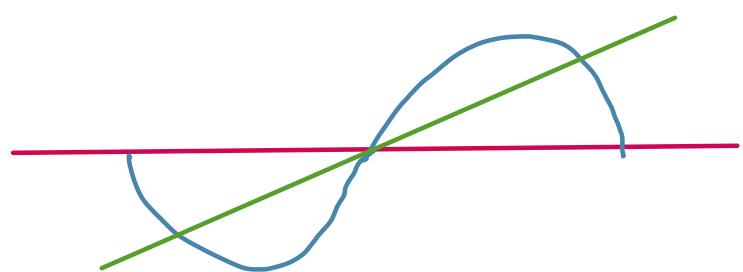
$$E_{out}(g) = E_x [\text{Var}(x) + \text{Bias}(x)]$$

$$f: [-1, +1] \rightarrow \mathbb{R} \quad f(x) = \sin \pi x$$

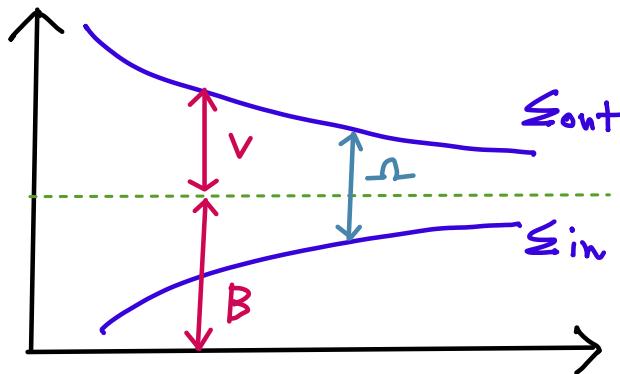
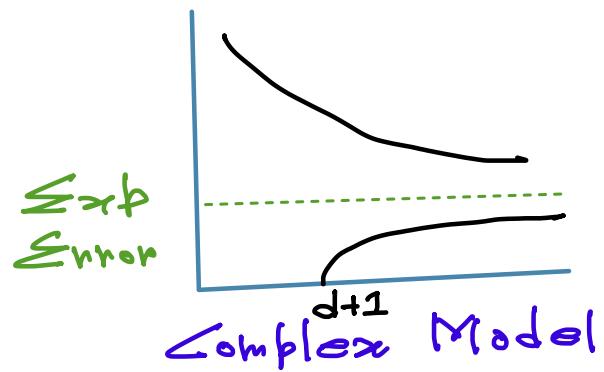
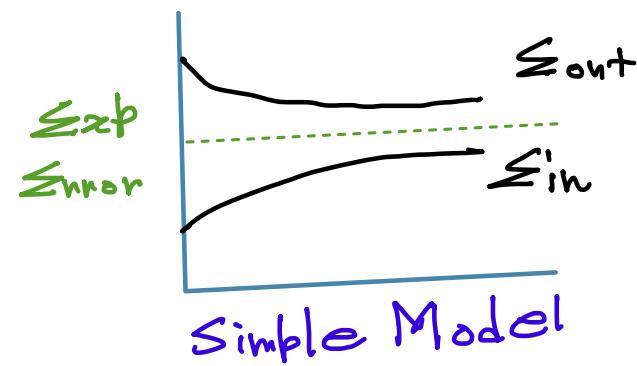


$$H_1: g(x) = b$$

$$H_2: g(x) = ax + b$$



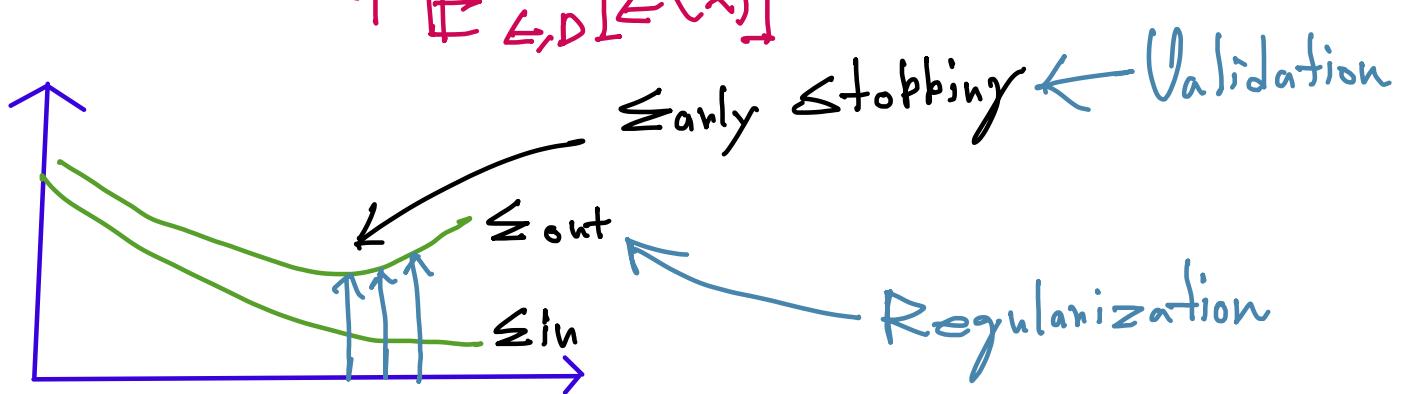
|D| = ≥ points



Overshooting $f(x) + \varepsilon(x) \rightsquigarrow \sigma^2$

$$\sum_{q=0}^{\infty} x^q a^q \quad \text{Normalized}$$

$$E_{out}(g) = E_D [(\bar{y}(x) - \bar{f}(x))^2] + (\bar{f}(x) - f(x))^2 + E_{\varepsilon, D} [\varepsilon(x)]$$



Regularization:

Legendre Poly $\gamma(\mathbf{z}) = \mathbf{w}^T \mathbf{z}$

$E_{in}(w) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{z}_n - y_n)^2$

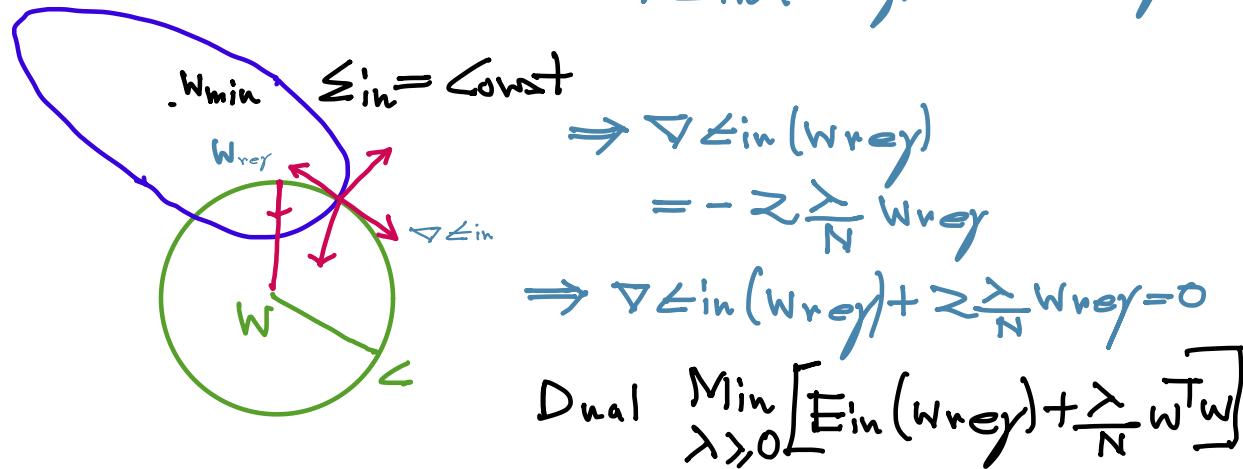
Minimize $= \frac{1}{N} (\mathbf{z} w - y)^T (\mathbf{z} w - y)$

$\frac{\partial E_{in}(w)}{\partial w} = 0 \Rightarrow w_{in} \leftarrow (\mathbf{z}^T \mathbf{z})^{-1} \mathbf{z}^T y$ s.t. $w^T w \leq C$

Hard Constraint Soft Constraint

$C_{\text{eff}}(g \geq z) \rightarrow 0 \quad w^T w \leq C$

$\nabla E_{in}(w_{\text{reg}}) \propto -w_{\text{reg}}$



$$\frac{\partial}{\partial w} \left[\frac{1}{N} (\mathbf{z} w - y)^T (\mathbf{z} w - y) + \frac{\lambda}{N} w^T w \right] = 0 \rightarrow E_{\text{reg}}(w)$$

$$\Rightarrow \mathbf{z}^T (\mathbf{z} w - y) + \lambda w = 0 \quad \lambda \uparrow \downarrow$$

$$w = (\mathbf{z}^T \mathbf{z} + \lambda N I)^{-1} \mathbf{z}^T y \quad \lambda \downarrow \uparrow$$

$$w^{t+1} \leftarrow w^t - \eta \nabla E_{\text{reg}}(w) \quad \text{Instead of } \nabla E_{in}(w) \text{ for reg}$$

$$= w^t - \eta \left(\nabla E_{in}(w) + \frac{\lambda}{N} w^t \right)$$

$$= w^t \left(1 - \frac{\lambda}{N} \right) - \eta \nabla E_{in}(w)$$

Weight Decay

$$E_{\text{out}}(y) = \mathbb{E}_x [e(y(x), y)] \quad \sigma^2 = \text{Var}[e(y(x), y)]$$

$(x_1, y_1) \dots (x_K, y_K)$

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E}_{x_k} [e(y(x_k), y)] = E_{\text{out}}(y)$$

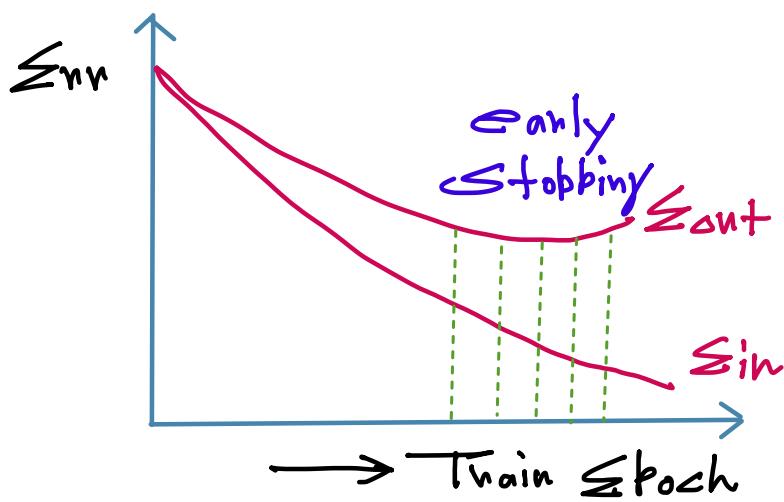
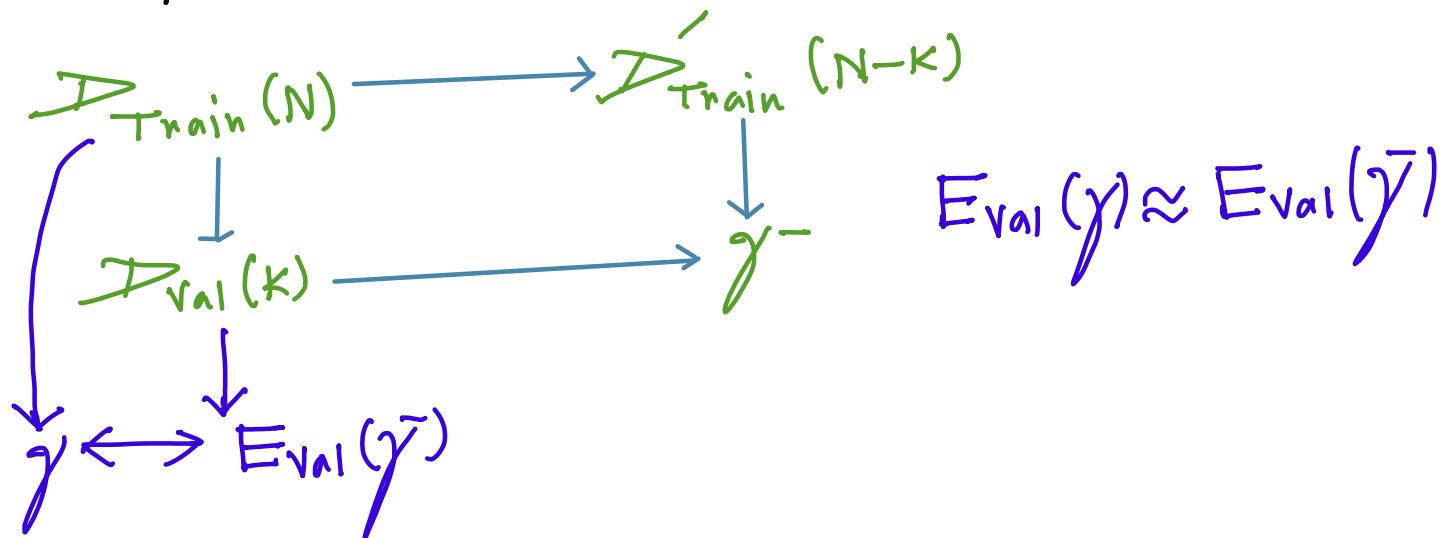
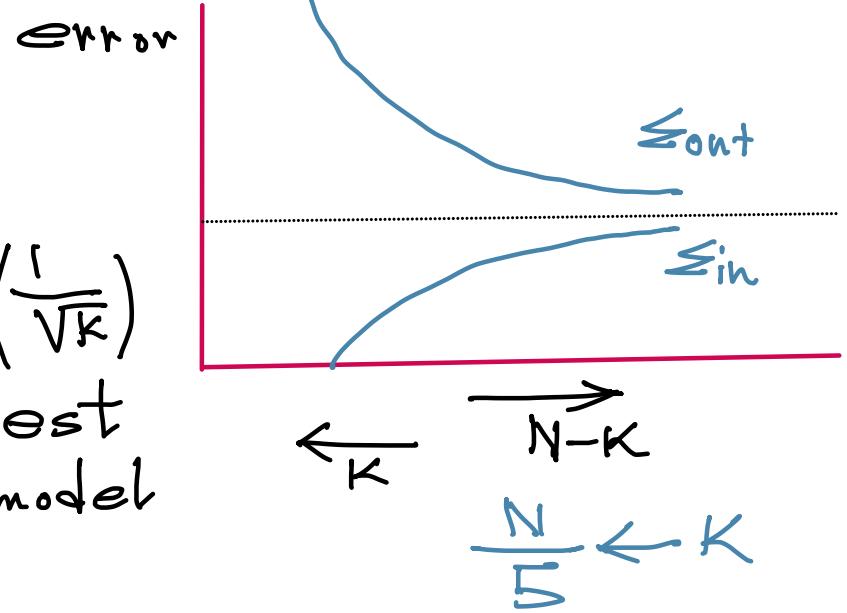
$$\text{Var} = \frac{1}{K^2} \sum_k \text{Var}(e(y(x_k), y)) = \frac{\sigma^2}{K}$$

$$D_{\text{Train}} \rightarrow N-K$$

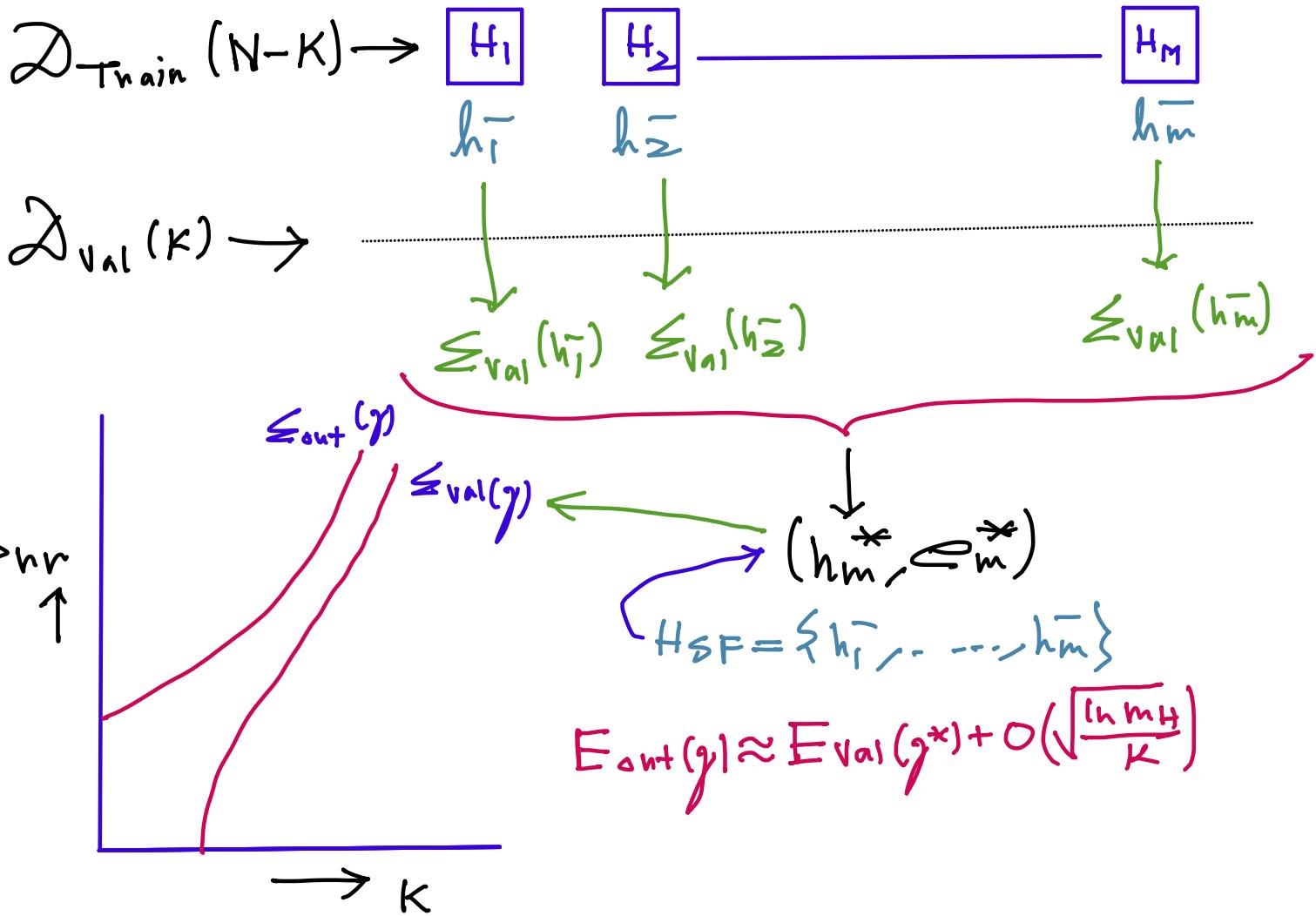
$$D_{\text{Val}} \rightarrow K$$

$$E_{\text{out}}(y) \approx E_{\text{val}}(y) \pm O\left(\frac{1}{\sqrt{K}}\right)$$

Small $K \rightarrow$ bad est
Large $K \rightarrow$ bad model



h_1, h_2
 $E_{\text{out}}(h_1) = E_{\text{out}}(h_2) = 0.5$
 $h_{\text{final}} = h_1 \text{ or } h_2$
 Consider $\epsilon = \min\{e_1, e_2\}$



$$\sum_{\text{out}}(y) \approx \sum_{\text{out}}(y^*) \approx \sum_{\text{val}}(y)$$

$\xleftarrow[\text{Small } K]$ $\xrightarrow[\text{Large } K]$

Leave One Out

$$(x_1, y_1) \dots (x_K, y_K) \dots (x_N, y_N)$$

$\underbrace{\quad \quad \quad}_{\text{Val}}$ $\underbrace{\quad \quad \quad}_{N-1}$

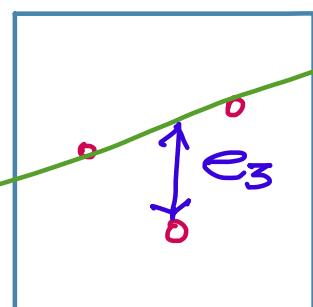
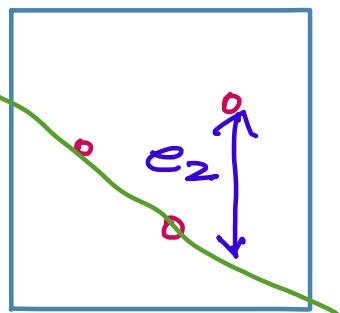
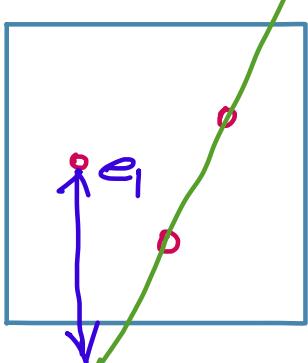
$N-1$ for train
 1 for Validation

$$E_{\text{val}} = \frac{1}{N} \sum_{n=1}^N E_{\text{val}}^n$$

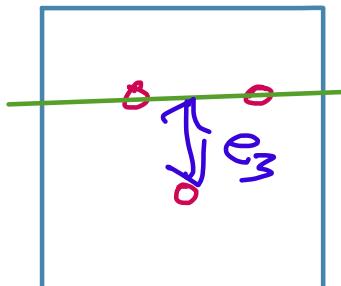
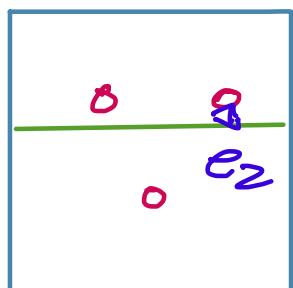
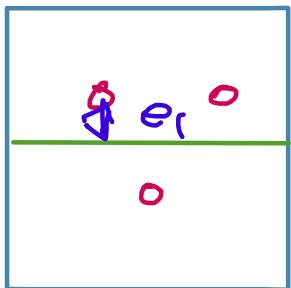
$$\sum_{\text{cv}}$$

$$H_1$$

$$\sum_{\text{cv}} = \frac{1}{3}(e_1 + e_2 + e_3)$$



H_2



$$\sum_{i=1}^3 \epsilon_i = \frac{1}{3} (\epsilon_1 + \epsilon_2 + \epsilon_3)$$

$\Rightarrow N$ can be large

$|P_1| |P_2| |P_3| \dots |D_K| \dots |P_n|$

Train on $\{D_j\} - D_K$
Validate only on D_K

$\underbrace{\qquad\qquad\qquad}_{K}$ n time

N/K Training Sessions // $(N-K)$ points to train

K-fold Cross Validation