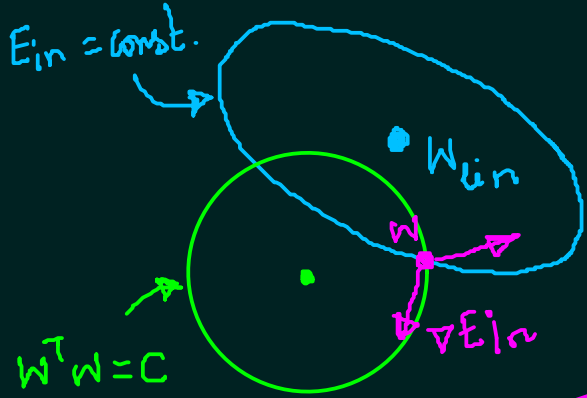


# Regularization:

unconstrained  $\rightarrow$  Constrained

**SUMMARY**

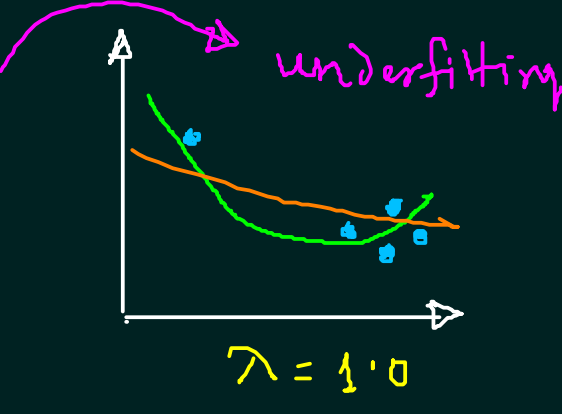
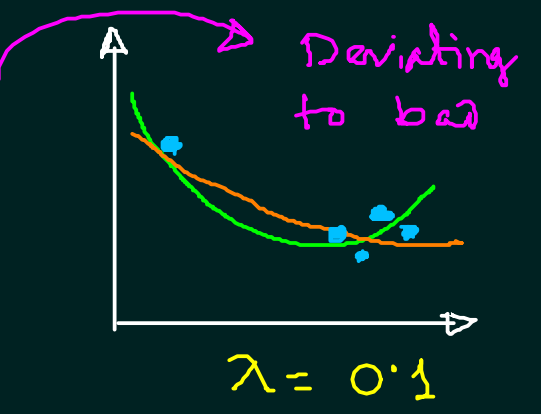
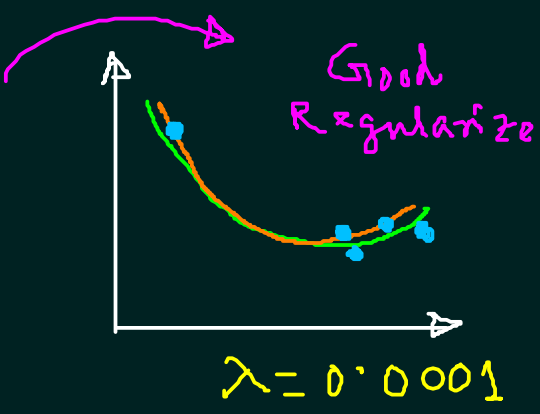
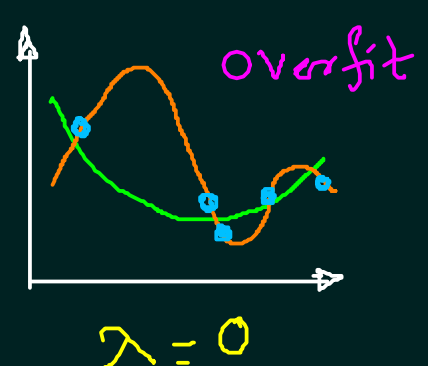
[Reduce overfitting Problem]



Minimize  $E_{in}(w)$   
 Minimize  $E_{aug}(w)$

$$E_{aug}(w) = E_{in}(w) + \frac{\lambda}{N} w^T w$$

$$E_{in}(w) \text{ subject to } w^T w \leq C$$



## Choosing a Regularizer:

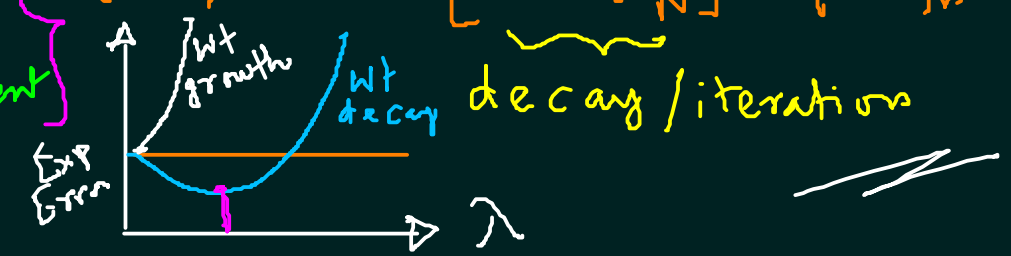
$$E_{aug}(h) = E_{in}(h) + \frac{\lambda}{N} \Omega(h)$$

$\Omega(h)$ : heuristic  $\rightarrow$  smooth, simple  $h$  functions

Popular: weight decay in stochastic gradient descent

$$w(t+1) \leftarrow w(t) \left[ 1 - 2\eta \frac{\lambda}{N} \right] - \eta \nabla E_{in}(w(t))$$

$\lambda$ : validation **TODAY**



$$E_{out}(h) = E_{in}(h) + \text{penalty}$$

→ overfit  
 → model  
 → noise  $e \rightarrow$  Stoch. / Det.

Regularization:

$$E_{avg}(h) = E_{in}(h) + \frac{\lambda}{N} \Omega(h)$$

Validation:  $E_{out}(h) \downarrow$

$$E_{in} \approx E_{out} \text{ (goal)}$$

▶ out-of-sample  $(x, y)$  → Error,  $e(h(x), y)$

$$E_{out}(h) = \mathbb{E}_x [e(h(x), y)]$$

$(h(x) - y)^2$   
 $\mathbb{I}[h(x) \neq y]$

$$\text{var}[e(h(x), y)] = \sigma^2 \quad (\text{ONE POINT})$$

▶ Set of points:  $(x_1, y_1) \dots (x_k, y_k)$

$$E_{val}(h) = \frac{1}{k} \sum_{k=1}^k e(h(x_k), y_k) \rightarrow \mathbb{E}[E_{val}(h)] = E_{out}(h)$$

$$E_{out}(h) = E_{in}(h) \pm O\left(\frac{1}{\sqrt{k}}\right) \quad \text{var}[E_{val}(h)] = \frac{1}{k^2} \sum_{k=1}^k \text{var}[e(h(x_k), y_k)] = \frac{\sigma^2}{k}$$

$$\mathcal{D} = (x_1, y_1) \dots (x_N, y_N)$$

$K$  point  $\Rightarrow \mathcal{D}_{val}$ ?  
 $(N-K)$  point  $\Rightarrow \mathcal{D}_{train}$  ✓

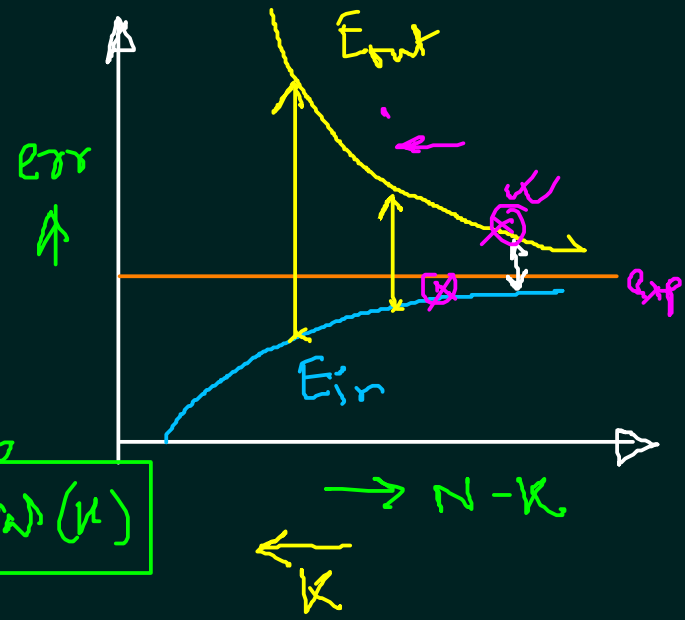
~~$$E_{out}(h) = E_{in}(h) \pm O\left(\frac{1}{\sqrt{K}}\right)$$~~ ✓

Small  $K \Rightarrow$  bad estimate

Large  $K \Rightarrow ?$

$$\mathcal{D} = \mathcal{D}_{train} \cup \mathcal{D}_{val}$$

$N \rightarrow N-K$  (train)  
 $N \rightarrow K$  (val)



$N$   $\mathcal{D}$

$\mathcal{D}_{train}^{(N-K)}$

$\mathcal{D}_{val}(K)$

Kislar

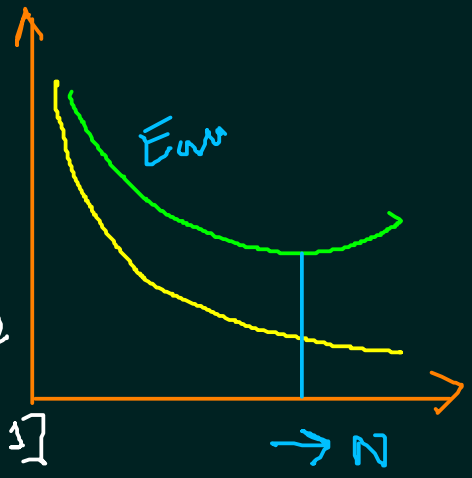
$g \neq g^-$   
 $g \rightarrow$  final

Eval?  $\rightarrow$   $g^-$   
 $E_{val}(g^-)$   
 $\uparrow$  est

$$K = \frac{N}{5} > \frac{N}{6}$$

$g^- \rightarrow E_{val}(g^-)$

(validation)



Optimistic bias

$$e = \min\{e_1, e_2\}$$

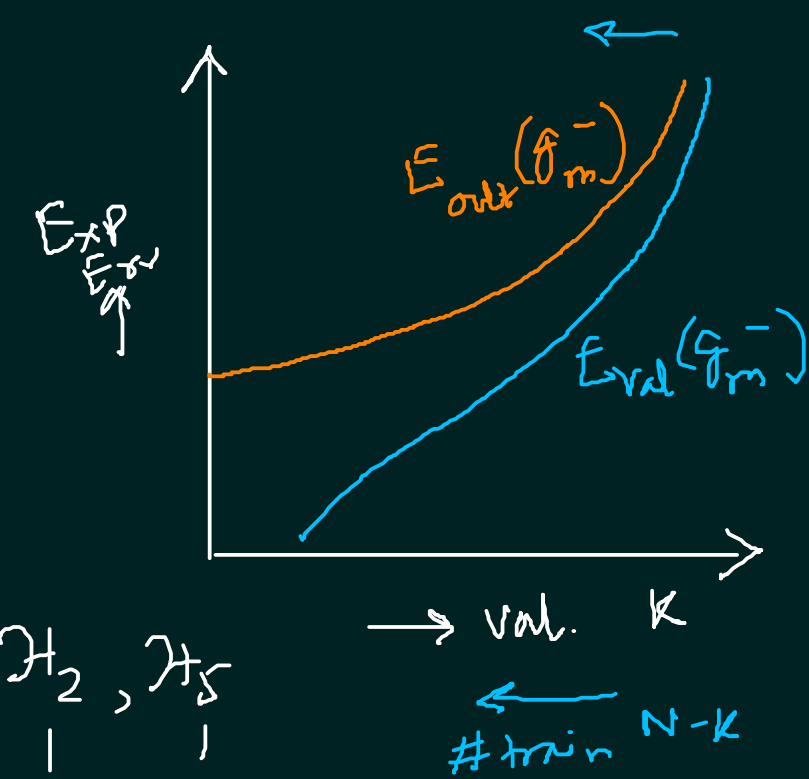
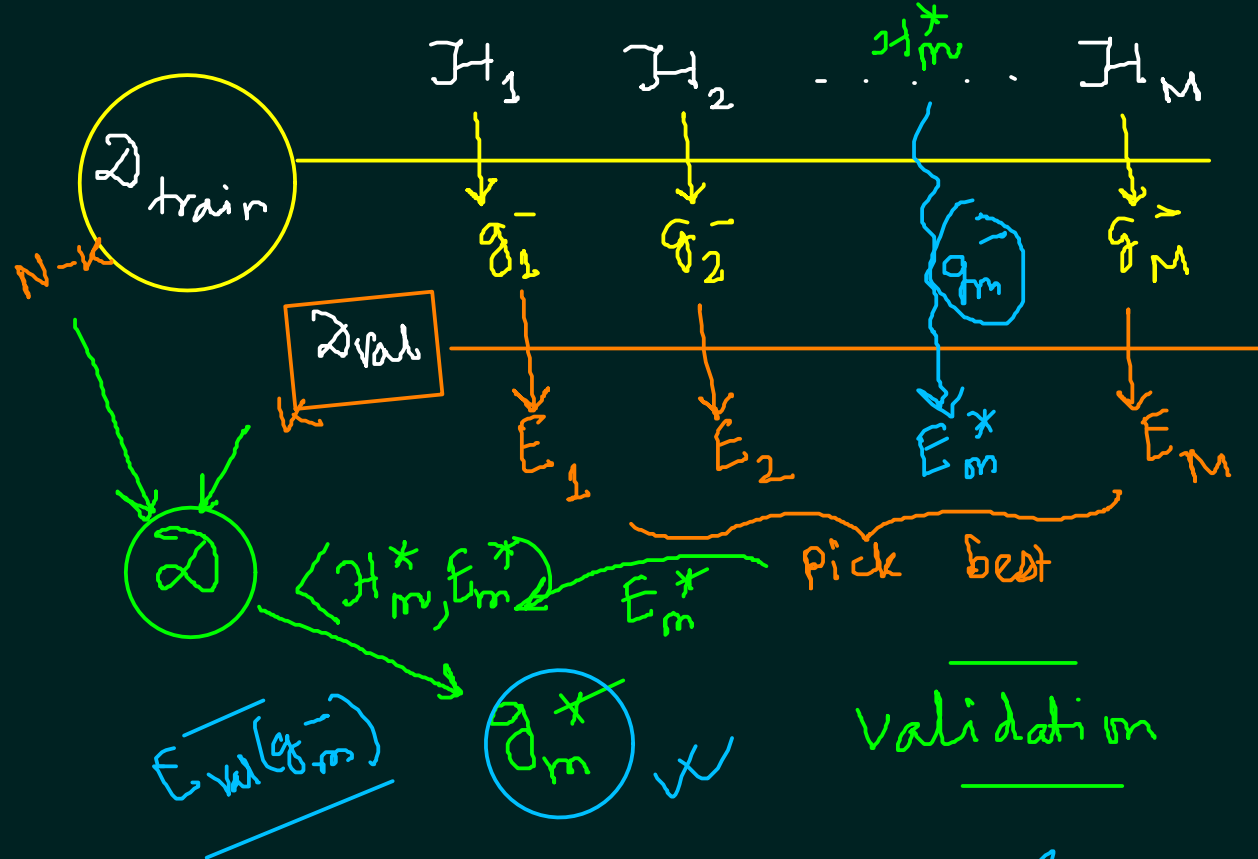
$$E_{out}(h_1) = E_{out}(h_2) = \frac{1}{2}$$

$$e_1 = e_2 \text{ Uniform } [0, 1]$$

$$\mathbb{E}(e) < 0.5$$

$K$  ✓

val



$$H_{val} = \{g_1^-, g_2^-, \dots, g_m^-\}$$

$$E_{out}(g_m^*) \leq E_{val}(g_m^*) + O\left(\sqrt{\frac{\ln M}{k}}\right)$$

$E_{in}$   
 $E_{val}$   
 $E_{out}$

with  $(1-\delta)$  prob

$$E_{out}(g) \approx E_{out}(g^-) \approx E_{val}(g^-) \left( \frac{1}{2k} \ln\left(\frac{2M}{\delta}\right) \right) \left( \frac{N}{k} \right)$$

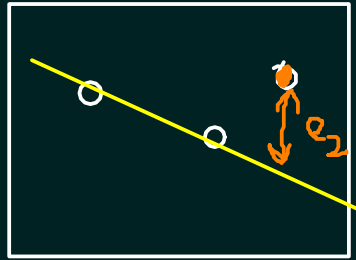
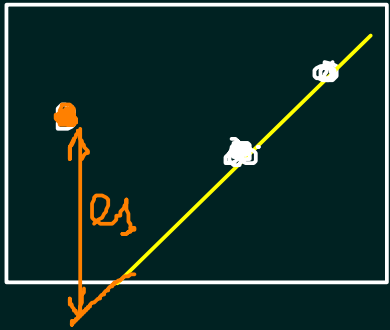
(small  $k$ ) (large  $k$ )

$$\mathcal{D}_n = (x_1, y_1) \text{ --- } (x_n, y_n) \text{ --- } (x_N, y_N)$$

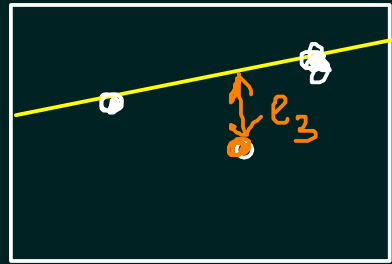
$N-1$   
for train

$$e_n = E_{\text{val}}(g_n^-) = e(g_n^-(x_n), y_n) \leftarrow \begin{matrix} \uparrow \\ \text{1 point} \\ \text{for validation} \end{matrix}$$

Cross validation error =  $\frac{1}{N} \sum_{n=1}^N e_n \leftarrow E_{\text{cv}}$



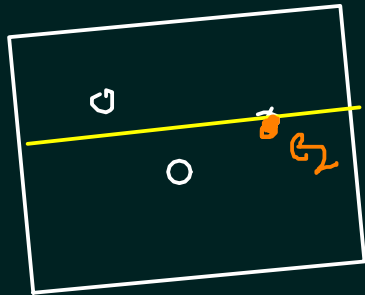
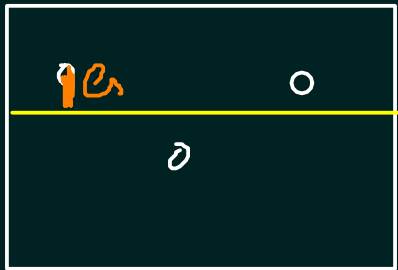
$H_1$



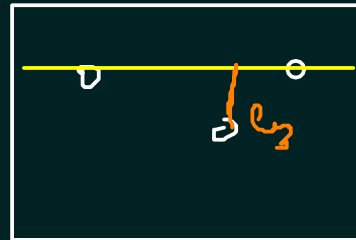
$$E_{\text{cv}} = \frac{1}{3} (e_1 + e_2 + e_3)$$

linear

1

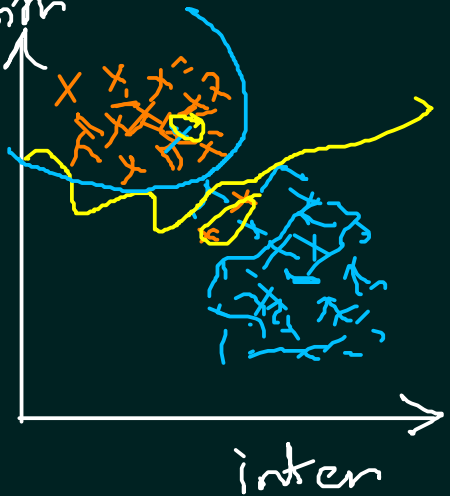


$H_2$



const

Symm



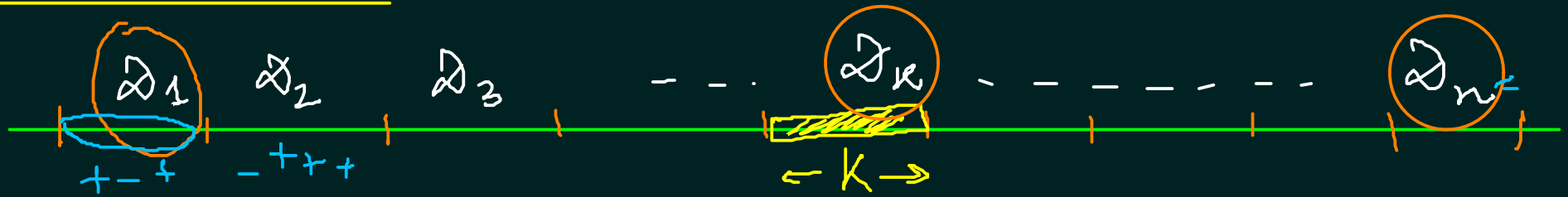
w/o val:  $E_{\text{in}} = 0\%$

$E_{\text{out}} = 2.5\%$

$$(1, x_1, x_2) \rightarrow (1, x_1, x_2, x_1^2, x_2^2, x_1^2 x_2, x_2^2, \dots)$$

w/ val:  $E_{\text{in}} = 0.8\%$ ,  $E_{\text{out}} = 1.5\%$

$N$  can be large



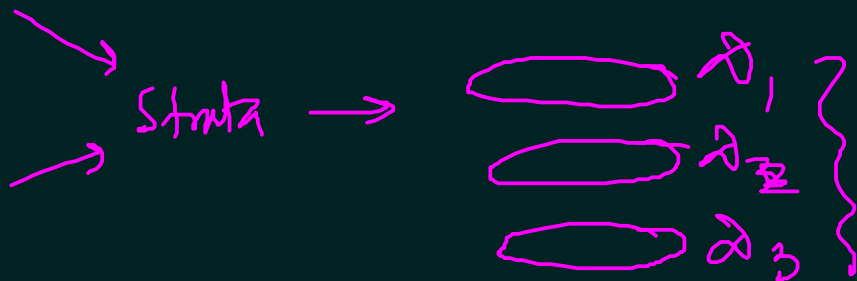
Train on  $\{x_i\} - x_k$   
Validate on only  $x_k$  } Repeat this  $n$  times

$\frac{N}{k}$  Training Sessions //  $(N-k)$  points to train

## $K$ -fold Cross-Validation

++++  
+++  
+++

-----  
-----  
-----



Stratified Approach ✓