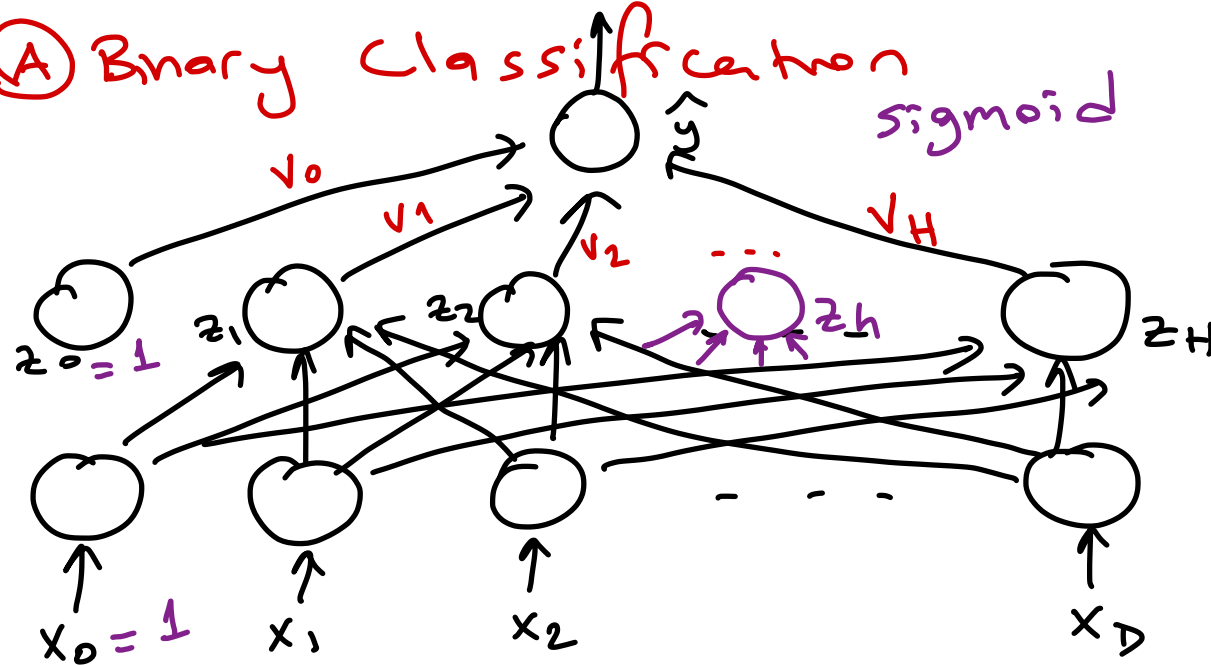


Multilayer Perceptrons

① Binary Classification



$$\hat{y}_i = \text{sigmoid}(v^T \cdot z_i)$$

$$z_{ih} = \text{sigmoid}(\underline{w_h^T} \cdot x_i)$$

\nwarrow weights for all incoming edges to z_h

\rightarrow binary cross-entropy

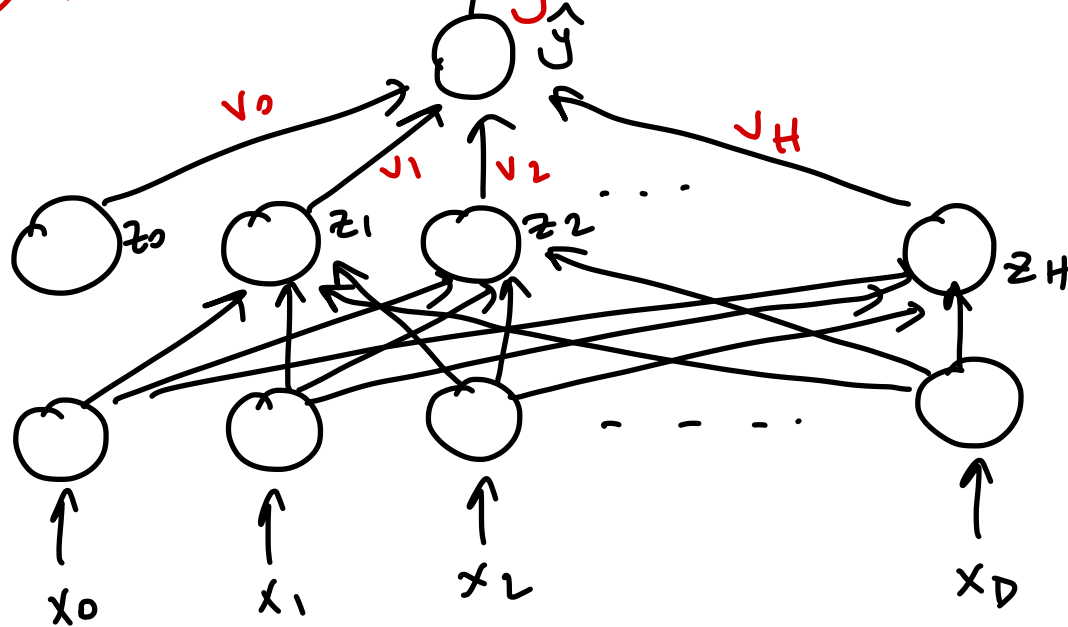
$$\text{Error}_i = - [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

$$\Delta v_h = \eta \cdot (y_i - \hat{y}_i) \cdot z_{ih}$$

$$\Delta w_{hd} = \eta (y_i - \hat{y}_i) v_h \cdot z_{ih} (1 - z_{ih}) \cdot x_{id}$$

y_i 's are either 0 or 1
 \hat{y}_i 's are between 0 and 1.

③ Nonlinear Regression



identity $\hat{y}_i = \underline{v^T \cdot z_i}$

sigmoid $z_{ih} = \text{sigmoid}(w_h^T \cdot x_i)$

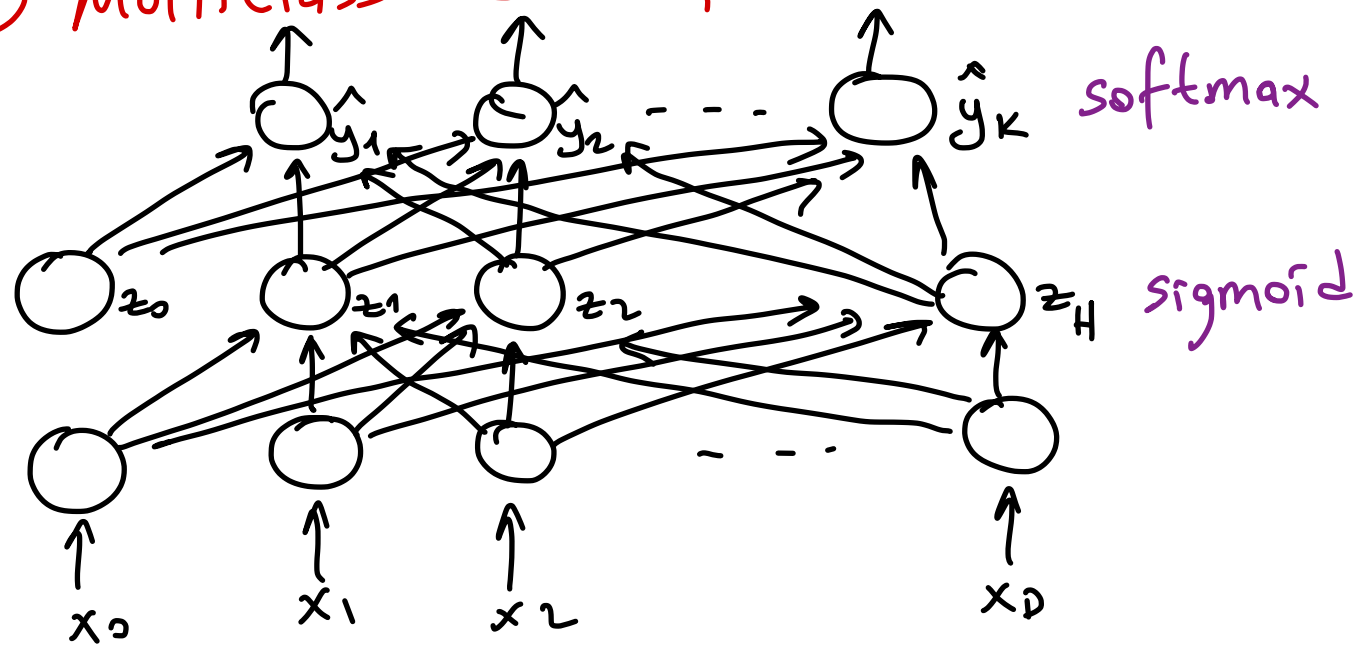
$$\text{Error}_i = \frac{1}{2} (y_i - \underline{\hat{y}_i})^2$$

$$\Delta v_h = \boxed{\eta \cdot (y_i - \hat{y}_i)} \cdot \boxed{z_{ih}}$$

$$\Delta w_{hd} = \boxed{\eta (y_i - \hat{y}_i)} \cdot v_h \cdot \boxed{z_{ih} (1 - z_{ih})} x_{id}$$

y_i 's and \hat{y}_i 's are real numbers.

③ Multiclass Classification



$$\hat{y}_{ic} = \frac{\exp[V_c^T \cdot z_i]}{\sum_{d=1}^K \exp[V_d^T \cdot z_i]}$$

$$z_{ih} = \text{sigmoid}(w_h^T \cdot x_i)$$

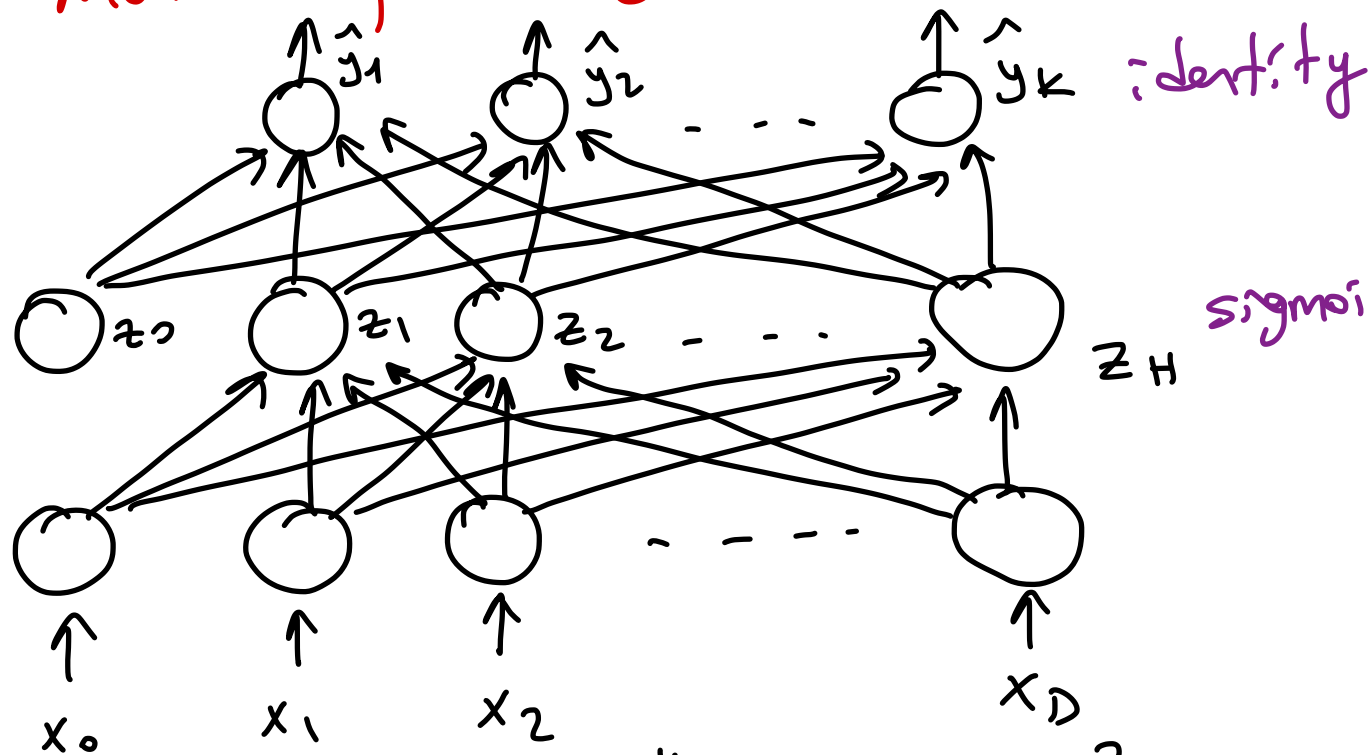
$$\text{Error}_i = - \sum_{c=1}^K y_{ic} \log(\hat{y}_{ic})$$

$$\Delta V_{ch} = \eta (y_{ic} - \hat{y}_{ic}) \cdot z_{ih}$$

$$\Delta W_{hd} = \eta \left[\sum_{c=1}^K (y_{ic} - \hat{y}_{ic}) \cdot V_{ch} \right] \cdot z_{ih} (1 - z_{ih}) \cdot x_{id}$$

y_{ic} 's are either 0 or 1.
 \hat{y}_{ic} 's are between 0 and 1.

① Multisoutput Regression



$$\hat{y}_{ic} = V_c^T \cdot z_i$$

$$z_{ih} = \text{sigmoid}(w_h^T \cdot x_i)$$

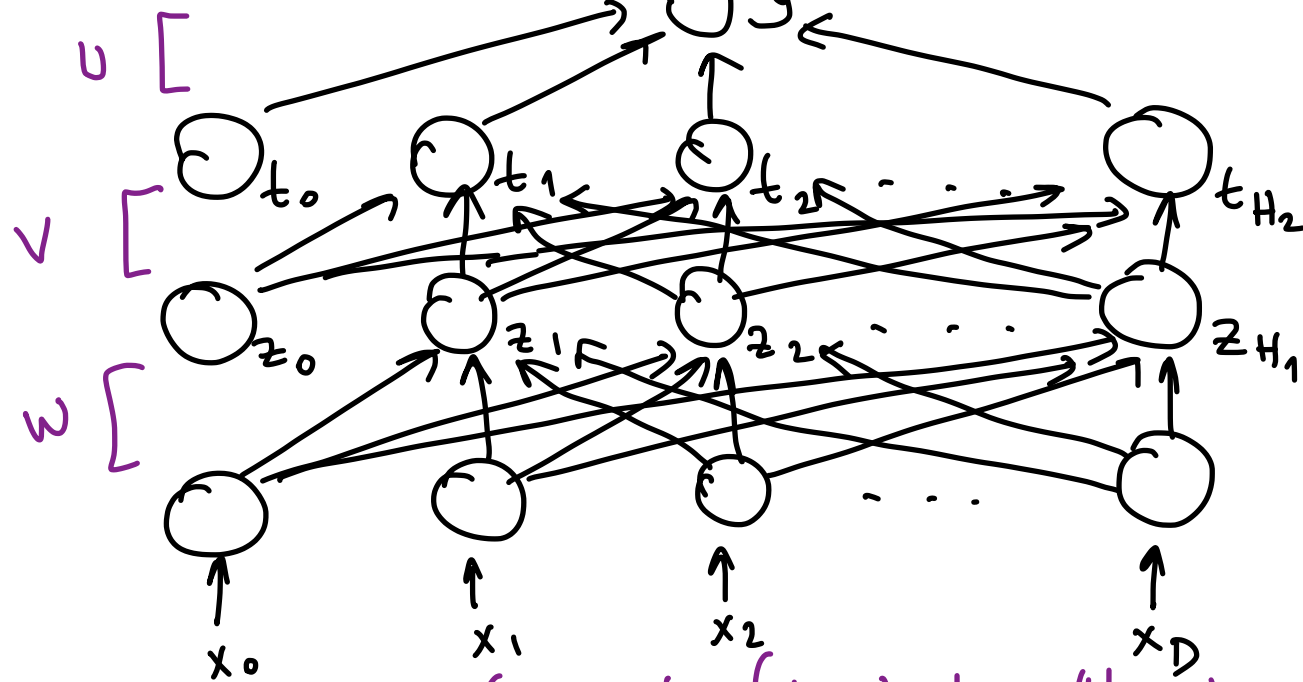
$$\text{Error}_i = \frac{1}{2} \sum_{c=1}^k (y_{ic} - \hat{y}_{ic})^2$$

$$\Delta V_{ch} = \eta (y_{ic} - \hat{y}_{ic}) \cdot z_{ih}$$

$$\Delta W_{hd} = \eta \left[\sum_{c=1}^k (y_{ic} - \hat{y}_{ic}) \cdot V_{ch} \right] \cdot z_{ih} (1 - z_{ih}) \cdot x_{id}$$

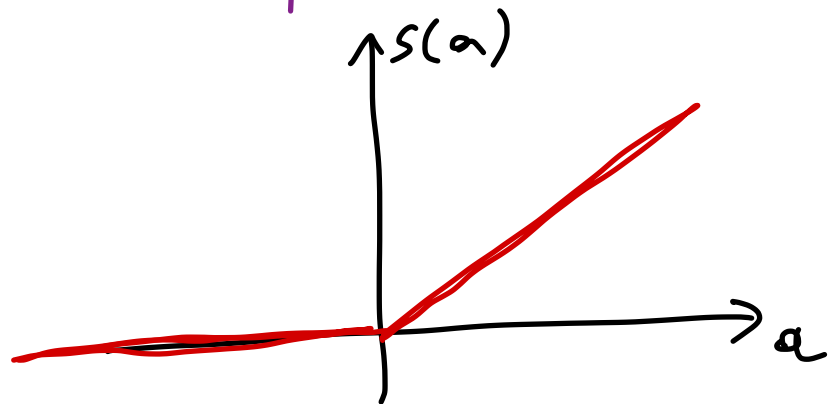
y_{ic} 's and \hat{y}_{ic} 's are real numbers.

Multiple Hidden Layers



$$\# \text{ of parameters} = \underbrace{(D+1) \times H_1}_W + \underbrace{(H_1+1) \times H_2}_V + \underbrace{(H_2+1)}_U$$

Rectified Linear Unit (ReLU)



$$s(a) = \begin{cases} a & \text{if } a > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$s'(a) = \begin{cases} 1 & \text{if } a > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\hat{y}_i = \text{sigmoid}(U^T \cdot t_i)$$

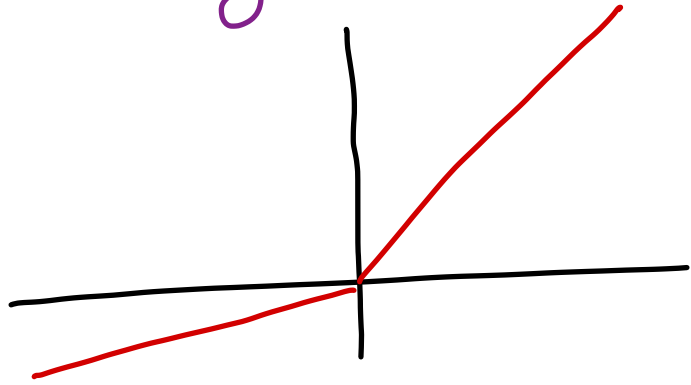
$$t_{ih} = \text{sigmoid}(V_h^T \cdot z_i)$$

$$z_{ih} = \text{sigmoid}(W_h^T \cdot x_i)$$

"vanishing gradients"

→ approaching to 0

Leaky ReLU



$$s(a) = \begin{cases} a & \text{if } a > 0 \\ \alpha a & \text{otherwise} \end{cases}$$

$$s'(a) = \begin{cases} 1 & \text{if } a > 0 \\ \alpha & \text{otherwise} \end{cases}$$

usually $\alpha = 0.01$

TRAINING PROCEDURES

Momentum:

$$s_h^{(t)} = \underbrace{\alpha s_h^{(t-1)}}_{\text{memory}} + (1-\alpha) \underbrace{\frac{\partial \text{Error}^{(t)}}{\partial w_h}}_{\text{current opinion}}$$

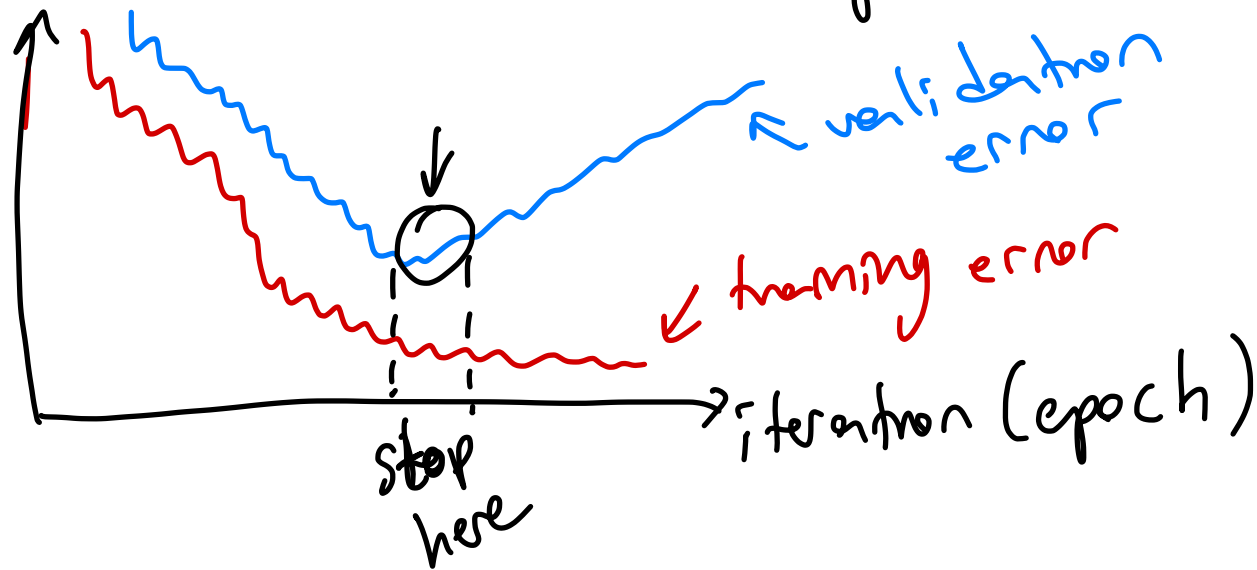
storage variable of the previous step *gradient of the current step.*

$$\Delta w_h^{(t)} = -\eta s_h^{(t)}$$

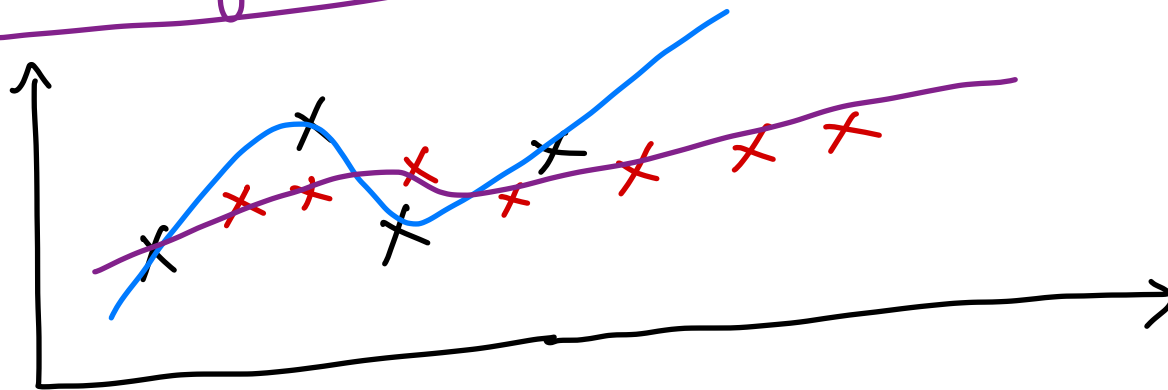
Adaptive Learning Rate: $\eta = ?$

→ starts with a large value
→ decrease if error increases

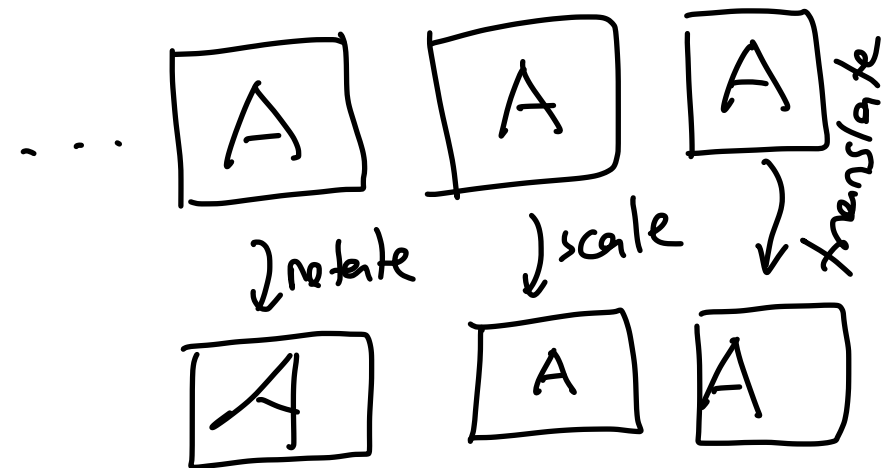
Early stopping: Stop if you think that your algorithm is overfitting.



Increasing Data Set Size:



Data augmentation



Weight Decay:

$$\text{Error}' = \text{Error} + \underbrace{\frac{\lambda}{2} \sum_{h=1}^H w_h^2}_{\text{weight decay}}$$

$$\frac{\partial \text{Error}'}{\partial w_h} = \frac{\partial \text{Error}}{\partial w_h} + \lambda \cdot w_h$$

ℓ_2 -norm
regularization
 $\sum_{h=1}^H w_h^2 = \|w\|_2^2$

$$\Delta w_h = -\eta \left[\frac{\partial \text{Error}}{\partial w_h} + \lambda w_h \right]$$