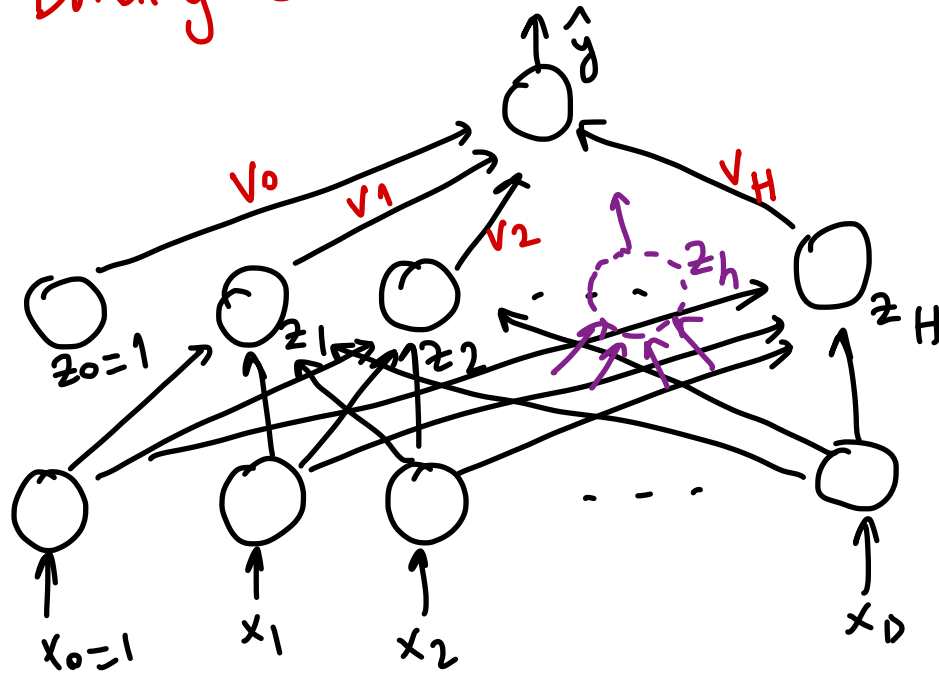


Multilayer Perceptrons

(A) Binary Classification



sigmoid
sigmoid

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

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

weights for all incoming edges to z_h

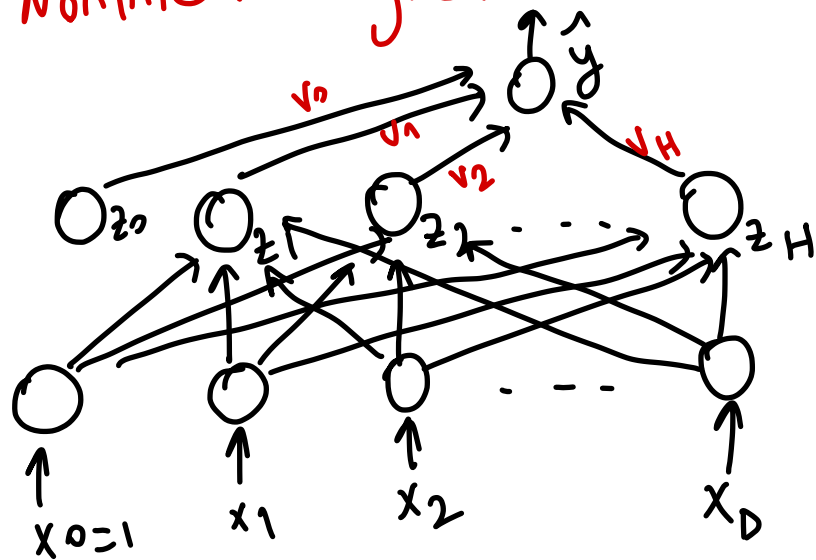
$$\text{Error}_i = - [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] \rightarrow \text{binary cross-entropy}$$

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

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

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

⑧ Nonlinear Regression



linear

sigmoid

$$\hat{y}_i = v^T \cdot z_i \Leftarrow \hat{y}_i \in \mathbb{R}$$

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

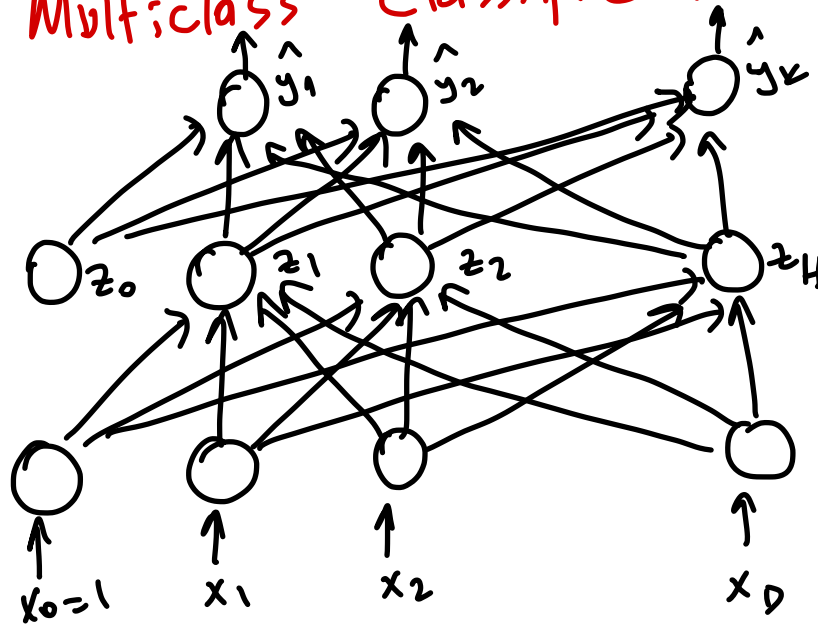
$$\text{Error}_i = \frac{1}{2} (y_i - \hat{y}_i)^2$$

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

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

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

(C) Multiclass Classification



softmax

sigmoid

$$\hat{y}_{ic} = \frac{\exp[v_c^T \cdot z_i]}{\sum_{d=1}^k \exp[v_d^T \cdot z_i]}$$

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

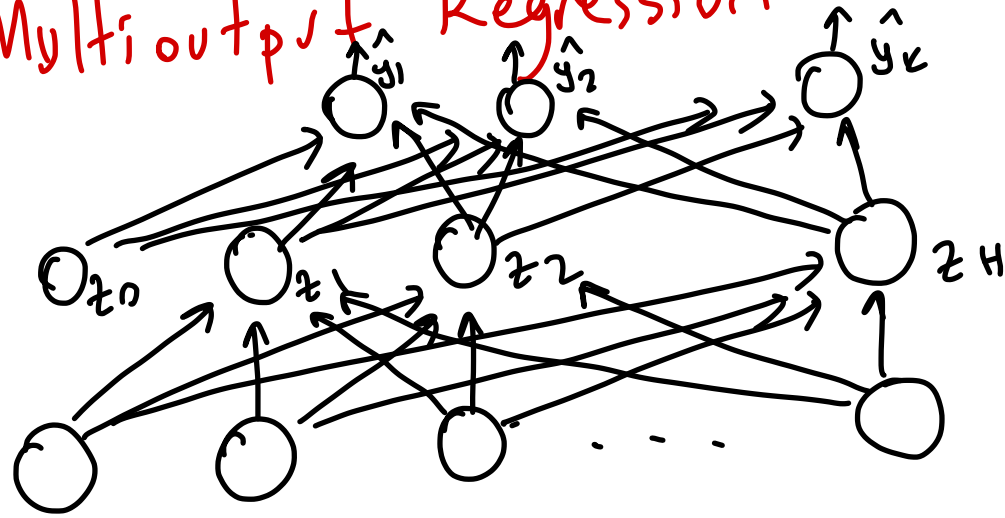
$$\text{Error}_i = - \sum_{c=1}^k y_{ic} \log(\hat{y}_{ic})$$

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

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

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

① Multioutput Regression



linear

sigmoid

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

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

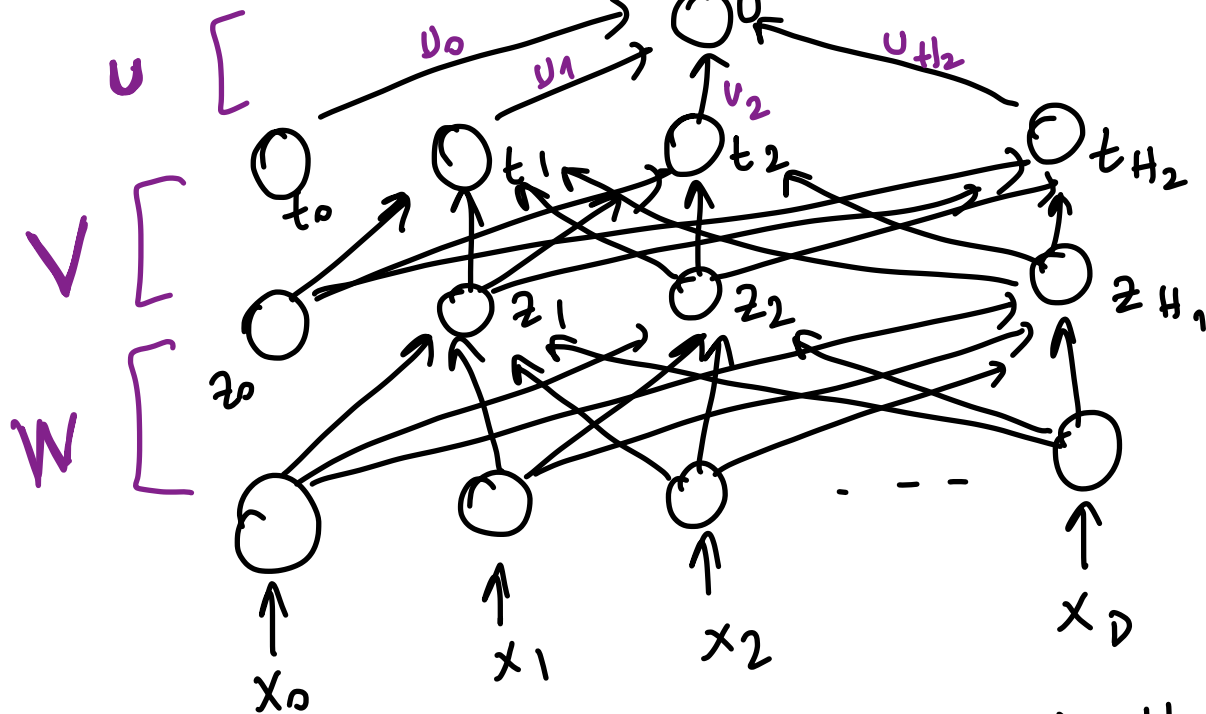
$$\text{Error}_i = \frac{1}{2} \sum_{c=1}^k (y_{ic} - \hat{y}_{ic})^2$$

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

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

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

Multiple Hidden Layers



sigmoid

sigmoid

sigmoid

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

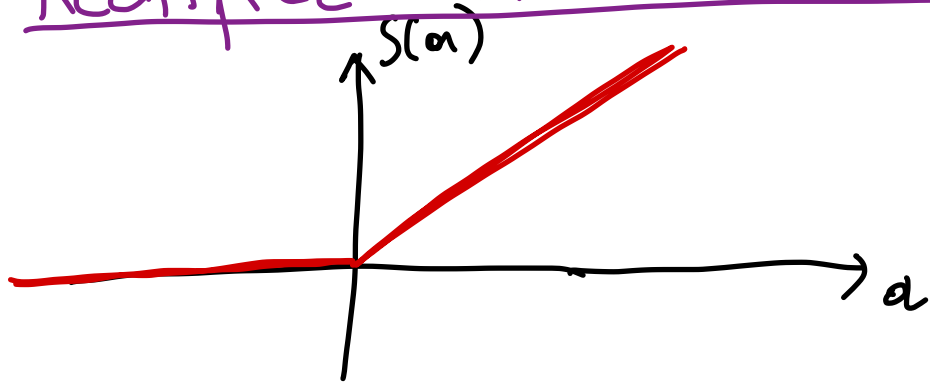
$$t_{ih} = \text{sigmoid}(v_h^T \cdot z_i)$$

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

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

"vanishing gradients" \rightarrow approaching to 0.

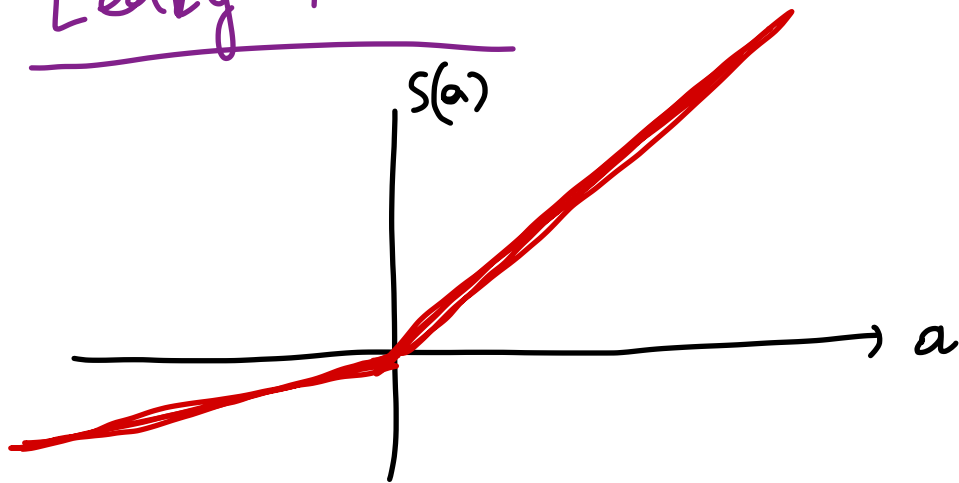
Rectified Linear Unit (ReLU)



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

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

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:

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

(state) variable of the previous step

gradient of the current step

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

Adaptive Learning Rate:

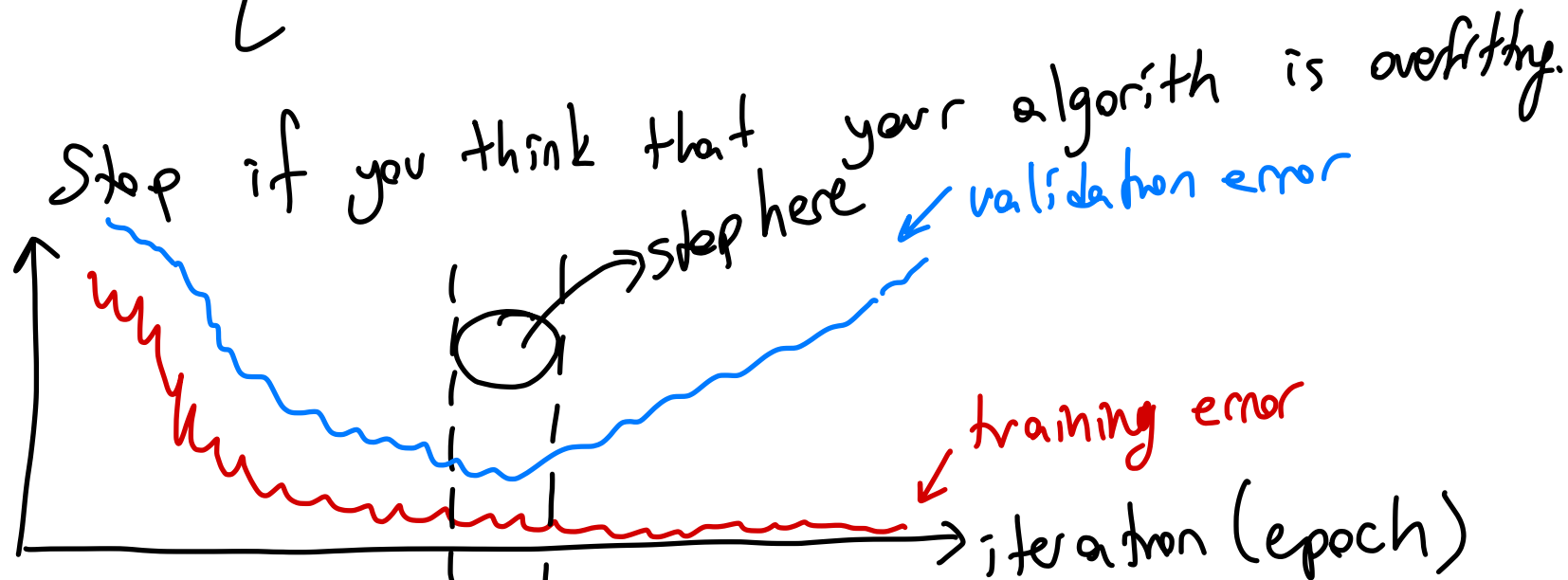
$$\eta^{(t)} = \eta^{(t-1)} \cdot 0.99$$

$\eta = ?$

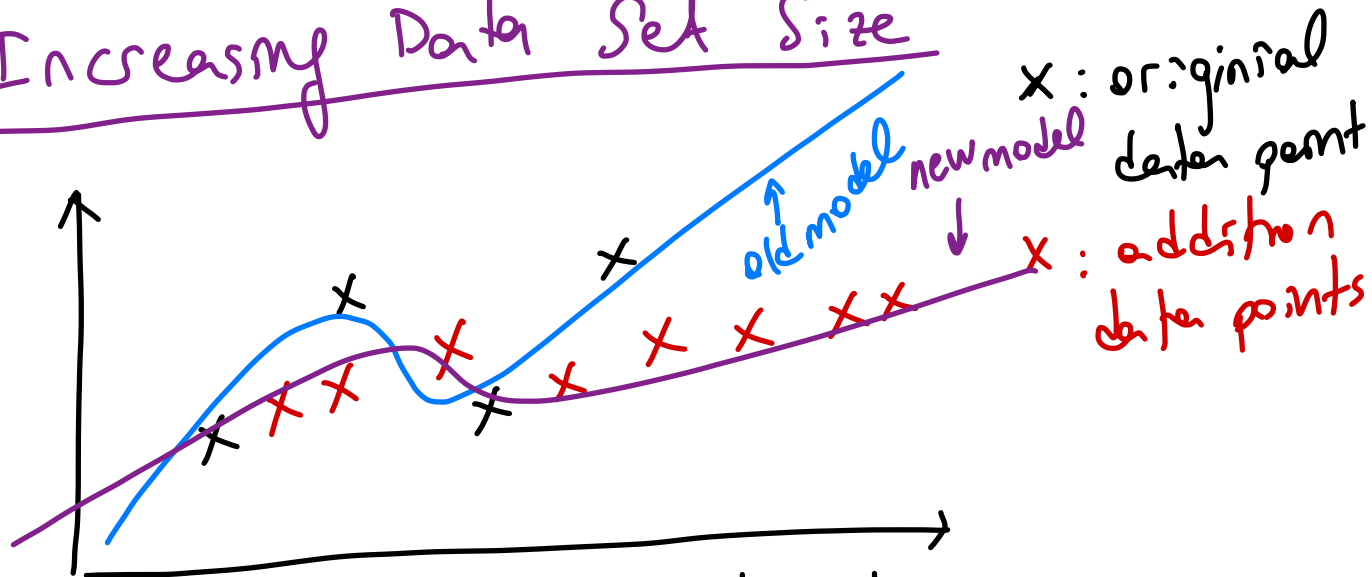
starts with a large value

decrease if error increases

Early Stopping:



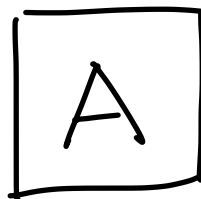
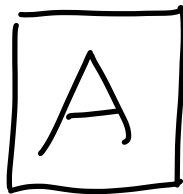
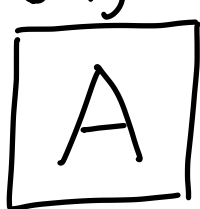
Increasing Data Set Size



test data points



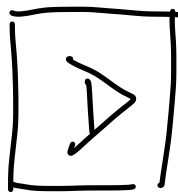
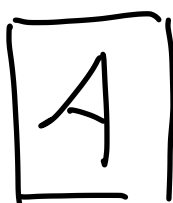
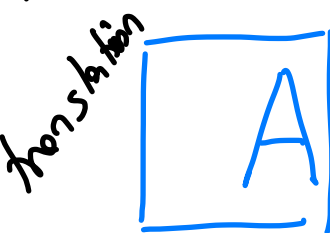
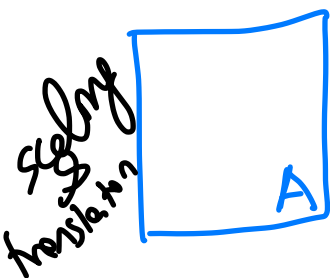
original pictures



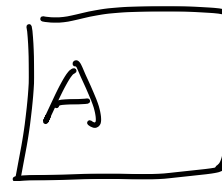
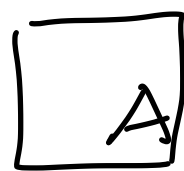
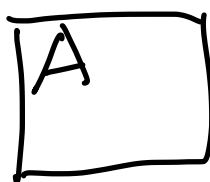
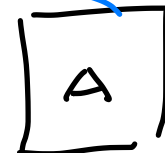
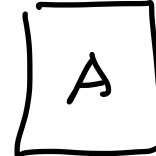
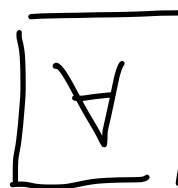
...

data augmentation

add these to the training set.



...



Weight Decay:

$$Error' = Error + \frac{\lambda}{2} \cdot \sum_{h=1}^H w_h^2$$

weight decay

l_2 -norm regularization

$$\frac{\partial Error'}{\partial w_h} = \frac{\partial Error}{\partial w_h} + \frac{\lambda}{2} \cdot 2 \cdot w_h$$

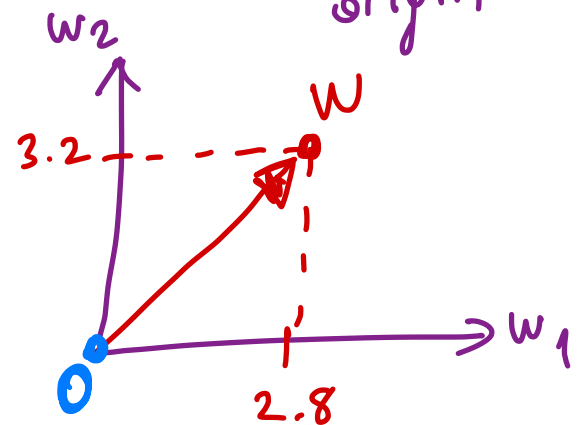
$$= \frac{\partial Error}{\partial w_h} + \lambda \cdot w_h$$

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

$$\sum_{h=1}^H w_h^2 = \|w\|_2^2$$

distance to the origin.

$$w = \begin{bmatrix} 2.8 \\ 3.2 \end{bmatrix}$$



Euclidean distance between w and O .

$$w_1^2 + w_2^2 = \sqrt{(w_1 - 0)^2 + (w_2 - 0)^2}^2$$