

1<sup>st</sup> Forward Propagation

$$net_{h1} = w_1 i_1 + w_2 i_2 + b_1 = 0.15 \times 0.05 + 0.2 \times 0.1 + 0.35$$

$$net_{h1} = 0.3775$$

$$outh_1 = \text{Sigmoid function} = \frac{1}{1 + e^{-0.3775}} =$$

$$outh_1 = 0.59326992$$

Similarly,

$$outh_2 = 0.596884370$$

$$net_{o1} = w_5 \times outh_1 + w_6 \times outh_2 = 0.4 \times 0.59326992 + 0.45 \times 0.596884370$$

$$net_{o1} = 1.105967$$

$$Output_{o1} = \frac{1}{1 + e^{-1.105967}} = 0.75136507$$

$$Output_{o2} = 0.772928465$$

### Backward Propagation

$$E_{01} = \frac{1}{2} (\text{target}_{01} - \text{output}_{01})^2$$
$$= \frac{1}{2} (0.01 - 0.7513)^2$$

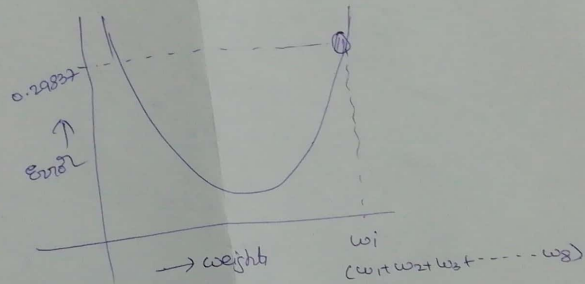
$$E_{01} = \frac{0.274811083}{0.023560026}$$

Similarly

$$E_{02} = 0.023560026$$

$$\text{Total Error } E_{\text{total}} = E_{01} + E_{02}$$

$$E_{\text{total}} = 0.298371109$$





$$w_5 = w_5 - \left( \frac{\partial E}{\partial u_5} \right) \times \eta \quad \text{(learning rule)}$$

$$\frac{\partial E_{\text{total}}}{\partial u_5} = \frac{\partial E_{\text{total}}}{\partial \text{out}_1} \times \frac{\partial \text{out}_1}{\partial \text{net}_1} \times \frac{\partial \text{net}_1}{\partial u_5}$$

$$\begin{aligned} \frac{\partial E_{\text{total}}}{\partial \text{out}_1} &= \frac{\partial}{\partial \text{out}_1} \left[ \frac{1}{2} (\text{target}_1 - \text{out}_1)^2 + \frac{1}{2} (\text{target}_2 - \text{out}_2)^2 \right] \\ &= \frac{\partial}{\partial \text{out}_1} \left[ \frac{1}{2} \times 2 (\text{target}_1 - \text{out}_1) \right] = \frac{\partial \text{out}_1}{\partial \text{out}_1} \\ &= (\text{target}_1 - \text{out}_1) \times -1 \\ &= \text{out}_1 - \text{target}_1 = 0.7513 - 0.01 \\ &= 0.7413 \end{aligned}$$

$$\boxed{\frac{\partial E_{\text{total}}}{\partial \text{out}_1} = 0.7413}$$

$$\frac{\partial \text{out}_1}{\partial \text{net}_1} = \frac{\partial}{\partial \text{net}_1} \quad \text{out}_1$$

$$= \frac{\partial}{\partial \text{net}_1} \frac{1}{1 + e^{-\text{net}_1}} =$$

$$= \text{out}_1 (1 - \text{out}_1) = 0.75(1 - 0.75)$$

$$= \boxed{0.186815662}$$

$$E_{\text{total}} = \frac{1}{2} [\text{target}_1 - \text{out}_1]^2 + \frac{1}{2} [\text{target}_2 - \text{out}_2]^2$$

$$\frac{\partial \text{net}_1}{\partial u_5} = w_5 \times \text{out}_4 + w_6 \times \text{out}_2 + 1 \times b_1$$

$$= \text{out}_4$$

$$\boxed{\frac{\partial \text{net}_1}{\partial u_5} = 0.59326992}$$

$$\begin{aligned} w_5^* &= w_5 - \left( \frac{\partial E_{\text{total}}}{\partial \text{out}_1} \times \frac{\partial \text{out}_1}{\partial \text{net}_1} \times \frac{\partial \text{net}_1}{\partial u_5} \right) \times \eta \\ &= 0.4 - (0.7413 \times 0.186815662 \times 0.59326992) \times 0.5 \end{aligned}$$

$$\boxed{w_5^* = 0.35892009}$$

→ activation function (Sigmoid, Tanh)  

$$\frac{e^{-x}}{1+e^{-x}}$$
  

$$\frac{e^x - e^{-x}}{e^x + e^{-x}}$$
  

$$1 - \tanh(x)$$
  

$$0 \leq \tanh(x) \leq 1$$
  
 ReLU,  $\begin{cases} 0 & x \leq 0 \\ x & x > 0 \end{cases}$   
 Rectified Linear Unit  

$$\begin{cases} 0 & x \leq 0 \\ x & x > 0 \end{cases}$$
  
 Softplus  

$$\begin{cases} 0 & x \leq 0 \\ 1 & x > 0 \end{cases}$$

→ Chain rule (Training MLP)

→ Back Propagation

→ Varying gradients

→ Hyper Parameter

→ Dropout (Random subset of features)

→ Regularization

→ Batch normalization

→ weight initialization (Uniform initialization)

(Sigmoid)

$$\left[ \frac{1}{\sqrt{f_{\text{norm}}}}, \frac{1}{\sqrt{f_{\text{norm}}}} \right], \left[ \frac{-\sqrt{6}}{\sqrt{f_{\text{norm}} + f_{\text{out}}}}, \frac{+\sqrt{6}}{\sqrt{f_{\text{norm}} + f_{\text{out}}}} \right]$$

(Xavier/Glorot)

He-init

$$\frac{\text{norm}}{\sqrt{f_{\text{norm}}}}, \text{norm} = \sqrt{\frac{2}{f_{\text{norm}}}}$$

$$\text{unif} = \sqrt{\frac{-\sqrt{6}}{f_{\text{norm}}}}, \sqrt{\frac{+\sqrt{6}}{f_{\text{norm}}}}$$

→ optimizers

↳ SGD, Adam, Adagrad

$$w_t = w_{t-1} - \eta \cdot g_t$$

$$\eta_t = \frac{\eta}{\sqrt{1 + \sum_{s=1}^t g_s^2}}$$

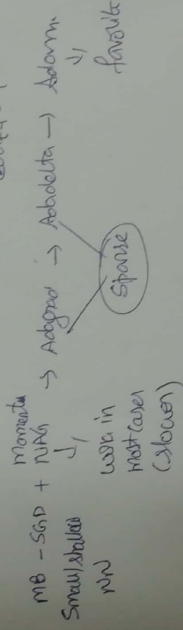
$$\eta_{t-1} = \frac{\sum_{i=1}^t g_i}{\sqrt{\sum_{i=1}^t g_i^2}} \rightarrow \left( \frac{\partial L}{\partial w} \right)_{(t-1)}$$

↳ Adadelta

$$\eta_t = \frac{\eta}{\sqrt{\sum_{s=1}^t g_s^2}}$$

$$\eta_t = \frac{\eta}{\sqrt{\sum_{s=1}^t g_s^2 + \epsilon}}$$

$$\text{adam}_t = \eta \cdot \text{adam}_{t-1} + (1 - \eta) g_t$$



→ loss (from loss function)

↳ mean-square error, mean square logarithm, mean Absolute Error

↳ hinge loss, square hinge

↳ categorical - cross

↳ binary - cross

$(1/(1+e^{-z})) \frac{e^z - e^{-z}}{e^z + e^{-z}} \xrightarrow{\text{d/dx}} 1 - \tanh^2(z)$   
 $0 \leq \frac{\tanh}{dz} \leq 1$ ,  $-1 \leq$   
 → activation function (Sigmoid, Tanh, ReLU,  $\begin{cases} 0 & y \leq 0 \\ 1 & y > 0 \end{cases}$ ,  $\begin{cases} 0 & y \leq 0 \\ 1 & y > 0 \end{cases}$ )  
 → memorisation  
 → chain rule (Training MLP)  
 → Back Propagation  
 → Vanishing gradient  
 → Hyper Parameter  
 → Dropout (Random subset of features)  
 → Regularization  
 → Batch normalization  
 → weight initialization (Sigmoid)  $\begin{bmatrix} \frac{1}{\sqrt{f_{\min}}} & \frac{1}{\sqrt{f_{\max}}} \end{bmatrix}$ , (Sigmoid) Xavier/Glorot  $\begin{bmatrix} -\sqrt{6} & \sqrt{6} \end{bmatrix}$ , He-Init  $\begin{bmatrix} -\sqrt{2} & \sqrt{2} \end{bmatrix}$ , uniform  $\begin{bmatrix} -\sqrt{6} & \sqrt{6} \end{bmatrix}$

→ optimizers

↳ SGD, AdaGrad

$$w_t = w_{t-1} - \eta^t g_t$$

$$\eta^t = \frac{\eta}{\sqrt{1 + \sum_{i=1}^t g_i^2}} \rightarrow \left( \frac{\partial L}{\partial w} \right)_{(t-1)}$$

$$\eta^t = \frac{\eta}{\sqrt{1 + \sum_{i=1}^t g_i^2}} \rightarrow \left( \frac{\partial L}{\partial w} \right)_{(t-1)}$$

↳ Adadelta

$$\eta^t = \frac{\eta}{\sqrt{1 + \sum_{i=1}^t g_i^2}}$$

$$e_{t+1} = \gamma e_{t+2} + (1-\gamma) g_{t+1}$$

↳ Adam

MB-SGD + Momentum  
Small/shallow NN  
works in most cases (stochastic)

→ Adagrad → Adadelta → Adam

Sparse

favourite

→ loss (from keras input losses)

- ↳ mean-square error, mean square logarithmic, mean Absolute error loss
- ↳ hinge loss, square hinge
- ↳ Categorical-Crossent
- ↳ binary-Cross