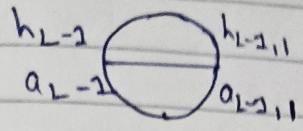


1

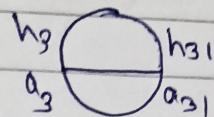
$$g = \text{softmax}(c_L)$$

$g \in \mathbb{R}^{K \times 1}$   
Pogelwerk  
 $K$  neurons.

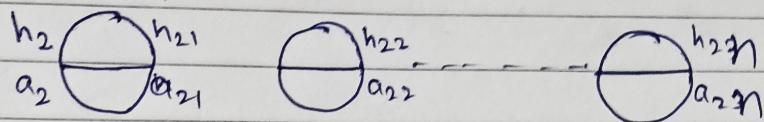
$w_L, b_L$



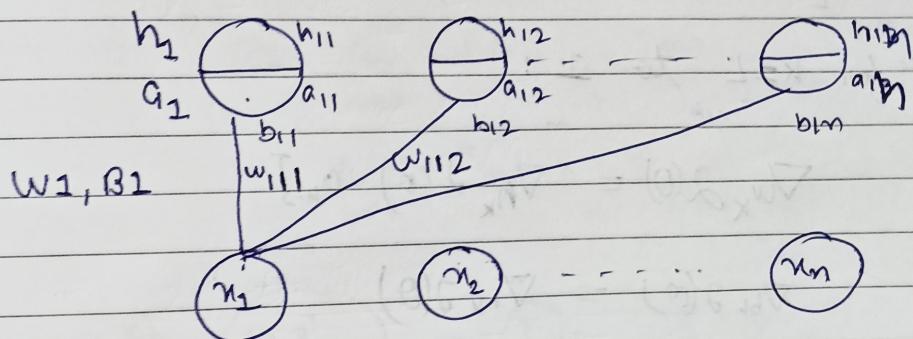
$w_{L-2}, b_{L-2}$



$w_3, b_3$



$w_2, b_2$



$w_1$  to  $w_{L-2} \equiv (n \times n)$  matrices having weights

$w_L = n \times K$  matrix

$b_1$  to  $b_{L-1} \equiv (1 \times n)$  vector having bias.

$b_L = 1 \times K$  vector.

$S_1$

$S \times 2$

$n \times 2 \times n \times 2$

$S \times 2$

$4 \times n \times n \times 2 \times n \times 2 +$

← (2) →

$$W_1 = \begin{bmatrix} w_{111} & w_{112} & w_{113} & \dots & w_{11n} \\ w_{121} & w_{122} & w_{123} & \dots & w_{12n} \\ w_{131} & w_{132} & w_{133} & \dots & w_{13n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ w_{1m1} & w_{1m2} & w_{1m3} & \dots & w_{1mn} \end{bmatrix}$$

$n \times n \equiv R^{n \times n}$

$$B_1 = \begin{bmatrix} b_{11} & b_{12} & b_{13} & \dots & b_{1n} \end{bmatrix}$$

$1 \times n \equiv R^n$

$a_1 = W_1 \cdot x^T + B_1^T$

$$x \in R^n \equiv [x_1 \ x_2 \ \dots \ x_n] \quad 1 \times n$$

We need to have,  
 $a_1$  same size as  $x$

$$\frac{n \times n}{R} \odot \frac{n \times 1}{R} + \frac{n \times 1}{R} = \frac{n \times 1}{R}$$

$$\Rightarrow a_1 \in R^{n \times 1}$$

$$h_1 = \text{sigmoid}(a_1) \text{ or } \text{act-f}^n(a_1)$$

$$h_1 \in R^{n \times 1} \equiv [h_{11} \ h_{12} \ \dots \ h_{1n}]^T$$

$$a_2 = W_2 \cdot h_1 + B_2^T$$

$$h_2 = \text{act-f}^n(a_2).$$

← ③

PageWork

~~Statement~~

(#) Weights & Bias Initialization

(\*) Initialise  $w_1, w_2, \dots, w_{L-1} \in \mathbb{R}^{n \times n}$

$w_L \in \mathbb{R}^{n \times k}$

$b_1, b_2, \dots, b_{L-1} \in \mathbb{R}^{1 \times n}$

$b_L \in \mathbb{R}^{1 \times k}$

$x \in \mathbb{R}^n$

(#) Forward propagation

$X = \text{list of } m \text{ data points.}$

for  $x$  in  $X$ :

$x \in \mathbb{R}^{1 \times n}$

(#) Input layer

$$a_1 = w_1 \cdot x^T + b_1^T$$

$$h_1 = \text{act-f}^n(a_1)$$

(#) Middle layers

$k=2 \text{ to } k=L-1$

$$a_k = w_k \cdot h_{k-1} + b_k^T$$

$$h_k = \text{act-f}^n(a_k)$$

(#) Output layer

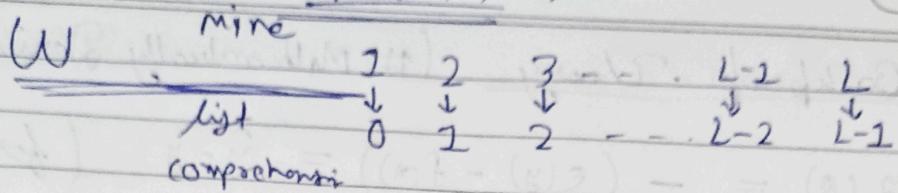
$$a_L = w_L^T \cdot h_{L-1} + b_L^T$$

$$\hat{y} = \text{Output-f}^n(a_L)$$

B	D	M	M	T	T	F

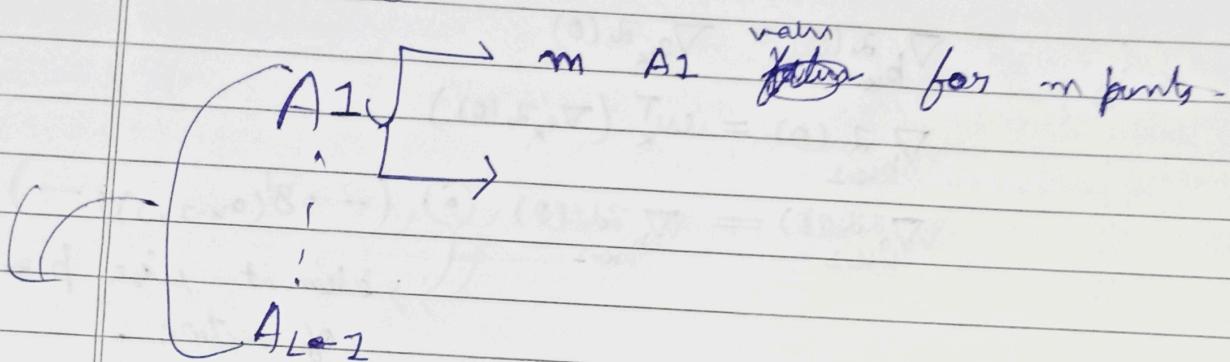
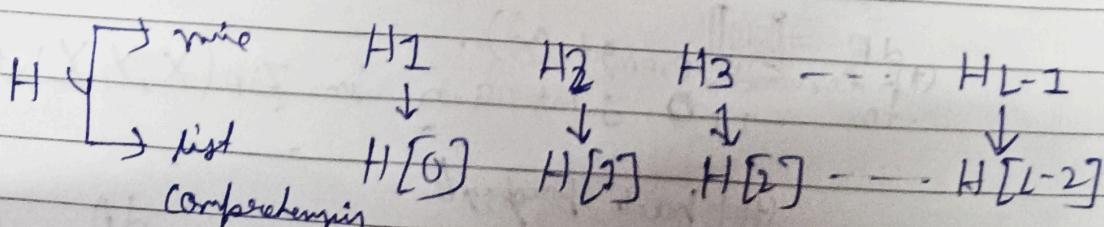
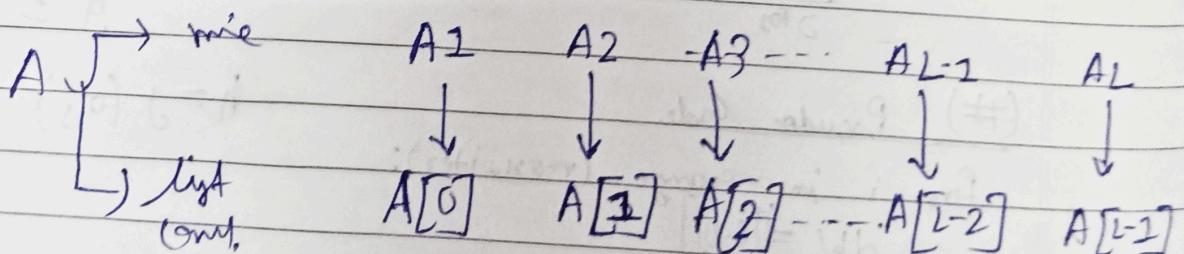
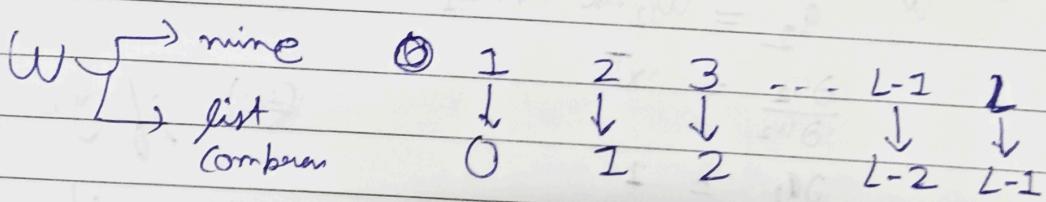
← 4

rough

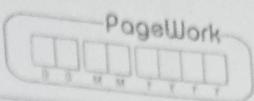


X → m points.

Y → m outputs for m points

A<sub>L</sub> = k A<sub>1</sub> ...H<sub>1</sub> → H<sub>L-1</sub> → m point.

← (5)



(#) Backprop Theory (# Mathematically Derived).

$$\nabla_{\theta} \mathcal{L}(\theta) = - (e(y) - f_m)) = \hat{y} - y \quad (\text{for softmax } \hat{y})$$

For  $k = L$  to  $k = 2$  :

$$\nabla_{w_k} \mathcal{L}(\theta) = \nabla_{a_k} \mathcal{L}(\theta) \cdot h_{k-1}^T$$

$$\nabla_{b_k} \mathcal{L}(\theta) = \nabla_{a_k} \mathcal{L}(\theta)$$

$$\nabla_{h_{k-1}} \mathcal{L}(\theta) = w_k^T (\nabla_{a_k} \mathcal{L}(\theta))$$

$$\nabla_{a_{k-1}} \mathcal{L}(\theta) = \nabla_{h_{k-1}} \mathcal{L}(\theta) \odot (\dots \delta'(a_{k-1}, j), \dots)$$

element wise product  
of vectors.

$$\nabla_{w_1} \mathcal{L}(\theta) = x^T$$

$$\nabla_{b_1} \mathcal{L}(\theta) = 1$$

Rough

$$q_1 = w_1 \cdot x^T + b_1^T$$

$$(\#) g = \text{act-f}^n$$

$$h_i = \text{act-f}^n(a_i)$$

$$\frac{\partial q_1}{\partial w_2} = x^T$$

$$(\#) \text{ if } g = \text{sigmoid}$$

$$\frac{\partial q_1}{\partial b_2} = 1$$

$$g' = g(a_{ij}) (1 - g(a_{ij}))$$

(#) Pseudo Code

for i in range (max-iter):

$$dW = [[0, 0]]$$

$$(\#) dR = [[0, 0]]$$

for forward-propag(X).  
for x, y, y-hat, a, h in zip(X, Y, Y-hat, A, H):

$$h = g(a_{ij})$$

(#) back propagate

- S<sub>Y</sub> u<sub>1</sub> → (##) update weights
- S<sub>u</sub> u<sub>1</sub> → (##) update biases

