

Lecture-9

* Feedforward Neural Networks.

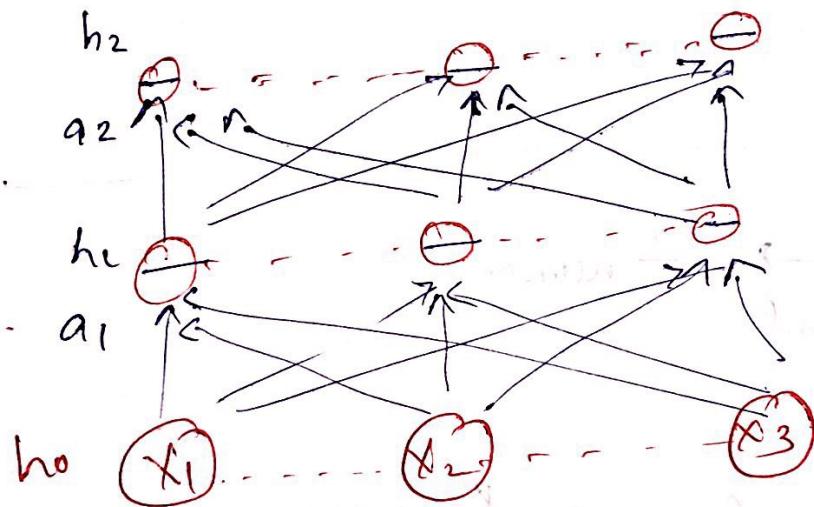
& Backpropagation

$$h_L = \vec{y} = f(\vec{x}) \in \mathbb{R}^k$$

\oplus \ominus

a₃.

K neurons
with
K classes



a_i - preactivation (aggregation)
 h_i - activation (non-linearity)

→ vectors

Pre-activation

$$a_i(x) = b_i + w_i h_{i-1}(x)$$

Ex: $a_1 = b_1 + w_1 h_0$

$$\begin{bmatrix} a_{11} \\ a_{12} \\ a_{13} \end{bmatrix} = \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix} + \begin{bmatrix} w_{111} & w_{112} & w_{113} \\ & w_{121} & \\ & & w_{133} \end{bmatrix} \begin{cases} h_0 = x_1 \\ h_0 = x_2 \\ h_0 = x_3 \end{cases}$$

$$= \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix} + \begin{bmatrix} w_{111}x_1 + w_{112}x_2 + w_{113}x_3 \\ w_{121}x_1 + w_{122}x_2 + w_{123}x_3 \\ w_{131}x_1 + w_{132}x_2 + w_{133}x_3 \end{bmatrix}$$

$$= \begin{bmatrix} \sum w_{11i}x_i + b_{11} \\ \sum w_{12i}x_i + b_{12} \\ \sum w_{13i}x_i + b_{13} \end{bmatrix}$$

$$h_i(x) = g(a_i(x))$$

$$\begin{bmatrix} h_{11} \\ h_{12} \\ h_{13} \end{bmatrix} = g\left(\begin{bmatrix} a_{11} \\ a_{12} \\ a_{13} \end{bmatrix}\right) = \begin{bmatrix} g(a_{11}) \\ g(a_{12}) \\ g(a_{13}) \end{bmatrix}$$

$$g(a_{13}) = \sigma(a_{13}) = \frac{1}{1+e^{-a_{13}}}$$

$$f(x) = h_L(x) = O(g_L(x))$$

$$y_i = f(x_i) = O\left(\omega^3 g(\omega^2 g(\omega^1 x + b_1) + b_2) + b_3\right)$$

$\overset{\longleftarrow R^n}{\underbrace{(n \times n) \times n \times n}}$

$$\theta = \omega_1, \dots, \omega_L, b_1, b_2, \dots, b_L$$

Lecture 4.2
 Learn Parameters of FeedForward
 Neural Networks

Algorithm : Gradient Descent.

```

 $t \leftarrow 0;$ 
max_iterations  $\leftarrow 1000;$ 
initialize  $\theta_0 = [\omega_0, b_0];$ 
while  $t++ < \text{max\_iterations}$  do
     $\theta_{t+1} \leftarrow \theta_t - \eta \nabla J(\theta_t);$ 
end
    
```

where $\nabla \theta_t = \left[\frac{\partial \mathcal{L}(\theta)}{\partial w_t}, \frac{\partial \mathcal{L}(\theta)}{\partial b_t} \right]^T$

$$\theta = [w, b]$$

$\begin{cases} \text{is converted} \\ \text{to} \end{cases} \rightarrow \theta = W_1, w_2, \dots, w_L, b_1, b_2, \dots, b_L$

Lecture 4.3

Output Functions & Loss functions

$$y_i \in \mathbb{R}^k$$

Loss Function

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^k (y_{ij} - \hat{y}_{ij})^2$$

$$f(x) = L = \theta a_L$$

$$= w_0 a_L + b_0$$

Entropy - measure of uncertainty

$$H(X) = H(p_1 \dots p_n)$$

Entropy

$$= - \sum_{i=1}^n p_i \log_2 p_i$$

Gross Entropy

Softmax Function

$$a_L = [a_{L1}, a_{L2}, a_{L3}]$$

$$\hat{y} = [\hat{y}_1, \hat{y}_2, \hat{y}_3]$$

$$\hat{y}_1 = \frac{e^{10}}{e^{10} + e^{-20} + e^{30}}$$

$$\text{Softmax Function} \Rightarrow a_L = w_L h_{L-1} + b_L$$

$$\hat{y}_j = O(a_L)_j$$

$$= \frac{a_{L,j}}{\sum_{i=1}^k a_{L,i}}$$

Principle way of computing two distribution events is Cross Entropy

Cross Entropy

$$L(\theta) = - \sum_{c=1}^K y_c \log \hat{y}_c$$

minimise

θ

$$L(\theta) = -\log \hat{y}_L \quad \begin{matrix} \text{(predicted correct)} \\ \text{label} \end{matrix}$$

maximise

θ

$$-L(\theta) = \log \hat{y}_L$$

$\log \hat{y}_L \rightarrow \log \text{likelihood of the data.}$

Outputs	
Output Activation	Real Values Probability
Linear	Softmax
Squared error	Cross-Entropy

Module 4.4

Backpropagation

$$L(0) = -\log \hat{y}_e$$

$$w(w_{112})$$

o/n layer

previous hidden layer

$$\frac{\partial L(0)}{\partial w_{111}} = \frac{\partial L(0)}{\partial y} \frac{\partial y}{\partial a_3} \frac{\partial a_3}{\partial h_2} \frac{\delta h_2}{\partial a_2}$$

$$\frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a_1} \frac{\partial a_1}{\partial w_{11}}$$

previous hidden layer

*the weights from the *i*th layer*

→ Compute the gradients with respect to the o/n units

→ Compute the gradients with respect to the hidden ~~8~~ layer unit

→ Compute the gradient w.r.t the weight & bias units

Module
4.5

Compute the gradients with respect
to the output units.

$$\frac{\partial L(\theta)}{\partial \hat{y}_L}$$

$$L(\theta) = -\log \hat{y}_l \quad (l = \text{true class label})$$

$$\frac{\partial}{\partial \hat{y}_i} (L(\theta)) = \frac{\partial}{\partial \hat{y}_i} (-\log \hat{y}_l)$$

$$\hookrightarrow_{\text{no. of classes}} = \begin{cases} 1 & \text{if } i=l \\ 0 & \text{otherwise} \end{cases}$$

More compactly

$$\frac{\partial}{\partial \hat{y}_i} (L(\theta)) = -\frac{\mathbb{I}_{(i=l)}}{\hat{y}_l}$$

$$\frac{\partial}{\partial \hat{y}_i} L(\alpha) = -\frac{T_{l=i}}{\hat{y}_e}$$

Gradient with respect to \hat{y}

$$\nabla_{\hat{y}} L(\alpha) = \begin{bmatrix} \frac{\partial L(\alpha)}{\partial \hat{y}_1} \\ \vdots \\ \frac{\partial L(\alpha)}{\partial \hat{y}_k} \end{bmatrix} = -\frac{1}{\hat{y}_e} \begin{bmatrix} T_{l=1} \\ \vdots \\ T_{l=k} \end{bmatrix}$$

$$= -\frac{1}{\hat{y}_e} e(l)$$

$$\frac{\partial L(\alpha)}{\partial a_{Li}} = \frac{\partial (-\log \hat{y}_e)}{\partial a_{Li}}$$

$$= \frac{\partial (-\log \hat{y}_e)}{\partial \hat{y}_e} \frac{\partial \hat{y}_e}{\partial a_{Li}}$$

$$\hat{y}_e = \frac{\exp(a_{Le})}{\sum_i \exp(a_{Li})}$$

$$\frac{\partial}{\partial a_{Li}} (-\log \hat{y}_e) = \frac{-1}{\hat{y}_e} \cdot \frac{\partial}{\partial a_{Li}} (\hat{y}_e)$$

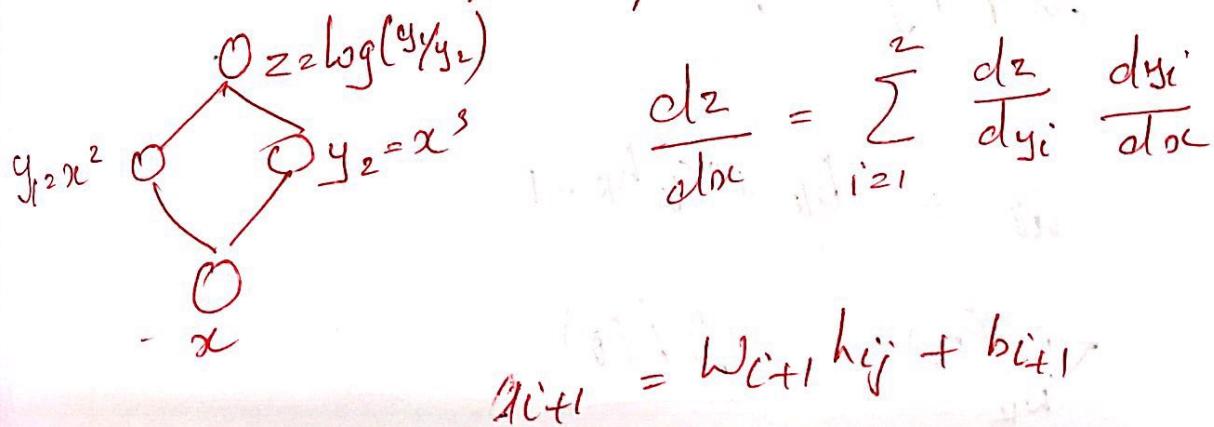
$$= \frac{-1}{\hat{y}_e} \cdot \frac{\partial}{\partial a_{Li}} \text{softmax}(a_L)_e$$

$\frac{\partial}{\partial x} \frac{v(x)}{v'(x)} = \frac{\partial v(x)}{\partial x} \frac{1}{v'(x)} - \frac{v(x)}{v'(x)^2} \frac{\partial v'(x)}{\partial x}$

$$= \frac{-1}{\hat{y}_e} \left(\frac{\partial}{\partial a_{Li}} \frac{\exp(a_L)_e}{\sum_i \exp(a_L)_i} \right)$$

expand

Module 4.6 Backpropagation : Computing Gradients
with respect to Hidden layers



$$(w_{i+1, \cdot, j})^T \nabla_{a_{i+1}} L(\theta) = \sum_{m=1}^k \frac{\partial L(\theta)}{\partial a_{i+1, m}} w_{i+1, m, j}$$

$$\frac{\partial L(\theta)}{\partial h_{ij}} = (w_{i+1, \cdot, j})^T \nabla_{a_{i+1}} L(\theta)$$

$$\nabla_{h_i} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial h_{i1}} \\ \vdots \\ \frac{\partial L(\theta)}{\partial h_{in}} \end{bmatrix} =$$

Module 4.7 Computing Gradients
 Backpropagation with respect Parameters

$$a_k = b_k + w_k h_{k-1}$$

$$\nabla_{w_k} L(\theta) = \frac{\partial L(\theta)}{\partial w_{kj}}$$

$$\frac{\partial \hat{a}_{ki}}{\partial w_{kij}} \Rightarrow h_{k-1,j}$$

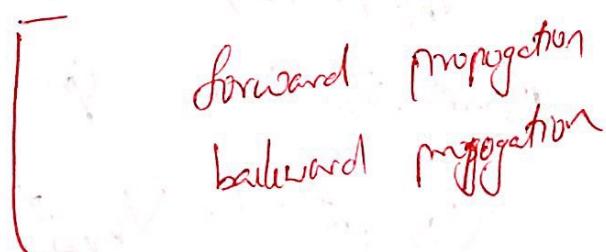
$$\frac{\partial L(o)}{\partial w_{kij}} = \frac{\partial L(o)}{\partial a_{ki}} \cdot \frac{\partial a_{ki}}{\partial w_{k-1,j}}$$

$$\Rightarrow \frac{\partial L(o)}{\partial a_{ki}} h_{k-1,j}$$

$$\nabla_{w_k} L(o) = \begin{bmatrix} \frac{\partial L(o)}{\partial w_{k00}} \\ \vdots \\ \frac{\partial L(o)}{\partial w_{k,n-1,n-1}} \end{bmatrix}$$

4.8

Algorithm



Forward propagation

for $k=1$ to $L-1$ do

$$\begin{cases} a_k = b_k + w_k h_{k-1}; \\ h_k = g(a_k) \end{cases}$$

end

$$a_L = b_L + w_L h_{L-1};$$

$$\hat{y} = O(a_L);$$

Backward

$$\nabla_{a_L} L(\theta) = - (e(y) - f(x)) / ;$$

for $k=L$ to 1 do

// wrt parameters

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

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

// wrt to layer below

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

// wrt layer below (one activation)

$$\nabla_{a_{k-1}} L(\theta) = \nabla_{h_{k-1}} L(\theta) \odot$$

$$[\dots, g'(a_{k-1}) \dots]$$

end

A.9 Derivative of the activation function

$$g(z) = \sigma(z)$$

$$= \frac{1}{1+e^{-z}}$$

$$g'(z) = -1 \left(\frac{1}{(1+e^{-z})^2} \right) \frac{d}{dz} (1+e^{-z})$$

$$= (-1) \frac{1}{(1+e^{-z})^2} (-e^{-z})$$

$$= \frac{1}{1+e^{-z}} \left(\frac{1+e^{-z}-1}{1+e^{-z}} \right)$$

$$= g(z)(1-g(z))$$

$$\overline{\tanh} \quad g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - (g(z))^2$$

4.10

Information content, Entropy, Cross Entropy

$$X \rightarrow A \quad B \quad C \quad D$$

$$\begin{bmatrix} 0.25 \\ p_1 \end{bmatrix} \quad \begin{bmatrix} 0.35 \\ p_2 \end{bmatrix} \quad \begin{bmatrix} 0.4 \\ p_3 \end{bmatrix} \quad P$$

estimate $\begin{bmatrix} 0.34 \\ q_1 \end{bmatrix} \quad \begin{bmatrix} 0.45 \\ q_2 \end{bmatrix} \quad \begin{bmatrix} 0.2 \\ q_3 \end{bmatrix}$ 2

One way
Square error $\sum_{i=1}^3 (p_i - q_i)^2$

$$X \rightarrow A \quad B \quad C \quad D$$

$$P \quad [0 \quad 1 \quad 0 \quad 0]$$

All the probability mass is forced on one item.

Prediction
 $q \quad [0.2 \quad 0.4 \quad 0.2 \quad 0.1]$

$$\sum_{i=1}^4 (p_i - q_i)^2$$

Classification

$$X \rightarrow A \quad B \quad C \quad D$$

$$\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

Prediction
 $[0.25 \quad 0.45 \quad 0.2 \quad 0.1]$

$$\mathcal{L}(o) = P // 2$$

Square Error $\Rightarrow \sum_{i=1}^4 (p_i - q_i)^2$

Expectation

X	A	B	C	D	
$P(X=i)$	0.3	0.4	0.2	0.1	
$V(X=i)$	10K	5K	10K	-30K	(gain)

Expected Reward \rightarrow

$$\text{Expected Gain} = \sum_{i \in \{A, B, C, D\}} P(X=i) V(X=i)$$

Information content

$$X = A \quad B \quad C \quad D$$

Information gain is high, when the event which happens is a surprising event

$$I_C(A) \propto \frac{1}{P(A)}$$

Event

$$f(P(A))$$

independent events

	X	Y
A	ON	
B		OFF
C		ON
D		OFF

$$I_C(X \wedge Y) = I_C(X) + I_C(Y)$$

⊗

$$f(P(x \cap y)) = f(P(x)) + f(P(y))$$

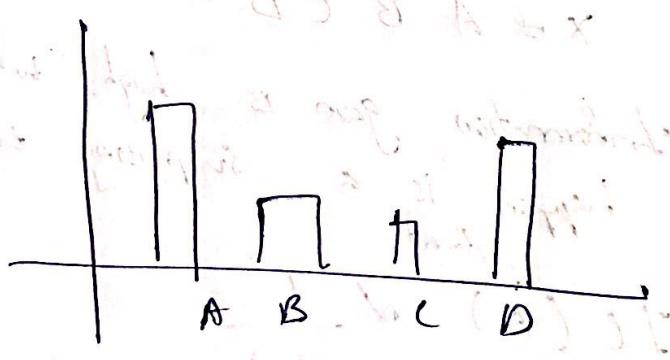
$$\underline{f(P(x))} =$$

$$f(P(x), P(y)) = f(P(x)) + f(P(y))$$

$$f(a \cdot b) = f(a) + f(b)$$

$$I(A) = -\log \left(\frac{1}{P(A)} \right)$$

$$I(A) = -\log P(A)$$



funktions
graph:

$$\sum_{i \in \{A, B, C, D\}} P(X=i) \log(P(X=i)) \times (P(X=A), P(X=B), P(X=C), P(X=D))$$

$$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$$

$$I(A), I(B), I(C), I(D)$$

$X \quad A \quad B \quad C \quad D$

A	0	0	$\frac{1}{4}$
B	0	1	$\frac{1}{4}$
C	1	0	$\frac{1}{4}$
D	1	1	$\frac{1}{4}$

$$I_C(\text{this}) = -\log_2\left(\frac{1}{4}\right)$$

$$= -\log_2(2^{-2})$$

$$= 2 \quad \text{no. of bits.}$$

$X \oplus A \quad B \oplus C \quad D \oplus C \quad A \oplus Y$

000

$$-\log_2\left(\frac{1}{8}\right)$$

001

$$-\log_2(2^{-3})$$

010

100

3

101

110

111

A $\frac{1}{2}$ 1

B $\frac{1}{4}$ 2

C $\frac{1}{8}$ 3

D $\frac{1}{8}$ 3

A B C D
P $\left[\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8}\right]$

1 2 3 3

Equation

$$-\sum_{i \in A, B, C, D} P(X=i) \log_2(P(X=i))$$

$$\text{Entropy} = 1.75$$

P $[P_1 \ P_2 \ P_3 \ P_4]$ - the

Q $[q_1 \ q_2 \ q_3 \ q_4]$ - definition

- $\log q_i$ bit for the i^{th} message

P $[P_1 \ P_2 \ P_3 \ P_4]$

Q $-\log q_1 \ -\log q_2 \ -\log q_3 \ -\log q_4$

How to compute expectation.

$$L(p, q) = - \sum p_i \log q_i \quad (\text{minimum value when } p = q)$$

P

We want Q such that

$$-\sum p_i \log q_i$$

Q

minimise

$$-\sum p_i \log q_i \text{ w.r.t } q_i$$

$$\min_x x^2$$

$$x=0$$

$$\min_{x,y} x^2 + y^2$$

$$\frac{\partial}{\partial q_2} (P_1 \log q_1 + P_2 \log q_2 + \dots + P_{10} \log q_{10})$$

$$\text{minimize}_{q_i} - \sum p_i \log q_i$$

$$\text{s.t. } \sum q_i = 1$$

$$\text{minimize}_{q_i} - \sum p_i \log q_i + \lambda (\sum q_i - 1)$$

Lagrange value.

↳ if violent this will be added
 ↳ if not correct, then the value will be zero.

$$-\frac{P_2}{q_2} + \lambda = 0$$

$$P_2 = \lambda q_2$$

$$p_i = q_i$$

Will be minimize of