

Sigmoid & Softmax

Notations.

Target $\rightarrow T$. (Ground truth, Labels)

obtained $\rightarrow Y$ (model-forward())
Model's output.

10 digit binary numbers.

$$\begin{array}{cccccccc} \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \frac{1}{2} & \frac{1}{2} & \dots & \dots & \dots & \dots & \dots & \dots \\ \downarrow & \downarrow & \downarrow & \dots & \dots & \dots & \dots & \downarrow \\ 2 & 2 & 2 & \dots & \dots & \dots & \dots & 2 \end{array} = 2^{10}.$$

#1's > #0's?

1010110101 $\rightarrow 1$

1111110000 $\rightarrow 1$

0000000011 $\rightarrow 0$

I/p \rightarrow Model \rightarrow o/p.

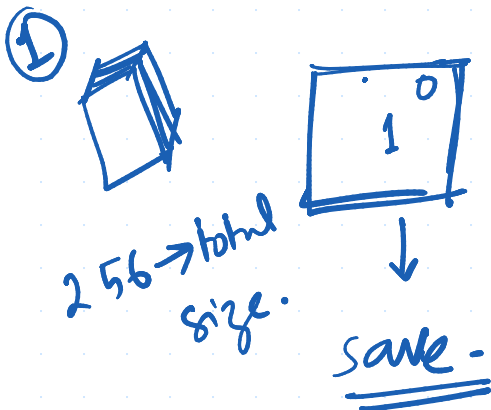
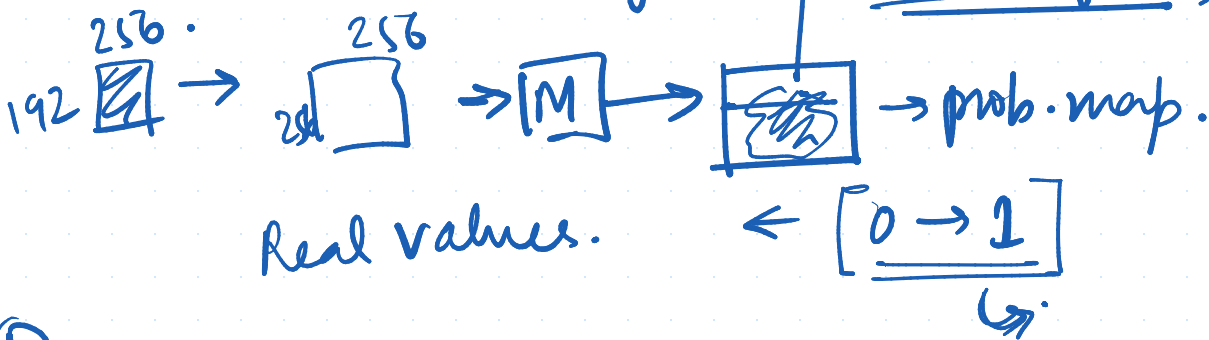
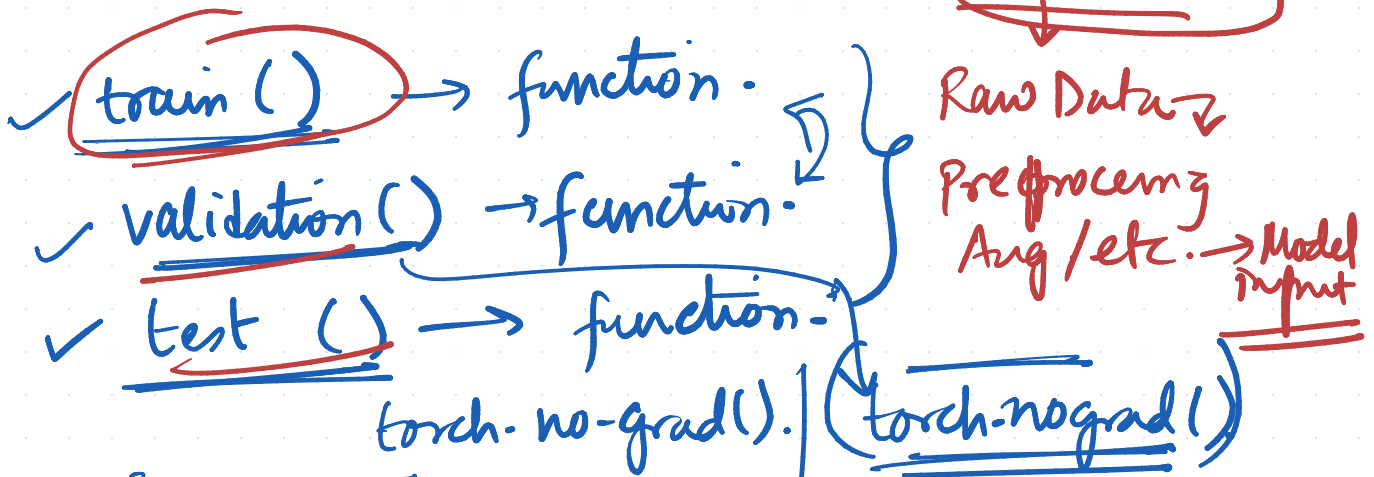
(1x10), (10x1)

$\left. \begin{array}{c} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{array} \right\} 10 \text{ neurons.}$

output $\rightarrow 0/1$ (Sigmoid). Activation.

Pytorch :- (pipeline)

- Data \rightarrow web/simulated/Collect.
- Model
- $Ip \rightarrow Model \rightarrow prediction.$
- Visualization.



round(output).

$\left\{ \begin{array}{l} 0.5 \downarrow \rightarrow 0 \\ 0.5 \uparrow \rightarrow 1 \end{array} \right\}$ threshold.

 (cv2.imwrite(—))



Dice → 70 %

Test function



Similar to validation function with some added functionality. Save. test

validation X



└───→ 2000 epochs

100 ↓



logs. etc.

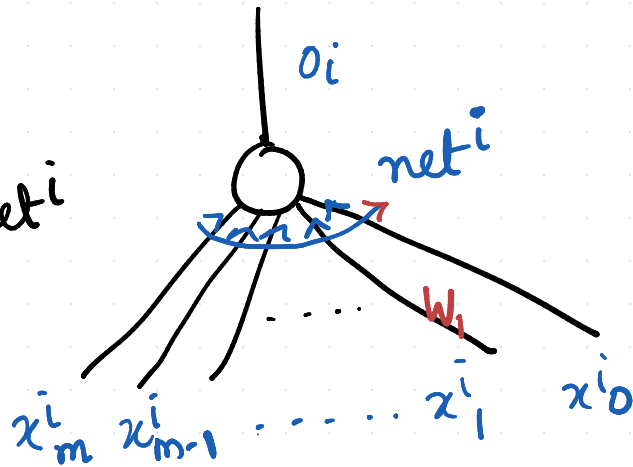
✓ Batch-wise → model → (cuda)
(Pipeline)

Assignment → (extend) ✓

Target $\rightarrow T$, Obtained $\rightarrow Y$.

Sigmoid.

$$o^i = \frac{1}{1 + e^{-x}} = \text{net}^i$$

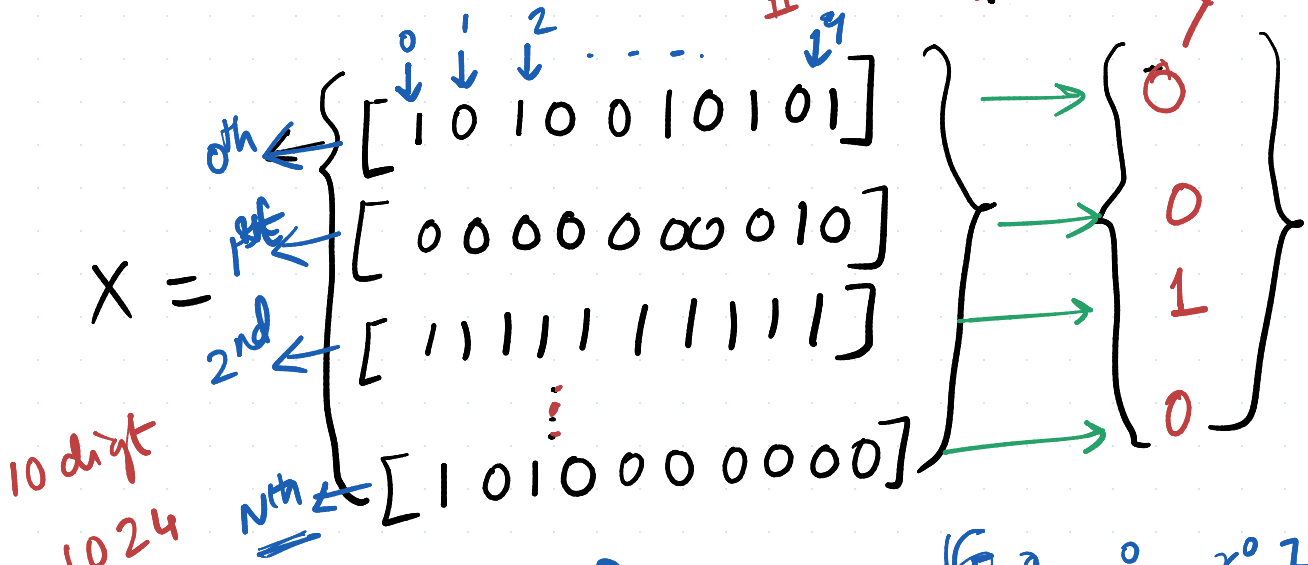


$$o^i = \frac{1}{1 + e^{-\text{net}^i}}$$

$x^i \rightarrow$ ith sample
 $x_j \rightarrow$ jth scalar.

$$\text{net}^i = W \cdot x^i = \sum_{j=0}^m w_j x_j^i$$

$! > 0 ?$

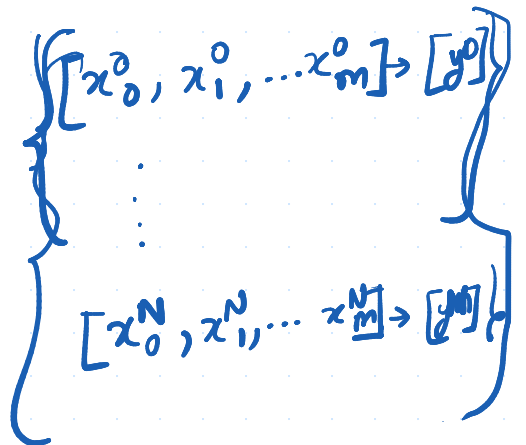


10 digit
 1024
 ip/op
 pairs.

$$x_7^0 = 1$$

$$x_7^1 = 0$$

(N+1)



Vectors \rightarrow capital letters.

scalars \rightarrow small letters.

$x^i \rightarrow$ i th input vector.

$o_i \rightarrow$ output scalar.

$net_i = W x^i$ (n -input/output

$W =$ weight vector. observations).

W & $x^i \rightarrow$ has m -components.

$$W = \langle w_m, w_{m-1}, \dots, w_2, w_1, w_0 \rangle$$

$$x^i = \langle x_m^i, x_{m-1}^i, \dots, x_2^i, x_0^i \rangle$$

Derivative of Sigmoid.

$$o_i = \frac{1}{1+e^{-net_i}} \rightarrow \text{for } i\text{th input.}$$

$$\ln o_i = -\ln(1+e^{-net_i})$$

$$\frac{1}{o_i} \frac{\partial o_i}{\partial net_i} = -\frac{1}{(1+e^{-net_i})} \cdot (-e^{-net_i})$$

$$= \frac{e^{-net_i}}{1+e^{-net_i}} = (1-o_i)$$

$$\frac{\partial o_i}{\partial net_i} = o_i(1-o_i)$$

$$x(1-x)$$

-10 in pytorch.

$-\infty \rightarrow 0$
 $+\infty \rightarrow 1$ } sigmoid.

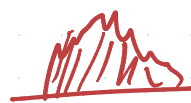
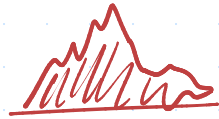
$+10 \rightarrow$ in pytorch.

$x = \text{torch.tanh}(t)$
 $t = \text{torch.sigmoid}(x)$
 $\text{print}(t)$
 $x(1-x)$

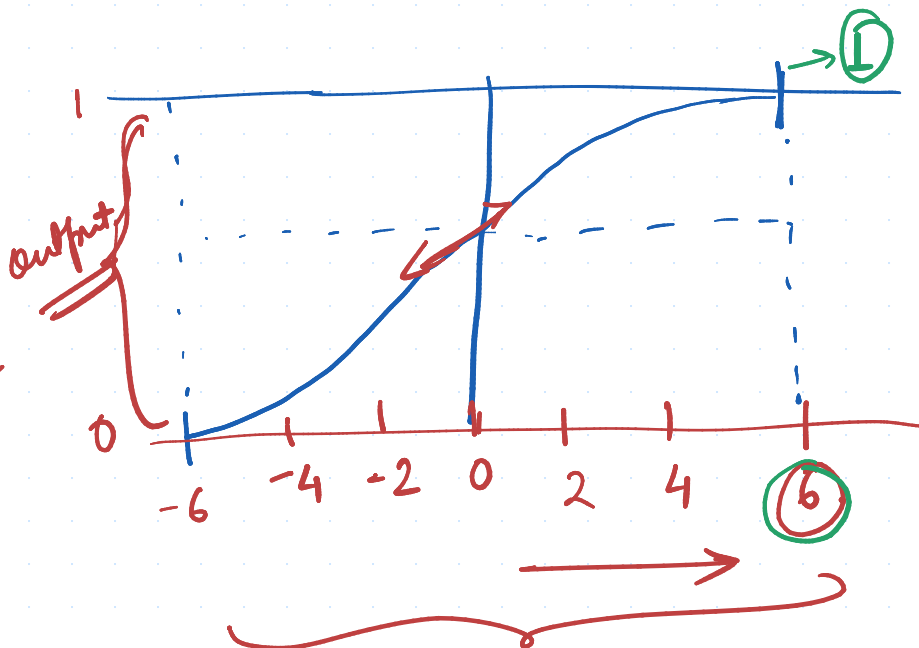
$x = [-10, -9, -8, \dots, 0, 1, 2, \dots, 9, 10]$

$[0, 1e^{-20}, \dots, \frac{1}{2}, \dots, 0.99, 1]$

\downarrow (torch.sigmoid) \downarrow



Derivative of sigmoid is maximum at where?



Maximum value of derivative:

(tune)

Normalize

Regularize

skts. \rightarrow $\text{overfit}(x)$

input

OOD. $\frac{1}{1000}$

$(-6, 6)$

$-100 \rightarrow 100$

\rightarrow skts.

$100 \rightarrow 1$

$6 \rightarrow 100$

1

$-100 \rightarrow -6$

$$\frac{\partial}{\partial x} (x(1-x)) \Rightarrow \frac{\partial}{\partial x} (x - x^2) \Rightarrow 1 - 2x$$

$$\left(\frac{x}{\|x\|} \right)$$

$$\frac{\partial}{\partial x} (x(1-x)) = 0$$

$$1 - 2x = 0$$

$$x = \frac{1}{2}$$

$$\frac{x - \mu}{\text{std}}$$

Maximum value of derivative of sigmoid. $\frac{1}{2} (1 - \frac{1}{2}) = \frac{1}{2} (\frac{1}{2}) = 0.25$.

$$o_i = \frac{1}{1 + e^{-net_i}}$$

$$net_i = \infty$$



$$net_i = -\infty$$

$$o_i = \frac{1}{1 + \frac{1}{e^{net_i}}}$$

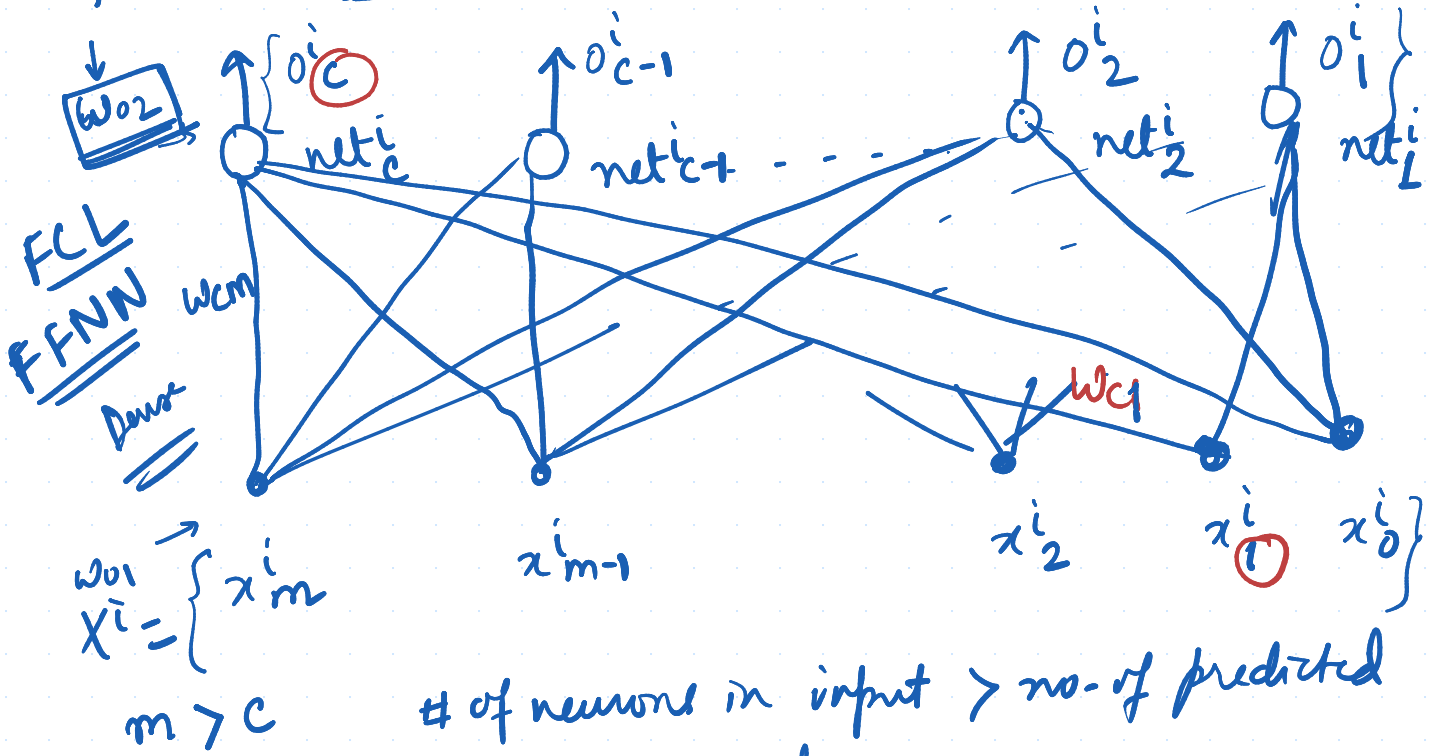
$$= \frac{1}{1 + e^{\infty}} = \frac{1}{\infty} = 0$$

$$o_i = \frac{1}{1 + \frac{1}{e^{net_i}}} = \frac{1}{1 + \frac{1}{e^{\infty}}}$$

$$= \frac{1}{1 + 0} = 1$$

$$y = mx + b$$

Softmax



of neurons in input > no. of predicted classes.

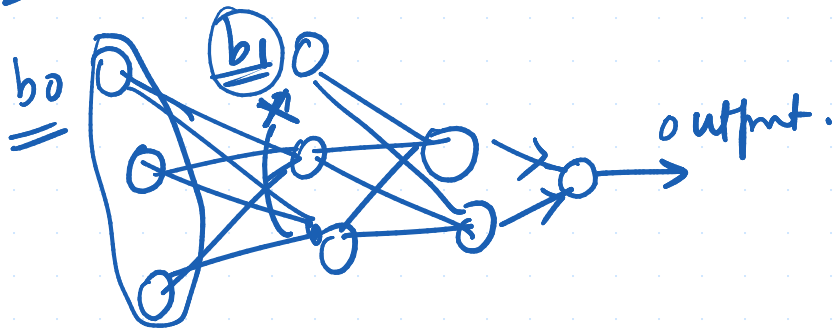
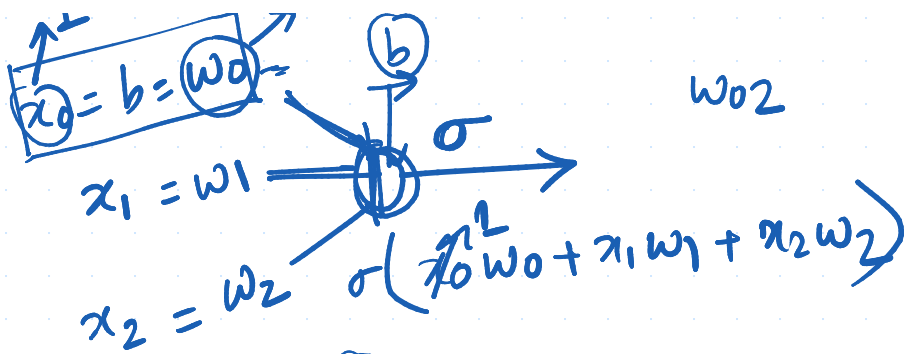
Output for class c , $c: 1$ to c .

w_{c1}
 $w_{(c-1)L}$

$$o_c^i = S(NET^i)_c = \frac{e^{net^i_c}}{\sum_{k=1}^c e^{net^i_k}}$$

Softmax.

Rough



Notations.

$i = 0, 1, 2, \dots, N$ $\{x_i^j\}_0^N$
 (N+1) ip/op pairs. , i-runs over training data.

$j = 0, \dots, m, \underline{(m+1)}$ → components.

also.

size of weight vector

input vector → size.

output $\rightarrow y$ (class.)

Rough

$$\mathcal{D} = \left\{ \begin{array}{l} [x_0^0, x_1^0, \dots, x_m^0] \rightarrow [o_0^0, o_1^0, \dots, o_c^0] \\ [x_0^1, x_1^1, \dots, x_m^1] \rightarrow [o_0^1, o_1^1, \dots, o_c^1] \\ \vdots \\ [x_0^N, x_1^N, \dots, x_m^N] \rightarrow [o_0^N, o_1^N, \dots, o_c^N] \end{array} \right.$$

$x \rightarrow (m+1)$

N (N+1) → o/p output pair.

$k = 0, \dots, C, C+1$ classes.

(C+1) → components in o/p. vector.

O^i, NET^i → are the vectors for the i^{th} input.

W_k → weight vector for the k^{th} output neuron.

(Actual output) / o/p.
Target vector (T) →

$C = 1, 2, \dots, C.$

$\left\{ \begin{array}{l} [\dots] \rightarrow [\cdot] \\ [\quad] \rightarrow [\cdot] \\ [\quad] \rightarrow [\cdot] \end{array} \right\}$
T

$T^i = \langle t_{c}^i, t_{c-1}^i, \dots, t_2^i, t_1^i \rangle$

i → for i^{th} input.

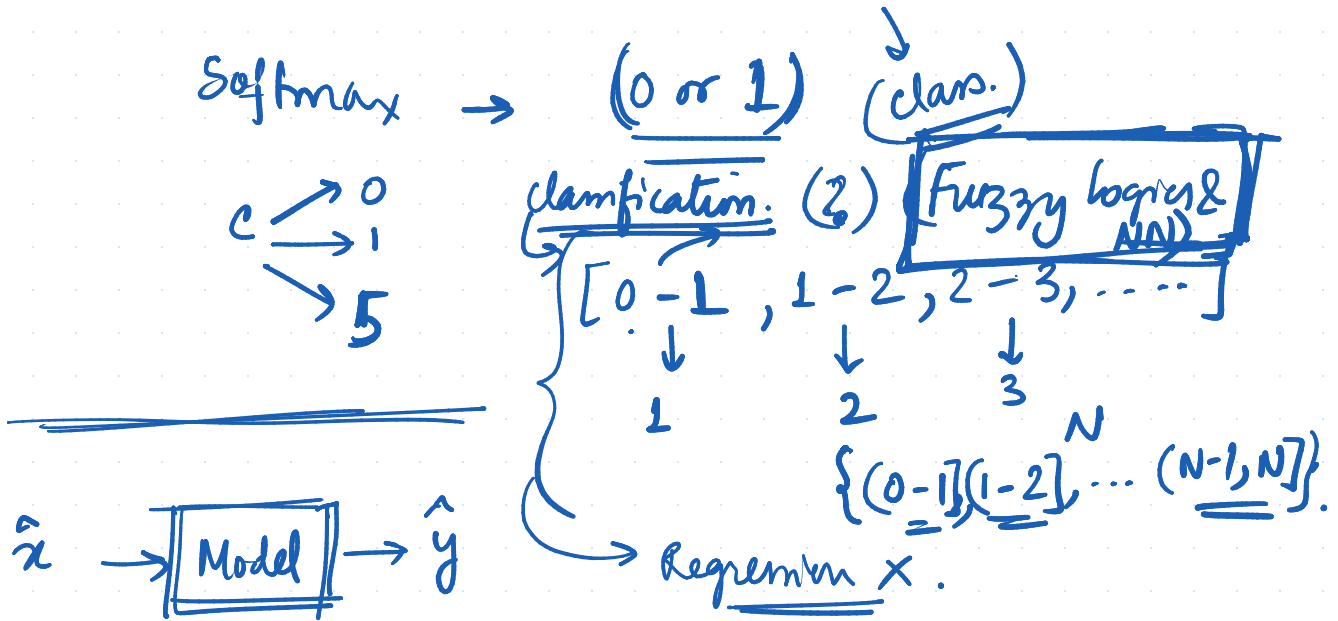
one of these c -component is 1, rest are 0.

$\begin{array}{ccc} 1 & 2 & 3 \\ \downarrow & \downarrow & \downarrow \\ \langle 0, 1, 0 \rangle & \rightarrow & 2 \end{array}$

$\langle 0, 0, 1 \rangle \rightarrow 3$

One hot vector?

(1, 2, 3) output.



Regression task \rightarrow continuous Real value.

$\hat{y} \rightarrow \underline{(0-100)} < \underline{1, 0, 0} >$

Softmax =

$$\left\{ \begin{matrix} -1 \\ 1 \\ 0 \end{matrix} \right\} \xrightarrow{\text{softmax}} \left\{ \begin{matrix} 0.02 \\ 0.5 \\ 0.48 \end{matrix} \right\} \xrightarrow{\text{argmax}} \underline{\underline{\left\{ \begin{matrix} 0 \\ 1 \\ 0 \end{matrix} \right\}}}$$

$$\left\{ \frac{e^{-1}}{e^{-1} + e^1 + e^0}, \frac{e^1}{e^{-1} + e^1 + e^0}, \frac{e^0}{e^{-1} + e^1 + e^0} \right\}$$

$$\downarrow$$

0 - 1 0 - 1 0 \rightarrow 1



Andrew Ng

Coursera
deeplearning.ai.

CS

1

Probability
distribution

Daphne Koller

$\{(0, 1], [1, 2], \dots, (N-1, N]\}$

↓
1

↓
2

↑

↓
N

N=5.

Dog Cat fish tree human

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

↓

↓

↓

↓

← 1

↔ 2

3

4

→ 5

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

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

class 1 is more close to class 2 than class 5.

Inherent ordering → (X)

~~1~~ → ~~2~~

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

Kilometers + Kg

(0, 1, 0, 0)

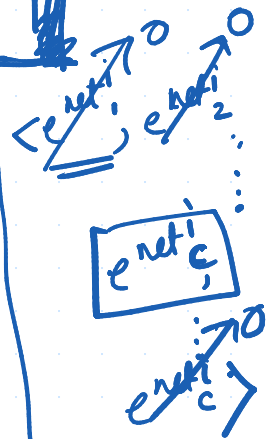


Softmax. $\left\{ \begin{aligned} o_c^i &= \frac{e^{\text{net}_c^i}}{\sum_{k=1}^c e^{\text{net}_k^i}} \end{aligned} \right.$, i^{th} input pattern.

$$\ln o_c^i = \text{net}_c^i - \ln \left(\sum_{k=1}^c e^{\text{net}_k^i} \right)$$

Derivative w.r.t. net_c^i :

$$\frac{1}{o_c^i} \cdot \frac{\partial o_c^i}{\partial \text{net}_c^i} = 1 - \frac{e^{\text{net}_c^i}}{\sum_{k=1}^c e^{\text{net}_k^i}}$$



Case 1 →
When class c for 0 & NET are the same.

$$= 1 - o_c^i$$

$$\frac{\partial o_c^i}{\partial \text{net}_c^i} = o_c^i (1 - o_c^i)$$

Case 2 → When class c' and net_c^i are different from c of 0.

$$\ln o_i^c = \underline{\text{net}_i^c} - \ln \left(\sum_{k=1}^c e^{\text{net}_i^k} \right)$$

Derivative w.r.t. net_i^c

$$\frac{1}{o_i^c} \frac{\partial o_i^c}{\partial \text{net}_i^c} = 0 - \frac{1 \times e^{\text{net}_i^c}}{\sum_{k=1}^c e^{\text{net}_i^k}} = -o_i^c.$$

$$\frac{\partial o_i^c}{\partial \text{net}_i^c} = -o_i^c o_i^c$$

when the classes are unequal.

Maximum Likelihood & Cross Entropy Loss.

Q) How did cross-entropy loss come to existence? (why this particular form?)



class value of input \rightarrow R.V.
(Random variable).



We are interested in modelling probability

$P(o^i | x^i)$, here $o^i \rightarrow$ output vector
 $x^i \rightarrow$ input vector.

$o_c^i \rightarrow$ component of $o^i \rightarrow$ probability of x^i belonging to the class c ($c=1, 2, \dots, C$)

c -components are redundant.

prob. of class $c \Rightarrow \underline{1 - \sum \text{prob. (class } \neq c \text{)}}.$

$$\langle 0.1, 0.2, 0.3, 0.4 \rangle = 1$$

↓

$$\langle \underbrace{0.1, 0.2, 0.3}_{0.6}, \boxed{0.4} \rangle = 1 - \underbrace{0.6}_{\boxed{0.4}}$$

Hence in case of 2-class, only one sigmoid neuron is required.

$o^i \rightarrow$ value between 0 & 1 \rightarrow interpreted as probability

$$\sigma\left(\frac{\square}{\square}\right) \rightarrow [0 \rightarrow 1]$$

$o^i \rightarrow$ probability of class being equal to 1.

$$P(\text{class} = 1 \text{ for } i\text{th input}) = o^i = \frac{1}{1 + e^{-net^i}}$$

Training data instance is labelled as 1 or 0.

Target value \rightarrow 1/0 $\rightarrow t^i$ for ith input.

Likelihood of observation for (two classes)

For N no. of $1/p - 0/p$ pairs.

$$L = \prod_{i=1}^N (o^i)^{t^i} (1 - o^i)^{(1 - t^i)}$$

$t^i = 1 \text{ or } 0.$

\downarrow

target

$0 \rightarrow$ output.

Bernoulli Distribution.

(2 \rightarrow Bernoulli)

Log-likelihood

maximize.

CE.

KL \rightarrow CE.

$$d = \prod_{i=1}^N (o_i)^{t_i} (1-o_i)^{(1-t_i)} \quad , t_i = 1 \text{ or } 0.$$

$$LL = \log \left(\prod_{i=1}^N (o_i)^{t_i} (1-o_i)^{(1-t_i)} \right)$$

$o_i \rightarrow \text{output.}$

Log-likelihood \rightarrow maximize

$$LL = \sum_{i=1}^N \{ t_i \log(o_i) + (1-t_i) \log(1-o_i) \}$$

Binary
cross entropy loss.

\rightarrow Negative of log-likelihood
 \rightarrow minimize.

$$\Rightarrow -LL = - \sum_{i=1}^N \{ t_i \log o_i + (1-t_i) \log(1-o_i) \}$$

$$\underline{\text{BCE loss}} = - \frac{1}{N} \sum_{i=1}^N \{ y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i) \}$$

L_1 , L_2 loss why not?



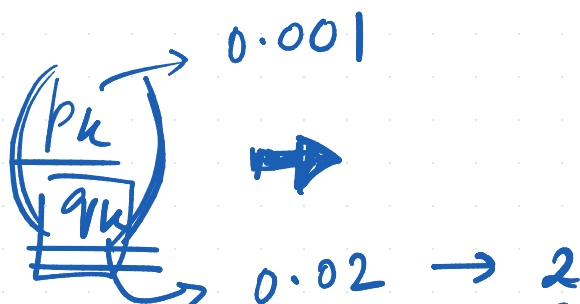
$$\left(\frac{1}{77}, \frac{1}{700000}\right) \rightarrow \left(\underline{\underline{0.001}}, \underline{\underline{0.0001}}\right)$$

$$\left(\underline{\underline{L1}}, \underline{\underline{L2}}\right)$$

Binary cross entropy?

$$c1 \rightarrow (1, 0, 0)$$

$$c2 \rightarrow (0, 1, 0)$$



$$\log\left(\frac{p_k}{q_k}\right)$$

$$\underline{\underline{(p_k - q_k)}}$$

$$\boxed{\log\left(\frac{p_k}{q_k}\right)}$$

X →

day 1

$$P(\text{chance of acci})$$

$$= \underline{\underline{0.00001}}$$

$$P(\text{chance of acci})$$

$$= 0.001$$

100%

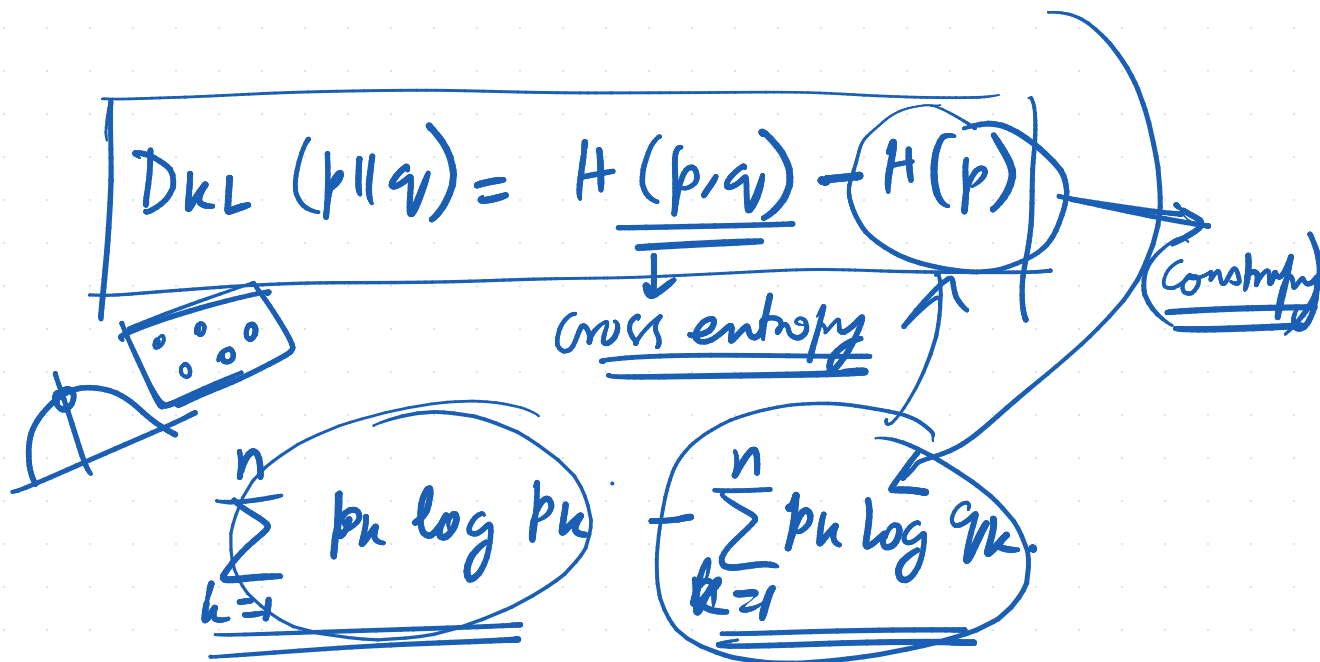
$$\underline{\underline{p_k \log\left(\frac{p_k}{q_k}\right)}}$$

→ Probability /
Distribution

(distance)

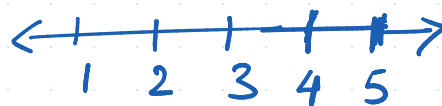
Distribution Difference

$$D_{KL}(p||q) = \sum_{k=1}^n p_k \log\left(\frac{p_k}{q_k}\right)$$



GT label $\rightarrow p$,

q \rightarrow model's output.



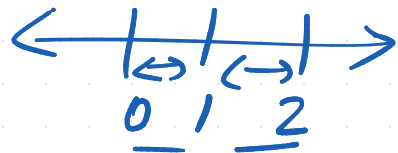
Entropy \uparrow Information Gain \downarrow

- LL \rightarrow negative of log-likelihood is called cross entropy, regarded as loss or error.

$\mathcal{E} \rightarrow$ notation for error.

Minimizing cross entropy brings o_i close to t^i ; Hence we established, equivalence b/w maximization of

likelihood \longleftrightarrow



Softmax

$o_c^i \rightarrow$ value b/w 0 and 1, interpreted

as probability, Multi-class situation.

\rightarrow value is the prob of the class being 'c' for the i th input.

$$P(\text{class of } i\text{th input} = c) = o_c^i.$$

Likelihood L of observations in case of softmax :-

For N no. of i/p - o/p pairs:-

$$L = \prod_{i=1}^N \prod_{k=1}^C (o_k^i)^{t_k^i}, \quad t_k^i = \begin{cases} 1 \\ 0 \end{cases}$$

For a pattern i , only one of t_k^i 's is 1,
rest are 0:-

$$L = \sum_{i=1}^N \sum_{k=1}^C t_k^i \log o_k^i.$$

$$-L = - \sum_{i=1}^N \sum_{k=1}^C t_k^i \log o_k^i \quad \text{Categorical Cross Entropy}$$

For softmax also maximize the
 likelihood = minimize the cross-entropy.

$-L =$ the cross-entropy, Loss or Error.

Backpropagation & Gradient Descent.

Greedy Algorithm.

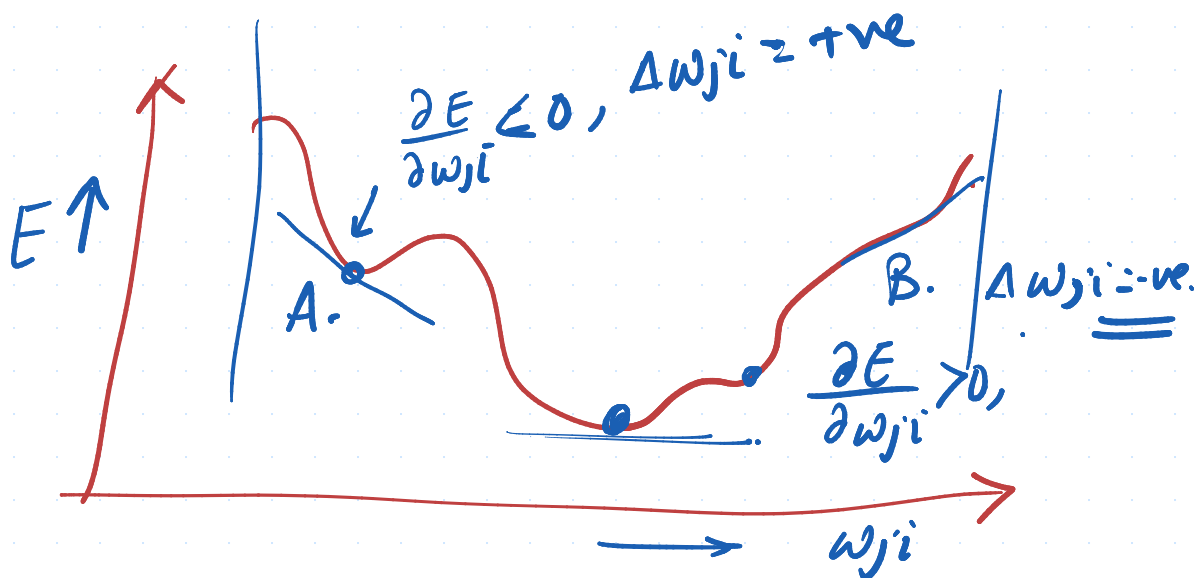
Always moves in the direction of reducing errors.

$$\Delta w_{ji} = -\eta \frac{\partial E}{\partial w_{ji}}$$

→ change in weight.
 η = learning rate.

E = Loss function.

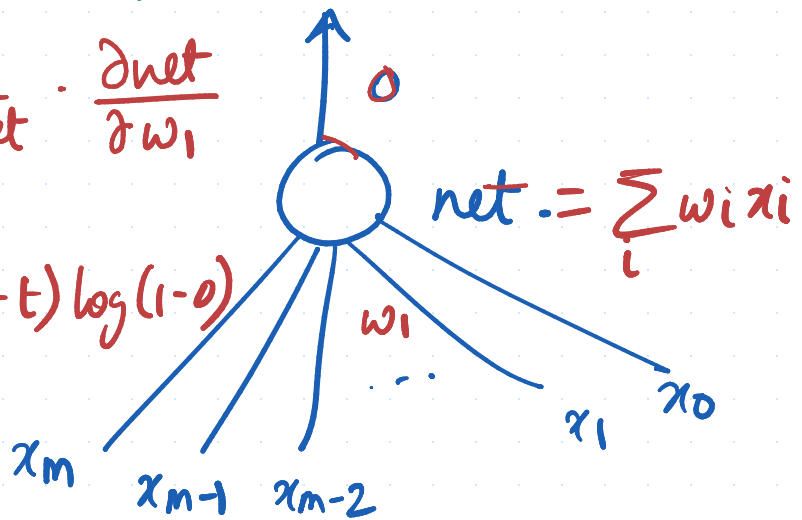
w_{ji} = weight of connection from the i^{th} neuron to the j^{th} neuron.



Single sigmoid Neuron & Cross Entropy Loss, derived for single data point, hence, dropping right suffix i .

$$\frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial o} \cdot \frac{\partial o}{\partial \text{net}} \cdot \frac{\partial \text{net}}{\partial w_1}$$

$$E = -t \log o - (1-t) \log(1-o)$$



$$\frac{\partial E}{\partial o} = -\frac{t}{o} - \frac{(1-t)(-1)}{(1-o)}$$

$$= \frac{-t(1-o) + (1-t)o}{o(1-o)} = \frac{-t + t\cancel{o} + 0 - t\cancel{o}}{o(1-o)}$$

$$\boxed{\frac{\partial E}{\partial o} = -\frac{t-o}{o(1-o)}}$$

Sigmoid Act.

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

$$\boxed{\frac{\partial o}{\partial \text{net}} = o(1-o)}$$

$$\frac{\partial \text{net}}{\partial w_1} = x_1$$

$$\frac{\partial E}{\partial w_1} = -\frac{(t-o)}{o(1-o)} \times \cancel{o(1-o)} \times x_1$$

$$\frac{\partial E}{\partial w_1} = -(t-o)x_1$$

$$\Delta w_1 = -\eta \frac{\partial E}{\partial w_1} = \eta (t-o)x_1$$

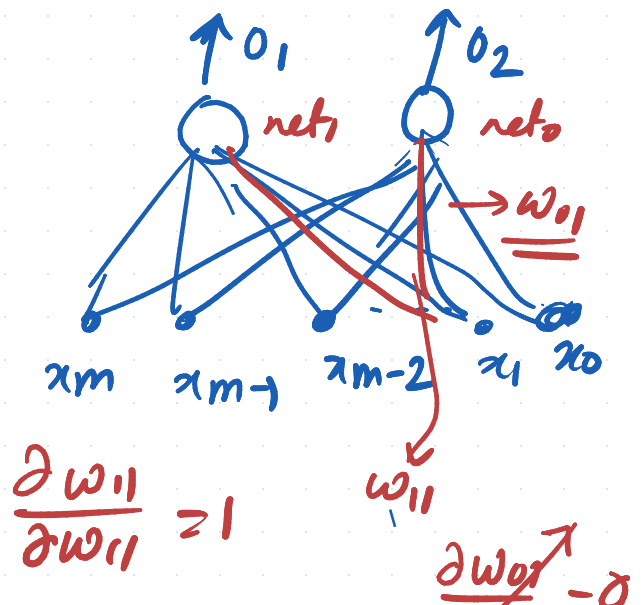
Multiple Neurons. → in output layer

Softmax & Cross Entropy loss → I illustrated

with 2 neurons.

$$O = \langle o_1, o_2 \rangle$$

$$\text{NET} = \langle \text{net}_1, \text{net}_2 \rangle$$



$$o_1 = \frac{e^{\text{net}_1}}{e^{\text{net}_1} + e^{\text{net}_0}}$$

$$o_0 = \frac{e^{\text{net}_0}}{e^{\text{net}_1} + e^{\text{net}_0}}$$

$$\frac{\partial O}{\partial \text{NET}} = \begin{bmatrix} \frac{\partial o_0}{\partial \text{net}_0} & \frac{\partial o_1}{\partial \text{net}_0} \\ \frac{\partial o_0}{\partial \text{net}_1} & \frac{\partial o_1}{\partial \text{net}_1} \end{bmatrix} = \begin{bmatrix} o_0(1-o_0) & -o_0 o_1 \\ -o_1 o_0 & o_1(1-o_1) \end{bmatrix}$$

Jacobian \rightarrow First order partial derivative

Hessian \rightarrow 2nd order partial derivative.

m-tensor \rightarrow higher order partial derivative.

ex. $x = \begin{bmatrix} \circ & \circ & \circ \\ \circ & \circ & \circ \\ \circ & \circ & \circ \end{bmatrix} \rightarrow$ 3-tensor
4-tensor

$$E = -t_1 \log o_1 - (1-t_1) \log (1-o_1)$$

$$\frac{\partial E}{\partial w_{11}} = -\frac{t_1}{o_1} \frac{\partial o_1}{\partial w_{11}} - \frac{t_0}{o_0} \cdot \frac{\partial o_0}{\partial w_{11}}$$

$$\begin{aligned} \frac{\partial o_1}{\partial w_{11}} &= \frac{\partial o_1}{\partial \text{net}_1} \cdot \frac{\partial \text{net}_1}{\partial w_{11}} + \frac{\partial o_1}{\partial \text{net}_0} \cdot \frac{\partial \text{net}_0}{\partial w_{11}} \\ &= 0_1(1-0_1)x_1 + \cancel{-0_1 \cdot 0} \cdot 0 \\ &= 0_1(1-0_1)x_1 \end{aligned}$$

$$\begin{aligned} \frac{\partial o_0}{\partial w_{11}} &= \frac{\partial o_0}{\partial \text{net}_1} \cdot \frac{\partial \text{net}_1}{\partial w_{11}} + \frac{\partial o_0}{\partial \text{net}_0} \cdot \frac{\partial \text{net}_0}{\partial w_{11}} \\ &= -0_0 \cdot 0_1 x_1 + 0 \end{aligned}$$

$$\begin{aligned} \frac{\partial E}{\partial w_{11}} &= -\frac{t_1}{o_1} \cdot \frac{\partial o_1}{\partial w_{11}} - \frac{t_0}{o_0} \cdot \frac{\partial o_0}{\partial w_{11}} \\ &= -\frac{t_1}{o_1} \cdot o_1(1-o_1)x_1 - \frac{t_0}{o_0} \cdot (-o_0 \cdot 0_1)x_1 \\ &= -t_1(1-o_1)x_1 + t_0 \cdot 0_1 x_1 \\ &= x_1(-t_1 + t_0) = -(t_1 - t_0)x_1 \end{aligned}$$

$\frac{\partial w_{01}}{\partial w_{11}} = 0$

$$\Delta w_{11} = -\eta \frac{\partial E}{\partial w_{11}} = \eta (t_1 - o_1) x_1$$