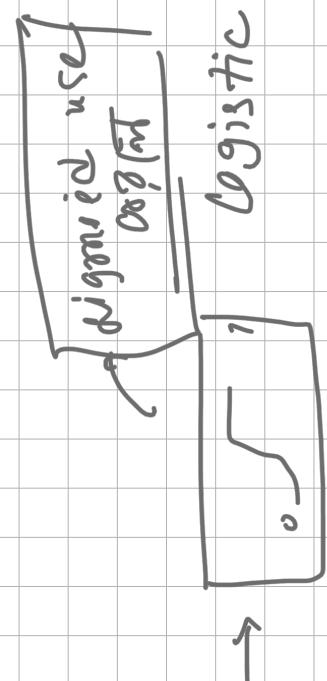
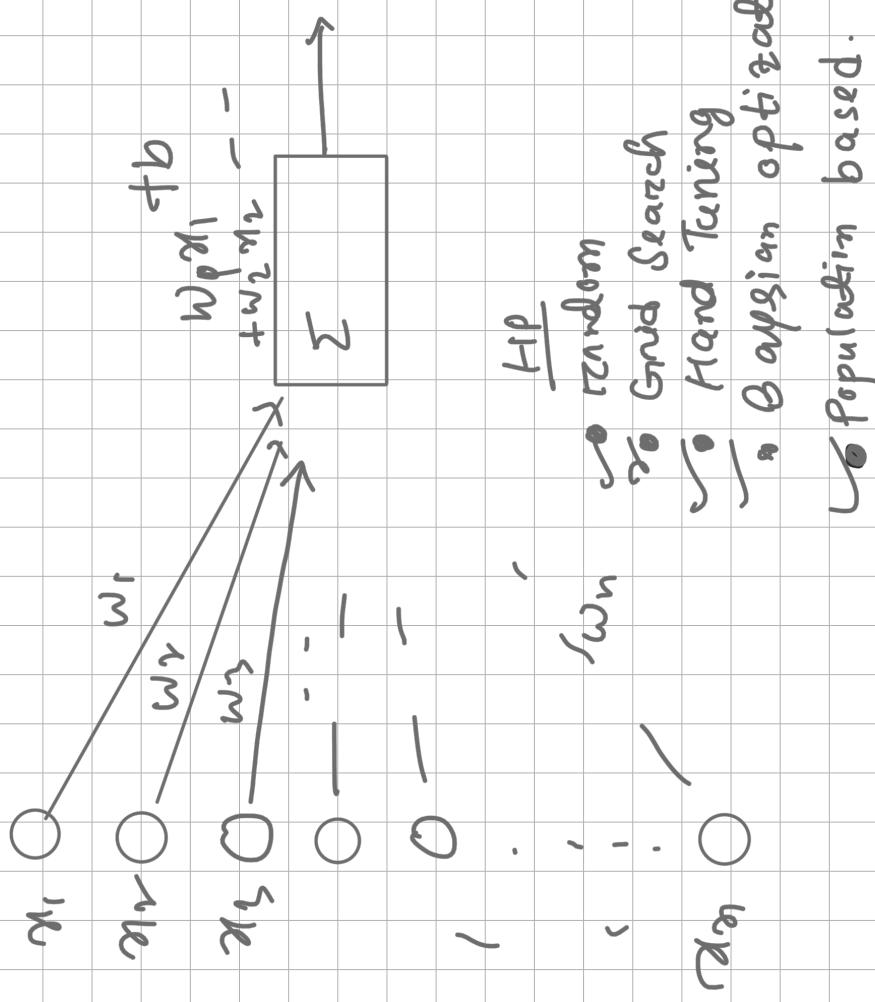


Machine learning



$o, 1$ use σ or softmax

linear.

- $o, 1$ use σ or softmax
- linear.
- \sum
- $\frac{1}{1 + e^{-x}}$ sigmoid
- softmax
- Grid Search
- Hand Tuning
- θ Bayesian optimization.
→ loss will be ordinal
→ λ, α will be continuous.
- population based.

Model Selection

$\text{err} \rightarrow []$

for i from 1 to ∞

$\text{err}[\text{size}] = \text{Cross Validation}(\text{training set}, K_i)$

if err increases significantly then,

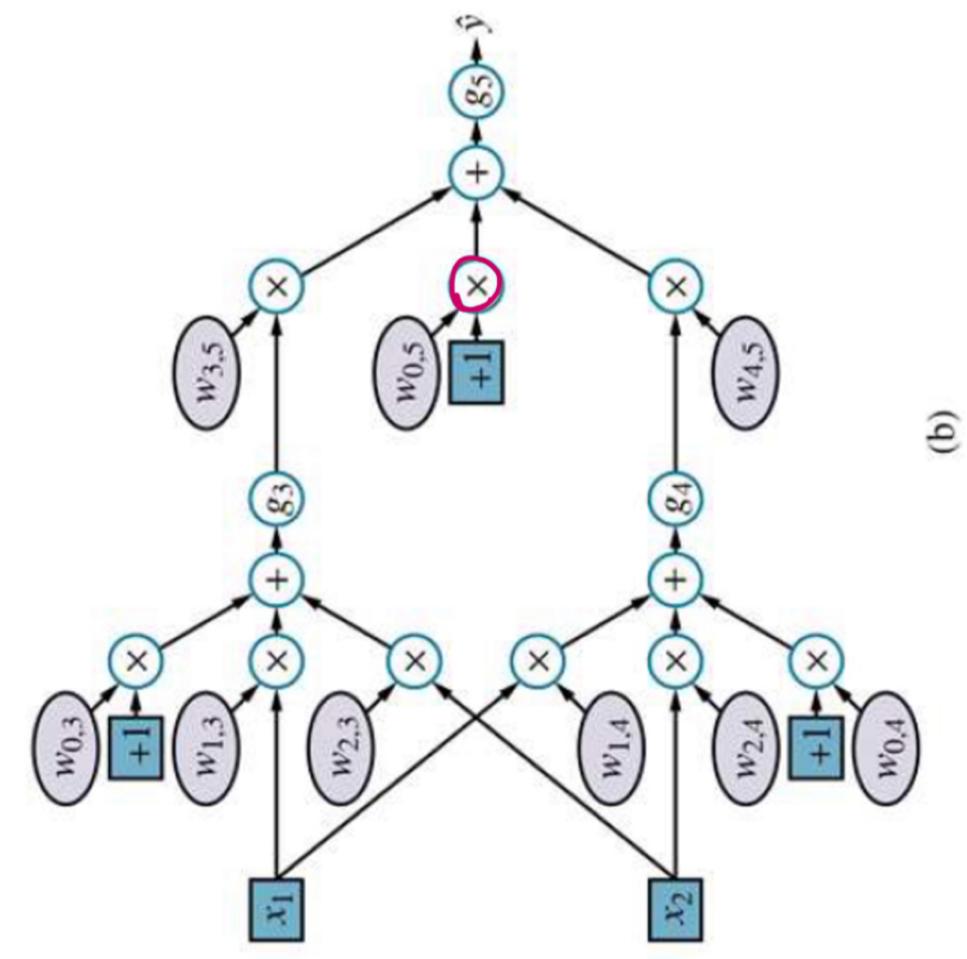
best size = $(\text{err}[\text{size}])$ again

$h = \text{learner}(\text{best size, training set})$

Return - $\text{ErronRate}(h, \text{test set})$



Now, for neural network implementation:



(b)

মুক্ত পরিমাণ

গুরুত্ব ক্রম

প্রক্রিয়া
প্রক্রিয়া
প্রক্রিয়া
প্রক্রিয়া
প্রক্রিয়া
প্রক্রিয়া

$$h(w) = a_5 = g(i_5)$$

$$= g(w_0 s + w_1 s_1 q_3 + w_4 s_4)$$

$$= g(w_0 s + w_1 s_1 g(i_3) + w_4 s_4)$$

অস্তর পদ্ধতি

অস্তর পদ্ধতি

forward

$$a = g(z)$$

$$= (z x_1)$$

$$= (b x_1) (h x_1)$$

$$z = w^T x$$

$$\frac{\partial \theta}{\partial \theta} \frac{\partial \theta}{\partial \theta} = \frac{m \theta}{n \theta}$$

layer 1

sum of features

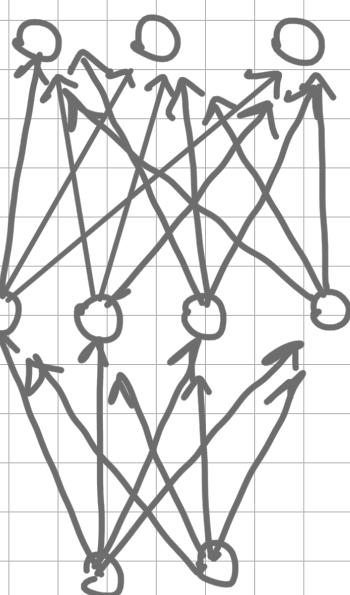
4 steps

$$w_{3,1}$$

$$w_{2,1}$$

$$w_{1,1}$$

$$w_{0,1}$$



backward pass

a

(da)

backward pass formula's :

$w_i(j) = \sum a_j w_i(j)$

$\rightarrow \delta_{01}$ in dense layer:

$$\delta^l = \frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} = \frac{\partial L}{\partial a} \cdot f'(z)$$

in the current layers

\hookrightarrow gradient

$$\boxed{\delta^{l-1} = \delta^l \cdot W^l}$$

derivative of
the activation
function

$$\frac{\partial L}{\partial w} = \delta^l \cdot \alpha^{l-1}$$

δ_{01} in backprop:

• grad out = activation. back (grad out)

$$(y_{01})$$

$(y_{01}) \rightarrow$ grad - out

- 0 grad weights = grad - out . self - input^T
 $\begin{pmatrix} y_{01} \end{pmatrix} = (y_{01}) \cdot q_{x_1} \rightarrow 1 \times 3$
- update weights = grad out^T . weight^T
 $L_3 \text{ from previous layer} = (1 \times 3) \cdot (1 \times 3)$
- grad input = grad out . weight^T
 $L_2 \text{ from previous layer} = (1 \times 3) \cdot (1 \times 3)$

z

$d\hat{w}$

Dropout

Adam optimieren

$$m_t = \beta_1 * m_{t-1} + (1 - \beta_1) * g_t$$

$$v_t = \beta_2 * v_{t-1} + (1 - \beta_2) * g_t^2$$

$$\hat{m}_t = m_t / (1 - \beta_1^t)$$

$$\hat{v}_t = v_t / (1 - \beta_2^t)$$

$$\theta = \theta - (\alpha * \hat{m}_t / \sqrt{(\hat{v}_t + \varepsilon)})$$

$\beta_1 \rightarrow$ Decay rate of momentum

$\frac{gt}{Ls \text{ grad}}$

$\beta_2 \rightarrow$ Decay rate of
Squared gradient

$\alpha \rightarrow$ learning rate / step size

Softmax

$$\Leftrightarrow f(x_i) = \frac{e^{x_i}}{\sum e^x} = \frac{e^{x_i}}{e^{x_1} + e^{x_2} + \dots + e^{x_n}}$$

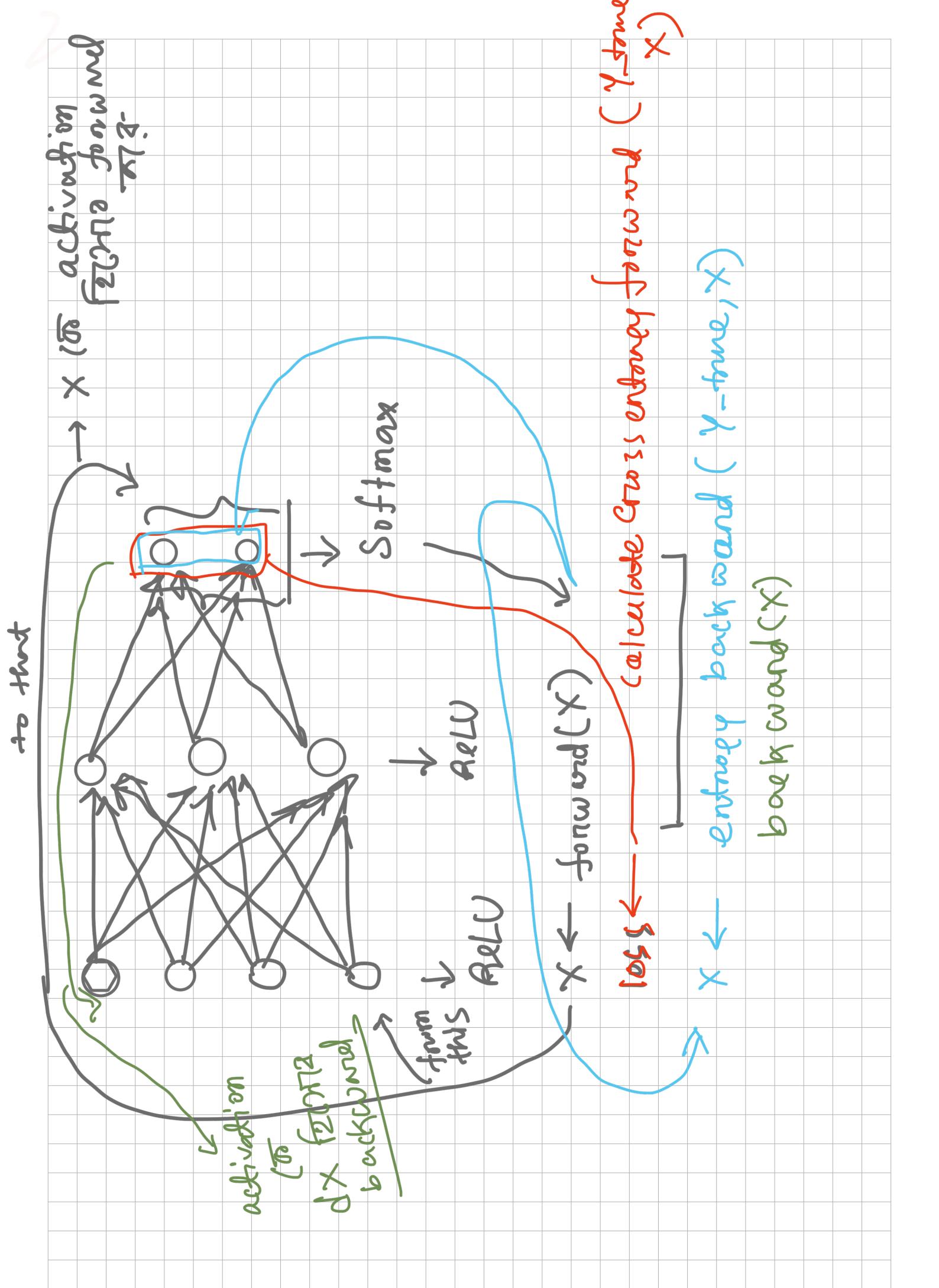
need to
check
for
overflow

$$\text{and, } f'(x_i) = \frac{\frac{d}{dx} e^{x_i} \cdot \sum e^x - e^{x_i} \frac{d}{dx} \sum e^x}{(\sum e^x)^2}$$

So, dom yahoo { } .

$$\frac{x_i (\sum e^x - e^{x_i})}{(\sum e^x)^2}$$

$$\frac{\partial}{\partial x_i}$$



für Softmax - detälin:

$$\frac{\partial L}{\partial z_i} = \frac{\partial f(z_i)}{\partial z_i} \cdot \frac{\partial z_i}{\partial z_i} = \Delta_i f(z_i) (1 - f(z_i))$$

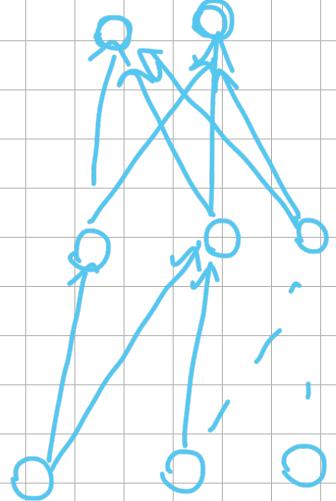
$$\begin{aligned}\frac{\partial L}{\partial z_j} &= \frac{\partial}{\partial z_i} f(z_j) \cdot \frac{\partial z_i}{\partial z_j} = \Delta_j f(z_i) f(z_j) (-1) \\ &= -\Delta_j f(z_i) f(z_j) f(z_i)\end{aligned}$$

Now i + ~~now~~: if $i \neq j$; $\Delta_i f(z_i) \neq 0 - \sum \Delta_j f(z_j)$

$$f(z_i) \left\{ \Delta_i - \sum \Delta_j f(z_j) \right\}$$
$$\frac{\partial L}{\partial z_i}$$

x input

Relu



$$z = xw$$

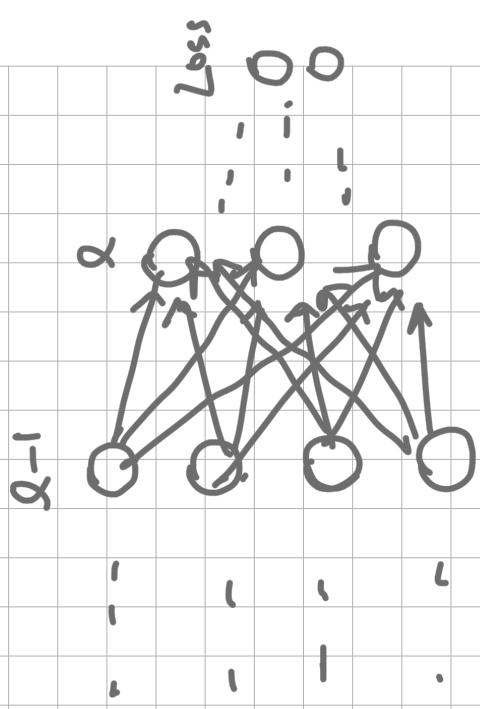
$$a = \delta(z)$$

loss function

$[x, z, a]$

$$-\rho \log p - (1-p) \log(1-p)$$

$$\begin{aligned} a &= \delta(z) \\ z &= xw \end{aligned}$$



zero activation in layer 0

$\alpha = \delta(z)$

no need of computing here

$$\frac{\partial}{\partial x} = w^T \alpha^{l-1}$$

$$L = -y \log \alpha - (1-y) \log(1-\alpha)$$

$$L = -y \log \delta(z) - (1-y) \log(1-\delta(z))$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial w}$$

$$= \frac{\partial L}{\partial \alpha^l} \cdot \frac{\partial \alpha^l}{\partial z} \cdot \frac{\partial z}{\partial w}$$

$$= \left(\frac{y}{\alpha} + \frac{1-y}{1-\alpha} \right) \alpha(1-\alpha) \frac{\partial z}{\partial w}$$

$$= \left(-\frac{y}{\alpha} + \frac{1-y}{1-\alpha} \right) \alpha(1-\alpha)$$

for loss activation in layer 1

we use it
initially
in layer 0
as layer 1

Square

in

set

com

see

co.

$$\frac{d^2x}{dt^2} = \frac{dp}{dq^{-1}} - \frac{d^2q}{dt^2}$$

actually,

$$M \cdot \frac{d^2q}{dt^2} =$$

$$1 - \frac{d^2p}{dt^2} \cdot \frac{d^2q}{dp} = \frac{-d^2p}{dp}$$

So,

$$Z_1' = M \cdot q^{-1} \text{ and}$$

Now, again $Z_1' = \frac{d^2q}{dt^2}$

* ignoring β for batch norm.

$$\mu = \frac{1}{m} \sum x_i$$

$$\sigma^2 = \frac{1}{m} \sum (x_i - \mu)^2$$

$$d_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$y_i = \gamma d_i \Leftrightarrow y_i = \gamma \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$



Suppose L is a function of y .

$$\frac{\partial L}{\partial \epsilon} = \frac{\partial L}{\partial y_i} \cdot \frac{\partial y_i}{\partial x_i} = \frac{\partial L}{\partial y_i} \cdot \frac{\partial}{\partial x_i} \left(\gamma \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \right) = \frac{\partial L}{\partial y_i} \cdot \gamma \frac{1}{\sqrt{\sigma^2 + \epsilon}} \cdot \frac{\partial}{\partial x_i} (x_i - \mu)$$

$$\Delta x_i = \Delta \hat{x}_i$$

$$\frac{\partial L}{\partial x_i} = \frac{\partial L}{\partial y_i^0} \cdot \frac{\partial y_i^0}{\partial x_i}$$

$$\hat{x}_i = \frac{x_i - \mu}{\sqrt{8v + \epsilon}}$$

$$= \sum \frac{\partial L}{\partial x_i} \cdot \frac{\partial y_i^0}{\partial x_i}$$

$$\frac{\partial y_i^0}{\partial x_i}$$

and,

$$= \sum \frac{\partial L}{\partial x_i} \cdot \frac{\partial y_i^0}{\partial x_i} \cdot \frac{\partial y_i^0}{\partial \mu} \cdot \frac{\partial \mu}{\partial v + \epsilon}$$

$$= -\frac{\partial L}{\partial x_i} \cdot \frac{1}{2} (x_i - \mu) (v + \epsilon)^{-3/2}$$

$$\hat{e} + \lambda \hat{e}$$

$$- \frac{\partial L}{\partial \dot{x}_i} (x_i - \bar{x})$$

$$+ \frac{\partial L}{\partial \dot{x}_i} \frac{(x_i - \bar{x})}{\sqrt{\delta^2 + \epsilon}} =$$

$$= \frac{\partial L}{\partial \dot{x}_i} \frac{(x_i - \bar{x})}{\sqrt{\delta^2 + \epsilon}}$$

$$\hat{e} + \lambda \hat{e}$$

$$= \frac{\partial L}{\partial \dot{x}_i} \frac{(x_i - \bar{x})}{\sqrt{\delta^2 + \epsilon}}$$

$$= \frac{\partial L}{\partial \dot{x}_i} \frac{(x_i - \bar{x})}{\sqrt{\delta^2 + \epsilon}}$$

$$\frac{\partial L}{\partial \dot{x}_i} \frac{(x_i - \bar{x})}{\sqrt{\delta^2 + \epsilon}} =$$

$$\begin{aligned} & \frac{\partial L}{\partial \dot{x}_i} \frac{(x_i - \bar{x})}{\sqrt{\delta^2 + \epsilon}} = \\ & = - \frac{\partial L}{\partial x_i} \frac{1}{\sqrt{\delta^2 + \epsilon}} + \frac{\partial L}{\partial x_i} (x_i - \bar{x}) (\delta^2 + \epsilon)^{-1/2} \cdot \frac{1}{m} \sum (x_i - \bar{x}) \\ & = - \frac{\partial L}{\partial x_i} \frac{1}{\sqrt{\delta^2 + \epsilon}} + \sum - \frac{1}{2} \frac{\partial L}{\partial x_i} (x_i - \bar{x}) (\delta^2 + \epsilon)^{-1/2} \cdot (-2) \frac{1}{m} (x_i - \bar{x}) \end{aligned}$$

$$= -\frac{\partial \varphi_i}{\partial x_1} \cdot \sqrt{g_{11}} + \frac{\partial \varphi_i}{\partial x_2} - 2 \frac{\partial}{\partial x_i} (\varphi_i - 1)$$