

$$\frac{\partial E}{\partial w_{ij}} =$$

~~$$\frac{\partial E}{\partial w_{jk}}$$~~

$$\frac{\partial E}{\partial w_{ik}} = \underbrace{\frac{\partial E}{\partial y_k} g'(z_k) y_j}_{\delta_k}$$

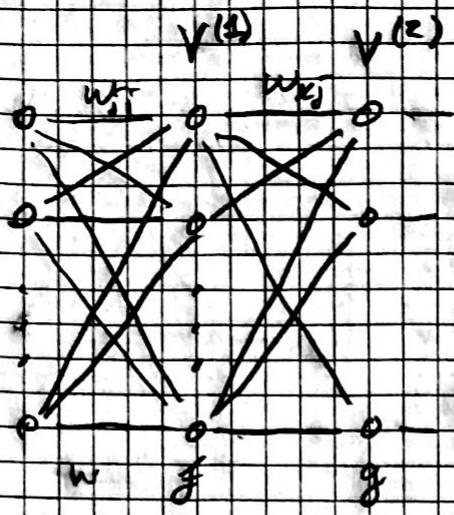
$$\begin{aligned}\frac{\partial E}{\partial w_{ij}} &= \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial w_{ij}} = \underbrace{\frac{\partial E}{\partial y_k} g'(z_k) w_{jk} g'(z_j) y_i}_{\delta_k} \\ &= \underbrace{\delta_k w_{jk} g'(z_j)}_{\delta_j} y_i = \delta_j y_i\end{aligned}$$

$$\frac{\partial E}{\partial w_x} = \frac{\partial E_{classical}}{\partial w_x} + \frac{\alpha}{2} \frac{\partial \|w\|^2}{\partial w_x}$$

$$E = E_{class} + \alpha \frac{1}{2} \|w\|^2$$

$$\frac{\alpha}{2} \frac{\partial \|w\|^2}{\partial w_x} = \alpha \|w\|$$

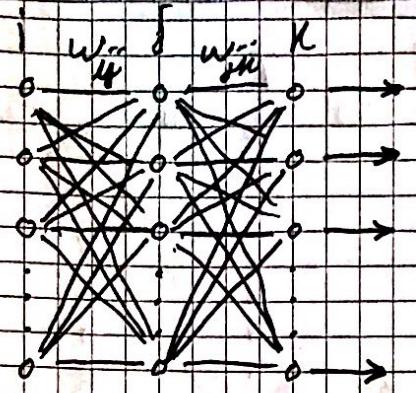
$$\|w\| = \sqrt{\sum (w_i)^2}$$



$v^{(1)} = \text{tan} \ \text{tanh}$ J
 $v^{(2)} = \text{tan} \ \text{out}$ K

$$v^{(1)} = f(h^{(1)}) \quad v^{(2)} = g(h^{(2)})$$

$$\frac{\partial E}{\partial w_{kj}} = C \cdot v^{(2)}(1 - v^{(2)}) \cdot v^{(1)}$$



$$w_{ij} = \text{input_id - hid}$$

$7 \times 256 \quad M$

$$w_{jk} = \text{hid_id - class}$$

$10 \times 7 \quad M$

$$\text{err} = \text{class_prob} - \text{dam_target}$$

$$y_j = \text{hid_output}$$

$$y_{jk} = \text{class_prob}$$

Wdm coefficient = learning rate weight decay

$$\text{model_to_them} \quad \begin{bmatrix} 1 & 2 & \dots \\ 2 & 1 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 \\ 1 & 3 \\ \vdots & \vdots \end{bmatrix}$$

grad should remain

$$\partial_i^m V_j^{m-1} \quad \text{in model form}$$

$$\left[\sum_k^m \partial_k^m \right] g'(h_k^m) V_j^m \underbrace{\quad}_{\text{class_prob}}$$

model.hid-to-dam

$$\Delta_n = \text{class_grad} - \partial_n - \text{temp}$$

$$\partial_n^m$$

$$\partial_j^m = g'(h_j^m) \sum_k W_{jk}^m \Delta_k$$

$$f(x)(1-f(x))$$

$$[\Sigma - \mathbf{0}^M] \text{ err}$$

g(

$$\text{err} = (\text{Class_Prob} - \text{data.target})$$

$$\text{cost} = \text{transpose}(\text{err}) * \text{model.wid_co_class}$$

2.4

g)

a) Validation data classification cost

0,4302

b) Validation data classification cost

0,3345

c)

10 22,6128

1 2,3026

0,0001 0,3483

0,001 0,2879

5 2,3026

d)

10 0,3708 0,4117

30 0,2868 0,3171

100 0,3021 0,3686

130 0,3169 0,3976

200 0,3345 0,4302

c)

18	0,3061
37	0,2652
83	0,3112
113	0,3137
236	0,3438

d)

$$n_hid = 37$$

$$\text{learning rate} = 0,35$$

$$\text{momentum} = 0,9$$

$$\text{early stop} = \text{true}$$

$$\text{mini-batch} = 100$$

$$Wd = 0$$

$$\text{test cost} \quad 0,2825$$

$$Wd = 0,001$$

$$\text{test cost} \quad 0,2571$$

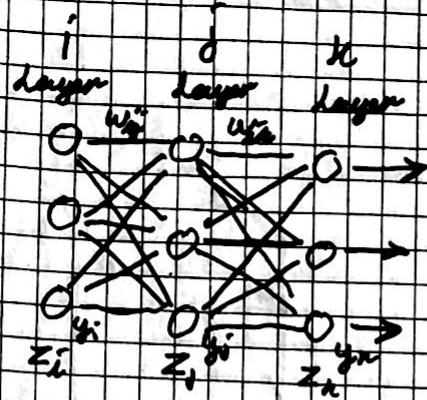
cross over derivation section
of

cross entropy

$$\begin{aligned}\frac{\partial H}{\partial \theta_i} &= - \sum_n y_n \frac{\partial \log(p_n)}{\partial \theta_i} \\ &= - \sum_n y_n \frac{\partial \log(p_n)}{\partial p_i} \times \frac{\partial p_i}{\partial \theta_i} = \\ &= - \sum_n y_n \frac{1}{p_n} \times \frac{\partial p_i}{\partial \theta_i}\end{aligned}$$

now seeing p_n as the soft max function we get

$$\begin{aligned}&-y_i(1-p_i) - \sum_{n \neq i} y_n \frac{1}{p_n} (-p_n p_i) \\ &= -y_i + y_i p_i - \sum_{n \neq i} y_n p_i \\ &= p_i \sum_n y_n - y_i \\ &\sum_n y_n = 1 \Rightarrow \\ &p_i - y_i\end{aligned}$$



$y_i = g(z_i)$

$$z_i = \sum_j w_{ij} y_j$$

$$y_j = g(z_j)$$

$$z_n = \sum_i w_{ni} y_i$$

$$y_n = g(z_n)$$

$g(z)$ is some differentiable function

incorrect error derivative

$$\frac{\partial E}{\partial y_k}$$

~~$$\frac{\partial E}{\partial z_n} = \frac{\partial E}{\partial y_n} \frac{\partial y_n}{\partial z_n} = \frac{\partial E}{\partial y_n} \frac{\partial y_n}{\partial z_n}$$~~

(a) $\frac{\partial E}{\partial z_n} =$

$$\frac{\partial E}{\partial y_n} = \frac{\partial E}{\partial z_n} \sum_j \frac{\partial z_n}{\partial w_{nj}} y_j$$

b) $\frac{\partial E}{\partial z_j}$

$$z_n = W_{jn}^T \cdot y_j$$

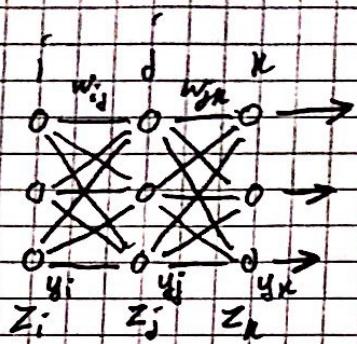
$$\frac{\partial E}{\partial w_{jn}} = \frac{\partial E}{\partial z_n} \frac{\partial z_n}{\partial W_{jn}^T}$$

c) $\frac{\partial E}{\partial w_{ji}}$

$$\frac{\partial z_n}{\partial w_{ji}} = \frac{\partial (W_{jn} y_j)}{\partial W_{ji}} = y_i$$

d) $\frac{\partial E}{\partial w_{ii}}$

~~$$\frac{\partial E}{\partial w_{ii}} = \frac{\partial E}{\partial z_n} = y_i \frac{\partial E}{\partial z_n}$$~~



$$y_i = g(z_i)$$

$$z_j = \sum_i w_{ij} y_i$$

$$y_j = g(z_j)$$

$$z_k = \sum_j w_{jk} y_j$$

$$y_k = g(z_k)$$

$$y_j = g(z_j)$$

$$z_j = \sum_i w_{ij} y_i = \sum_i w_{ij} g(z_i)$$

$$y_j = g(z_j) = g\left(\sum_i w_{ij} g(z_i)\right)$$

$$z_k = \sum_j w_{jk} y_j = \sum_j w_{jk} g\left(\sum_i w_{ij} g(z_i)\right)$$

$$y_k = g\left(\sum_j w_{jk} g\left(\sum_i w_{ij} g(z_i)\right)\right)$$

$$\frac{\partial E}{\partial y_k} = C$$

$$\frac{\partial E}{\partial z_k} = \frac{\partial E}{\partial y_k} \frac{\partial y_k}{\partial z_k}$$

$$\frac{\partial E}{\partial w_{jk}} = \sum_i \mu$$

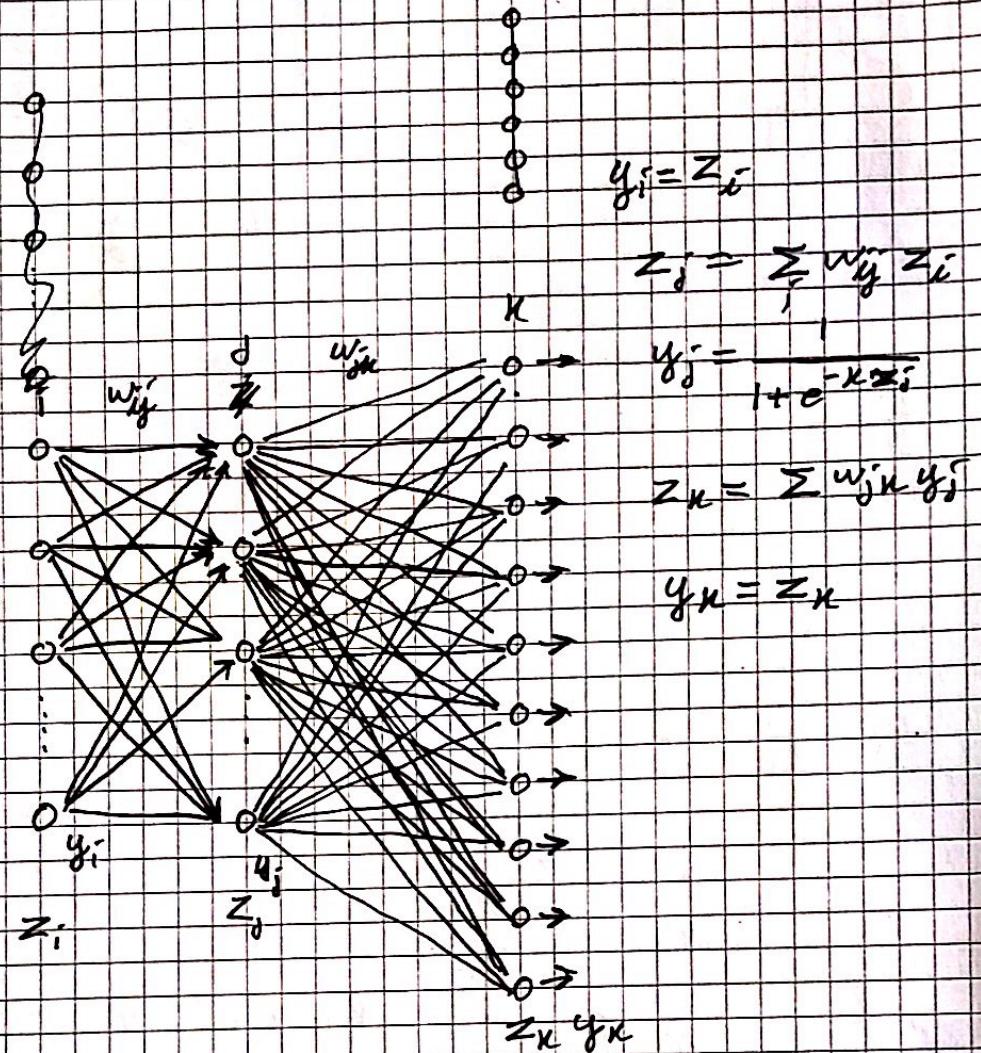
$$\frac{\partial E}{\partial z_j} = \frac{\partial E}{\partial y_k} \frac{\partial y_k}{\partial z_k} \frac{\partial z_k}{\partial y_j} \frac{\partial y_j}{\partial z_j}$$

$$\frac{\partial E}{\partial z_j} = C g'(z_k) \sum_j w_{jk} g'(z_j)$$

$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial E}{\partial y_k} g'(z_k) \cdot y_j$$

$$\frac{\partial E}{\partial w_{ij}} + \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial w_{jk}}$$

$$\frac{\partial E}{\partial y_j} = \frac{\partial E}{\partial z_j}$$



$$E = E_{\text{classification}} + \alpha E_{\text{weight decay}}$$

$$\frac{\partial E_{\text{classification}}}{\partial z}$$

$$\frac{\partial E}{\partial z} = \frac{\partial E_{\text{classification}}}{\partial z} + \alpha \left| \frac{\partial E_{\text{classification}}}{\partial z} \right|$$

$$y_j = \frac{1}{1 + e^{-kz_j}}$$

$$y'_j = \frac{k e^{kz_j}}{(1 + e^{kz_j})^2}$$

$$\delta_i^M = g'(h_i^M) (\xi_i^M - v_i^M)$$

$$\frac{\partial E}{\partial z}$$

$$\frac{\partial E}{\partial y_k} = g'(z_k) [\xi_k^M - v_k]$$

$$\frac{\partial E}{\partial y_j} = g'(z_i) \sum w_{jk} \frac{\partial E}{\partial y_k}$$