

$$\textcircled{1} \text{ P.T. } \sigma(x) = \sigma(x+c).$$

$$\sigma(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}} \rightarrow \textcircled{1}.$$

$$\sigma(x_i+c) = \frac{e^{x_i+c}}{\sum_j e^{x_j+c}} = \frac{e^{x_i} \cdot e^c}{\sum_j e^{x_j} \cdot e^c}$$

$$= \frac{e^c \cdot e^{x_i}}{e^c \sum_j e^{x_j}} = \frac{e^{x_i}}{\sum_j e^{x_j}} \rightarrow \textcircled{2}.$$

$\textcircled{1} = \textcircled{2}$. Hence proved.

$$\frac{d}{dx} \sigma(x) = ? \quad \frac{d}{dx} \sigma(x) = \frac{d}{dx} \left(\frac{1}{1+e^{-x}} \right)$$

$$= \frac{d}{dx} (1+e^{-x})^{-1} = -1 (1+e^{-x})^{-2} \left[\frac{d}{dx} (1+e^{-x}) \right]$$

$$= - (1+e^{-x})^{-2} \cdot [0 + (-e^{-x})] = \frac{e^{-x}}{(1+e^{-x})^2}$$

$$= \underbrace{\frac{e^{-x}}{1+e^{-x}}}_{1-\sigma(x)} \cdot \underbrace{\frac{1}{(1+e^{-x})}}_{\sigma(x)}$$

$$= \sigma(x) \cdot [1 - \sigma(x)]$$

(2b)

$$\hat{y} = \text{softmax}(\theta)$$

$$CE(y, \hat{y}) = - \sum_i y_i \cdot \log(\hat{y}_i)$$

$$\frac{\partial}{\partial \theta} \cdot CE(y, \hat{y}) = ?$$

Sol:

$$\frac{\partial}{\partial \theta} \left[e - \sum_i y_i \cdot \log(\hat{y}_i) \right] = ?$$

$$= - \frac{\partial}{\partial \theta} \left[\sum_i y_i \log \left(\frac{\text{softmax}(\theta_i)}{\text{sigmoid}(\theta_i)} \right) \right] = ?$$

y is one-hot-encoded vector.

\therefore choose k , s.t. $y_k \neq 0$

$$= - \frac{\partial}{\partial \theta} \left[y_k \cdot \log \left(\frac{\text{softmax}(\theta_k)}{\text{sigmoid}(\theta_k)} \right) \right] = ?$$

$$= - \frac{\partial}{\partial \theta} \cdot \left[y_k \cdot \log \frac{e^{\theta_k}}{\sum_{j=1}^d e^{\theta_j}} \right].$$

$$= - \frac{\partial}{\partial \theta} \left[y_k \cdot \left(\log e^{\theta_k} - \log \sum_{j=1}^d e^{\theta_j} \right) \right]$$

$$= - \frac{\partial}{\partial \theta} \left[y_k \cdot \left(\theta_k - \log \sum_{j=1}^d e^{\theta_j} \right) \right]$$

$\textcircled{*}$ choose i , s.t. $i = k$,

$$= - \frac{\partial}{\partial \theta_i} \left[y_k \cdot \left(\theta_k - \log \sum_{j=1}^d e^{\theta_j} \right) \right]$$

$$= - \left[y_{ik} \cdot \left(1 - \frac{1}{\sum_{j=1}^d e^{\theta_j}} e^{\theta_k} \right) \right]$$

$$= - \left[y_{ik} \cdot (1 - \text{softmax}(\theta_k)) \right]$$

$$= - \left[y_{ik} \cdot (1 - \hat{y}_k) \right] = (\hat{y}_k - 1) \cdot y_{ik}.$$

But $y_{ik} = 1$

$$\Rightarrow \frac{\partial}{\partial \theta_k} CE(y, \hat{y}) = (\hat{y}_k - 1), \rightarrow \textcircled{1}$$

choose \hat{y}_g s.t. $\hat{y} \neq \hat{y}_k$.

$$\therefore - \frac{\partial}{\partial \theta_g} \left[y_{ik} \left(\theta_k - \log \sum_{j=1}^d e^{\theta_j} \right) \right]$$

$$= - \left[y_{ik} \cdot \left(0 - \frac{1}{\sum_{j=1}^d e^{\theta_j}} \cdot e^{\theta_g} \right) \right]$$

$$= - \left[y_{ik} \cdot (-\text{softmax}(\theta_g)) \right]$$

$$= y_{ik} \cdot \text{softmax}(\theta_g) = y_{ik} \cdot \hat{y}_g = \hat{y}_g. \rightarrow \textcircled{2}$$

From $\textcircled{1} + \textcircled{2}$,

$$\frac{\partial}{\partial \theta_k} CE(y, \hat{y}) = \hat{y}_{ik} - 1 \rightarrow \textcircled{3}$$

$$\frac{\partial}{\partial \theta_g} CE(y, \hat{y}) = \hat{y}_g - 0 \rightarrow \textcircled{4}$$

But, $y_k = 1$ & $y_g = 0$. in ~~without~~ in 1-hot encoding format
 $\therefore \textcircled{3} + \textcircled{4} \Rightarrow$

$$\frac{\partial}{\partial \theta_k} CE(y, \hat{y}) = \hat{y}_k - y_k$$

$$\frac{\partial}{\partial \theta_g} CE(y, \hat{y}) = \hat{y}_g - y_g$$

(generalize)

$$\boxed{\frac{\partial}{\partial \theta} CE(y, \hat{y}) = \hat{y} - y}$$

$$Q) J = CE = - \sum_{i=1}^N y_i \cdot \log(\hat{y}_i)$$

$$\hat{y}_i = \text{softmax}(h_i w_2 + b_2)$$

$$h = \text{sigmoid}(x w_1 + b_1)$$

$$\frac{\partial J}{\partial x} = ?$$

$$J = - \sum_{i=1}^N y_i \log \left(\frac{\text{softmax}(\theta_i)}{\text{sigmoid}(\theta_i)} \right)$$

$$\therefore J(\theta).$$

$$\theta = h w_2 + b_2$$

$$\therefore \theta(h).$$

$$h = \text{sigmoid}(z); z = x w_1 + b_1$$

$$\therefore \theta(h(z)).$$

$$z = x \cdot w_1 + b,$$

$\therefore z$ (Ans).

$$\therefore \frac{d J}{d x} = \frac{d J(\theta)}{d \theta} \cdot \frac{d \theta(h)}{d h} \cdot \frac{d h(z)}{d z} \cdot \frac{d z(x)}{d x}.$$

From previous question, $\frac{d J(\theta)}{d \theta} = (\hat{y} - y) \rightarrow ①$.

$$\frac{d \theta(h)}{d h} = w_2 \rightarrow ②$$

$$\frac{d h(z)}{d z} = \sigma(z) \circ \sigma(1-z) \rightarrow ③$$

element-wise prod.

$$\frac{d z(x)}{d x} = w_1 \rightarrow ④$$

$$\frac{d J}{d x} = (\hat{y} - y) * (w_2) \circ (\sigma(z) \circ \sigma(1-z)) * w_1$$

\downarrow
 $1 \times D_y \cdot D_y \times D_h \quad \cdot \quad \cancel{1 \times D_h} \quad 1 \times D_h \quad \cdot \quad D_h \times D_x$
 $1 \times D_h \quad \cdot \quad 1 \times D_h \quad \cdot \quad D_h \times D_x$
 $= 1 \times D_h \times D_h \times D_x$
 $\sim \underline{\underline{1 \times D_x}}$

② $D_x \quad H \quad D_y$

$$\text{No of params} = ((1+D_x) * H) + ((H+1) * D_y)$$

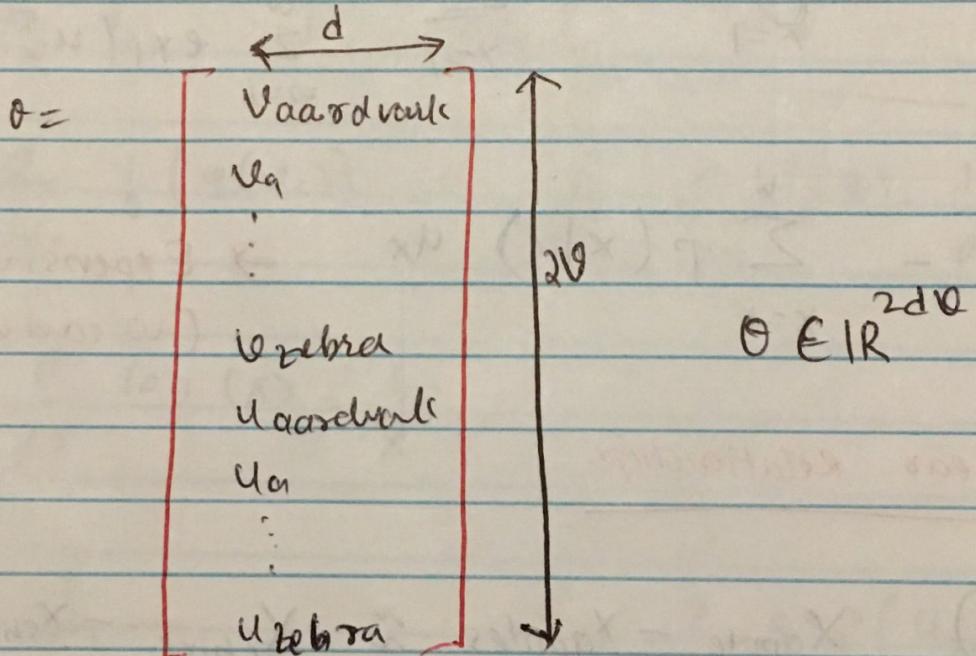
CNF

$$w=1$$

Today I am giving a lecture

v_{today} , v_I , v_{am}

$$p(\text{old}) = \frac{\exp(v_{\text{today}}^T \cdot v_I)}{\sum} \quad p(\text{old}) = \frac{\exp(v_{\text{am}}^T \cdot v_I)}{\sum}$$



$$\theta^{\text{new}} = \theta^{\text{old}} - \nabla_{\theta} J(\theta)$$

- Use SGD instead of batch gradient descent
- Use small random numbers as initial values.

Stochastic Gradient Descent with word vectors:

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} 0 \\ \vdots \\ \nabla_{U_{like}} \\ \vdots \\ 0 \\ \nabla_{U_I} \\ \vdots \\ \vdots \\ \nabla_{U_{learning}} \\ \vdots \\ -0 \end{bmatrix} \Rightarrow \text{very sparse.}$$

$\sigma(x) = \frac{1}{1+e^{-x}}$

IIP: $-\infty \leq x \leq +\infty$

OIP: $0 \leq \sigma(x) \leq 1$

DL CS
DL NLP

Pset 1: Skip-Gram Model

+

Negative Sampling

skip-gram:

CBOW: Avg all outside words &
predict center word.

→ Continuous Bag of words.

Stochastic Gradient Descent with word vectors:

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} 0 \\ \vdots \\ \nabla_{V_{like}} \\ \vdots \\ 0 \\ \nabla_{U_I} \\ \vdots \\ \nabla_{V_{learning}} \\ \vdots \\ -0 \end{bmatrix} \Rightarrow \text{very sparse.}$$
$$\sigma(x) = \frac{1}{1+e^{-x}}$$

IIP: $-\infty \leq x \leq +\infty$
OIP: $0 \leq \sigma(x) \leq 1$

DL CS
DL NLP

Pset 1: Skip-Gram Model

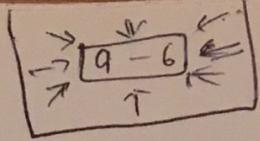
+
Negative Sampling

skip-gram:

CBOW: Avg all outside words &
predict center word.

→ Continuous Bag of words.

SVD + skip - gram,

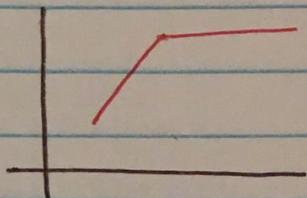


GloVe:

for every pair
of words.

$$\hat{J}(\theta) = \frac{1}{2} \sum_{i,j=1}^w f(p_{ij}) \left(\underbrace{u_i^T v_j}_{\text{Predicted co-occ stat}} - \underbrace{p_{ij}}_{\text{Actual co-occurrence statistics}} \right)^2$$

→ scalable.



Hyperparameters

① $d \approx 300$

② $w \approx 8$

③ time

Softmax:

Sigmoid : 2 class classification

Softmax : Multiclass classification.

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

$$\sigma(\theta^T x) = \frac{1}{1+e^{-\theta^T x}}$$

$$W = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix}$$

$$x = [x_1 \ x_2 \ x_3]$$

→ Prep
 → Adverb
 →

$$P(y|x) = \frac{\exp(w_y x)}{\sum_{c=1}^C \exp(w_c \cdot x)}$$

$$\frac{a}{\text{count}} \frac{b}{\sum} \frac{c}{\epsilon}$$

$$w \in \mathbb{R}^{C \times d}$$

$$w \in \mathbb{R}^{C \times d}$$

→ How often?
 → prep.
 → Bias?

Max	$p(y x)$	→ How
or Min	$-p(y x)$	→ Kappa value
or Min	$-\log(p(y x))$	→ Pmi

Cross-entropy:

$$p = [0, 1, 0]$$

$$q = [0.01, 0.9, 0.09]$$

$$H(p, q) = -\sum_c p_c * \log(q_c)$$

$$\sigma(x) = \sigma(x + c)$$

[1 2]

$$\sigma[1 2] \quad \sigma[2 3]$$

at 6

$$w = 1, 1$$

$$\frac{\exp(3)}{\exp(3) + \exp(6)}$$

$$\frac{\exp(5)}{\exp(5) + \exp(10)}$$

$$\frac{\exp(1)}{\exp(1) + \exp(2)}$$

$$\frac{\exp(2)}{\exp(1) + \exp(3)}$$

$$\frac{1}{1+2} \quad \frac{2}{1+2}$$

$$2^{2+3} = 2^2 \cdot 2^3$$

$$2^5 = 2^2 \cdot 2^3$$

$$\frac{\exp(x_1)}{\exp(x_1) + \exp(x_2)}$$

$$\frac{\exp(x_1+c)}{\exp(x_1+c) + \exp(x_2+c)}$$

$$= \frac{\exp(x_1) \cdot \exp(c)}{\exp(c) [\exp(x_1) + \exp(x_2)]}$$

~~$$\sigma(x) \quad \sigma(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}} \rightarrow ①$$~~

$$\sigma(x_i+c) = \frac{e^{(x_i+c)}}{\sum_j e^{(x_i+c)}} = \frac{e^{x_i} \cdot e^c}{\sum_j e^{x_j} \cdot e^c}$$

∴

$$= \frac{e^{x_i} / e^c}{e^c \cdot \sum_j e^{x_j}} = \frac{e^{x_i}}{\sum_j e^{x_j}} \rightarrow ②$$

① = ② Hence proved.

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$$

$$\frac{e^{x_i}}{\sum_j e^{x_j}}$$

$$x = [1 \ 2]$$

$$\begin{matrix} 1 \times 2 \\ 2 \times 1 \\ 2 \times 2 \end{matrix}$$

np.sum

$$x = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

$$a = [1 \ 2]$$

$$b = [3 \ 4] \text{ axis}$$

RNN, CNN, LSTM.

$$x = np.exp(x)$$

x.reshape()

$$[1 \ 2]$$

$$[[1, 2]]$$

$$(2, 1)$$

$$(1, 2)$$

$$a^b$$

$$-1, \text{len}(x)$$

$$e^x$$

$$e^{\frac{e^a}{e^a + e^b}}$$

$$x(0)$$

$$[1, 2]$$

$$[-1, 0] = \left[\frac{1}{e}, 1 \right]$$

$$1 \ 2 \quad 2 \ 4 \quad -1 \quad 0$$

$$3 \ 4 \quad 4 \ 2 \ 4 \quad -2 \quad 0$$

$$1 \ 3 \quad 2 \ 4$$

$$2 \ 4 \quad 2 \ 4$$

$$1 \ 2 \quad 2$$

$$3 \ 5 \quad 5$$

$$1 \ 3 \quad 2 \ 4$$

$$2 \ 5 \quad 2 \ 4 \quad -1 \quad 0$$

$$-2 \quad 0$$

$$1 \ 3 \quad 2 \ 5 \quad -1 \quad 1 \quad -1$$

~~2 5~~
$$2 \ 5 \quad 2 \ 5 \quad -2 \quad 0 \quad 1 \quad -2$$

$$x.T$$

$$-1 \ -2 \ -1 \ -2$$

$$0 \ 0 \ -1 \ -2$$

$$\sigma(x) = \frac{1}{1+e^{-x}} \cdot \frac{1}{e^x} \cdot \frac{\partial (x^{B_{xy}})^2}{\partial x} =$$

$$\frac{\partial}{\partial x} \sigma(x) = \frac{\partial}{\partial x} (1+e^{-x})^{-1}$$

$$\frac{\partial}{\partial x} e^{\frac{x}{g(x)}} \cdot g'(x) = -x$$

$$= -1 (1+e^{-x})^{-2} \cdot [0 + e^{-x} \cdot (-1)]$$

$$\begin{aligned} \frac{\partial (\sigma(x))}{\partial x} &= \frac{-1 - e^{-x}}{(1+e^{-x})^2} = \frac{-1(1+e^{-x})}{(1+e^{-x})^2} & y = e^{-x} \\ &= \frac{-1}{1+e^{-x}} = -\sigma(x). & \ln y = -x \end{aligned}$$

$$\therefore \frac{\partial}{\partial x} r(x) = r(x) (1 - \sigma(x)).$$

$$\frac{\partial}{\partial x} r(x) = \frac{\partial}{\partial x} (1+e^{-x})^{-1}$$

$$= -1 (1+e^{-x})^{-2} \cdot \frac{\partial}{\partial x} (1+e^{-x})$$

$$= -1 (1+e^{-x})^{-2} \cdot [0 + \frac{\partial}{\partial x} e^{-x}]$$

$$= -1 (1+e^{-x})^{-2} \frac{\partial}{\partial x} e^{-x} (-e^{-x}).$$

$$= \frac{e^x}{(1+e^{-x})^2} \frac{e^x}{1+e^{-x}} \cdot \frac{1}{1+e^{-x}} \cdot (1 - \sigma(x)) \sigma'(x).$$

$$\frac{e^{-x}}{g(x)} \quad f(x) = e^{g(x)} \quad g(x) = -x$$

$$\frac{\partial}{\partial x} f(x) = e^z$$

$$f'(x) =$$

$$\frac{d}{dx} f(x) = \frac{d}{dz} \cdot e^z \cdot \frac{d}{dx} -x$$

$$= e^z \cdot (-1) = -e^{-x}$$

$$\frac{\partial}{\partial x} e^{-x} = -e^{-x}$$

$$= f(z) = e^z \quad z = g(x) = -x.$$

$$\therefore [f(g(x))]$$

$$\Rightarrow \frac{d}{dz} f(z) \cdot \frac{d}{dx} g(x)$$

$$\Rightarrow e^z \cdot \frac{d}{dx} (-x) = -e^z = -e^{-x}.$$

$$\{(x_i, y_i)\}_{i=1}^N$$

$$\theta = \begin{bmatrix} w_{:,1} \\ \vdots \\ w_{:,d} \end{bmatrix} = w(:) \in \mathbb{R}^{cd}$$

$$\nabla_\theta J(\theta) = \begin{bmatrix} \nabla_{w,:1} \\ \vdots \\ \nabla_{w,:d} \end{bmatrix} \in \mathbb{R}^{cd}$$

$$L = \begin{bmatrix} 0 & 0 & \rho \\ 0 & 0 & b \\ 0 & b & 0 \\ \vdots & \ddots & \vdots \\ 0 & 0 & b \end{bmatrix}$$

$\leftarrow |V| \rightarrow$
 d z
 standvark zebra

... marmots in Paris are amazing..

$$x_{\text{window}} = [x_m \quad x_E \quad x_p \quad x_a \quad x_g]^T$$

y = location | Not location

$$x \in \mathbb{R}^{5d}$$

col vector.

$$0 \leq x \leq 1$$

$$-\infty \leq \log(x) \leq 0$$

$$x \rightarrow 0 \quad \log(x) \rightarrow -\infty$$

$$x \rightarrow 1 \quad \log(x) \rightarrow 0$$

$$\hat{y} = \text{softmax}(\theta)$$

$$y = [0, 1, 0]$$

$$\hat{y} = [0.01, 0.95, 0.04]$$

$$CE(y, \hat{y}) = - \sum_i y_i \log(\hat{y}_i)$$

$$\frac{d}{d\theta} \text{softmax}(\theta) = ? \quad r(\theta) \cdot [1 - r(\theta)]$$

$$= - \sum_i y_i [\log(r(\theta))]$$

$$= - \sum_i y_i \log($$

$$CE(y, \hat{y}) = - y_x \log(\hat{y}_x)$$

$$= - y_x \log(r(\theta_x))$$

$$\frac{\partial}{\partial \theta} CE(y, \hat{y}) = - \frac{\partial}{\partial \theta_x} y_x \log[r(\theta_x)]$$

$$= - \left[y_x \left[\log(r(\theta_x)) \right]' + \log(r(\theta_x))' y_x' \right]$$

$$= - \left[y_x \frac{1}{r(\theta_x)} \cdot r(\theta_x)' \cdot (1 - r(\theta_x)) \right]$$

$$= - \left[y_x \left[(1 - r(\theta_x)) \right] \right]$$

$$= \left[\frac{r(\theta_x)}{(r(\theta_x)) - 1} \right] \quad x; y_x = 1$$

$$x = [x_1 \ x_2] \quad A = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} \cdot \frac{d}{dx} \sigma(v) =$$

$\frac{1}{1+e^{-v}} - \frac{(1+e^{-v})}{(1+e^{-v})^2}$
 $\frac{-1}{(1+e^{-v})^2} [a_1 + a_2]$

$$x^T \cdot A = [x_1 \ x_2] \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \sum_i x_i A_i$$

$$\frac{d}{dx} x^T \cdot A = \left[\begin{array}{l} \frac{d}{dx_1} \sum_i x_i A_i = a_1 \\ \frac{d}{dx_2} \sum_i x_i A_i = a_2 \end{array} \right]$$

$$CE(y, \hat{y}) = - \sum_i y_i \cdot \log(\hat{y}_i)$$

$$\frac{d}{d\theta} CE(y, \hat{y}) = ? \left[\begin{array}{l} -\frac{d}{d\theta_1} \sum_i y_i \cdot \log(\hat{y}_i) \\ \vdots \\ -\frac{d}{d\theta_n} \sum_i y_i \cdot \log(\hat{y}_i) \end{array} \right]$$

$$[0 \ 1] = -\frac{d}{d\theta_1} \sum_i y_i \log(\sigma(\theta_1))$$

$$[\sigma(\theta_1) \ \sigma(\theta_2)]$$

$$= \begin{bmatrix} -y_1 (1 - \sigma(\theta_1)) \\ \vdots \\ -y_n (1 - \sigma(\theta_n)) \end{bmatrix}$$

$y_1 \cdot \frac{1}{\sigma(\theta_1)} \cdot \frac{\sigma'(\theta_1)}{\sigma(\theta_1)^2}$
 $y_1 \cdot \frac{1}{\sigma(\theta_1)}$

$$= \begin{bmatrix} -y_1 & (1 - \sigma \hat{y}_1) \\ \vdots & \vdots \\ -y_n & (1 - \sigma \hat{y}_n) \end{bmatrix} \quad \hat{y} - y$$

$$= \begin{bmatrix} y_1 & (\hat{y}_1 - 1) \\ \vdots & \vdots \\ y_n & (\hat{y}_n - 1) \end{bmatrix} \quad \hat{y} - y.$$

when ~~y_j~~ $j = i$

$$\begin{aligned} &= \hat{y} - y \\ &= \hat{y} - 1 \end{aligned}$$

$$\Rightarrow [0 \ 0 \ \dots \ \hat{y}_i - 1 \ \dots \ 0] \quad \hat{y} - y.$$

$$\text{when } j \neq i \quad \frac{\partial}{\partial \theta_j} C_{\theta}(y, \hat{y}) = \hat{y}_j - 1$$

$$\Rightarrow \boxed{y \neq (\hat{y} - 1)}$$

$$\text{when } i = k; \quad \frac{\partial}{\partial \theta_k} C_{\theta}(y, \hat{y}) = \hat{y}_i - 1$$

$$i \neq k; \quad \frac{\partial}{\partial \theta_i} C_{\theta}(y, \hat{y}) = 0$$

$\ln \theta$

$$\hat{y} = \text{softmax}(\theta)$$
$$CE = -\frac{1}{n} \sum_i y_i \cdot \log(\hat{y}_i).$$

~~$\frac{\partial}{\partial \theta}$~~ $\frac{\partial}{\partial \theta} CE = ?$

$$\frac{\partial}{\partial \theta} (CE(y, \hat{y})) \leftarrow \hat{y} = \text{softmax}(\theta)$$

$$\hat{y}_k = \frac{e^{\theta_k}}{\sum_{l=1}^L e^{\theta_l}}$$

$$\hat{y} = \text{softmax}(\theta)$$

$$CE = -\sum_i y_i \log(\hat{y}_i)$$

$\frac{\partial}{\partial \theta} (CE(y, \hat{y})) = ?$

$$\frac{\partial}{\partial \theta} \left[-\sum_i y_i \log(\hat{y}_i) \right] = ?$$

$$= \frac{\partial}{\partial \theta} \left[-\sum_i y_i \cdot \log(\text{softmax}(\theta_i)) \right] = ?$$

$$\Rightarrow \frac{\partial}{\partial \theta} \left[-\sum_i y_i \cdot \log \left(\frac{e^{\theta_i}}{\sum_{j=1}^L e^{\theta_j}} \right) \right] = ?$$

$$\log(\tilde{P}(\theta_{\text{fin}}))$$

$$= \left[\frac{\partial}{\partial \theta_1} \left[- \sum_i y_i \cdot \log \left(\frac{e^{\theta_i}}{\sum_{j=1}^d e^{\theta_j}} \right) \right] \right. \\ \vdots \\ \left. \frac{\partial}{\partial \theta_d} \left[- \sum_i y_i \cdot \log \left(\frac{e^{\theta_i}}{\sum_{j=1}^d e^{\theta_j}} \right) \right] \right]$$

$$= \frac{\partial}{\partial \theta_1} \left[- \sum_i y_i \cdot \log \left(\frac{e^{\theta_i}}{\sum_{j=1}^d e^{\theta_j}} \right) \right]$$

$$\frac{\partial}{\partial \theta_i} \log \frac{e^{\theta_i}}{\sum_{j=1}^d e^{\theta_j}} = \left[\log \exp(\theta_i) - \log \sum_{j=1}^d \exp(\theta_j) \right] \quad \frac{\partial}{\partial X} e^X = e^X$$

$$= \theta_i - \log \sum_{j=1}^d \exp(\theta_j) \quad [1 - \log]$$

$$\frac{\partial}{\partial \theta_i} \left(\theta_i - \log \sum_{j=1}^d \exp(\theta_j) \right) \quad \frac{\partial}{\partial \theta_i} \exp(\theta_i)$$

$$= 1 - \frac{\exp(\theta_i)}{\sum_{j=1}^d \exp(\theta_j)} \quad 1 - \text{softmax}(\theta_i)$$

$$\therefore \frac{\partial}{\partial \theta_X} \log(\text{softmax}(\theta_X)) = 1 - s(\theta_X)$$

$$\frac{\partial}{\partial \theta_1} \left(\theta_2 - \log \sum_{i=1}^d \exp(\theta_i) \right)$$

~~$\frac{\partial}{\partial \theta_1} = 0 - \frac{0}{\sum_{i=1}^d \exp(\theta_i)}$~~

$$y_1 \cdot \log(s(\theta_1)) + y_2 \cdot \log(s(\theta_2))$$

$$\frac{\partial}{\partial \theta_1} \left[- \sum_i y_i \cdot \log [\text{softmax}(\theta_i)] \right] = ?$$

$$\Rightarrow \frac{\partial}{\partial \theta_1} \left[- y_1 \cdot \log [\text{softmax}(\theta_1)] \right] = ?$$

$$\Rightarrow - \left[y_1 \cdot (1 - \text{softmax}(\theta_1)) \right] = - \left[y_1 \cdot (1 - g_1) \right]$$

$$= \begin{bmatrix} y_1 (\hat{y}_1 - 1) \\ y_2 (\hat{y}_2 - 1) \\ \vdots \\ y_d (\hat{y}_d - 1) \end{bmatrix}$$

when $i = k$

$$\begin{aligned} \frac{\partial}{\partial \theta_i} (F(y_i, \hat{y}_i)) &= y_i (\hat{y}_i - 1) & \frac{\partial y}{\partial \theta} &= 1 \\ &= 1 \cdot (\hat{y}_i - 1) & & \\ &= \hat{y}_i - 1 & & \end{aligned}$$

$$\begin{array}{c} i \neq k \\ \frac{\partial}{\partial \theta_i} (F_k) = \end{array} \quad x + y + z$$

$$\hat{y} = \text{softmax}(\theta)$$

$$E(y, \hat{y}) = - \sum_i y_i \cdot \log(\hat{y}_i),$$

$$\frac{\partial}{\partial \theta} E(y, \hat{y}) = ?$$

$$\frac{\partial}{\partial \theta} E(y, \hat{y}) = \frac{\partial}{\partial \theta} \cdot - \sum_i y_i \cdot \log(\text{softmax}(\theta_i))$$

choose k , s.t. $y_k \neq 0 \Rightarrow$

$$= \cancel{\frac{\partial}{\partial \theta} - \sum_i y_i \cdot \frac{\partial}{\partial \theta} \cdot - y_k \cdot \log(\text{softmax}(\theta_i))}$$

$$= \left[\begin{array}{l} \frac{\partial}{\partial \theta_1} - y_k \cdot \log(\text{softmax}(\theta_k)) \\ \vdots \\ \frac{\partial}{\partial \theta_d} - y_k \cdot \log(\text{softmax}(\theta_k)) \end{array} \right]$$

$$\text{when } i=k; \quad \frac{\partial}{\partial \theta_i} = ?$$

$$\Rightarrow \frac{\partial}{\partial \theta_k} \left[+ y_k \cdot \log(\text{softmax}(\theta_k)) \right]$$

$$= - \left[y_k \cdot (1 - \text{softmax}(\theta_k)) \right] = - [y_k (1 - \hat{y}_k)]$$

$$= y_k (\hat{y}_k - 1) = (\cancel{y_k}) \hat{y}_k - \cancel{1}$$

when $i \neq k$

$$= \frac{\partial}{\partial \theta_i} \left[y_k \cdot \log(\text{softmax}(\theta_k)) \right]$$

=

$$\log(xy) = \frac{1}{xy}$$

$$\frac{\partial}{\partial \theta_k} \left[\log \left(\frac{e^{\theta_i}}{\sum_{j=1}^d e^{\theta_j}} \right) \right] =$$

$$= \frac{\partial}{\partial \theta_i} \left[\theta_i - \log \sum_{j=1}^d e^{\theta_j} \right]$$

$$= \frac{\partial}{\partial \theta_i} 0 - \frac{1}{\sum_{j=1}^d e^{\theta_j}} \cdot \frac{\partial}{\partial \theta_k} \cdot \left(\sum_{j=1}^d e^{\theta_j} \right)$$

$$= - \frac{e^{\theta_k}}{\sum_{j=1}^d e^{\theta_j}} = - \text{softmax}(\theta_k)$$

$$= - \frac{\partial}{\partial \theta_i} \left[y_k \cdot \text{f}(\theta_i) \right]$$

$$= - \left[y_k \cdot (-s(\theta_i)) + s(\theta_k) \cdot 1 \right]$$

$$= y_k \cdot \text{softmax}(\theta_i)$$

$$= y_k \cdot \hat{y}_i = 1 \cdot \hat{y}_c = y_i$$

$$\therefore \frac{\partial}{\partial \theta_j} CE(y, \hat{y}) = \hat{y}_j - 1 ; \text{ when } j = k$$

$$= \hat{y}_j ; \text{ when } j \neq k$$

$$= [\hat{y}_1, \hat{y}_2, \dots, \hat{y}_{j-1}, \hat{y}_{j+1}, \dots, \hat{y}_n]$$

$$\therefore \frac{\partial}{\partial \theta_j} CE(y, \hat{y}) = \hat{y}_j - y_j$$

$$= \hat{y} - y$$

$$\{x_i, y_i\}_{i=1}^n$$

$x_i \rightarrow$ words | sentences (vectors)

$y_i \rightarrow$ sentiment | words | sentences.

$$P(y|x) = \frac{e^{w_y \cdot x}}{\sum_{c=1}^C e^{w_c \cdot x}}$$

$$\{x_i, y_i\}_{i=1}^n$$

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N -\log \left(\frac{e^{f_{y_i}}}{\sum_{c=1}^C e^{f_c}} \right)$$

for each $\{x_i, y_i\}$,

$$f_y = f_y(x) = w_y \cdot x = \sum_{j=1}^d w_{y,j} \cdot x_j$$

$$f = w \cdot x$$

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N -\log \left(\frac{e^{f_y_i}}{\sum_{c=1}^C e^{f_{yc}}} \right)$$

$$f_y = w_y \cdot x$$

$$\{x_i, y_i\}; \quad f_y' = w_y \cdot x = \sum_{j=1}^d w_{yj} \cdot x_j$$

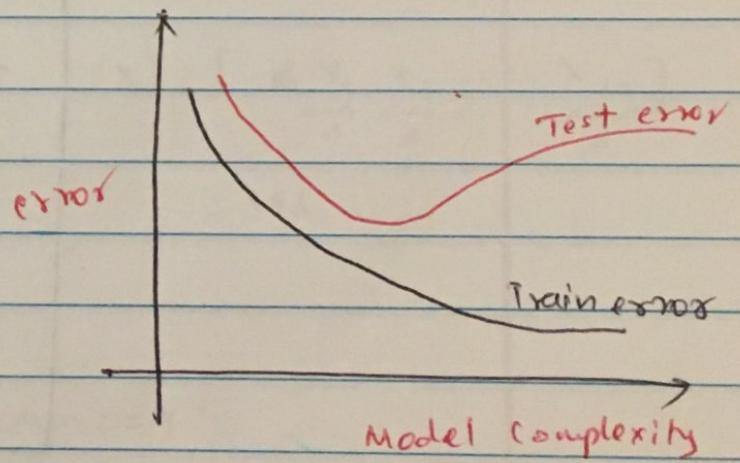
$$J(\theta) = \frac{1}{N} \sum_{i=1}^N -\log \left(\frac{e^{f_y_i}}{\sum_{c=1}^C e^{f_c}} \right)$$

$$\begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \quad [x_1 \ x_2]$$

- Why cross entropy over MSOS error.
- Squared error doesn't work well with neural networks. ∴ use cross-entropy.

Regularized: (prevents over-fitting)

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N -\log \left(\frac{e^{f_y_i}}{\sum_{c=1}^C e^{f_c}} \right) + \lambda \cdot \sum_{k=1}^d \theta_k^2$$



$$\theta = \begin{bmatrix} w_{1,1} \\ \vdots \\ w_{1,d} \end{bmatrix} \quad \theta \in \mathbb{R}^{cd}$$

$$\nabla_{\theta} f(\theta) = \begin{bmatrix} \nabla_{w_{1,1}} \\ \vdots \\ \nabla_{w_{1,d}} \end{bmatrix} \quad \nabla_{\theta} f(\theta) \in \mathbb{R}^{cd}$$

Deep learning!

Learn both w and x .

$$\nabla_{\theta} f(\theta) = \begin{bmatrix} \nabla_{w_{1,1}} \\ \vdots \\ \nabla_{w_{1,d}} \\ \nabla_{x_{\text{ardvark}}} \\ \vdots \\ \nabla_{x_{\text{zebra}}} \end{bmatrix} \quad \nabla_{\theta} f(\theta) \in \mathbb{R}^{fd+vd}$$

Large
fd+vd
over-fitting
danger

Word vec = word embeddings = word representations

|v|

$$d = d \begin{bmatrix} 0 & 0 & & \\ 0 & 0 & & \\ 0 & 0 & \dots & 0 \\ 0 & 0 & & 0 \\ 0 & 0 & & 0 \end{bmatrix}$$

Xaardvark. Xzebra

ML / NLP,

NLP

→ Rarely training single words work.
∴ use window classification.

Ex! 'in' says something about location.

merleau in Paris are amazing

$$x_{\text{window}} = [x_m \ x_n \ x_p \ x_a \ x_m]^T$$

$x_{\text{window}} = x \in \mathbb{R}^{5d}$, col vector!

$$x = \begin{bmatrix} x_1 & x_2 & \dots & x_{10} \end{bmatrix}^T$$

$\xrightarrow{d} \xrightarrow{d} \dots \xrightarrow{d}$

= 15d

→ dimensionality

→ variables

→ der. w.r.t fc when c=9

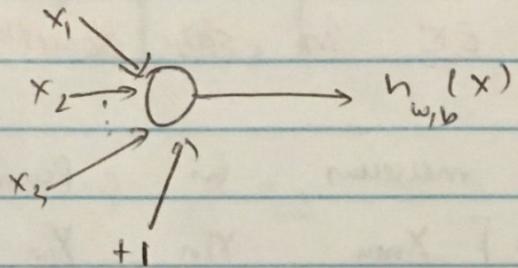
fc when c=9

→ Matrix notations

$$x_{\text{window}} = [x_{mu} \ x_{in} \ x_{pa} \ x_{aa} \ x_{am}]$$

$$\delta_{\text{window}} = \begin{bmatrix} \nabla x_{mu} \\ \nabla x_{in} \\ \vdots \\ \vdots \\ \nabla x_{am} \end{bmatrix} \in \mathbb{R}^{5d}$$

Neural Nets \Rightarrow

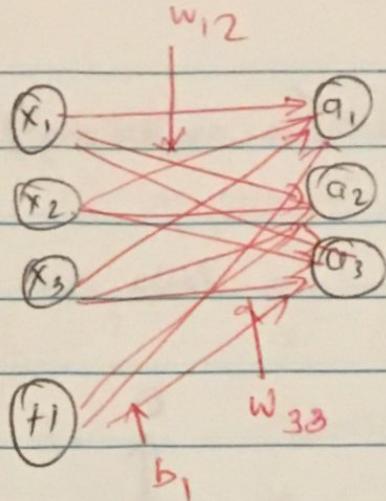


$$h_{w,b}(x) = f(w^T x + b).$$

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

$W \in \mathbb{R}^{a \times c}$

a : no of activation dimensions
 c : no of dimensions



$$a_1 = f(w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + b_1)$$

$$a_2 = f(w_{21}x_1 + w_{22}x_2 + w_{23}x_3 + b_2)$$

etc.

$$\begin{array}{c} W \quad X \quad B \\ \left[\begin{matrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{matrix} \right] \left[\begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} \right] \left[\begin{matrix} b_1 \\ b_2 \\ b_3 \end{matrix} \right] \\ 3 \times 3 \quad 3 \times 1 \quad 3 \times 1 \end{array}$$

$$Z = \dots = 3 \times 1$$

$$Z = W X + B$$

$$Z = W z + b$$

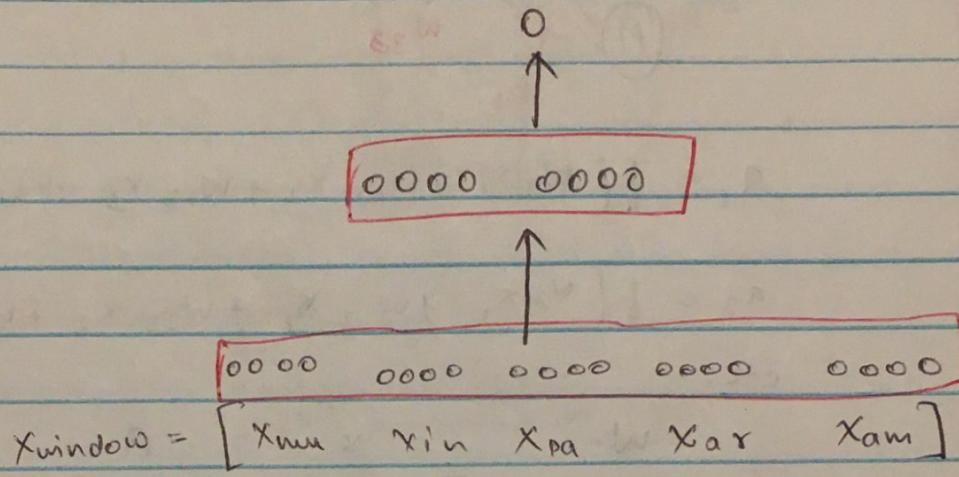
$$Z = [z_1 \ z_2 \ z_3]$$

$$a = f(z)$$

$$f(z) = f([z_1 \ z_2 \ z_3]) = [f(z_1), f(z_2), f(z_3)]$$

$s = \text{score}("Paris \text{ is amazing}")$

$s = \text{score}(\text{museums in Paris are amazing})$



$$s = v^T a$$

$$a = f(z)$$

$$z = Wx + b$$

$$s = v^T f(Wx + b)$$

A red box encloses three equations defining the dimensions of matrices x , w , and v .
 $x \in \mathbb{R}^{20 \times 1}$
 $w \in \mathbb{R}^{8 \times 20}$
 $v \in \mathbb{R}^{8 \times 1}$

$$\boxed{\frac{d}{d(\text{vec})} (\text{scalar}) = \text{vector}}$$

Task: Complete a few questions
in 1st assignment.

$$\frac{\partial J}{\partial x} = ? \quad \tau @ \text{hidden}$$

softmax output

$y = 1\text{-hot}$

cross entropy.

$$h = \text{sigmoid}(xw_1 + b_1) \quad \hat{y} = \text{softmax}(hw_2 + b_2)$$

→ what's objective function? Reg / not reg?

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N -\log \left(\frac{e^{b y_i}}{\sum_{c=1}^C e^{b c}} \right) + \lambda \sum_{k=1}^d \theta_k^2$$

$$\frac{\partial J(\theta)}{\partial \theta} = ? \quad b y_i = w_{y_i} \cdot x = \sum_{j=1}^d (w_{y_i})_j x_j$$

$$b c = w_c \cdot x = \sum_{j=1}^d (w_c)_j x_j$$

$$\therefore J(\theta) = \frac{1}{N} \sum_{i=1}^N -\log \left(\frac{\exp \left(\sum_{j=1}^d (w_{y_i})_j \cdot x_j \right)}{\sum_{c=1}^C \exp \left(\sum_{j=1}^d (w_c)_j \cdot x_j \right)} \right) + \lambda \sum_{k=1}^d \theta_k^2$$

$$\frac{d}{d \theta} x^T a = \frac{d}{d \theta} \sum_i x_i a_i = \frac{d}{d \theta} \sum_i$$

σ @ hidden
 s @ output.

$$\frac{\partial}{\partial w_i} \sum_i x_i a_i$$

$$\sum_i \frac{\partial}{\partial x_i} x_i a_i$$

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N \lambda \sum_{k=1}^d \theta_{ik}^2 - \log \left(\frac{\exp \left(\sum_{j=1}^d (w_{y_i})_j \cdot x_j \right)}{\sum_{c=1}^C \exp \left(\sum_{j=1}^d (w_c)_j \cdot x_j \right)} \right)$$

$$\frac{\partial J(\theta)}{\partial x} = \left[\frac{\partial}{\partial x_1} J(\theta), \frac{\partial}{\partial x_2} J(\theta), \dots, \frac{\partial}{\partial x_d} J(\theta) \right]$$

$$\frac{\partial}{\partial w_i} J(\theta) = \frac{1}{N} \sum_{i=1}^N - \frac{\partial}{\partial w_i} \cdot \log \left(\frac{\exp \left(\sum_{j=1}^d (w_{y_i})_j \cdot x_j \right)}{\sum_{c=1}^C \exp \left(\sum_{j=1}^d (w_c)_j \cdot x_j \right)} \right)$$

$$\cancel{\frac{\partial}{\partial w_i}} \cdot \log \left(\frac{a}{b} \right) = \log a - \log b$$

$$= \sum_{j=1}^d (w_{y_i})_j \cdot x_j - \sum_{z=1}^C (w_z)_i \cdot x_z$$

$$= \frac{\partial}{\partial w_i} \sum_{j=1}^d (w_{y_i})_j \cdot x_j - \frac{\partial}{\partial w_i} \sum_{z=1}^C (w_z)_i \cdot x_z$$

$$\cancel{\frac{\partial}{\partial w_i}} = (w_{y_i})_i - (w_z)_i$$

$$\frac{\partial}{\partial t} *^2$$

$$\frac{\partial}{\partial x} e^{-x} = -e^{-x}, \quad \frac{\partial}{\partial z} e^z; \quad \frac{d}{dx} -x$$

$$\log \left(\frac{\exp \left(\sum_{j=1}^d (w_{y_i})_j \cdot x_j \right)}{\sum_{c=1}^C \exp \left(\sum_{z=1}^d (w_c)_z \cdot x_z \right)} \right)$$

$$a = \sum_{j=1}^d (w_{y_i})_j \cdot x_j - \log \sum_{c=1}^C \exp \left(\sum_{z=1}^d (w_c)_z \cdot x_z \right)$$

$$\frac{\partial a}{\partial x_1} = (w_{y_i})_1 - \frac{1}{\sum_{c=1}^C \exp \left(\sum_{z=1}^d (w_c)_z \cdot x_z \right)} \cdot \frac{\partial}{\partial x_1} \left(\sum_{c=1}^C \exp \left(\sum_{z=1}^d (w_c)_z \cdot x_z \right) \right)$$

$$= (w_{y_i})_1 - \frac{1}{\sum_{c=1}^C \exp \left(\sum_{z=1}^d (w_c)_z \cdot x_z \right)} \cdot \sum_{c=1}^C \frac{\partial}{\partial x_1} \exp \left(\sum_{z=1}^d (w_c)_z \cdot x_z \right)$$

$$= (w_{y_i})_1 - \frac{\sum_{c=1}^C \cancel{\exp \left(\sum_{z=1}^d (w_c)_z \cdot x_z \right)} \cdot (w_c)_1}{\sum_{c=1}^C \cancel{\exp \left(\sum_{z=1}^d (w_c)_z \cdot x_z \right)} \cdot}$$

$$= w_{\text{y}_i}$$

$$\frac{1}{N} \sum_{i=1}^N \frac{\sum_{c=1}^C \exp \left(\sum_{z=1}^d (w_c)_z \cdot x_z \right) \cdot (w_c)_1}{\sum_{c=1}^C \exp \left(\sum_{z=1}^d (w_c)_z \cdot x_z \right)}$$

$$f(g(x))$$

$$J(\theta) = CE(y, \hat{y}) = - \sum_{i=1}^N y_i \log(\hat{y}_i)$$

$$\frac{\partial J(\theta)}{\partial \theta} = ?$$

$$\left[f\left(\frac{z}{g(x)} \right) \right]' = \frac{d}{dz} f(z) \cdot \frac{d}{dx} g(x).$$

$$h = \text{sigmoid}(xw_1 + b_1) \quad \hat{y} = \text{softmax}(hw_2 + b_2)$$

$$J(\theta) \rightarrow \theta \rightarrow h \rightarrow x.$$

$$\therefore \frac{\partial J(\theta)}{\partial \theta} = \frac{\partial J(\theta)}{\partial \theta} \cdot \frac{d}{dh} \hat{y} \cdot \frac{d}{dx} h$$

$$= \frac{\partial J}{\partial \theta} \cdot \frac{d}{dh} \hat{y} \cdot \frac{d}{dx} h$$

$$J = - \sum_{i=1}^N y_i \log \left(\text{softmax}(hw_2 + b_2) \right)$$

$$\partial \theta \quad \frac{\partial J}{\partial \theta}$$

$$J = - \sum_{i=1}^N y_i \log \left(\text{softmax}(hw_2 + b_2) \right)$$

$$\theta = h w_2 + b_2 \quad \frac{\partial \theta}{\partial h}.$$

$$h = \sigma(xw_1 + b_1) \quad \frac{\partial h}{\partial x}.$$

$$\text{Let } J = - \sum_{i=1}^N y_i \log (\hat{y}_i)$$

$$= - \sum_{i=1}^N y_i \log (\text{softmax}(\theta_i))$$

$\therefore J(\theta)$.

$$\theta = h w_2 + b_2$$

$\therefore \theta(h)$.

$$h = \text{sigmoid}(\frac{z}{x w_1 + b_1})$$

$$h = \sigma(z)$$

$\therefore h(z)$.

$$\therefore \frac{\partial J}{\partial x} = \frac{\partial J(\theta)}{\partial \theta} \cdot \frac{\partial \theta(h)}{\partial h} \cdot \frac{\partial h(x)}{\partial x}$$

$$= (\hat{y} - y) \cdot w_2 \cdot \frac{d}{dx} \cdot \sigma(x w_1 + b_1)$$

$$= (\hat{y} - y) \cdot w_2 \cdot \frac{d}{dx} \cdot [1 + e^{-(x w_1 + b_1)}]^{-1}$$

$$= (\hat{y} - y) \cdot w_2 \cdot -1 \cdot (1 + e^{-(x w_1 + b_1)})^{-2} \cdot \frac{d}{dx} e^{-(x w_1 + b_1)}$$

$$= - \frac{1}{1 + e^{-(x w_1 + b_1)^2}} \cdot (-1) \cdot e^{-(x w_1 + b_1)}$$

$$\text{row} = \text{#n}$$

$$\text{col} = \text{x dim}$$

$$D_h \times D_x = n \times x$$

$$x \times n$$

$$z = x w_1 + b,$$

$$\therefore \frac{\partial z}{\partial x}.$$

$$\frac{\partial E}{\partial x} = \frac{\partial E(\theta)}{\partial \theta} \cdot \frac{\partial \theta(b)}{\partial b} \cdot \frac{\partial b(z)}{\partial z} \cdot \frac{\partial z(x)}{\partial x}.$$

$$= [\hat{y} - y] \cdot w_2 \sigma(x w_1 + b_1) \odot (1 - \sigma(x w_1 + b_1))$$

$$\Rightarrow 1 * C \odot (h).$$

$(\ast h \oplus h)$

$$(\hat{y} - y) \frac{d}{dz} \text{sig}(x w_1^T + b_1)$$

$D_h \times D_x$

$$1 * D_x \cdot D_x * D_h$$

$$1 * D_h.$$

$$1 * D_h.$$

$$\frac{d}{dt} \sigma(z) = r(z) \cdot 1 - \sigma(z).$$

$$\textcircled{D}_Y \quad \boxed{H} \quad \textcircled{D}_Y \Rightarrow (H * D_X + 1) \\ + (D_Y * H + 1)$$

$$r(x) = \frac{1}{1 + \exp(-x)}$$

$$f = r(x) \quad r^{-1}(x) = r(x) \cdot (1 - r(x))$$

$$f = \frac{1}{1 + \exp(-x)} \quad x \in \mathbb{R}^{20 \times 1}$$

$$z = w_x + b \quad x \in \mathbb{R}^{20 \times 1} \\ a = f(z) \quad w \in \mathbb{R}^{8 \times 20} \\ \text{score} = v^T a \quad v \in \mathbb{R}^{1 \times 8}$$

$$\xi = v^T f(w_x + b) \quad \rightarrow \text{Max-Margin} \\ \rightarrow \text{Minimize } \xi \\ \rightarrow \xi = \max(0, 1 - s + s_c)$$

~~scribble~~

$$\frac{\partial}{\partial x}$$

$$\frac{\partial}{\partial w_i} \sum_i v_i x_i$$

x_1, x_2, \dots, x_n

$$s = v^T f(w^T x + b)$$

$$s_c = v^T f(w^T x_c + b),$$

$$\delta = \max(0, 1 - s + s_c),$$

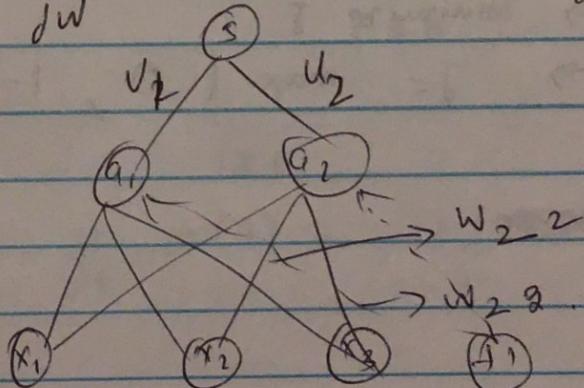
Training with Back-prop:

Compute derivatives of s & s_c w.r.t.

v, w, x, b .

$$\textcircled{1} \quad \frac{ds}{dv} = \frac{d}{du} v^T f(w^T x + b) \cdot = f'(w^T x + b) = \alpha$$

$$\textcircled{2} \quad \frac{ds}{dw} = \frac{d}{dw} \cdot v^T f(w^T x + b) = \frac{\partial}{\partial w} v^T \cdot \alpha.$$



a_i depends only on w_i

to

$$a_i \rightarrow w_i$$

$$v_i \rightarrow a_i \rightarrow w_i$$

$$[w, v_i] \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} \quad \frac{\partial}{\partial w} v^T a = \delta_{w_{11}} v^T a$$

$$\frac{\partial}{\partial w_{ij}} (v^T a) \rightarrow$$

$$[v_1 \ v_2] \begin{bmatrix} \sigma(w_{11}x_1 + w_{12}x_2 + \\ w_{13}x_3) \\ \sigma(w_{21}x_1 + w_{22}x_2 + \\ w_{23}x_3) \end{bmatrix}$$

$$\frac{\partial}{\partial w} v^T a = \begin{bmatrix} \delta_{w_{11}}(v^T a) & \delta_{w_{12}}(v^T a) & \delta_{w_{13}}(v^T a) \\ \vdots & \vdots & \vdots \\ \delta_{w_{21}}(v^T a) & \delta_{w_{22}}(v^T a) & \delta_{w_{23}}(v^T a) \end{bmatrix}$$

$$\delta_{w_{11}}(v^T a) \rightarrow \delta_{w_{11}} v_1 a_1$$

$$\delta_{w_{ij}}(v^T a) \rightarrow \delta_{w_{ij}} v_i a_i$$

$$\therefore \frac{\partial}{\partial w_{ij}} v^T a \rightarrow \frac{\partial}{\partial w_{ij}} v_i a_i.$$

$$= v_i \frac{\partial}{\partial w_{ij}} a_i = v_i \cdot \frac{\partial}{\partial w_{ij}} \sigma(z)$$

$$z_i = \sum_j w_{ij} x_j + b$$

$$\therefore \frac{\partial}{\partial z_i} \sigma(z)$$

$$\therefore = v_i \cdot \frac{\partial}{\partial z_i} \sigma(z) \cdot \frac{\partial}{\partial w_{ij}} \sigma(z)$$

$$\boxed{\frac{\partial}{\partial w_{ij}} (v^T a) = v_i \cdot \sigma'(z) \cdot x_j}$$

$$\frac{\partial}{\partial w_{ij}} s = \delta_i \cdot x_j$$

$$\begin{bmatrix} \delta_1 & \delta_2 & \delta_3 \\ \delta_1 \cancel{x}_1 & \cancel{x}_2 & \cancel{x}_3 \end{bmatrix} \quad \begin{bmatrix} \delta_1 x_1 & \delta_1 x_2 & \delta_1 x_3 \\ \delta_2 x_1 & \delta_2 x_2 & \delta_2 x_3 \end{bmatrix}$$

$$\therefore \frac{\partial}{\partial W} s = \underbrace{\delta x^T}_{\text{error msg}} \quad ; \quad \frac{\partial}{\partial W} s_c = \delta_c x_c^T.$$

$$\begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} \delta_1 \\ \delta_2 \end{bmatrix}$$

$1 \times 3 \times 2 \times 1 \quad 3 \times 2$ $3 \times 1 \times 1 \times 2 = 3 \times 2$

\Rightarrow Outer-product

$$\gamma = \max(0, 1 - s + s_c)$$

$$\therefore \frac{\partial \gamma}{\partial W} = \frac{\partial \gamma(s)}{\partial W}$$

$$\therefore \frac{\partial \gamma}{\partial W} = -\delta x^T + \delta_c x_c^T.$$

①

$$\frac{\partial s}{\partial b} = \frac{\partial v^T a}{\partial b} = \frac{\partial v^T f(wx + b)}{\partial b}$$

$$\frac{\partial}{\partial b_i} \frac{\partial s}{\partial b} = \begin{bmatrix} \frac{\partial s}{\partial b_1} & \frac{\partial s}{\partial b_2} \end{bmatrix}$$

$$\therefore \frac{\partial s}{\partial b_i} = v^T f(wx + b_i) \\ = \sum_i v_i \cdot f(w_i x + b_i)$$

$$\therefore s(a) = a(z) \cdot z(b)$$

$$\therefore \frac{\partial c}{\partial b_i} = \frac{\partial s(a_i)}{\partial a_i} \cdot \frac{\partial a(z)}{\partial z} \cdot \frac{\partial z(b)}{\partial b}$$

$$= \frac{\partial}{\partial a} v_i^T a \cdot \frac{\partial}{\partial b_i} f(z) \cdot \frac{\partial}{\partial b_i} w_i x + b_i$$

$$= v_i \sigma'(z) \cdot (1)$$

$$= s_i$$

$$\frac{\partial s}{\partial x} = \begin{bmatrix} \frac{\partial s}{\partial x_1} & \frac{\partial s}{\partial x_2} & \frac{\partial s}{\partial x_3} \end{bmatrix}$$

$$[u_1 \ u_2] \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}.$$

$$s = v^T f(wx + b).$$

$$= v^T a; \quad a = f(\underline{z}) \quad ; \quad z = wx + b \\ s(a); \quad a(z); \quad z(x).$$

$$\therefore \frac{ds}{dx_1} = \frac{ds(a)}{da}, \quad \frac{da(z)}{dz} \cdot \frac{dz(x)}{dx_1}.$$

$$= \quad s = \sum_i v_i a_i; \quad a_i = f(z_i); \quad z_i = w_i x + b_i \\ s(a_i) \quad a_i(z_i) \quad z_i(x).$$

$$= v_i \cdot \sigma'(z_i) \cdot \left[\frac{dz_i}{dx_1} \right]$$

$$\frac{\partial s}{\partial x_1} = \left[\frac{\partial s}{\partial x_1} \quad \frac{\partial s}{\partial x_2} \right]$$

$$\frac{\partial}{\partial x_1} v^T a$$

$$\sigma'(z)$$

$$a = f(w)$$

$$\frac{\partial}{\partial x_1} f$$

$$a = [a_1 \ a_2]$$

$$= \begin{bmatrix} f(w_1 x + b_1) & f(w_2 x + b_2) \end{bmatrix}$$

$$= [f(z_1) \ f(z_2)]$$

$$\frac{\partial}{\partial x} b(\vec{x}) =$$

$$z_1 = \sum_i w_i \quad ; \quad z_1 = \sum_i w_{ij} x_j \quad ; \quad [w_1^T, s] \quad ; \quad [w_1^T, w_2^T, s]$$

$$z_1 = w_1 x + b_1 \quad ; \quad z_2 = w_2 x + b_2 \quad ; \quad w^T s.$$

\therefore ~~s(a)~~; $s(a, a_1)$; $a_1(z_1)$; $a_2(z_2)$
 $z_1(x)$; $z_2(x)$.

$$\therefore \frac{\partial}{\partial x} s(a, a_1) = \frac{\partial}{\partial a_1} s(a_1) + \frac{\partial}{\partial a_2} s(a_2).$$

$$\Rightarrow \text{If } \frac{\partial}{\partial x} s(a_1) = \cancel{\frac{\partial}{\partial a_1} s(a_1)} \cdot \cancel{\frac{\partial}{\partial z_1} z_1(x)}$$

$$= \sigma'(z_1).$$

$$s(a) = \cancel{s(a)}^T \quad s(a) = \sum_i v_i a_i$$

$$\therefore \frac{\partial}{\partial x} s(a_1) \cdot \frac{\partial}{\partial a_1} a_1(z_1) \cdot \frac{\partial}{\partial z_1} z_1(x).$$

$$= v_1 \cdot \sigma'(z_1), w_1$$

$$\therefore \frac{\partial}{\partial x} s = s_1 w_1 + \frac{\partial}{\partial x} s = s_2 w_2.$$

$$= s_1 w_1 + s_2 w_2.$$

$$= \sum_{i=1}^2 s_i w_i$$

$$J = \max(0, 1-s+s_c)$$

→ Minimize J .

$$\frac{\partial}{\partial(x,y)} f(x,y) = \frac{\partial}{\partial x} f(x,y) + \frac{\partial}{\partial y} f(x,y)$$

$$\frac{\partial J}{\partial(s, s_c)} = \frac{\partial J(s, s_c)}{\partial s} + \frac{\partial J(s, s_c)}{\partial s_c}$$

$$s \underset{a.}{=} \frac{\partial (1-s+s_c)}{\partial s} + \frac{\partial (1-s+s_c)}{\partial s_c}$$

$$J(v, w, x, b) = \max(0, 1-s+s_c)$$

$$\frac{\partial J}{\partial(v, w, x, b)} = \frac{\partial (1-s+s_c)}{\partial(v, w, x, b)}$$

$$J = 1-s+s_c$$

∴ $J(s, s_c)$

$$s = v^T a$$

$$\therefore s(a)$$

$$a = f(wx+b) - f(z)$$

$$\therefore a(z)$$

$$z = wx+b$$

∴ $z(\cdot)$

$$\frac{d}{dx} \left[f(\theta(s, s_c)) \right] = \frac{\partial f}{\partial z}(z) \cdot \frac{\partial \theta}{\partial x}(s, s_c).$$

$\begin{array}{l} z \mapsto t \mapsto x \\ s_c \rightarrow t \rightarrow x \\ s_c \rightarrow x \end{array}$

$$= \cancel{\frac{\partial \theta}{\partial x}(s, s_c)} + \cancel{\frac{\partial \theta}{\partial x}(s, s_c)}.$$

$$= \cancel{\frac{\partial \theta}{\partial x}(s, s_c)}.$$

$$= \theta(z) = \max(0, 1 - s + s_c)$$

$$\therefore \theta(z).$$

$$z = 1 - s + s_c,$$

. . . $\theta(s_c, s_c)$

$$\cancel{\theta} \quad s(x)$$

$$\frac{\partial \theta}{\partial z} \cdot \left[\frac{\partial z(s)}{\partial s}, \cancel{\frac{\partial z}{\partial x}}, + \frac{\partial z(s_c)}{\partial s_c} \right]$$

$$= 1 \cdot \left[\frac{\partial z(s)}{\partial s}, \frac{\partial s(x)}{\partial x} + \frac{\partial z(s_c)}{\partial s_c} + \frac{\partial s_c(x)}{\partial x} \right]$$

$$\frac{\partial}{\partial x} f(x) + g(x),$$

$$\frac{\partial}{\partial x} b(x) + \frac{\partial}{\partial x} g(x).$$

$$J = \max(0, 1 - s + s_c)$$

Minimize J .

$$\nabla J(u, w, x, b) = ?$$

$$\textcircled{1} \quad \frac{\partial}{\partial(x,y)} f(x,y) = \frac{\partial}{\partial x} f(x,y) + \frac{\partial}{\partial y} f(x,y),$$

$$\textcircled{2} \quad \frac{\partial}{\partial x} f(g(x)) = \frac{\partial}{\partial z} f(z) \cdot \frac{\partial}{\partial x} g(x).$$

$$\frac{\partial J}{\partial x} = \frac{\partial}{\partial x} (1 - s + s_c),$$

$$= 0 - \frac{\partial}{\partial x} s + \frac{\partial}{\partial x} s_c.$$

$$s = v^T a \quad \therefore s(a)$$

$$a = f(z); z = w^T x + b \quad \therefore a(z)$$

$$z = w^T x + b \quad \therefore z(x).$$

$$\therefore \frac{\partial s}{\partial x} = \frac{\partial s(a)}{\partial a} \cdot \frac{\partial a(z)}{\partial z} \cdot \frac{\partial z(x)}{\partial x},$$

$$= v \cdot \sigma'(z) \cdot w^T.$$

$$\boxed{\frac{\partial J}{\partial x} = 0 - \frac{\partial}{\partial(u,w,x,b)} (+ \frac{\partial}{\partial(u,w,x,b)} s_c)}$$

$$\frac{d}{d(v,w,x,b)} s = \frac{d}{dv} s + \frac{d}{dw} s + \frac{d}{dx} s + \frac{d}{db} s$$

$$\frac{d}{d(v,w,x,b)} s_c = \frac{d}{dv} s_c + \frac{d}{dw} s_c + \frac{d}{dx} s_c + \frac{d}{db} s_c$$

$$\boxed{\frac{d}{d(v,w,x,b)} J(v,w,x,b) = -\frac{d}{d(v,w,x,b)} s + \frac{d}{d(v,w,x,b)} s_c}$$

$$s = v^T a = v^T f(z) = v^T f(wx+b).$$

$$\frac{d}{dv} s = \frac{d}{dv} v^T a = \frac{d}{dv} \sum_i v_i a_i = [a_1, \dots, a_n] = a$$

$$\therefore \boxed{\frac{d}{dv} s = a} \rightarrow ①$$

$2 \times 3 \times 3 \times 1$

2×1

~~$$\frac{d}{dw} s = \cancel{\frac{d}{dw} \frac{d}{da} s(a)} \cdot \cancel{\frac{d}{dz} f(z)} \cdot \cancel{\frac{d}{dx} f'(wx+b)}$$~~

2×3

~~$$= v \cdot f'(z) \cdot x$$~~

$$: 2 \times 1, 2 \times 1, \begin{matrix} 1 \times 3 \\ 3 \times 1 \\ 1 \times 2 \end{matrix}, 1 \times 2$$

$$\frac{\partial}{\partial w} s = \frac{\partial}{\partial w} v^T a = \frac{\partial}{\partial w} v^T f(z) = \frac{\partial}{\partial w} v^T f(wx + b)$$

∂ derivative w.r.t. w_{ij}

$$\frac{\partial}{\partial w_{ij}} v^T f(wx + b)$$

~~2x2~~ $1 \times 2 \times 2 \times 1$

$$= \frac{\partial}{\partial w_{ij}} v_i \cdot f(z)$$

$$= \frac{\partial}{\partial w_{ij}} \sum_i v_i \cdot f(z_i)$$

$$z = (w_x + b)$$

$$z_i = (w_i x + b)$$

$$= \sum_i v_i \cdot \frac{\partial}{\partial w_{ij}} f(z_i)$$

$$x = 1, 2$$

$$= \sum_i v_i \cdot \frac{\partial}{\partial w_{ij}} f(z_x).$$

$$= \sum_i v_i \cdot \frac{\partial}{\partial w_{ij}} f(w_i x + b)$$

$$\begin{aligned} \vec{W} &\in \mathbb{R}^{2 \times 3}, \quad x \in \mathbb{R}^{3 \times 1} \\ w_{j1} &\in \mathbb{R}^{2 \times 1} \end{aligned}$$

$$\frac{\partial}{\partial w_{ij}} v_i^T f(w_i x + b).$$

$i = 2$

$$= \frac{\partial}{\partial w_{ij}} \sum_p v_p^T f(z_p).$$

$$\begin{aligned} &= v_1 f(z_1) + v_2 f(z_2) \\ &= v_1 f(w_1 x + b) + v_2 f(w_2 x + b) \\ &= 0 + \end{aligned}$$

$$= \cancel{\frac{\partial}{\partial w_{ij}}} \quad \text{when } p \neq i;$$

$$\frac{\partial s}{\partial w_{ij}} = \cancel{\frac{\partial}{\partial w_{ij}}} v_i f(z_i).$$

$$\frac{\partial}{\partial w_{ij}} v_i \cdot f(z_i) = \cancel{\frac{\partial}{\partial w_{ij}}} v_i.$$

$$= v_i \cdot \frac{\partial}{\partial w_{ij}} f(z_i) \quad ; \quad z_i = w_i x + b$$

$$= v_i \cdot \frac{\partial}{\partial z_i} f(z_i) \cdot \frac{\partial}{\partial w_{ij}} (w_i x + b)$$

$$= v_i \cdot \underbrace{f'(z_i)}_{s_i} \cdot x_j$$

$\frac{\partial}{\partial w_{ij}} s = s_i x_j$

from $\frac{\partial}{\partial w_{ij}} s$ to $\frac{\partial}{\partial w} s$

$$\boxed{\begin{aligned}\frac{\partial}{\partial w_{ij}} s &= s_i \cdot x_j \\ &= \underbrace{s_i x_j^\top}_{\text{outer product}}\end{aligned}} \rightarrow \textcircled{2}$$

$$[s_1 \ s_2] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$s \in \mathbb{R}^{2 \times 1}$: responsibility / error msg from each activation
 i.e., $s_1 x_1 \ s_2 x_2 \ s_1 x_1 \ s_1 x_2 \ s_2 x_3$.

$$\frac{\partial}{\partial b} s = \frac{\partial}{\partial b} v^T a = \frac{\partial}{\partial b} v^T f(z) = \frac{\partial}{\partial b} v^T f(wx+b)$$

$$= \frac{\partial}{\partial b} \sum_i v_i \cdot f(z_i)$$

$$\bullet \frac{\partial s}{\partial b_i} = \frac{\partial}{\partial b_i} v_i \cdot f(z_i)$$

$$= v_i \cdot \frac{\partial}{\partial b_i} f(z_i)$$

$$= v_i \cdot \frac{\partial}{\partial z_i} f(z_i) \cdot \frac{\partial}{\partial b_i} (w_i x + b_i)$$

$$1 \times 2 \times 2 \times 1 \\ = 1 \times 1$$

$$\delta_1 w_1 + \delta_2 w_2$$

$$= v_i \cdot f'(z_i) \cdot 1 \\ = \delta_i$$

$1 \times 1 \times$

$\frac{d}{dx} s = 8$

→ ⑤

$$\frac{d}{dx} s = \frac{\partial}{\partial x} \sum_i v_i \cdot f(z_i)$$

$$\begin{bmatrix} \delta_1 \delta_2 \\ \vdots \end{bmatrix} \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} & w_{23} \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$$= \sum_i v_i \cdot \frac{\partial}{\partial x} f(z_i) ; z_i = w_i x + b_i \quad 3 \times 2$$

$$= \sum_i v_i \cdot \frac{\partial}{\partial z_i} f(z_i) \cdot \frac{\partial}{\partial x} z_i(x) \quad \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix} \begin{bmatrix} \delta_1 \\ \delta_2 \\ \delta_3 \end{bmatrix}$$

$$= \sum_i v_i \cdot f'(z_i) \cdot \frac{\partial}{\partial x} w_i x + b_i \quad \begin{matrix} 3 \times 2 \times 2 \times 1 \\ s_1, w_1 + \delta_2 w_2 \\ s_1 \end{matrix}$$

$$= \sum_i \delta_i \cdot w_i = w_{11} \delta_1 + w_{21} \delta_2 = s_1$$

$$w_{11} \delta_1 + w_{21} \delta_2$$

$$w_{13} \delta_1 + w_{23} \delta_2$$

$$\begin{bmatrix} \delta_1 w_1 + \delta_2 w_2 + \delta_3 w_3 \\ \delta_2 w_2 + \delta_3 w_3 \end{bmatrix} = w^T \cdot \underline{s} \\ \underline{s} =$$

$\frac{d}{dx} s = w^T \underline{s}$

$$\sum_i u_i f_i(x)$$

$$\frac{d}{dx} v^T a$$

$$= \frac{d}{dx} \sum_i u_i f(z_i)$$

$$= \frac{d}{dx} \sum_i u_i \cdot f(w_i x + b_i)$$

$$\frac{d}{dx_j} \sum_i u_i \cdot f(w_i x + b_i)$$

$$= \sum_i u_i \cdot \frac{d}{dx_j} f(z_i)$$

$$= \sum_i u_i \cdot \frac{d}{dz_i} f(z_i) \cdot \frac{d}{dx_j} (w_i x + b_i)$$

$$= \sum_i u_i \cdot f'(z_i) \cdot \frac{d}{dx_j} \sum_p w_{ip} x_p$$

$$= \sum_i u_i \cdot f'(z_i) \cdot w_{ij}$$

$$\frac{d}{dx_j} \sum_i$$

$$= \sum_i s_i \cdot w_{ij}$$

$$s_1 w_{12} + s_2 w_{22}$$

$$= \delta \cdot w_{\cdot j}^T$$

$$\frac{\partial}{\partial x} s = \begin{bmatrix} s \cdot w_1^T & s \cdot w_2^T & s \cdot w_3^T \end{bmatrix}$$

$$\boxed{\frac{\partial}{\partial x} s = s w^T} \rightarrow \textcircled{4}$$

$$\boxed{\frac{\partial}{\partial v} s = a}$$

$$\frac{\partial}{\partial w} s = s x^T$$

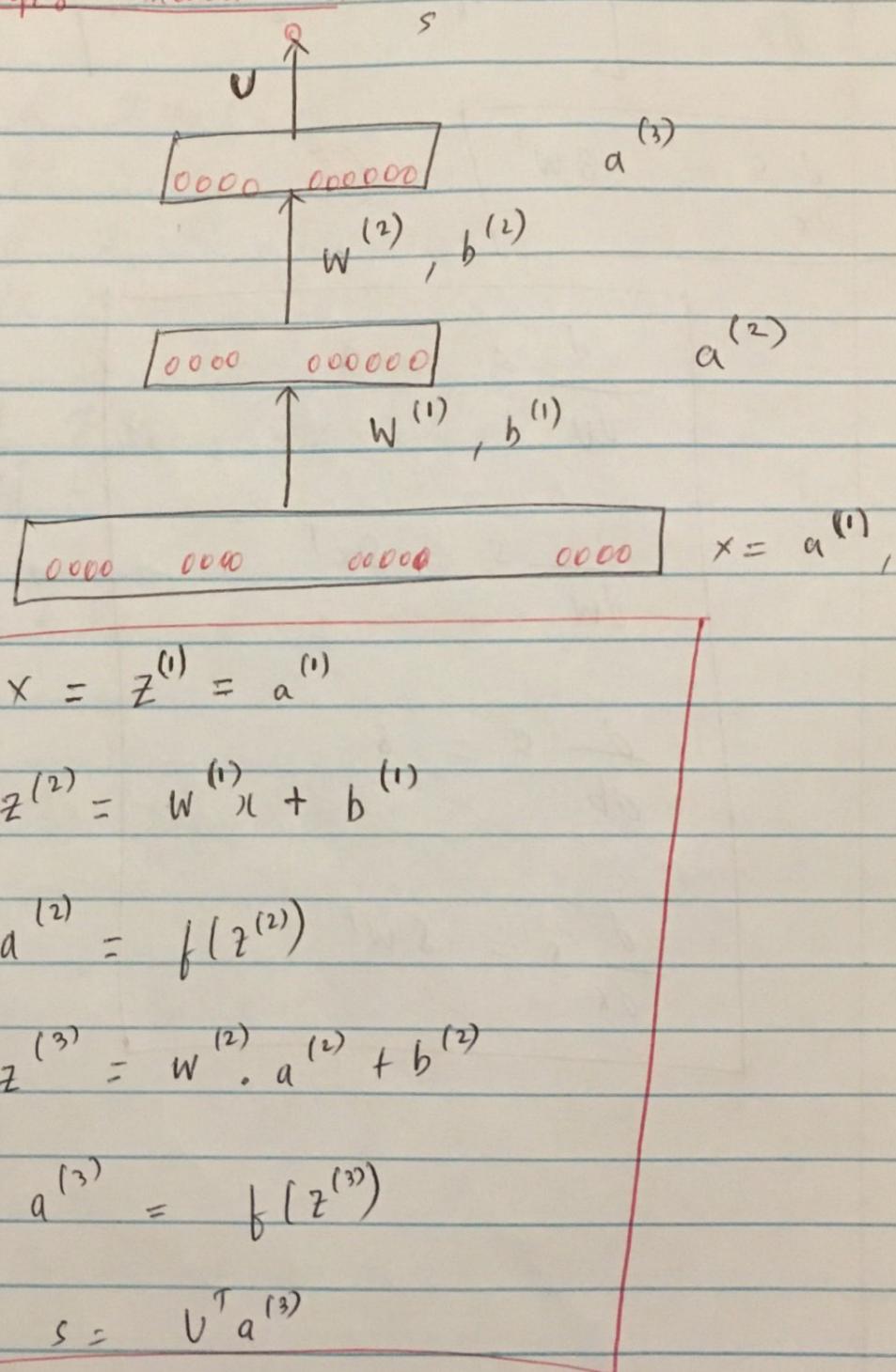
$$; \quad \delta = [s_1 \quad s_2]$$

$$= [v_1 f'(z_1) \quad v_2 f'(z_2)]$$

$$\frac{\partial}{\partial b} s = s$$

$$\boxed{\frac{\partial}{\partial x} s = s w^T}$$

2-layer neural net:



$$s = v^T \cdot a$$

$$\frac{\partial s}{\partial w_{ij}} = v^T \cdot b(z_i)$$

$$\frac{\partial s}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} v_i \cdot b(z_i)$$

$$s = v^T \cdot a^{(3)}$$

$$= v^T \cdot f(z^{(3)})$$

$$= v^T \cdot f(w^{(2)} \cdot a^{(2)} + b^{(2)})$$

$$= v^T \cdot f(w^2 \cdot f(w^{(1)} \cdot a^{(1)} + b^{(1)}) + b^{(2)})$$

$$= v^T \cdot f(w^{(2)} \cdot f(w^{(1)} \cdot x + b^{(1)}) + b^{(2)})$$

$$\frac{\partial s}{\partial w_{ij}} = v_i \cdot f'(z_i) \cdot x_j$$

$$\frac{\partial s}{\partial w_{ij}} = \delta_i \cdot x_j$$

$$\begin{bmatrix} \delta_1 x_1 & \delta_1 x_2 & \delta_1 x_3 \\ \delta_2 x_1 & \delta_2 x_2 & \delta_2 x_3 \end{bmatrix}$$

$$\frac{\partial s}{\partial w_{ij}^{(2)}} = \delta_i^{(2)} \cancel{\cdot a_j^{(2)}} = \delta_i^{(2)} \cdot a_j^{(2)}$$

$$\begin{bmatrix} \delta_1 & \delta_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\therefore \frac{\partial s}{\partial w^{(2)}} = \underbrace{\delta^{(3)} \cdot a^{(2)^T}}_{\text{outer prod}}$$

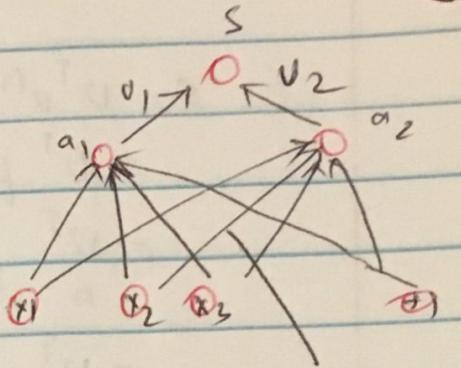
$$\delta_1^{(3)} \quad \delta_2^{(3)} \quad \cancel{a^{(3)}} \quad f'(z_1^{(3)})$$

$$\delta^{(3)} = v \circ f'(z^{(3)})$$

$$\frac{\partial s}{\partial w^{(1)}} = \delta^{(2)} \underbrace{\cdot a^{(1)^T}}_{\text{outer prod.}}$$

$$\delta^{(2)} = \cdot f'(z^{(2)})$$

$$\frac{\partial}{\partial w_{ij}} s = \frac{\partial}{\partial w_{ij}} v^T a$$



$$= \frac{\partial}{\partial w_{ij}} \sum_y v_y a_y$$

$$= \frac{\partial}{\partial w_{ij}} \sum_y v_y f(w_y x + b)$$

w_{ij} appears only in w_i ,
i.e., when $i=y$.

$$\therefore \frac{\partial}{\partial w_{ij}} s = \frac{\partial}{\partial w_{ij}} v_i \cdot f(w_i x + b)$$

$$= v_i \cdot \frac{\partial}{\partial w_{ij}} f(w_i x + b)$$

$$= v_i \cdot f'(z_i) \cdot \frac{\partial}{\partial w_{ij}} (w_i x + b)$$

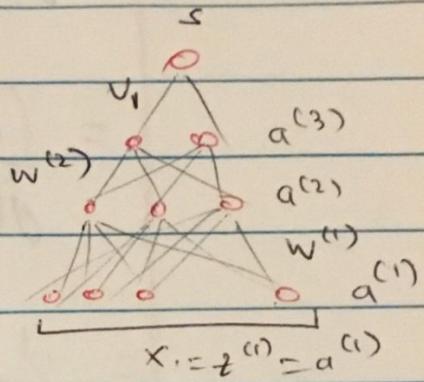
$$= v_i \cdot f'(z_i) \cdot \cancel{w_j} x_j$$

$$\boxed{\frac{\partial}{\partial w_{ij}} s = s_i \cancel{w_{ij}}} \quad \text{Since } s_i \cdot x_j$$

$$\therefore \frac{d}{d w_{ij}^{(1)}} s = \delta_i \circ x_j$$

For 2 layer Neural Net;

$$\frac{\partial}{\partial w_{ij}^{(2)}} s = \delta_i^{(3)} \circ a_j^{(2)}$$



$$\frac{d}{d w} s = f^{(3)} a^{(2)}$$

outer product.

$$\frac{\partial}{\partial w_{ij}^{(1)}} s = \delta_i^{(2)} \circ a_j^{(1)}$$

$$\delta_i^{(2)} = ?$$

$$\frac{d}{d}$$

$$\delta_i^{(3)} = u_i \circ f'(z_i^{(3)})$$

$$\delta_i^{(3)} = v_i \cdot f'(z_i^{(3)})$$

$$\delta^{(3)} = \underbrace{u \circ f'(z^{(3)})}_{\text{Hadamard product}}$$

$$x_{np}^m$$

1

$i \rightarrow$ layer \times no
 $x \rightarrow$ node name
 $m \rightarrow$ layer no
 $n \rightarrow$ row
 $p \rightarrow$ col

$$\frac{\partial}{\partial w_{ij}^{(1)}} \cdot f = \frac{\partial}{\partial w_{ij}^{(1)}} \cdot v^T a$$

$$= \frac{\partial}{\partial w_{ij}^{(1)}} + \sum_p v_p a_p^{(3)}$$

$$= \frac{\partial}{\partial w_{ij}} \sum_p v_p \cdot \frac{d}{\partial w_{ij}}^{(1)} \cdot a_p^{(3)}$$

$$= U_1 \cdot \frac{\partial}{\partial w_{ij}^{(1)}} a_1^{(3)} + U_2 \cdot \frac{\partial}{\partial w_{ij}^{(1)}} a_2^{(3)}$$

$$= v_1 \cdot \frac{d}{dw_{ij}^{(1)}} f(z_1^{(3)}) + v_2 \cdot \frac{d}{dw_{ij}^{(2)}} f(z_2^{(3)})$$

$$= v_1 \frac{d}{dw_i^{(1)}} f(w_1^{(2)} a^{(2)} + b_1^{(2)})$$

$$\frac{\partial f \cup_2}{\partial w_{ij}^{(1)}} b \left(w_2^{(2)} a^{(2)} + b_2^{(2)} \right)$$

$$\begin{pmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}$$

$w_{11}x_1 + w_{12}x_2 + w_{13}x_3$

$$2 \times 3 \times 3 \times 1$$

2+1

$$= v_1 f'(z_1^{(3)}) \frac{d}{d w_{ij}^{(1)}} \left(w_1^{(2)} a^{(2)} + b_1^{(2)} \right)$$

$$+ v_2 f'(z_2^{(3)}) \frac{d}{d w_{ij}^{(1)}} \left(w_2^{(2)} a^{(2)} + b_2^{(2)} \right)$$

$$= v_1 f'(z_1^{(3)}) \cdot \frac{d}{d w_{ij}^{(1)}} w_1^{(2)} a^{(2)}$$

$$+ v_2 f'(z_2^{(3)}) \frac{d}{d w_{ij}^{(1)}} \cdot w_2^{(2)} a^{(2)}$$

$$= v_1 f'(z_1^{(3)}) \cdot w_1^{(2)} \frac{d}{d w_{ij}^{(1)}} a^2 + v_2 f'(z_1^{(3)}) \frac{d}{d w_{ij}^{(1)}} w_2^{(2)} a^{(2)}$$

$$a^2 = f(z^{(2)}) = f(w^{(1)}x + b^{(1)})$$

$$\frac{d}{d w_{ij}^{(1)}} a^{(2)} = \frac{d}{d w_{ij}} f(w^{(1)}x + b^{(1)})$$

$\begin{bmatrix} \bullet & x_j \\ x_j & 0 \end{bmatrix}$

$$= f'(z^{(2)}) \cdot \frac{d}{d w_{ij}} (w^{(1)}x + b^{(1)})$$

$$= f'(z^{(2)}) \begin{bmatrix} x_j \\ 0 \end{bmatrix}$$

$$= b'(z^{(2)}) x_j$$

$$\partial \tilde{f} / \partial (w_p^{(2)} a_j^{(2)} b)$$

$$= \sum_P v_P \cdot \frac{\partial}{\partial w_{ij}^{(1)}} a_P^{(3)} \frac{z_P^{(3)}}{z_P}$$

$$= \sum_P v_P \cdot \frac{\partial}{\partial w_{ij}^{(1)}} f(w_P^{(2)} a^{(2)} + b_P^{(2)})$$

$$= \sum_P v_P f'(z_P^{(3)}) \cdot \frac{\partial}{\partial w_{ij}^{(1)}} a^{(2)}$$

$$= \sum_P v_P f'(z_P^{(3)}) w_P^{(2)} \cdot \frac{\partial}{\partial w_{ij}^{(1)}} f(z^{(2)})$$

$$= \sum_P v_P f'(z_P^{(3)}) w_P^{(2)} \cdot f'(z^{(2)}) \frac{\partial}{\partial w_{ij}^{(1)}} [w^{(2)} a^{(2)}]$$

$$= \sum_{P=1}^n v_P \underbrace{f'(z_P^{(3)}) w_P^{(2)}}_{\delta_P^{(2)}} \cdot f'(z^{(2)}) \underbrace{x_j}_{a_j^{(1)}}$$

$$= \sum_P \underbrace{\delta_P^{(3)}}_{\delta_P^{(2)}} \cdot w_P^{(2)} \cdot f'(z^{(2)}) \underbrace{x_j}_{a_j^{(1)}}$$

$$2 \times 3 \times 3 \times 1 \quad \frac{d}{d w_{ij}^{(1)}} \begin{bmatrix} w_{11}^T x_1 + w_{12}^T x_2 + w_{13}^T x_3 \\ w_{21}^T x_1 + w_{22}^T x_2 + w_{23}^T x_3 \end{bmatrix}$$

$$\begin{bmatrix} 0 \\ x_j \\ 0 \end{bmatrix} \odot \begin{bmatrix} x_j \\ 0 \end{bmatrix}$$

$$\therefore \frac{d}{d w_{ij}^{(1)}} s = \sum_p \delta_p^{(3)} \cdot w_p^{(2)} \cdot f'(z^{(2)}) x_j$$

$$\begin{array}{c} 2 \times 3 \\ \hline \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix} \\ \hline \begin{bmatrix} \delta_1^{(3)} & \delta_2^{(3)} & \delta_3^{(3)} \end{bmatrix} \end{array}$$

$$= f'(z^{(2)}) \cdot x_j \left[\sum_p \delta_p^{(3)} \cdot w_p^{(2)} \right]$$

$$= f'(z^{(2)}) \cdot x_j \left[\delta^{(3)} w^{(2)} \right]$$

$$= \underbrace{w^{(2)T}}_{\delta_i^{(2)}} \underbrace{\delta^{(3)} \circ f'(z^{(2)})}_{a_j^{(1)}} \cdot \underbrace{x_j}_{1 \times 1}$$

$$2 \times 3, 2 \times 1$$

$$3 \times 2, 2 \times 1$$

$$f'(z^{(2)}) x_j$$

$$[x_j, z^{(1)}]$$

$$\underbrace{\left(w^{(2)T} \cdot \delta^{(3)} \right) \circ f'(z^{(2)})}_{3 \times 1} x_j$$

~~w^{(2)T} \cdot \delta^{(3)} \circ f'(z^{(2)})~~

element-wise hadamard.

$$3 \times 1$$

$$\frac{\partial}{\partial w_{ij}^{(1)}} s = \delta^{(2)} x_j \cdot \underbrace{\delta^{(2)} \cdot x_j^T}_{1 \times 1}$$

$$\begin{array}{c} \delta_1^{(3)} \quad \delta_2^{(3)} \\ \hline \begin{bmatrix} w_{11}^{(2)} & w_{12}^{(2)} & w_{13}^{(2)} \\ w_{21}^{(2)} & w_{22}^{(2)} & w_{23}^{(2)} \end{bmatrix} \end{array}$$

$$1 \times 2 \cdot x_2 \cdot 3$$

$$x_1, x_2, x_3, \delta^{(2)} x_1, \delta^{(2)} x_2, \delta^{(2)} x_3$$

$$\delta^{(2)} x_1, \delta^{(2)} x_2, \delta^{(2)} x_3$$

$$g^{(l)} = ((w^{(l)})^T s^{(l+1)}) \circ f'(z^{(l)})$$

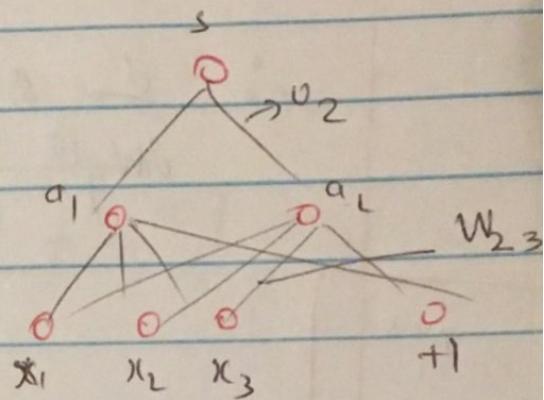
$$\delta = \max(0, 1 - s + s_c)$$

$$s = v^T a$$

$$= v^T f(z)$$

$$= v^T \cdot f(wx + b)$$

$\min \delta$.



$$\therefore \nabla \delta_{(v, w, x, b)} = ?$$

$$\nabla \delta = -\nabla s + \nabla s_c$$

$$\nabla s = \nabla v^T a$$

$$\left[\frac{\partial s}{\partial v} = \frac{\partial}{\partial v} v^T a = \frac{\partial}{\partial v} \sum_i v_i a_i = a \right] \rightarrow ①$$

$$\frac{\partial s}{\partial w} \cdot 0 \quad \frac{\partial s}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \sum_p v_p f(w_p x + b)$$

$$= \sum_p v_p \cdot \frac{\partial}{\partial w_{ij}} f(w_p x + b)$$

$$= \sum_p v_p \cdot f'(z_p) \frac{\partial}{\partial w_{ij}} (w_p x + b)$$

$$= \sum_p \text{sp. } \frac{\partial}{\partial w_j} (w_p x + b)$$

$$\Rightarrow \text{when } i \neq p \quad \frac{\partial}{\partial w_{ij}} (w_i x + b) = 0$$

$$i = p \quad \frac{\partial}{\partial w_{ij}} (w_i x + b) = x_j$$

$$\therefore \frac{\partial}{\partial w_{ij}} s = 0 + \delta_i x_j = \delta_i x_j$$

$\frac{\partial}{\partial w} s = \underline{\delta x^T}$
outer prod

→ ②

δ_1	δ_2	x_1
x_2	y_3	

$$\frac{\partial s}{\partial b_j} = \frac{\partial}{\partial b_j} \sum_i v_i \cdot f(w_i x + b_j)$$

$$= \sum_i v_i \cdot \frac{\partial}{\partial b_j} f(z_i)$$

$$= \sum_i v_i \cdot f'(z_i) \frac{\partial}{\partial b_j} (w_i x + b_j)$$

$$= \cancel{v_i \cdot f'(z_i)} \quad \text{when } j \neq i; \quad \frac{\partial}{\partial b_j} (z_i) = 0$$

else 1

$$= v_i \cdot f'(z_i)$$

$$\frac{\partial s}{\partial b_i} = \delta_i$$

$$\therefore \frac{\partial s}{\partial b} = \delta$$

$$\boxed{\frac{d}{db} s = \delta} \rightarrow ③$$

$$\frac{d}{dx_j} s = \frac{d}{dx_j} v^T a = \frac{d}{dx_j} \sum_i v_i a_i$$

$$= \sum_i v_i \cdot \frac{d}{dx_j} f(z_i)$$

w_1, w_{11}, w_{12}
 w_{11}, w_{22}, w_{23}

$$= \sum_i v_i \cdot f'(z_i) \frac{d}{dx_j} (w_i x + b_i)$$

$$= \sum_i v_i \cdot f'(z_i) \cdot w_{ij}$$

$$= \sum_i \delta_i w_{ij}$$

$\delta_1 w_{1j}$
 $+ \delta_2 w_{2j}$

$$= \delta w_{.j}^T = w_{.j} \cdot \delta^T$$

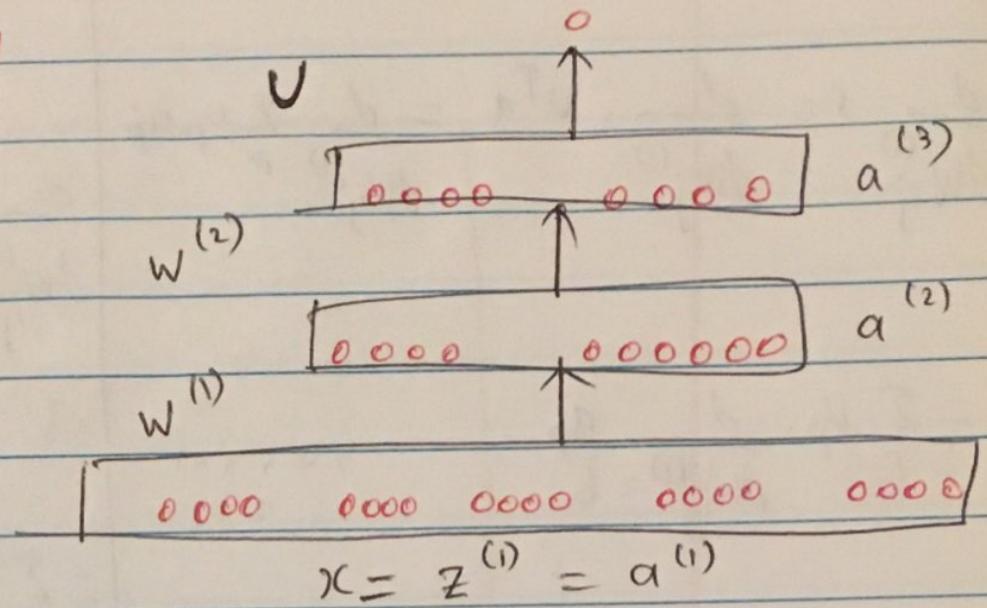
$$\therefore \frac{d}{dx} s = [w_{.1} \delta^T \quad w_{.2} \delta^T \quad w_{.3} \delta^T]$$

$$\begin{bmatrix} \delta_1 & \delta_2 \end{bmatrix} \downarrow \begin{bmatrix} w_{1j} \\ w_{2j} \end{bmatrix}$$

$$= \delta \cdot w^T$$

$$\boxed{\frac{d}{dx} s = w^T \delta} \rightarrow ④$$

$$\begin{aligned}
 x &= z^{(1)} = a^{(1)} \\
 z^{(2)} &= w^{(1)}x + b^{(1)} \\
 a^{(2)} &= f(z^{(2)}) \\
 z^{(3)} &= w^{(2)}a^{(2)} + b^{(2)} \\
 a^{(3)} &= f(z^{(3)}) \\
 s &= v^T a^{(3)}
 \end{aligned}$$



$$\frac{\partial}{\partial w^{(2)}} s = \delta^{(3)} a^{(2)T}$$

$$\delta^{(3)} = v \circ f'(z^{(3)})$$

$$\frac{\partial}{\partial w^{(1)}} f = \delta^{(2)} a^{(1)T}$$

$$\delta^{(2)} = (\quad) \circ f'(z^{(2)})$$

$$\begin{matrix} w^{(2)} & \overset{\rho}{\underset{\circ}{\circ}} \\ w^{(2)} & \overset{\circ}{\underset{\circ}{\circ}} & \overset{\circ}{\underset{\circ}{\circ}} \\ \overset{\circ}{\underset{\circ}{\circ}} & \overset{\circ}{\underset{\circ}{\circ}} & \overset{\circ}{\underset{\circ}{\circ}} \end{matrix} \quad p=2.$$

$$\frac{d}{dw_{ij}^{(1)}} s = \frac{d}{dw_{ij}^{(1)}} \cdot v^T a = \frac{d}{dw_{ij}^{(1)}} \sum_p v_p a_p$$

$$\frac{d}{dw_{ij}^{(1)}} w^{(1)} x = \begin{pmatrix} x_j & 0 \\ 0 & x_j \end{pmatrix}$$

$$= x_j$$

$$= \sum_p v_p \cdot f'(z_p) \cdot w^{(2)} x_j$$

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

$$1 \times 1 \quad 1 \times 1 \quad 2 \times 3 \quad 1 \times 1$$

$$= \sum_p \delta_p^{(3)} \cdot w^{(2)} \cdot x_j$$

$$\delta_1^{(2)} x_j + \delta_p^{(3)} w^{(2)} x_j + \delta_2^{(3)} w^{(2)} x_j$$

$$= x_j \bullet \delta^{(3)} \cdot w^{(2)}$$

$$= (a_j^{(1)} \circ \delta^{(3)}) \cdot w^{(2)}$$

$$\frac{d}{dw_{ij}^{(1)}} s = \frac{d}{dw_{ij}^{(1)}} \cdot v^T a = \frac{d}{dw_{ij}^{(1)}} \sum_p v_p a_p$$

$$= \sum_p v_p \cdot \frac{d}{dw_{ij}^{(1)}} \cdot a_p$$

$$= \sum_p v_p \cdot \frac{d}{dw_{ij}^{(1)}} \cdot f(z_p)$$

$$= \sum_p v_p \cdot f'(z_p) \cdot \frac{d}{dw_{ij}^{(1)}} \cdot z_p$$

$$= \sum_p v_p \cdot f'(z_p) \cdot \frac{d}{dw_{ij}^{(1)}} \cdot (w^{(2)} a^{(2)} + b^{(2)})$$

$$= \sum_p v_p \cdot f'(z_p) \cdot w^{(2)} \cdot \frac{d}{dw_{ij}^{(1)}} a^{(2)}$$

$$= \sum_p v_p \cdot f'(z_p) \cdot w^{(2)} \cdot \frac{d}{dw_{ij}^{(1)}} \cdot (w^{(1)} x + b^{(1)})$$

④

$$\frac{\partial s}{\partial v} = a \checkmark$$

$$\frac{\partial s}{\partial w} = \delta x^T$$

$$\delta = \begin{bmatrix} \delta_1 & \delta_2 \end{bmatrix}$$

$$\frac{\partial s}{\partial b} = \delta$$

$$\frac{\partial s}{\partial x} = w^T \delta$$

2×3 3×1
 2×1

$$\begin{bmatrix} w_{11}x_1 + w_{12}x_2 + w_{13}x_3 \\ w_{21}x_1 + w_{22}x_2 + w_{23}x_3 \end{bmatrix}$$

$$\frac{\partial}{\partial w_{ij}^{(1)}} s = \frac{\partial}{\partial w_{ij}^{(1)}} v^T f(z^{(3)})$$

$$= \frac{\partial}{\partial w_{ij}^{(1)}} \sum_{p=1}^2 v_p \cdot f(z_p^{(3)})$$

$$= \sum_{p=1}^2 v_p \cdot \frac{\partial}{\partial w_{ij}^{(1)}} f(z_p^{(3)})$$

$$= \sum_{p=1}^2 v_p \cdot f'(z_p^{(3)}) \cdot \frac{\partial}{\partial w_{ij}^{(1)}} z_p^{(3)}$$

$$= \sum_{p=1}^2 \delta_p^{(3)} \frac{\partial}{\partial w_{ij}^{(1)}} w^{(2)} a^{(2)} + b^{(2)}$$

$$= \sum_{p=1}^2 \delta_p^{(3)} w^{(2)} \frac{\partial}{\partial w_{ij}^{(1)}} w^{(1)} a^{(1)} + b^{(1)}$$

$$= \sum_{p=1}^2 \delta_p^{(3)} w^{(2)} \underbrace{\frac{\partial}{\partial w_{ij}^{(1)}}}_{f'(z^{(2)})} f(z^{(2)})$$

$$= \sum_{p=1}^2 \delta_p^{(3)} w^{(2)} f'(z^{(2)}) \frac{\partial}{\partial w_{ij}^{(1)}} w^{(1)} a^{(1)} + b^{(1)}$$

$$= \sum_{p=1}^2 \delta_p^{(3)} w^{(2)} f'(z^{(2)}) x_j$$

$\delta^{(2)}$

$$= \sum_{p=1}^2 \delta_p^{(3)} w^{(2)} f'(z^{(2)}) a_j^{(1)} \quad \delta^{(2)} = v \circ f'(z^{(2)})$$

|x1| 2x3 2x1 |x1|
 1x1 2x3 2x1 1x1
 $\frac{d}{\partial w_i^{(1)}} s = \delta^{(2)} a^{(1)T}$

 \Leftarrow $\delta^{(2)} =$

$$= f'(z^{(2)}) a_j^{(1)} \left[\sum_{p=1}^2 \delta_p^{(3)} w^{(2)} \right]$$

$$= f'(z^{(2)}) a_j^{(1)} \circ \delta^{(3)} w^{(2)T}$$

$$= \underbrace{\left(w^{(2)T} \delta^{(3)} \right)}_{\delta^{(2)}} \circ f'(z^{(2)}) a_j^{(1)}$$

 $\delta^{(2)}$

$$\delta^{(2)} = \left(w^{(2)T} \delta^{(3)} \right) \circ f'(z^{(2)})$$

$$\frac{\partial}{\partial w_{ij}^{(1)}} \delta = ?$$

$$z_p^{(3)} = w_p^{(2)} a^{(2)} + b^{(2)}$$

$$\frac{\partial}{\partial w_{ij}^{(1)}} v^T a = \frac{\partial}{\partial w_{ij}^{(1)}} \sum_p v_p a_i$$

$$= \sum_p v_p \cdot \frac{\partial}{\partial w_{ij}^{(1)}} a_p^{(2)}$$

$$= \sum_p v_p \cdot \frac{\partial}{\partial w_{ij}^{(1)}} f(z_p^{(3)})$$

$$= \sum_p v_p \cdot f'(z_p^{(3)}) \cdot \frac{\partial}{\partial w_{ij}^{(1)}} z_p^{(3)}$$

$$= \sum_p v_p \cdot f'(z_p^{(3)}) \cdot \left[w_p^{(2)} \cdot \frac{\partial}{\partial w_{ij}^{(1)}} \cdot f(w_{ij}^{(1)} a^{(1)} + b^{(1)}) \right]$$

$$= \sum_p v_p \cdot f'(z_p^{(3)}) \cdot \left[w^{(2)} \cdot f'(z^{(2)}) \cdot \frac{\partial}{\partial z^{(2)}} z^{(2)} \right]$$

$$= \sum_p s_p^{(3)} \cdot \left[w^{(2)} \cdot \left(f'(z^{(2)}) \cdot \frac{\partial}{\partial w_{ij}^{(1)}} (w^{(1)} a^{(1)} + b^{(1)}) \right) \right]$$

$$= \sum_p s_p^{(3)} \cdot \left[w^{(2)} \cdot \left(f'(z^{(2)}) \cdot x_j \right) \right]$$

$$= \sum_p s_p^{(3)} \cdot \left[w_p^{(2)} \cdot \left(f'(z^{(2)}) \cdot a_j^{(1)} \right) \right] \quad f^{(3)} \cdot x$$

$$\begin{matrix} 1 \times 1 & 1 \times 2 & 1 \times 2 & @1 \\ & \delta^{(2)} & & a^{(1)} \end{matrix}$$

2x1

$$\delta^{(2)} = \text{Uo} f'(z^{(2)}) \quad 2 \times 1$$

$$\delta^{(1)} = \text{Uo} \quad \delta^{(3)} w^{(2)T} \circ f'(z^{(2)}) \quad a_j^{(1)}$$

$$\frac{d}{d w^{(1)}} s = \underbrace{\left(w^{(1)T} \delta^{(3)} \right) \circ f'(z^{(2)})}_{\delta^{(2)}} \quad a_j^{(1)}$$

$$= \delta^{(2)} a^{(1)} \quad (2 \times 1) \quad 3 \times 1$$

$2 \times 1 \times 1 \times 5$

$$\delta^{(3)} = \text{Uo} f'(z^{(3)})$$

$$\delta^{(2)} = \left(w^{(2)} \delta^{(3)} \right) \circ f'(z^{(2)})$$

$$\therefore \delta^{(1)} = \boxed{w^{(1)} \delta^{(2)}} \circ f'(z^{(1)})$$

\therefore For last hidden layer, $\delta^{(n)} = \text{Uo} f'(z^{(n)})$

$$\delta^{(n)} = \left(w^{(n)T} \delta^{(n+1)} \right) \circ f'(z^{(n)})$$

$$\frac{d}{d w^{(2)}} s = \delta^{(3)} a^{(1)T}$$

$$\frac{d}{d w^{(1)}} s = \delta^{(2)} a^{(1)T}$$

$$\begin{aligned}
 & \leq \sum_p \delta_p^{(3)} \left[W_p^{(2)} \cdot \left(f'(z^{(2)}) \circ a_j^{(1)} \right) \right] \\
 & \quad \begin{matrix} 3 \times 2 & 1 \times 2 \\ 1 \times 2 & 2 \times 1 \\ 1 & 1 \end{matrix} \\
 & \quad \begin{matrix} 1 \times 1 & 2 \times 1 & 2 \times 1 & 1 \times 1 & 2 \times 1 \\ 1 \times 1 & 2 \times 2 & 2 \times 1 & 1 & 2 \times 1 \\ 1 \times 2 & 1 \times 1 & 2 \times 1 & 1 & 2 \times 1 \end{matrix} \\
 & = \sum_p \delta_p^{(3)} \left[\underbrace{W_p^{(2)} \cdot f'(z^{(2)})}_{\begin{matrix} 1 \times 2 \\ 1 \times 2 \end{matrix}} \right] \circ \underbrace{a_j^{(1)}}_{\begin{matrix} 1 \times 1 \\ 1 \times 1 \end{matrix}}
 \end{aligned}$$

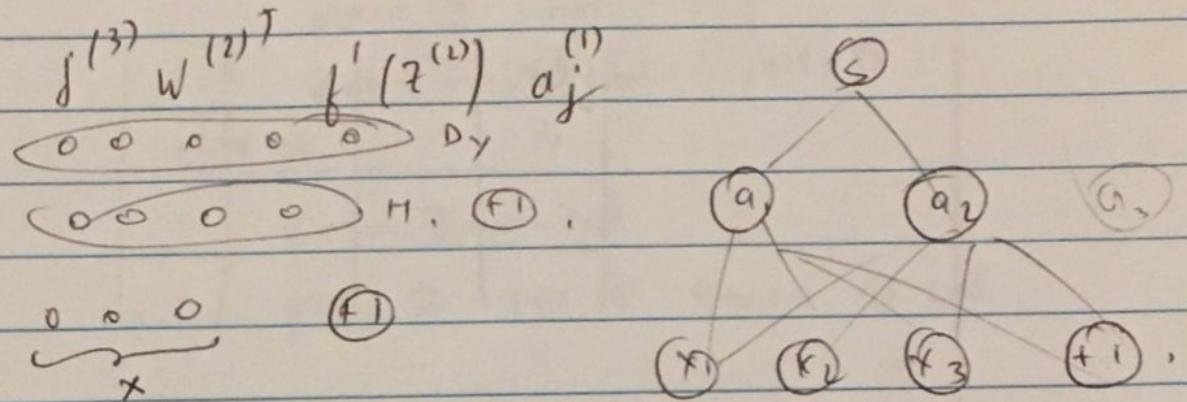
$$= s^{(3)} \alpha$$

$$s^{(3)} = [s_1^{(3)}, s_2^{(3)}]$$

$$\begin{aligned}
 a^{(1)} &= [a_1^{(1)}, a_2^{(1)}, a_3^{(1)}] \\
 &\quad \begin{matrix} 1 \times 1 & 2 \times 1 \\ 1 & 1 \end{matrix} \\
 & \quad \begin{matrix} f'(z^{(2)}) \circ a_j^{(1)} \end{matrix} \cdot \left[\sum_p \delta_p^{(3)} W_p^{(2)T} \right] \\
 & \quad \begin{matrix} f'(z^{(2)}) \circ a_j^{(1)} \end{matrix} \cdot \left[\delta^{(3)} \cdot \underbrace{W^{(2)T}}_{1 \times 1} \right] \\
 & \quad \left[\delta^{(3)} \cdot \underbrace{W^{(2)T}}_{1 \times 2} \right] \circ \underbrace{f'(z^{(2)}) \circ a_j^{(1)}}_{1 \times 1} \\
 & \quad \left[\delta^{(3)} \cdot \underbrace{W^{(2)T}}_{1 \times 2} \right]
 \end{aligned}$$

$\theta^{(1)} \text{ no. } \times \text{ no. } n$

$$\left[\delta^{(3)}, w^{(2)^T} \right]_{1 \times 2}, f'(z^{(2)})_{1 \times 2}, a_j^{(1)}_{1 \times 1}$$



$$(D_x * H) + (H * D_y) + \cancel{\star}$$

~~px~~

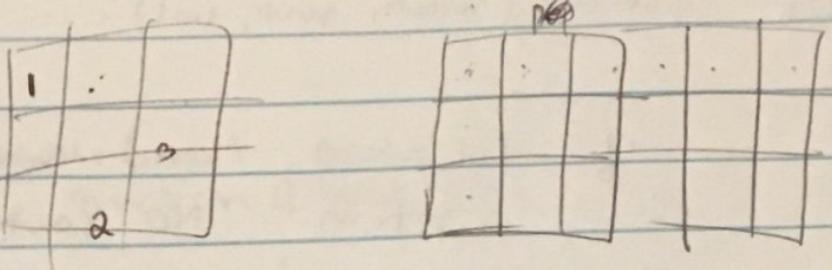
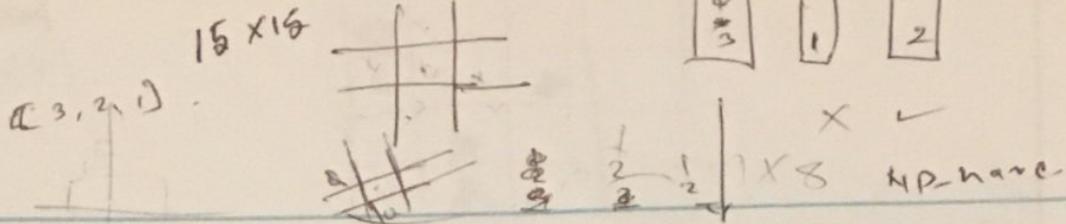
$$[(D_x + 1) * H]$$

~~px~~

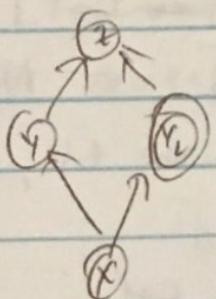
$$D_H \times \cancel{(x+1)}$$

$$+ [(H + 1) * D_y]$$

7 4 8x8
8x8
deb nancil
1
2
3



→ Compute ∇ loss example wise.



$$\frac{\partial Z}{\partial X} = \frac{\partial Z}{\partial Y_1} \frac{\partial Y_1}{\partial X} + \frac{\partial Z}{\partial Y_2} \frac{\partial Y_2}{\partial X}$$

$$\boxed{\frac{\partial Z}{\partial X} = \sum_{i=1}^n \frac{\partial Z}{\partial Y_i} \frac{\partial Y_i}{\partial X}}$$

canPlace

→ check (board, pos, num)

'Y' if occ

'N' if not occ.

→ complete getpos

X

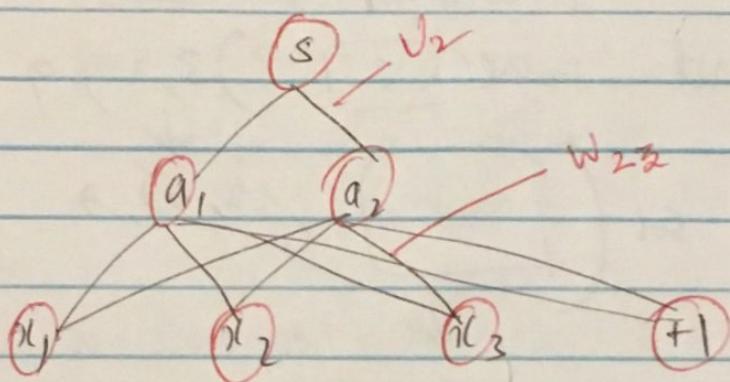
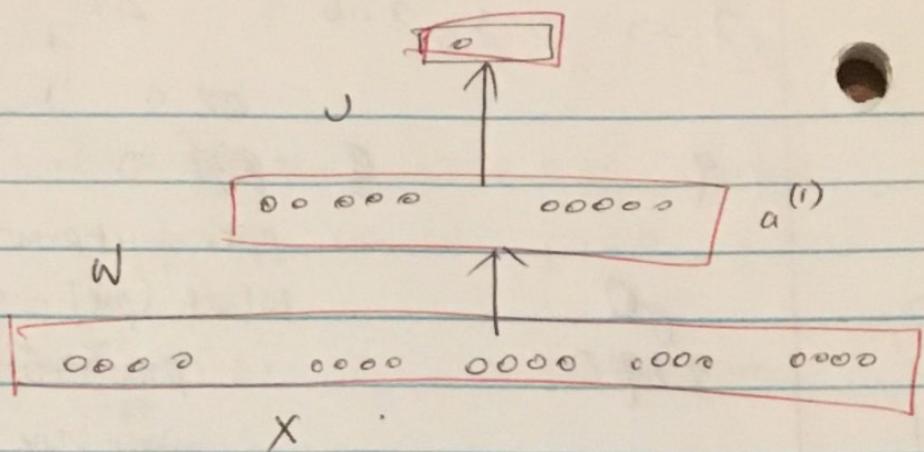
1-3, 1-6

1-3, 1-6

8 poss

1-3, 1-6

1-3, 1-6



$$s = v^T a = v^T f(z) = v^T f(w^T x + b)$$

$$s \in \mathbb{R}^{1*1}; v^T \in \mathbb{R}^{1*2}; z \in \mathbb{R}^{2*1}$$

$$w \in \mathbb{R}^{2*3}; b \in \mathbb{R}^{2*1}$$

$$\begin{matrix} 1*2 & 2*3, 3*1 & 1*1 \\ 1*2 & 2*1 & 2*1 \end{matrix}$$

$$\begin{matrix} 1*2, 2*1 \\ 1*1 \end{matrix}$$

$$F = \max(0, 1 - s + s_c)$$

Minimize \bar{F} .

$$\therefore \nabla \bar{F}_{(v, w, x, \rho)} = ?$$

$$\begin{aligned} s > s_c &\Rightarrow 0 \\ s \leq s_c &\Rightarrow P \\ s = s_c &\Rightarrow 1 \end{aligned}$$

$$\frac{d\bar{F}}{dx} = -\frac{ds}{da} + \frac{\partial s_c}{da}$$

$$\therefore \frac{ds}{du} = \frac{\partial v^T a}{\partial u} = \frac{\partial}{\partial u} \sum_i v_i a_i = a$$

$$\therefore \boxed{\frac{ds}{du} = a} \rightarrow \textcircled{1}$$

$$\frac{ds}{dw} = \frac{\partial v^T a}{\partial w} \quad \cancel{\text{if } \frac{\partial a}{\partial w}}$$

$$\text{Try for } w_{ij}; \quad \frac{\partial v^T a}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \sum_i v_i a_i$$

$$= v_i \frac{\partial}{\partial w_{ij}} a_i = v_i \frac{\partial}{\partial w_{ij}} f(z_i)$$

$$= v_i \frac{\partial}{\partial w_{ij}} f(w_i x + b_i) = v_i \cdot \cancel{\frac{\partial}{\partial z_i}} f(z_i) \cdot \frac{\partial}{\partial w_{ij}} (w_i x + b_i)$$

$$= v_i \cdot f'(z_i) \cdot x_i$$

$$= \cancel{v_i} \cdot \cancel{f'}$$

$$v_1 f(z_1) x_1 \quad v_1 f(z_2) x_2 \quad v_1 f(z_3) x_3 \\ v_2 f'(z_1) x_1 \quad \dots \quad v_2 f(z_2) x_3$$

$$\therefore \frac{\partial}{\partial w_{ij}} s = \underbrace{v_i \cdot f'(z_i)}_{s_i} \cdot x_j = s_i x_j$$

From $w_{ij} \rightarrow w$:

$$\frac{\partial}{\partial w} s = ?$$

$$\text{Let } v = [v_1 \ v_2] \\ f'(z) = [f'(z_1) \ f'(z_2)]$$

$$x = [x_1 \ x_2 \ x_3].$$

$$s = [s_1 \ s_2 \ s_3].$$

$$\frac{\partial}{\partial w} s = \underbrace{s x^T}_{\text{outer product}} \rightarrow \textcircled{2}$$

$$\frac{\partial}{\partial x_j} s = \frac{\partial}{\partial x_j} v^T a = \frac{\partial}{\partial x_j} \# \sum_i v_i \cdot a_i$$

$$= \sum_i v_i \cdot \frac{\partial}{\partial x_j} a_i$$

$$= \sum_i v_i \cdot \frac{\partial}{\partial x_j} f(z_i).$$

$$= \sum_i v_i \cdot \frac{\partial}{\partial x_j} f(w_i x + b_i)$$

$$= \sum_i v_i \frac{d}{dx_i} f(z_i) \cdot \frac{d}{dx_j} (w_i x + b_i)$$

$$= \sum_i v_i \cdot b'(z_i) \cdot w_{ij}$$

$$= \sum_i s_i \cdot w_{ij} = s \cdot w_{ij}^T$$

$$\therefore \boxed{\frac{d}{dx_j} s = s w_{ij}^T}$$

From $x_j \rightarrow x$

$$\boxed{\frac{d}{dx} s = s w^T} \rightarrow ③$$

$$\frac{d}{db} s = ?$$

$$\frac{d}{db_i} v^T a = \frac{d}{db_i} \sum_i v_i \cdot a_i$$

$$= v_i \cdot \frac{d}{db_i} a_i = v_i \cdot \frac{d}{db_i} f(w_i x + b_i)$$

$$= v_i \cdot \frac{d}{dx_i} f(z_i) \frac{d}{db_i} (w_i x + b_i)$$

$$= v_i \cdot f'(z_i) \cdot 1$$

$$= s_i$$

$$\therefore \frac{d s}{d b_i} = \delta_i$$

$$\boxed{\frac{d s}{d b} = \delta} \rightarrow ④$$

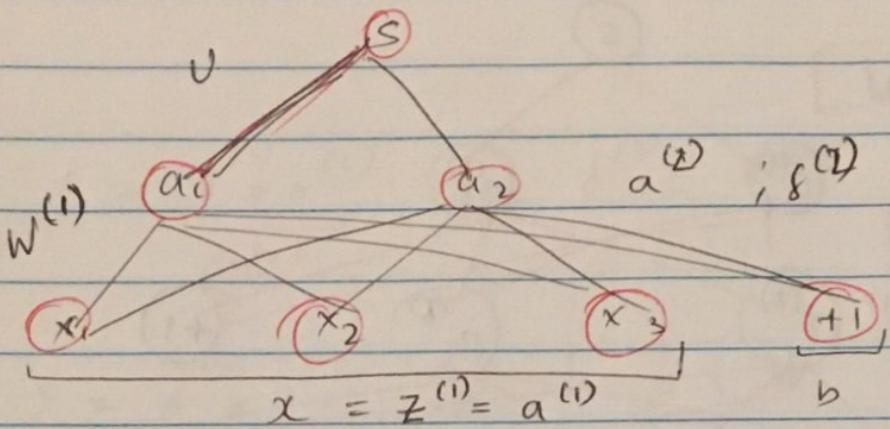
1, 2, 3, 4 \Rightarrow

$$\frac{d s}{d v} = a$$

$$\frac{d s}{d w} = \delta x^T ; \quad s = [\delta_1, \delta_2] \\ \text{Howard} \quad = [v_1 \sigma(z_1) \quad v_2 \sigma(z_2)]$$

$$\frac{d s}{d x} = \delta w^T$$

$$\frac{d s}{d b} = \delta$$



$$\frac{\partial s}{\partial v} = \alpha ; \quad \frac{\partial s}{\partial w} = \delta x^T ; \quad \delta = \begin{bmatrix} \delta_1 & \delta_2 \\ v_1 \sigma' f(z_1) & v_2 \sigma' f(z_2) \end{bmatrix}$$

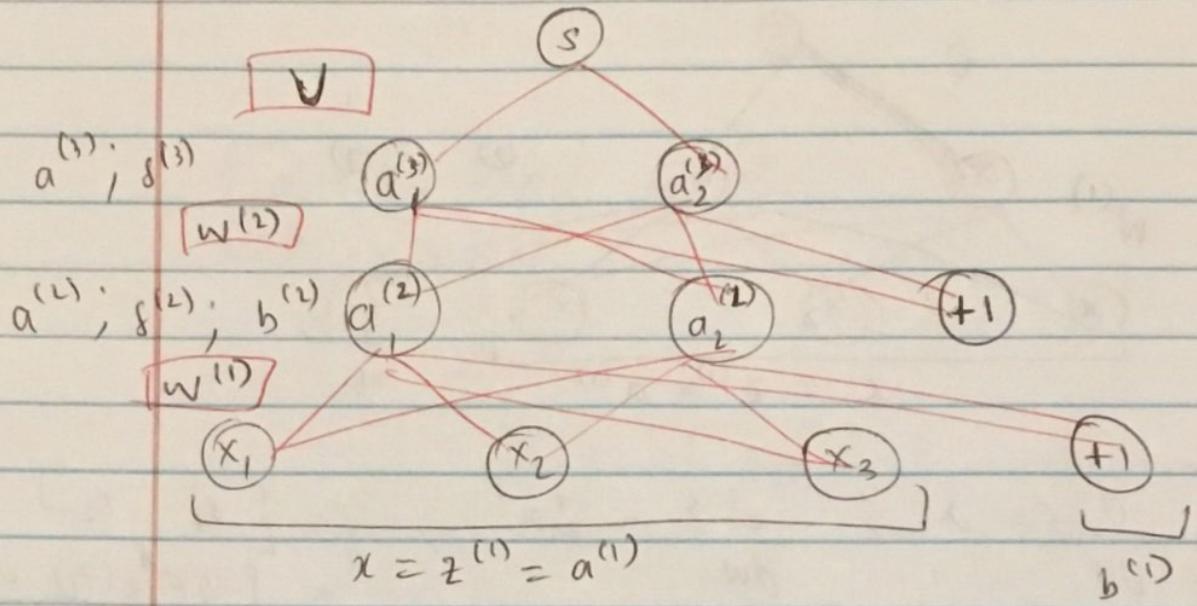
$$\frac{\partial s}{\partial x} = \delta w^T ; \quad \frac{\partial s}{\partial b} = \delta$$

Renaming:

$$\frac{\partial s}{\partial v} = \alpha^{(1)} ; \quad \frac{\partial s}{\partial w^{(1)}} = \delta^{(2)} \cdot \alpha^{(2)} ; \quad \frac{\partial s}{\partial a^{(1)}} = \delta^{(2)} \cdot a^{(1)}$$

$$\frac{\partial}{\partial b^{(1)}} = \delta^{(1)}$$

2 layer Neural Network:



To find:

$$\frac{d}{dV} S$$

$$\frac{d}{db^{(1)}} S$$

$$\frac{d}{w^{(1)}} S$$

$$\frac{d}{db^{(2)}} S$$

$$\frac{d}{dW^{(1)}} S$$

$$\frac{d}{da^{(2)}} S$$

$$\frac{d}{da^{(1)}} S$$

$250 \xrightarrow{0} 25\%$

$x \rightarrow 100$

$$\frac{d}{dx} s = a^{(3)} \quad 25$$

$\circ \Rightarrow$ element
 $\bullet \Rightarrow$ dot
 \Rightarrow outer

$2 \times$

$1 \times 2, 2 \times$

$2 \times 2 \times$

$3 \times 2, 2 \times 1$

$$\frac{d}{d w^{(2)}} s = \underbrace{\delta^{(3)} a^{(2)}}_{\text{outer}} ; \quad \delta^{(3)} = \cancel{a^{(3)}} \underbrace{v \circ b'(z^{(3)})}_{\text{element wise}}$$

$$\frac{d}{d w^{(1)}} s = \underbrace{\delta^{(2)} a^{(1)}}_{\text{outer}} ; \quad \delta^{(2)} = (W^{(2)} \delta^{(3)}) \circ b'(z^{(1)})$$

2×1

$2 \times 1 \quad \cancel{3 \times 1}$

$$\frac{d}{d x} s = \frac{d}{d a^{(1)}} s = \delta^{(1)} a^{(1)}$$

$$\therefore \boxed{\delta^{(1)} = (W^{(1)}^T \cdot \delta^{(1+1)}) \circ b'(z^{(1)})}$$

i) cr limit; acc. our bal, open dche.

30

30×25

$30 \times$

9

$$J(\theta)$$

$$\nabla J(\theta) \underset{\epsilon \rightarrow 0}{\approx} \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

choose $\epsilon \approx 10^{-9}$

\rightarrow If FP is correct

\rightarrow use gradient check to
check BP.

$$J(v, v, w, x, b)$$

$\frac{1}{2} ($

$$J(\theta) = \frac{1}{2} \theta^2 - \frac{3}{5}$$

$$\nabla J(\theta) = \frac{1}{2} .$$

$$g(\theta) = \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

~~J~~

$$J(\theta + \epsilon) = \frac{1}{2} \theta + \frac{1}{2} \epsilon - \frac{3}{5}$$

$$= \frac{1}{2} \cancel{\theta} + \frac{10^{-9}}{2} - \cancel{\frac{3}{5}}$$

$$J(\theta - \epsilon) = \cancel{\frac{1}{2} \theta} + \cancel{\frac{10^{-9}}{2}} + \cancel{\frac{3}{5}}$$

$$= \frac{10^{-9}}{2} * = \frac{1}{2} \frac{10^{-9}}{2} \times \frac{1}{2} \cancel{10^{-9}} \\ = \cancel{\frac{1}{2}} \cdot \cancel{10^{-9}} = \frac{1}{4} .$$

$$\theta = [\theta_1 \quad \theta_2].$$

$$J(\theta) \approx -\frac{1}{2}\theta - \frac{1}{3}.$$

$$\nabla J(\theta) =$$

$$\frac{\partial J(\theta)}{\partial \theta_i} = \frac{1}{2} \xrightarrow{\text{cancel } 2} \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

$$[\theta_1 \quad \theta_2]$$

$$\frac{\partial}{\partial \theta_i} [\epsilon \quad \theta] = c_i$$

$$\epsilon * e_i$$

$$\therefore \cancel{\frac{\partial}{\partial \theta_i}} \theta^{(i+)} = \theta + \epsilon \vec{e}_i$$

$$\frac{\partial J(\theta)}{\partial \theta_i} = \frac{J(\theta^{(i+)}) - J(\theta^{(i)})}{2\epsilon}$$

$$\theta^{(i-)} = \theta - \epsilon \vec{e}_i$$

$$\frac{\partial J(\theta)}{\partial \theta_i} = \frac{J(\theta^{(i+)}) - J(\theta^{(i)})}{2\epsilon}$$

$$\delta(\theta)$$

$$\frac{\partial}{\partial \theta} \delta(\theta) = \approx \theta^+$$

$$\frac{\partial J(\theta)}{\partial \theta} = \lambda t \frac{\delta(\theta + \epsilon) - \delta(\theta - \epsilon)}{2\epsilon}$$

$$\theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_n]$$

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

$$\frac{\partial}{\partial \theta_i} \delta(\theta) =$$

$$\theta^{i+} = \theta + \epsilon_i * \epsilon$$

$$\theta^{i-} = \theta - \epsilon_i * \epsilon$$

$$\frac{\partial J(\theta)}{\partial \theta_i} = \lambda t \frac{\delta(\theta^{i+}) - \delta(\theta^{i-})}{2\epsilon}$$

choose $\epsilon \approx 10^{-4}$

$$\frac{\partial J(\theta)}{\partial \theta_i} = \frac{\delta(\theta^{i+}) - \delta(\theta^{i-})}{2\epsilon}$$

$$\frac{\partial \delta(\theta)}{\partial \theta_i} \approx \frac{\delta(\theta)}{\delta \theta_i}$$

$$\delta(\theta) = x^2$$

$$\nabla \delta(\theta) = 2x$$

$$\frac{\partial \delta(\theta)}{\partial x} = 2x = 2(1) = 2$$

$$\delta(\theta) = h_\theta(x)$$

$$J(\theta) = \sum_{i=1}^m (h_\theta(x^i) - y^i)^2$$

[1 2]

$$\Rightarrow \delta(\theta) = \text{cost} = \text{sum}([1 4])$$

- 5.

$$\nabla \delta(\theta) - \text{grad} = \text{sum}([1 2] \text{ err}^i)$$

$$= \frac{1}{2} [1 2]$$

$$= 2 \times 3 = 6$$

$$= [2 4]$$

5, [2 4])

[1, 2)

↓

b

$$(5, [2 4])$$

$$g(\theta) = \theta^2 ; \quad \frac{\partial g(\theta + c) - g(\theta)}{2c}$$

$\begin{bmatrix} 1 & 2 \end{bmatrix}$
 \downarrow
 b
 \downarrow
 $(2, [1 \ 4])$

$\begin{bmatrix} 1 & 2 \end{bmatrix}$
 \downarrow
 ∇
 \downarrow
 $\overbrace{2 \cdot 10^{-4}}$

123.456

$$Y = x^2 ; \quad \frac{\partial Y}{\partial x} = 2x$$

$$Y = 15,241.383 ; \quad \nabla Y = 246.912.$$

$$x = [2 \dots 3]^\top$$

$$x = [123.456].$$

$$\text{cost} = f(x) \quad \text{cost, grad} = f(x)$$

~~grad~~

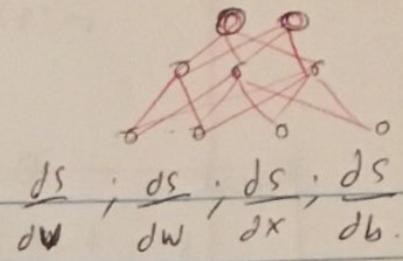
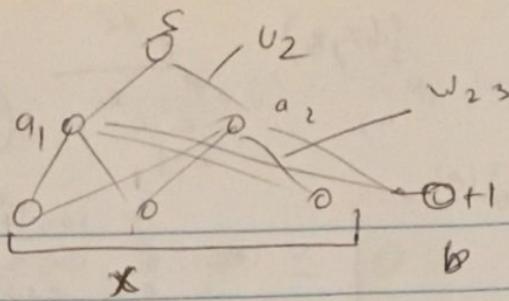
$$f(x) \quad \downarrow \quad \nabla$$

$$\text{cost} = 15,241.383 \quad \text{grad} = [246.912]$$

$$[1, 2] \quad \downarrow$$

\downarrow
 $f(x)$

\swarrow \searrow
 $\text{cost} = 2$ $\text{grad} = [2, 4]$
 \nwarrow \nearrow
 $g(x)$



$$s = v^T a = v^T f(z) = v^T f(z_1, z_2)$$

$$s = v^T (f(z_1)^{t_1} \quad f(z_2)^{t_2})$$

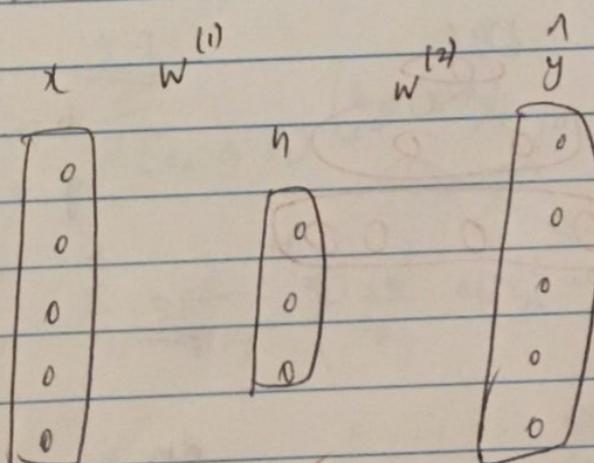
$$s = v^T (f(w_1^T x)^{t_1} \quad f(w_2^T x)^{t_2})$$

$$s = v^T (\sigma(w_1^T x)^{t_1} \quad \sigma(w_2^T x)^{t_2})$$

$$\sigma(z) = 1/(1+e^{-z})$$

$$\begin{aligned} &(\text{D}_x + 1 \neq \text{H}) \\ &+ (\text{D}_y + 1 \neq \text{D}_Y) \end{aligned}$$

$$s = v^T a$$



$$n = \text{sigmoid}(w^{(1)} x + b^{(1)})$$

$$\begin{aligned} &D_H * D_X = D_X \times 1 \\ &D_H * 1 + D_H * 1 \\ &D_H * 1 \end{aligned}$$

$$\hat{y} = \text{softmax}(w^{(2)} n + b^{(2)})$$

$$\begin{aligned} &D_Y * D_H = D_H \times 1 \\ &D_Y * 1 + D_Y * 1 \\ &D_Y * 1 \end{aligned}$$

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

$$f(x_i)$$

$$\frac{e^{x_i}}{\sum_{i=1}^n e^{x_i}}$$

(+1)

Repeat

$$\theta = \theta + \eta \nabla_\theta J(\theta)$$

$$\theta = \theta + \eta \frac{\partial J(\theta)}{\partial \theta}$$

$$y_1 \quad y_2 \\ 1 \quad 0 \quad \approx 100$$

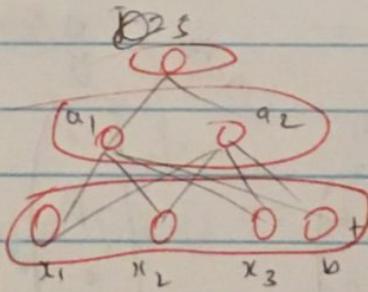
$$\begin{matrix} \phi & 1 & \approx 0 \\ 0 & 0 & \approx \\ 0 & 1 & \approx \end{matrix}$$

$$\frac{\partial}{\partial (x,y)} f(x,y) = \frac{\partial}{\partial x} f(x,y) + \frac{\partial}{\partial y} f(x,y)$$

$$\frac{\partial S}{\partial (v, \theta, w, x, b)} = \frac{\partial S}{\partial v} + \frac{\partial S}{\partial w} + \frac{\partial S}{\partial x} + \frac{\partial S}{\partial b}.$$

Forward: \rightarrow , δ .

Backprop:



FP:

$$h = \text{Sigmoid}(w^{(1)}x + b^{(1)})$$

$$y = \text{Softmax}(w^{(2)}h + b^{(2)})$$

BP:

$$\nabla(w^{(1)}, w^{(2)}, b^{(1)}, b^{(2)}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(\hat{y}^{(i)})$$

$$CE = - \sum_{i=1}^m y_i \log(\hat{y}_i)$$

$$-\infty \leq \log(x) \leq 0$$

$$0 \leq x \leq 1$$

$x \rightarrow 1$
for

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 = 0$$

$$\log(x) \rightarrow 0 ; x \rightarrow 1$$

$$\log(x) \rightarrow -\infty ; x \rightarrow 0$$

Needed: Cost, grad (w_1, w_2, b_1, b_2).

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^i) - y^i)^2$$

$\frac{\partial}{\partial \theta} \text{ cost}$
 $2x \frac{\partial}{\partial \theta} (h_\theta y)$

$$\nabla J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial \theta} (h_\theta(x^i) - y^i)$$

$$= \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial \theta} (h_\theta(x^i)) - 0$$

$\theta_0 x_0$
 $\theta_j x_j$

$$= \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial [\theta_0, \theta_1, \theta_2]} (h_\theta(x^i))$$

$$= \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^d \frac{\partial}{\partial \theta_j} \left(\sum_{k=1}^d \theta_k x_{ik}^j \right)$$

$k = j ; \quad x_j^i$

$$= \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^d x_j^i$$

$$\frac{\partial}{\partial(x,y)} f(x,y) \quad \frac{\partial}{\partial x} f(\cancel{x},y) + \frac{\partial}{\partial y} f(x,\cancel{y}).$$

Repeat

{

$$\theta = \theta + \alpha \nabla J(\theta)$$

3

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_\theta(x^i) - y^i)^2$$

$$\begin{aligned} \theta \cdot \frac{\partial}{\partial \theta} J(\theta) &= \frac{1}{2} \sum_{i=1}^m \cancel{\theta} \frac{\partial}{\partial \theta} (h_\theta(x^i) - y^i) \\ &= \sum_{i=1}^m \frac{\partial}{\partial \theta} (h_\theta(x^i)) (h_\theta(x^i) - y^i) \end{aligned}$$

$$= \sum_{i=1}^m \frac{\partial}{\partial \theta} \sum_{j=1}^d \theta_j x_j^i$$

$$= \sum_{i=1}^m \left[\frac{\partial}{\partial \theta_1} \sum_{j=1}^d \theta_j x_j^i \quad \frac{\partial}{\partial \theta_2} \sum_{j=1}^d \theta_j x_j^i \quad \dots \quad \frac{\partial}{\partial \theta_d} \sum_{j=1}^d \theta_j x_j^i \right]$$

$$= \sum_{i=1}^m [x_1^i \ x_2^i \ \dots \ x_d^i]$$

$$= \sum_{i=1}^m n^i (h_\theta(x^i) - y^{(i)})$$

$$\therefore \theta = \theta + \alpha \sum_{i=1}^m x^i (h_\theta(x^i) - y^{(i)})$$

$$\theta_j = \theta_j + \alpha \nabla_{\theta_j} J(\theta)$$

$$\theta_j = \theta_j + \alpha \sum_{i=1}^m x_j^i (h_{\theta}(x^i) - y^i)$$

$$\theta = \theta + \alpha \sum_{i=1}^m x^i (h_{\theta}(x^i) - y^i)$$

$$\theta = \theta + \underline{\alpha}$$

$$\begin{bmatrix} x_1^i & x_2^i \\ 1 & 5 \end{bmatrix}$$

$$\delta(w, w_2, b_1, b_2) =$$

Repeat .

$$[w, w_2, b_1, b_2] = [w, w_2, b_1, b_2] + \eta \nabla J(w, w_2, b_1, b_2)$$

?

$$\nabla J(w, w_2, b_1, b_2) = \nabla_w J(w, w_2, b_1, b_2) + \nabla_{w_2} J(\text{parus})$$

+ ∇

$$\frac{\partial}{\partial w^{(1)}} \delta = \delta^{(1+1)} \cdot a^{(1)^T}$$

$$\delta^{(1)} = (w^{(1)} \cdot \delta^{(1+1)}) \circ f^{-1}(z^{(1)})$$

$$\frac{\partial}{\partial b^{(1)}} \delta = \delta^{(1)}$$

Needed: $\frac{d}{dW^{(1)}} \quad \frac{d}{dW^{(2)}} \quad \frac{d}{dB^{(1)}} \quad \frac{d}{dB^{(2)}} = ?$

(cost = ?)

$$\begin{matrix} & 0 & & 0 \\ D_X & 0 & D_H & D_Y \\ & 0 & 0 & 0 \\ & 0 & 0 & 0 \\ & 0 & & 0 \end{matrix}$$

No of params: $((D_X+1) * D_H) + (D_Y * (D_H+1))$

$$\text{cost } (y, \hat{y}) = - \sum_{i=1}^m y_i \log(\hat{y}_i)$$

Can:

Forward: T , softmax

Backprop:

$$\frac{d}{dW^{(1)}} \quad \frac{d}{dW^{(2)}}$$

Cannot:

Cost:

y_i

$$1, 2, 3. \quad \begin{matrix} n \\ n \\ n=1 \end{matrix}$$

$$x + 2x + 3x = n$$

Abstract Solution:

`lweights, = backprop();`
`ly biases`

$$Cost = - \sum_i y_i \cdot \log(\hat{y}_i) \rightarrow$$

return (cost,eweights,bbiases)

$$\hat{y}_i = \text{Softmax}(\mathbf{w}_2^T \mathbf{h} + b_2)$$

↓
 $D_y \times D_H \rightarrow D_H \times 1$
 $D_y \times 1$ $D_y \times 1$
 $D_y \times 1$

$$\hat{y}_t = \text{softmax} \left(\frac{\beta}{D_{Y \times 1}} \right)$$

$$\sigma \left(w^{(1)} x + b^{(1)} \right)$$

$$a^{(i)} = \sigma(w^{(i)T}x + b^{(i)})$$

$$\hat{y}_i = [0.009 \quad \cancel{0.01}] \quad (t \leftarrow \text{small.})$$

0 + \log(0.01)

initialize network weights

do

for Each ex: train example:

pred = forward(nn, ex)

actual = label(ex)

$$\text{error} = \text{actual} - \text{pred}$$

Complete A Who

Compute Δw_{IH}

update n/w weights

until params wont change

rehm nn .

$$W^{(1)} \xrightarrow{\alpha^{(2)}} W^{(2)}$$

$$\underbrace{0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0}_{x = a^{(1)}} \quad] \quad \begin{matrix} 0 \\ b^{(1)} \end{matrix}$$

$$\frac{a^{(1)} \cdot (1)}{a}, a$$

$\frac{d}{dw} (z)$