

Derivatives for paper

Backpropagation formulas:

$$\frac{\partial \mathcal{L}}{\partial w_{ij}} = \delta_j \sigma_i$$

$$\delta_j = \phi'(\text{net}_j) \sum_{k \in \text{output nodes}} \delta_k w_{jk}$$

$$\text{where } \delta_j = \frac{\partial \mathcal{L}}{\partial \sigma_j} \frac{\partial \sigma_j}{\partial \text{net}_j}$$

1) $\sigma_i = \text{softmax error at node } i$

$$x_i \rightarrow y_i \text{ output}$$

$$y_i = \text{softmax}(W x_i)$$

$$\Rightarrow \delta_i = \underbrace{\phi'(\text{net}_i)}_{f'(x_i)} \otimes \underbrace{\delta_i - t_i}_{g_i - t_i} \underbrace{\sum_j w_{ji} \delta_j}_{\frac{\partial (W x_i)}{\partial x_i}}$$

$$\Rightarrow \delta_i = f'(x_i) \otimes (y_i - t_i) W_s$$

$$= W_s^T (y_i - t_i) \otimes f'(x_i)$$

Derivative with respect to V at node p_2 :

$$\text{Using } \frac{\partial \mathcal{L}}{\partial w_{ij}} = \delta_j \sigma_i = \delta_j \frac{\partial \text{net}_j}{\partial w_{ij}}$$

$$\text{Let } j = p_2; i \in \{a, p_1\}$$

$$\text{Let } w_{ij} = V$$

$$\Rightarrow \frac{\partial \mathcal{L}_{p_2}}{\partial V_{app}} = \delta_{p_2, \text{com}} \frac{\partial \Gamma(a, p_a)}{\partial V_{app}}$$

(2)

where $\Gamma(a, p_1) = \begin{bmatrix} a \\ p_1 \end{bmatrix}^T V \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} a \\ p_1 \end{bmatrix} + W \begin{bmatrix} a \\ p_1 \end{bmatrix}$
 (Let $v = a \oplus p_1$)

$$\text{or } \Gamma(a, p_1)^T = v^T V_{\alpha\beta} v^\beta + W_{\gamma\lambda} v^\lambda$$

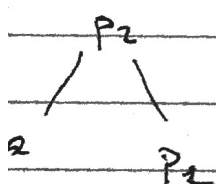
$$\Rightarrow \frac{\partial \Gamma(a, p_1)}{\partial v_{\alpha\beta}} \delta = \delta_{\alpha}^{\lambda} v^{\lambda} v^{\beta}$$

$$\Rightarrow \frac{\partial \delta p_2}{\partial v_{\alpha\beta}} = \delta p_{2, \text{com}} \delta \delta_{\alpha}^{\lambda} v^{\lambda} v^{\beta}$$

$$= \delta p_{2, \text{com}} v^{\lambda} v^{\beta}$$

($v = a \oplus p_1$)
 i.e. $v^{\lambda} = a^{\lambda} + p_1^{\lambda}$

Error manages for the two children
 of p_2



By the backpropagation formula,

$$\delta^{a, \text{com}} = f'(a) \sum_{\text{to output node}} \delta_k \frac{\partial \text{net}_k}{\partial \text{out}_j}$$

$$\Rightarrow \delta^{a, \text{com}} = f'(a) \delta p_{2, \text{com}} \frac{\partial \Gamma(a, p_1)}{\partial a} + f'(a) \delta^{a, s}$$

Similarly,

$$\delta p_{2, \text{com}} = f'(p_1) \delta p_{2, \text{com}} \frac{\partial \Gamma(a, p_1)}{\partial p_1}$$

$$+ f'(p_1) \delta p_{2, s}$$

$$\Gamma(a, p_1)^T = v^T V_{\alpha\beta} v^\beta + W_{\gamma\lambda} v^\lambda$$

($v = a \oplus p_1$)

$$= (a^{\alpha} + p_1^{\alpha})(a^{\beta} + p_1^{\beta}) V_{\alpha\beta}$$

$$+ W_{\gamma\lambda} (a^{\gamma} + p_1^{\gamma})$$

③

$$\Rightarrow \frac{\partial \Gamma(a, p_{\pm})}{\partial a} = W a + U^T V a + U^T V b$$

(same for $\frac{\partial \Gamma(a, p_{\pm})}{\partial p_{\pm}}$)

$$= W + (U^T V + V^T U)^T \begin{bmatrix} a \\ p_{\pm} \end{bmatrix}$$

(in paper's notation)

Simplifying,

$$g_{a, com} = f'(a) \otimes \left\{ g_{p_{\pm}, com} \frac{\partial \Gamma(a, p_{\pm})}{\partial a} + g_{a, s} \right\}$$

$$g_{p_{\pm}, com} = f'(p_{\pm}) \otimes \left\{ g_{p_{\pm}, com} \frac{\partial \Gamma(a, p_{\pm})}{\partial p_{\pm}} + g_{p_{\pm}, s} \right\}$$

Or in the paper's notation,

$$g_{i, s} \rightarrow f'(k_i) \otimes g_{i, s}$$

$$g_{a, com} = f'(a) \otimes g_{p_{\pm}, com} \frac{\partial \Gamma(a, p_{\pm})}{\partial a} + g_{a, s}$$

$$g_{p_{\pm}, com} = f'(p_{\pm}) \otimes g_{p_{\pm}, com} \frac{\partial \Gamma(a, p_{\pm})}{\partial p_{\pm}} + g_{p_{\pm}, s}$$

where $\frac{\partial \Gamma(a, p_{\pm})}{\partial a} = \frac{\partial \Gamma(a, p_{\pm})}{\partial p_{\pm}}$

$$= W + (U^T V + V^T U)^T \begin{bmatrix} a \\ p_{\pm} \end{bmatrix}$$

(4)

Full derivative for \mathcal{L} w.r.t V

By the chain rule,

$$\frac{\partial \mathcal{L}}{\partial V} = \sum_{\text{nodes } i} \frac{\partial \mathcal{L}}{\partial x_i} \frac{\partial x_i}{\partial \text{net}_i} \frac{\partial \text{net}_i}{\partial V}$$

$$= \sum_{\text{nodes } i} g_{i, \text{com}} \frac{\partial \text{net}_i}{\partial V}$$

 $\frac{\partial \text{net}_i}{\partial V} \neq 0$ only for $i \in \{P_1, P_2\}$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial V} = g_{P_1, \text{com}} \frac{\partial \text{net}_{P_1}}{\partial V} + g_{P_2, \text{com}} \frac{\partial \text{net}_{P_2}}{\partial V}$$

$$\frac{\partial \text{net}_{P_1}}{\partial V} = \frac{\partial P(b, c)}{\partial V} \quad \frac{\partial \text{net}_{P_2}}{\partial V} = \frac{\partial \mathcal{L}_{P_2}}{\partial V}$$

$$= \begin{bmatrix} b \\ c \end{bmatrix} \begin{bmatrix} b \\ c \end{bmatrix}^T \quad \text{(computed earlier)}$$

(computed earlier)

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial V \propto b \wedge} = g_{P_1, \text{com}} \propto (b \downarrow + c \downarrow) (b \uparrow + c \uparrow)$$

$$+ \frac{\partial \mathcal{L}_{P_2}}{\partial V \propto b \wedge}$$