

1. Derivative of Softmax Function

$$t \neq j \quad r_{a_j}^{(L)} = \text{soft}(\mathbf{r}_z^{(L)}, j) = \frac{e^{r_{z_j}^{(L)}}}{\sum_s e^{r_{z_s}^{(L)}}} \quad \text{soft}(\mathbf{r}_{z_{\neq t}}^{(L)}, j) = \frac{e^{r_{z_j}^{(L)}}}{\sum_{s \neq t} e^{r_{z_s}^{(L)}}}$$

$$\ln \text{soft}(\mathbf{r}_z^{(L)}, j) = \ln \frac{e^{r_{z_j}^{(L)}}}{\sum_s e^{r_{z_s}^{(L)}}} = r_{z_j}^{(L)} - \ln \sum_s e^{r_{z_s}^{(L)}}$$

$$\ln(1 - \text{soft}(\mathbf{r}_z^{(L)}, j)) = \ln \frac{\sum_{s \neq j} e^{r_{z_s}^{(L)}}}{\sum_s e^{r_{z_s}^{(L)}}} = \ln \sum_{s \neq j} e^{r_{z_s}^{(L)}} - \ln \sum_s e^{r_{z_s}^{(L)}}$$

$$\ln \text{soft}(\mathbf{r}_z^{(L)}, t) = \ln \frac{e^{r_{z_t}^{(L)}}}{\sum_s e^{r_{z_s}^{(L)}}} = r_{z_t}^{(L)} - \ln \sum_s e^{r_{z_s}^{(L)}}$$

$$\ln(1 - \text{soft}(\mathbf{r}_z^{(L)}, t)) = \ln \frac{\sum_{s \neq t} e^{r_{z_s}^{(L)}}}{\sum_s e^{r_{z_s}^{(L)}}} = \ln \sum_{s \neq t} e^{r_{z_s}^{(L)}} - \ln \sum_s e^{r_{z_s}^{(L)}}$$

$$\frac{\partial}{\partial r_{z_j}^{(L)}} \ln \text{soft}(\mathbf{r}_z^{(L)}, j) = \frac{\partial}{\partial r_{z_j}^{(L)}} \left(r_{z_j}^{(L)} - \ln \sum_s e^{r_{z_s}^{(L)}} \right) = 1 - \frac{e^{r_{z_j}^{(L)}}}{\sum_s e^{r_{z_s}^{(L)}}} = 1 - \text{soft}(\mathbf{r}_z^{(L)}, j)$$

$$\frac{\partial}{\partial r_{z_j}^{(L)}} \ln(1 - \text{soft}(\mathbf{r}_z^{(L)}, j)) = \frac{\partial}{\partial r_{z_j}^{(L)}} \left(\ln \sum_{s \neq j} e^{r_{z_s}^{(L)}} - \ln \sum_s e^{r_{z_s}^{(L)}} \right) = -\frac{e^{r_{z_j}^{(L)}}}{\sum_s e^{r_{z_s}^{(L)}}} = -\text{soft}(\mathbf{r}_z^{(L)}, j)$$

$$\frac{\partial}{\partial r_{z_j}^{(L)}} \ln \text{soft}(\mathbf{r}_z^{(L)}, t) = \frac{\partial}{\partial r_{z_j}^{(L)}} \left(r_{z_t}^{(L)} - \ln \sum_s e^{r_{z_s}^{(L)}} \right) = -\frac{e^{r_{z_j}^{(L)}}}{\sum_s e^{r_{z_s}^{(L)}}} = -\text{soft}(\mathbf{r}_z^{(L)}, j)$$

$$\frac{\partial}{\partial r_{z_j}^{(L)}} \ln(1 - \text{soft}(\mathbf{r}_z^{(L)}, t)) = \frac{\partial}{\partial r_{z_j}^{(L)}} \left(\ln \sum_{s \neq t} e^{r_{z_s}^{(L)}} - \ln \sum_s e^{r_{z_s}^{(L)}} \right) = \frac{e^{r_{z_j}^{(L)}}}{\sum_{s \neq t} e^{r_{z_s}^{(L)}}} - \frac{e^{r_{z_j}^{(L)}}}{\sum_s e^{r_{z_s}^{(L)}}} = \text{soft}(\mathbf{r}_{z_{\neq t}}^{(L)}, j) - \text{soft}(\mathbf{r}_z^{(L)}, j)$$

2. Derivative of Softmax Cross Entropy Function

$$C = - \sum_{t=1}^{m^{(L)}} \sum_r ({}^r y_t \ln \text{soft}({}^r \mathbf{z}^{(L)}, t) + (1 - {}^r y_t) \ln(1 - \text{soft}({}^r \mathbf{z}^{(L)}, t)))$$

$$= - \sum_r ({}^r y_j \ln \text{soft}({}^r \mathbf{z}^{(L)}, j) + (1 - {}^r y_j) \ln(1 - \text{soft}({}^r \mathbf{z}^{(L)}, j))) - \sum_{\substack{t=1 \\ t \neq j}}^{m^{(L)}} \sum_r ({}^r y_t \ln \text{soft}({}^r \mathbf{z}^{(L)}, t) + (1 - {}^r y_t) \ln(1 - \text{soft}({}^r \mathbf{z}^{(L)}, t)))$$

$$\frac{\partial C}{\partial {}^r z_j^{(L)}} = - \left({}^r y_j (1 - \text{soft}({}^r \mathbf{z}^{(L)}, j)) + (1 - {}^r y_j) (-\text{soft}({}^r \mathbf{z}^{(L)}, j)) \right) - \sum_{\substack{t=1 \\ t \neq j}}^{m^{(L)}} \left({}^r y_t (-\text{soft}({}^r \mathbf{z}^{(L)}, j)) + (1 - {}^r y_t) (\text{soft}({}^r \mathbf{z}_{\neq t}^{(L)}, j) - \text{soft}({}^r \mathbf{z}^{(L)}, j)) \right)$$

$$= - \left({}^r y_j - \text{soft}({}^r \mathbf{z}^{(L)}, j) \right) - \sum_{\substack{t=1 \\ t \neq j}}^{m^{(L)}} \left((1 - {}^r y_t) \text{soft}({}^r \mathbf{z}_{\neq t}^{(L)}, j) - \text{soft}({}^r \mathbf{z}^{(L)}, j) \right)$$

$$= - \left({}^r y_j - \text{soft}({}^r \mathbf{z}^{(L)}, j) \right) - \sum_{\substack{t=1 \\ t \neq j}}^{m^{(L)}} (1 - {}^r y_t) \text{soft}({}^r \mathbf{z}_{\neq t}^{(L)}, j) + \sum_{\substack{t=1 \\ t \neq j}}^{m^{(L)}} \text{soft}({}^r \mathbf{z}^{(L)}, j)$$

$$= \text{soft}({}^r \mathbf{z}^{(L)}, j) - {}^r y_j + (m^{(L)} - 1) \text{soft}({}^r \mathbf{z}^{(L)}, j) - \sum_{\substack{t=1 \\ t \neq j}}^{m^{(L)}} (1 - {}^r y_t) \text{soft}({}^r \mathbf{z}_{\neq t}^{(L)}, j)$$

$$= m^{(L)} \text{soft}({}^r \mathbf{z}^{(L)}, j) - {}^r y_j - \sum_{\substack{t=1 \\ t \neq j}}^{m^{(L)}} (1 - {}^r y_t) \text{soft}({}^r \mathbf{z}_{\neq t}^{(L)}, j)$$

3. Backpropagation

$$C = - \sum_{t=1}^{m^{(L)}} \sum_r (r_{y_t} \ln \text{soft}(\mathbf{r}^{(L)}, t) + (1 - r_{y_t}) \ln(1 - \text{soft}(\mathbf{r}^{(L)}, t)))$$

$$r_{z_j}^{(L)} = b_j^{(L)} + \sum_i w_{i,j}^{(L)} r_{a_i}^{(L-1)}$$

$$\frac{\partial C}{\partial b_j^{(L)}} = \sum_r \frac{\partial C}{\partial r_{z_j}^{(L)}} \frac{\partial r_{z_j}^{(L)}}{\partial b_j^{(L)}} = \sum_r \frac{\partial C}{\partial r_{z_j}^{(L)}}$$

$$\frac{\partial C}{\partial w_{i,j}^{(L)}} = \sum_r \frac{\partial C}{\partial r_{z_j}^{(L)}} \frac{\partial r_{z_j}^{(L)}}{\partial w_{i,j}^{(L)}} = \sum_r \frac{\partial C}{\partial r_{z_j}^{(L)}} r_{a_i}^{(L-1)}$$

$$\frac{\partial C}{\partial r_{a_i}^{(L-1)}} = \sum_j \frac{\partial C}{\partial r_{z_j}^{(L)}} \frac{\partial r_{z_j}^{(L)}}{\partial r_{a_i}^{(L-1)}} = \sum_j \frac{\partial C}{\partial r_{z_j}^{(L)}} w_{i,j}^{(L)}$$

$$r_{a_j}^{(l)} = \text{act } r_{z_j}^{(l)} \quad r_{z_j}^{(l)} = b_j^{(l)} + \sum_i w_{i,j}^{(l)} r_{a_i}^{(l-1)}$$

$$r_{a_j}^{(1)} = \text{act } r_{z_j}^{(1)} \quad r_{z_j}^{(1)} = b_j^{(1)} + \sum_i w_{i,j}^{(1)} r_{x_i}$$

$$\frac{\partial C}{\partial b_j^{(l)}} = \sum_r \frac{\partial C}{\partial r_{a_j}^{(l)}} \frac{\partial r_{a_j}^{(l)}}{\partial r_{z_j}^{(l)}} \frac{\partial r_{z_j}^{(l)}}{\partial b_j^{(l)}} = \sum_r \frac{\partial C}{\partial r_{a_j}^{(l)}} \frac{\partial r_{a_j}^{(l)}}{\partial r_{z_j}^{(l)}}$$

$$\frac{\partial C}{\partial b_j^{(1)}} = \sum_r \frac{\partial C}{\partial r_{a_j}^{(1)}} \frac{\partial r_{a_j}^{(1)}}{\partial r_{z_j}^{(1)}} \frac{\partial r_{z_j}^{(1)}}{\partial b_j^{(1)}} = \sum_r \frac{\partial C}{\partial r_{a_j}^{(1)}} \frac{\partial r_{a_j}^{(1)}}{\partial r_{z_j}^{(1)}}$$

$$\frac{\partial C}{\partial w_{i,j}^{(l)}} = \sum_r \frac{\partial C}{\partial r_{a_j}^{(l)}} \frac{\partial r_{a_j}^{(l)}}{\partial r_{z_j}^{(l)}} \frac{\partial r_{z_j}^{(l)}}{\partial w_{i,j}^{(l)}} = \sum_r \frac{\partial C}{\partial r_{a_j}^{(l)}} \frac{\partial r_{a_j}^{(l)}}{\partial r_{z_j}^{(l)}} r_{a_i}^{(l-1)}$$

$$\frac{\partial C}{\partial w_{i,j}^{(1)}} = \sum_r \frac{\partial C}{\partial r_{a_j}^{(1)}} \frac{\partial r_{a_j}^{(1)}}{\partial r_{z_j}^{(1)}} \frac{\partial r_{z_j}^{(1)}}{\partial w_{i,j}^{(1)}} = \sum_r \frac{\partial C}{\partial r_{a_j}^{(1)}} \frac{\partial r_{a_j}^{(1)}}{\partial r_{z_j}^{(1)}} r_{x_i}$$

$$\frac{\partial C}{\partial r_{a_i}^{(l-1)}} = \sum_j \frac{\partial C}{\partial r_{a_j}^{(l)}} \frac{\partial r_{a_j}^{(l)}}{\partial r_{z_j}^{(l)}} \frac{\partial r_{z_j}^{(l)}}{\partial r_{a_i}^{(l-1)}} = \sum_j \frac{\partial C}{\partial r_{a_j}^{(l)}} \frac{\partial r_{a_j}^{(l)}}{\partial r_{z_j}^{(l)}} w_{i,j}^{(l)}$$