

# Custom Neural Networks

## 1 Vanilla 2-layers NN

### Forward

$$X \in \mathcal{R}^{n \times m}, W_1 \in \mathcal{R}^{d \times m}, W_2 \in \mathcal{R}^{k \times dn}, B_1 \in \mathcal{R}^{d \times n}, B_2 \in \mathcal{R}^k, y \in \mathcal{R}^k$$

$$pre_{a_1} = \frac{X - \mathbb{E}(X)}{\sqrt{\mathbb{E}(X - \mathbb{E}X)^2}} \in \mathcal{R}^{n \times m} \quad (1)$$

$$a_1 = W_1 pre_{a_1}^\top + B_1 \in \mathcal{R}^{d \times n} \quad (2)$$

$$h_1 = \text{ReLU}(a_1) \in \mathcal{R}^{d \times n} \quad (3)$$

$$pre_{a_2} = Flatten(h_1) \in \mathcal{R}^{dn} \quad (4)$$

$$a_2 = (W_2 pre_{a_2}^\top)^\top + B_2 \in \mathcal{R}^k \quad (5)$$

$$h_2 = a_2 I_k \in \mathcal{R}^k \quad (6)$$

$$s = Softmax(h_2) \in \mathcal{R}^k \quad (7)$$

$$\mathcal{L} = -(y \odot \ln s) \mathbf{1}_k^\top \in \mathcal{R} \quad (8)$$

$$\mathcal{L}(X, W_1, W_2, y) = \text{CrossEntropy}(y, Softmax(W_2 Flatten(\text{ReLU}(W_1 X + B1)) + B2))$$

### Backward

We have

$$\left( \frac{d\mathcal{L}}{dW_2} \right)_{k \times dn} = \left( \frac{\partial \mathcal{L}}{\partial s} \right)_k \left( \frac{\partial s}{\partial h_2} \right)_{k \times k} \left( \frac{\partial h_2}{\partial a_2} \right)_{k \times k} \left( \frac{\partial a_2}{\partial W_2} \right)_{k \times k \times dn} \quad (1)$$

or, equivalently,

$$\left( \frac{d\mathcal{L}}{dW_2} \right)_{k \times dn} = \left( \frac{\partial \mathcal{L}}{\partial h_2} \right)_k \left( \frac{\partial h_2}{\partial a_2} \right)_{k \times k} \left( \frac{\partial a_2}{\partial W_2} \right)_{k \times k \times dn} \quad (2)$$

In a simplified vectorized form, we have

$$\frac{d\mathcal{L}}{dW_2} = ((s - \mathbf{y}) \odot \mathbf{1}_k)^\top \mathbf{pre}_{\mathbf{a}_2} \quad (3)$$

Meanwhile

$$\begin{aligned} \left( \frac{d\mathcal{L}}{dW_1} \right)_{d \times m} &= \left( \frac{\partial \mathcal{L}}{\partial h_2} \right)_k \left( \frac{\partial h_2}{\partial a_2} \right)_{k \times k} \left( \frac{\partial a_2}{\partial pre_{a_2}} \right)_{k \times dn} \left( \frac{\partial pre_{a_2}}{\partial h_1} \right)_{dn \times d \times n} \\ &\quad \left( \frac{\partial h_1}{\partial a_1} \right)_{d \times n \times d \times n} \left( \frac{\partial a_1}{\partial W_1} \right)_{d \times n \times d \times m} \end{aligned} \quad (4)$$

Therefore,

$$\frac{d\mathcal{L}}{dW_1} = ((s - \mathbf{y}) \odot \mathbf{1}_k \mathbf{W}_2)_{dn \rightarrow d \times n} \odot \left( \frac{\partial \text{ReLU}(\mathbf{a}_1)}{\partial \mathbf{a}_1} \right)_{d \times n} \mathbf{pre}_{\mathbf{a}_1} \quad (5)$$

Biases gradients are computed as follows:

$$\frac{d\mathcal{L}}{dB_2} = (s - \mathbf{y}) \odot \mathbf{1}_k \quad (6)$$

$$\frac{d\mathcal{L}}{dB_1} = ((s - \mathbf{y}) \odot \mathbf{1}_k \mathbf{W}_2)_{dn \rightarrow d \times n} \odot \left( \frac{\partial \text{ReLU}(\mathbf{a}_1)}{\partial \mathbf{a}_1} \right)_{d \times n} \quad (7)$$

## 2 1-layer Transformer Encoder

### Forward

$$X \in \mathcal{R}^{n \times d_0}, W_1 \in \mathcal{R}^{d_0 \times d_1}, \text{cls\_tok} \in \mathcal{R}^{1 \times d_1}, W_{\{q,k,v\}} \in \mathcal{R}^{d_1 \times d_2}, W_2 \in \mathcal{R}^{d_2 \times K}, \\ y \in \mathcal{R}^K$$

$$h_0 = \frac{X - \mathbb{E}(X)}{\sqrt{\mathbb{E}(X - \mathbb{E}X)^2}} \in \mathcal{R}^{n \times d_0} \quad (1)$$

$$h_1 = h_0 W_1 \in \mathcal{R}^{n \times d_1} \quad (2)$$

$$h_{pre_2} = \text{concat}(h_1, \text{cls\_tok}) \in \mathcal{R}^{(n+1) \times d_1} \quad (3)$$

$$\{Q, K, V\} = h_{pre_2} W_{\{q,k,v\}} \in \mathcal{R}^{(n+1) \times d_2} \quad (4)$$

$$S = \text{Softmax}(QK^\top) \in \mathcal{R}^{(n+1) \times (n+1)} \quad (5)$$

$$h_2 = SV + h_{pre_2} \in \mathcal{R}^{(n+1) \times d_2} \quad (6)$$

$$z = h_{2(n+1)} W_2 \in \mathcal{R}^K \quad (7)$$

$$s = \text{Softmax}(z) \in \mathcal{R}^K \quad (8)$$

$$\mathcal{L} = -\frac{1}{K} \mathbf{1}_{d_2} (y \odot \ln s)^\top \in \mathcal{R} \quad (9)$$

$$\mathcal{L}(X, W_1, \text{cls\_tok}, W_{\{q,k,v\}}, W_2, y) = \\ \text{CrossEntropy}(y, \text{Softmax}(\text{SelfAttention}(\text{concat}((XW_1), \text{cls\_tok}), W_{\{q,k,v\}})[:, \text{seq\_len+1}, :], W_2)))$$

### Backward

$$\frac{d\mathcal{L}}{dW_2} = ((z - y)^\top h_{2(n+1)})^\top = h_{2(n+1)}^\top (z - y) \quad (1)$$

$$\frac{d\mathcal{L}}{dW_v} = (((z - y)W_2^\top)^\top S_{(n+1)} h_{pre_2})^\top = h_{pre_2}^\top S_{(n+1)}^\top (z - y)W_2^\top \quad (2)$$

Now, with

$$\alpha = \text{diag}(S_{(n+1)}) - (Q_{(n+1)} K^\top)^\top (Q_{(n+1)} K^\top) \quad (3)$$

$$= \text{diag}(S_{(n+1)}) - K Q_{(n+1)}^\top Q_{(n+1)} K^\top \quad (4)$$

$$= \frac{\partial S_{(n+1)}}{\partial Q_{(n+1)} K^\top} \quad (5)$$

we have

$$\frac{d\mathcal{L}}{dW_k} = (((z - y)\mathbf{W}_2^\top)^\top \mathbf{V}^\top \boldsymbol{\alpha})^\top \mathbf{Q}_{(n+1)} \mathbf{h}_{pre_2})^\top \quad (6)$$

$$= \mathbf{h}_{pre_2}^\top \mathbf{Q}_{(n+1)}^\top ((z - y)\mathbf{W}_2^\top)^\top \mathbf{V}^\top \boldsymbol{\alpha} \quad (7)$$

$$\frac{d\mathcal{L}}{dW_q} = \left( (((z - y)\mathbf{W}_2^\top)^\top \mathbf{V}^\top \boldsymbol{\alpha} \mathbf{K})^\top \mathbf{h}_{pre_{2(n+1)}} \right)^\top \quad (8)$$

$$= \mathbf{h}_{pre_{2(n+1)}}^\top ((z - y)\mathbf{W}_2^\top)^\top \mathbf{V}^\top \boldsymbol{\alpha} \mathbf{K} \quad (9)$$

And with

$$\boldsymbol{\beta} = \mathbf{Q}_{(n+1)} \mathbf{W}_k^\top + \mathbf{W}_q \mathbf{K}^\top \quad (10)$$

we have

$$\frac{d\mathcal{L}}{d\text{cls\_tok}} = (((z - y)\mathbf{W}_2^\top)^\top (\mathbf{S}_{(n+1)} \mathbf{W}_v + \boldsymbol{\alpha} \boldsymbol{\beta} \mathbf{V} + \mathbf{1}))^\top \quad (11)$$

$$\frac{d\mathcal{L}}{dW_1} = (((z - y)\mathbf{W}_2^\top)^\top (\mathbf{S}_{(n+1)} \mathbf{W}_v \mathbf{h}_0 + \boldsymbol{\alpha} \boldsymbol{\beta} \mathbf{h}_0 \mathbf{V} + \mathbf{h}_0))^\top \quad (12)$$