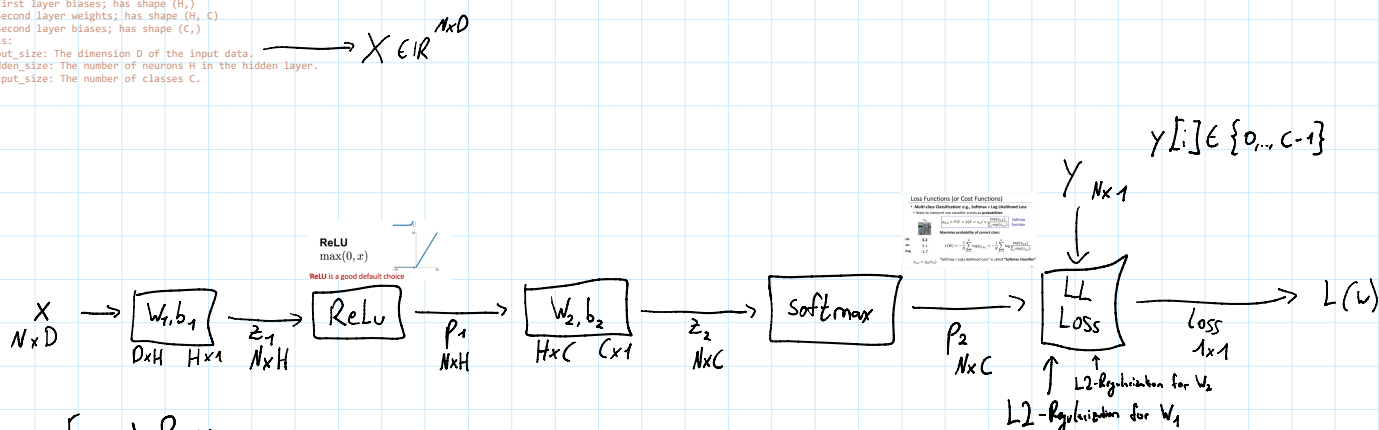


input - fully connected layer - ReLU - fully connected layer - softmax

- W1: First layer weights; has shape (D, H)
- b1: First layer biases; has shape (H,)
- W2: Second layer weights; has shape (H, C)
- b2: Second layer biases; has shape (C,)

Inputs:

- input_size: The dimension D of the input
- hidden_size: The number of neurons H in
- output_size: The number of classes C


$$z_1 = \overset{N \times D}{X} \cdot \overset{D \times H}{W_1} + \left(\overset{N \times 1}{\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}} \cdot \overset{1 \times H}{b_1} \right)^T \in \mathbb{R}^{N \times H}$$

$$p_1 = \max(0, z_1^{NH}) \in \mathbb{R}^{N \times H}$$

$$Z_2 = \overset{N \times H}{p_1} \cdot \overset{H \times C}{W_2} + \left(\overset{N \times 1}{\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}} \cdot \overset{1 \times C}{b_2} \right) \in \mathbb{R}^{N \times C}$$

$$\rho_2[i:j][c] = \frac{\exp(z_2[i:c])}{\sum_{j=1}^c \exp(z_2[i:j])} \quad c \in \{1, \dots, C\} \quad \left| \quad \rho_2 = \frac{\exp(z_2)}{\text{sum}(\exp(z_2), \text{axis}=1, \text{keepdims}=True)} \in \mathbb{R}^{N \times C}$$

$$L(w) = -\frac{1}{N} \cdot \sum_{i=1}^N \log(p_i[i, y[i]]) + L_2(w_1) + L_2(w_2) \in \mathbb{R}^{1 \times 1}$$

with $L2(V_1) = \text{reg} \cdot \text{sum}(V_1 \ast \ast 2)$, $L2(W_2) = \text{reg} \cdot \text{sum}(W_2 \ast \ast 2)$

Backpropagation, Gradients: (biases still have to be adjusted to batch-average)

$$\begin{aligned} \frac{\partial L}{\partial b_2} &= \frac{\partial L}{\partial z_2} = \frac{p_2}{\partial z_2} \cdot \frac{\partial L}{\partial p_2} \\ &= p_2 - \left(1 \text{ where } [i, y [i]] \right) \in \mathbb{R}^{N \times C} \quad (\text{average over } N) \end{aligned}$$

$$\frac{\partial L}{\partial W_2} = \left(\frac{\partial L}{\partial W_{2,1}}, \frac{\partial L}{\partial W_{2,2}}, \dots, \frac{\partial L}{\partial W_{2,C}} \right)^T = \left(\frac{\partial z_2}{\partial W_{2,1}} \cdot \frac{\partial L}{\partial z_2}, \frac{\partial z_2}{\partial W_{2,2}} \cdot \frac{\partial L}{\partial z_2}, \dots, \frac{\partial z_2}{\partial W_{2,C}} \cdot \frac{\partial L}{\partial z_2} \right)^T$$

$$= \begin{matrix} H \times N \\ P_1^T \end{matrix} \cdot \begin{matrix} N \times C \\ \frac{\partial L}{\partial z_1} \end{matrix} \quad \text{EIR}^{H \times C}$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial z_1}{\partial b_1} \cdot \frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial z_1} = \frac{\partial p_1}{\partial z_1} \cdot \frac{\partial L}{\partial p_1} = \frac{\partial p_1}{\partial z_1} \cdot \frac{\partial z_2}{\partial p_1} \cdot \frac{\partial L}{\partial z_2}$$

$$= \frac{\partial L}{\partial z_1} \cdot \overbrace{w_2^T}^{N \times C \cdot \frac{\partial z_2}{\partial p_1}} \cdot \underbrace{(1 \text{ if } z_2^T [i] > 0, \text{ else } 0)}_{\text{element-wise}} \in \mathbb{R}^{N \times H} \quad (\text{average over } N)$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial z_1}{\partial w_1} \cdot \frac{\partial L}{\partial z_1} = \frac{D_{xN}}{X^T} \cdot \left(\frac{\partial L}{\partial b_1} \right)^T \in \mathbb{R}^{D_{xH}}$$