

Feed-Forward Neural Network

Forward Propagation

Single Data Point

$$\begin{aligned} m &= W^{(1)}x + b^{(1)} \\ h &= \text{ReLU}(m) \\ z_2 &= W^{(2)}h + b^{(2)} \\ y &= \text{softmax}(z_2) \\ \mathcal{L} &= \mathcal{L}_{\text{CE}}(y, t) \end{aligned}$$

Batch

$$\begin{aligned} M &= W^{(1)}x + 1(b^{(1)})^T \\ \mathcal{H} &= \text{ReLU}(M) \\ \mathcal{Z} &= W^{(2)}\mathcal{H} + 1(b^{(2)})^T \\ \mathcal{Y} &= \text{softmax}(\mathcal{Z}) \\ \mathcal{L} &= \frac{1}{N} \sum_i \mathcal{L}_{\text{CE}}(y^{(i)}, t^{(i)}) \end{aligned}$$

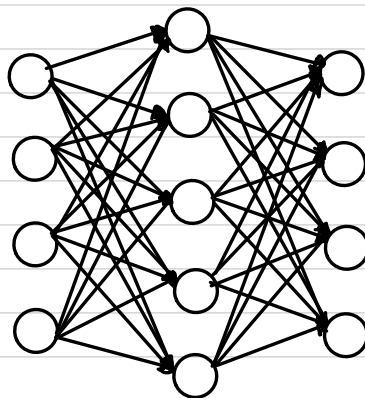
Backpropagation

$$\begin{aligned} \bar{\mathcal{L}} &= 1 \\ \bar{z}_2 &= y - t \\ \bar{h} &= \frac{\partial \mathcal{L}}{\partial z_2} = \bar{z}_2 \\ \bar{W}^{(2)} &= \bar{z}_2 (h)^T \\ \bar{b}^{(2)} &= \bar{z}_2 \\ \bar{m} &= \bar{h} \circ \text{ReLU}'(m) \\ \bar{W}^{(1)} &= \bar{m} (x)^T \\ \bar{b}^{(1)} &= \bar{m} \end{aligned}$$

$$\begin{aligned} \bar{\mathcal{L}} &= 1 \\ \bar{\mathcal{Z}} &= \frac{1}{N} (Y - \mathcal{Y}) \\ \bar{\mathcal{H}} &= \frac{\partial \mathcal{L}}{\partial \mathcal{Z}} W^{(2)} \\ \bar{W}^{(2)} &= (\bar{\mathcal{Z}})^T \mathcal{H} \\ \bar{b}^{(2)} &= (\bar{\mathcal{Z}})^T 1 \\ \bar{M} &= \bar{\mathcal{H}} \circ \text{ReLU}'(M) \\ \bar{W}^{(1)} &= (\bar{M})^T \mathcal{X} \\ \bar{b}^{(1)} &= (\bar{M})^T 1 \end{aligned}$$

Adam (Adaptive Moment Estimation)

$$\begin{aligned} m_t &= \beta_1 m_{t-1} + (1 - \beta_1) \frac{\partial \mathcal{L}}{\partial w_i} \\ v_t &= \beta_2 v_{t-1} + (1 - \beta_2) \left(\frac{\partial \mathcal{L}}{\partial w_i} \right)^2 \\ \hat{m}_t &= \frac{m_t}{1 - \beta_1^t} \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t} \\ w_{t+1} &= w_t + \alpha \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \end{aligned}$$



Gradient Clipping

$$g = \eta \frac{g}{\|g\|}$$

Weight Initialization

Xavier Initialization

$$x = \sqrt{\frac{6}{n_{\text{inputs}} + n_{\text{outputs}}}}$$

sigmoid or softmax

$$W \sim \text{Uniform}(-x, x)$$

He / Kaiming Initialization

$$W \sim \text{Uniform}\left(-\sqrt{\frac{6}{n_{\text{inputs}}}}, \sqrt{\frac{6}{n_{\text{outputs}}}}\right) \text{ ReLU}$$

Dropout

$$h' = \begin{cases} 0 & \text{with probability } p \\ \frac{h}{1-p} & \text{otherwise} \end{cases}$$

Note:

* for code get rid of 1 and transpose \rightarrow let numpy work