

# The core of Neural Network

This is to give a very brief summary of essence of Neural Network.

## Model

Simple Neural Network

Forward Neural Network

$$[z]_{1 \times r} = [h]_{1 \times n} [w]_{n \times r} + [b]_{1 \times r}$$

$$[\hat{y}]_{1 \times r} = \text{soft max}([z]_{1 \times r})$$

$$L = [y]_{1 \times r} \circ \log[\hat{y}]_{1 \times r}$$

Back Propagation

$$\frac{\partial L}{\partial z} = [\hat{y}]_{1 \times r} - [y]_{1 \times r}$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial w} = ([\hat{y}]_{1 \times r} - [y]_{1 \times r})_{n \times 1} \circ \text{transpose}([h]_{r \times 1})$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial b} = [\hat{y}]_{1 \times r} - [y]_{1 \times r}$$

$$\frac{\partial L}{\partial h} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial h} = ([\hat{y}]_{1 \times r} - [y]_{1 \times r}) [w]_{n \times r}^T$$

Variant with dropout

Forward

$$[r]_{1 \times n} \sim \text{Ber}(p)$$

$$[\tilde{h}]_{1 \times n} = [r]_{1 \times n} \circ [h]_{1 \times n}$$

$$[z]_{1 \times r} = [\tilde{h}]_{1 \times n} [w]_{n \times r} + [b]_{1 \times r}$$

$$[\hat{y}]_{1 \times r} = \text{soft max}([z]_{1 \times r})$$

$$L = [y]_{1 \times r} \circ \log[\hat{y}]_{1 \times r}$$

Backward

$$\frac{\partial L}{\partial z} = [\hat{y}]_{1 \times r} - [y]_{1 \times r}$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial w} = ([\hat{y}]_{1 \times r} - [y]_{1 \times r})_{n \times 1} \circ \text{transpose}([h]_{r \times 1})$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial b} = [\hat{y}]_{1 \times r} - [y]_{1 \times r}$$

$$\frac{\partial L}{\partial \tilde{h}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial \tilde{h}} = ([\hat{y}]_{1 \times r} - [y]_{1 \times r}) [w]_{n \times r}^T$$

$$\frac{\partial L}{\partial h} = \frac{\partial L}{\partial \tilde{h}} \frac{\partial \tilde{h}}{\partial h} = ([\hat{y}]_{1 \times r} - [y]_{1 \times r}) [w]_{n \times r}^T \circ [r]_{1 \times n}$$

## **Solver**

Stochastic Gradient Descent

Momentum

$$v_t = \mathcal{W}_{t-1} + \eta \nabla_{\theta} J(\theta)$$

$$\theta = \theta - v_t$$

NAG

$$v_t = \mathcal{W}_{t-1} + \eta \nabla_{\theta} J(\theta - \mathcal{W}_{t-1})$$

$$\theta = \theta - v_t$$

## Conjugate Gradient

$$\min \frac{1}{2} x^T Q x + b^T x + c$$

choose  $x^1$ ,  $p_1 = -\nabla f(x^1)$

if  $\nabla f(x^k) = 0$ , stop;

otherwise  $x^{k+1} = x^k + a_k p_k$

$$a_k = -\frac{p_k^T \nabla f(x^k)}{p_k^T Q p_k}$$

$$p_{k+1} = -\nabla f(x^{k+1}) + \lambda_k p_k$$

$$\lambda_k = \frac{p_k^T Q \nabla f(x^k)}{p_k^T Q p_k}$$

## BFGS

$$d_k = H_k \nabla f(x^k)$$

$$x^{k+1} = x^k - a_k d_k$$

$$f(x) = f(x^{k+1}) + \nabla f(x^{k+1})(x - x^{k+1}) + \frac{1}{2}(x - x^{k+1})^T \nabla^2 f(x^{k+1})(x - x^{k+1})$$

$$\nabla f(x) - \nabla f(x^{k+1}) = \nabla^2 f(x^{k+1})(x - x^{k+1})$$

$$\nabla f(x^{k+1}) - \nabla f(x^k) = \nabla^2 f(x^{k+1})(x^{k+1} - x^k)$$

$$g_{k+1} - g_k = \nabla^2 f(x^{k+1})(x^{k+1} - x^k)$$

$$s_k = x^{k+1} - x^k$$

$$y_k = g_{k+1} - g_k$$

$$H_{k+1} = (\nabla^2 f(x^{k+1}))^{-1} = B_{k+1}^{-1}$$

$$B_{k+1} = B_k + E_k$$

$$E_k = a_k U_k U_k^T + b_k V_k V_k^T$$

$$y_k = (B_k + a_k U_k U_k^T + b_k V_k V_k^T) s_k$$

$$y_k - B_k s_k = a_k U_k U_k^T s_k + b_k V_k V_k^T s_k$$

$$U_k = y_k \Rightarrow a_k = \frac{1}{y_k^T s_k}$$

$$V_k = -B_k s_k \Rightarrow b_k = -\frac{1}{s_k^T B_k s_k}$$

$$B_{k+1} = B_k + \frac{y_k y_k^T}{y_k^T s_k} - \frac{B_k s_k s_k^T B_k}{s_k^T B_k s_k}$$

$$H_{k+1} = \left( I - \frac{s_k y_k^T}{y_k^T s_k} \right) H_k \left( I - \frac{y_k s_k^T}{y_k^T s_k} \right) + \frac{s_k s_k^T}{y_k^T s_k}$$

## L-BFGS

$$H_{k+1} = \left( I - \frac{s_k y_k^T}{y_k^T s_k} \right) H_k \left( I - \frac{y_k s_k^T}{y_k^T s_k} \right) + \frac{s_k s_k^T}{y_k^T s_k}$$

$$\rho_k = \frac{1}{y_k^T s_k}, V_k = I - \rho_k y_k s_k^T$$

$$H_{k+1} = V_k^T H_k V_k + \rho_k s_k s_k^T$$

$$H_k = V_{k-1}^T H_{k-1} V_{k-1} + \rho_{k-1} s_{k-1} s_{k-1}^T$$

⋮

$$H_1 = V_0^T H_0 V_0 + \rho_0 s_0 s_0^T$$

$$H_k = (V_{k-1}^T \cdots V_0^T) H_0 (V_0 \cdots V_{k-1}) + \\ (V_{k-1}^T \cdots V_1^T) \rho_0 s_0 s_0^T (V_1 \cdots V_{k-1}) + \\ (V_{k-1}^T \cdots V_2^T) \rho_1 s_1 s_1^T (V_2 \cdots V_{k-1}) + \\ \dots$$

$$V_{k-1}^T \rho_{k-2} s_{k-2} s_{k-2}^T V_{k-1} +$$

$$\rho_{k-1} s_{k-1} s_{k-1}^T$$