

NN BP Algorithm

Model: $f_{t+1}(x) = g(W_t f_t(x) + b_t)$

when $t=0$ $f_1(x) = g(W_0 x + b_0)$

when $t=T-1$ $f_T(x) = g(W_{T-1} f_{T-1}(x) + b_{T-1})$

Vector-form BP

$\Phi(\theta) = \text{Loss}(f_T(x), y)$

our interest is: $\nabla_{W_t} \Phi(\theta^k)$ for $t=0, 1, \dots, T-1$

① $\nabla_{W_{T-1}} \Phi(\theta^k) = \nabla_{W_{T-1}} f_T \nabla_{f_T} \Phi(\theta^k)$

② $\nabla_{W_{T-2}} \Phi(\theta^k) = \nabla_{W_{T-2}} f_{T-1} \nabla_{f_{T-1}} f_T \nabla_{f_T} \Phi(\theta^k)$

$P_{T-1} = \nabla_{f_{T-1}} f_T P_T$ this term will be very complicated

where $f_T = g(W_{T-1} f_{T-1} + b_{T-1})$

③ we have $P_T \rightarrow P_{T-1} \rightarrow \dots \rightarrow P_k = \nabla_{f_{T-1}} u \cdot \nabla_u f_T$
 $= W_{T-1}^T \cdot g'(W_{T-1} f_{T-1} + b_{T-1})$

a) then $P_{k+1} = \nabla_{f_k} f_{k+1} P_k$

b) $\nabla_{W_{k-2}} \Phi(\theta) = \nabla_{W_{k-2}} f_{k-1} \nabla_{f_{k-1}} \Phi(\theta)$

$= \nabla_{W_{k-2}} f_{k-1} \cdot P_{k-1}$

where $f_{k-1} = g(W_{k-2} f_{k-2} + b_{k-2})$

c) then we have:

$P_T \rightarrow \dots \rightarrow P_k \rightarrow P_{k-1}$
 $\downarrow \quad \downarrow$
 $\nabla_{W_{k-1}} \Phi \quad \nabla_{W_{k-2}} \Phi$

In vector-BP algo

① $\frac{\partial \Phi}{\partial W_{ij}^{T-1}} = \frac{\partial \Phi}{\partial f_i^{T-1}} \cdot \frac{\partial f_i^{T-1}}{\partial h_i^{T-1}} \cdot \frac{\partial h_i^{T-1}}{\partial W_{ij}^{T-1}}$
 $\delta_i^{T-1} \quad f_j^{T-1}(x)$

② $\frac{\partial \Phi}{\partial W_{ij}^{T-2}} = \sum_s \left[\frac{\partial \Phi}{\partial h_s^{T-1}} \cdot \frac{\partial h_s^{T-1}}{\partial f_i^{T-1}} \cdot \frac{\partial f_i^{T-1}}{\partial h_i^{T-1}} \cdot \frac{\partial h_i^{T-1}}{\partial W_{ij}^{T-2}} \right]$
 $= \left(\sum_s \delta_s^{T-1} \cdot W_{si}^{T-1} \right) \cdot g'(h_i^{T-1}) \cdot f_j^{T-2}(x)$
 $\delta_i^{T-1} \quad f_j^{T-2}(x)$

$\frac{\partial \Phi}{\partial f_i^{T-1}} = \frac{\partial \Phi}{\partial f_i^{T-1}} \cdot \frac{\partial f_i^{T-1}}{\partial h_i^{T-1}}$

Scalar-form BP

$$W_k \Rightarrow W_{d_{k+1}, d_k}^k \Rightarrow f_i^{k+1}(x) = g\left(\sum_{j=1}^{d_k} w_{i,j}^k f_j^k(x) + b^k\right)$$

$$k = 0, 1, 2, \dots, T-1$$

$$i = 1, 2, \dots, d_{k+1}$$

when $k=0$ $f_i^1(x) = g\left(\sum_{j=1}^{d_0} w_{i,j}^0 \underbrace{f_j^0(x)}_{x_j} + b^0\right)$

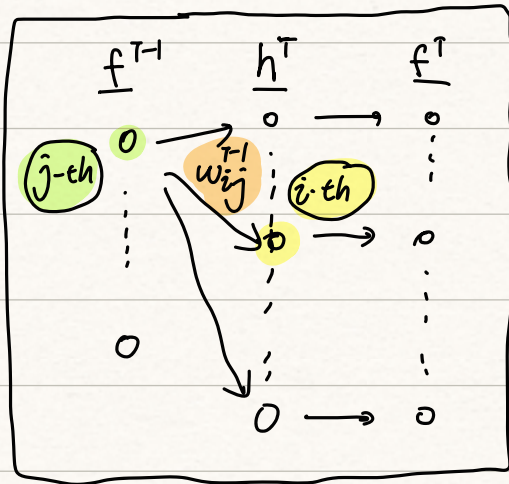
when $k=T-1$ $f_i^T(x) = g\left(\sum_{j=1}^{d_{T-1}} w_{i,j}^{T-1} f_j^{T-1}(x) + b^{T-1}\right)$

$$\Phi(\theta) = \text{Loss}(f^T(x), y)$$

$$\textcircled{1} \frac{\partial \Phi}{\partial w_{ij}^{T-1}} = \frac{\partial \Phi}{\partial f_i^T} \cdot \frac{\partial f_i^T}{\partial h_i^T} \cdot \frac{\partial h_i^T}{\partial w_{ij}^{T-1}}$$

$$= \underbrace{\frac{\partial \Phi}{\partial f_i^T}}_{\left(\frac{\partial \Phi}{\partial h_i^T}\right)} \cdot g'(h_i^T) \cdot f_j^{T-1}(x)$$

$$h_i^T := \sum_{j=1}^{d_{T-1}} w_{i,j}^{T-1} f_j^{T-1}(x) + b^{T-1}$$



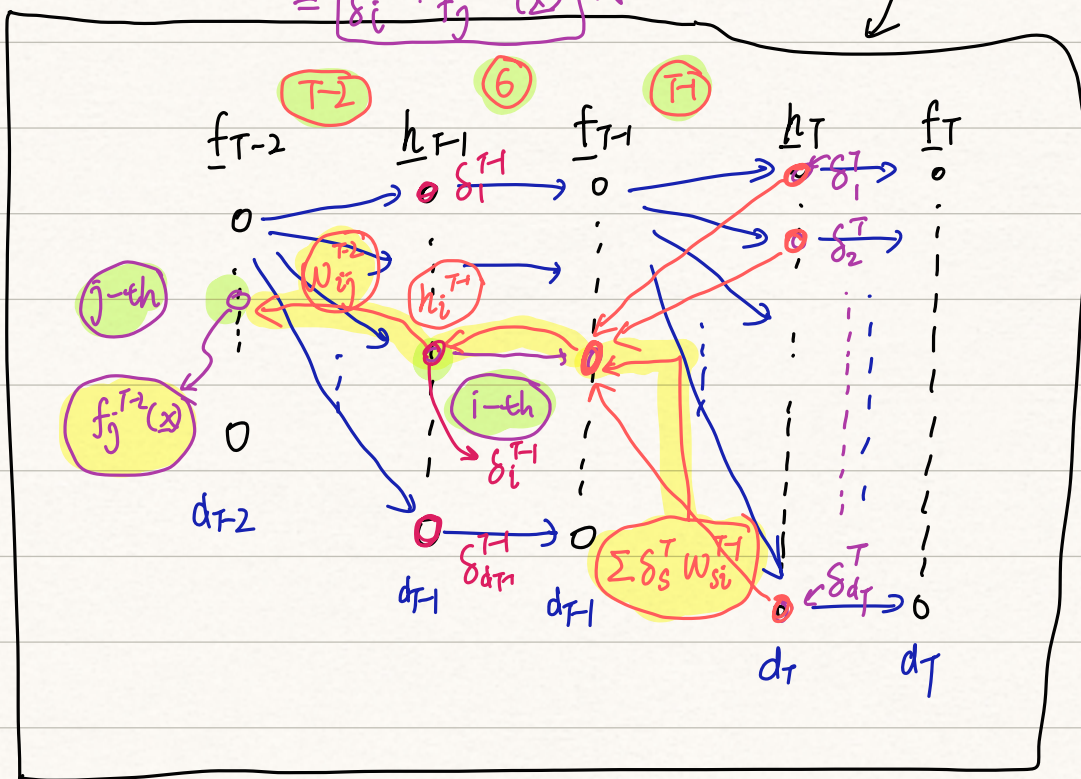
$$\begin{aligned} \textcircled{2} \frac{\partial \Phi}{\partial w_{ij}^{T-2}} &= \sum_{s=1}^{d_T} \frac{\partial \Phi}{\partial h_s^T} \cdot \frac{\partial h_s^T}{\partial w_{ij}^{T-2}} \\ &= \sum_{s=1}^{d_T} \frac{\partial \Phi}{\partial h_s^T} \cdot \frac{\partial h_s^T}{\partial f_i^{T-1}} \cdot \frac{\partial f_i^{T-1}}{\partial w_{ij}^{T-2}} \\ &= \left(\sum_{s=1}^{d_T} \frac{\partial \Phi}{\partial h_s^T} \cdot \frac{\partial h_s^T}{\partial f_i^{T-1}} \right) g'(h_i^{T-1}) \cdot f_j^{T-2}(x) \end{aligned}$$

$$f_i^{T-1}(x) = g\left[\sum_{j=1}^{d_{T-2}} w_{ij}^{T-2} f_j^{T-2}(x) + b^{T-2}\right]$$

$$h_s^T = \sum_{j=1}^{d_{T-1}} w_{sj}^{T-1} \cdot f_j^{T-1}(x) + b^{T-1}$$

$$= \left(\sum_{s=1}^{d_T} \frac{\partial \Phi}{\partial h_s^T} \cdot w_{si}^{T-1} \right) b'(h_i^{T-1}) f_j^{T-2}(x)$$

$$= \delta_i^{T-1} \cdot f_j^{T-2}(x) \star$$



Algo

Initialization

$$\Rightarrow \delta_j^T := \frac{\partial \Phi}{\partial f_j^T} \cdot b'(h_j^T)$$

↓

Update of error $\delta_i^k = \left[\sum_{s=1}^{d_{k+1}} \delta_s^{k+1} \cdot w_{si}^k \right] \star b'(h_i^k)$

Update of gradient $\Rightarrow \frac{\partial \Phi}{\partial w_{ij}^e} = \delta_i^{t+1} \cdot f_j^e(x)$