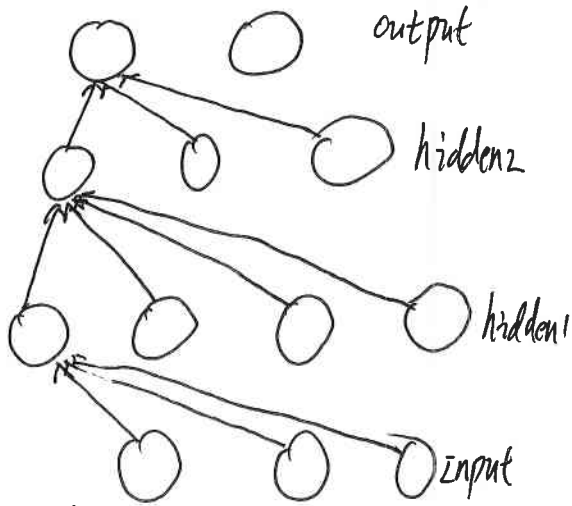


# 多层感知机



前向传播

$$a_h = \sum_{i=1}^I w_{hi} x_i \quad \text{输入} \rightarrow \text{隐层}$$

$$b_h = \theta_h(a_h)$$

隐层  $\rightarrow$  隐层 (L: 隐层数量) 反向传播

$$a_h = \sum_{h \in H_{L-1}} w_{hh} b_h$$

$$b_h = \theta_h(a_h)$$

隐层  $\rightarrow$  输出层

$$a_k = \sum_{h \in H_L} w_{hk} b_h$$

对于多分类:

$$P(c_k | x) = y_k = \frac{e^{a_k}}{\sum_{k=1}^K e^{a_k}}$$

$$P(z | x) = \prod_{k=1}^K y_k^{z_k}$$

$$z_1 = (1, 0, 0, 0) \quad \text{4类}$$

$$z_2 = (0, 1, 0, 0)$$

$$\text{对于一个样本 } \sum_{k=1}^K z_k = 1$$

$$\mathcal{L}(x, z) = - \sum_{k=1}^K z_k \ln y_k$$

$$\text{其实: } \mathcal{L}(s) = - \ln \prod_{(x, z) \in S} P(z | x)$$

$$= - \sum_{(x, z) \in S} (\ln P(z | x))$$

$$\text{令 } \mathcal{L}(x, z) = - \ln P(z | x)$$

$$\mathcal{L}(s) = \sum_{(x, z) \in S} \mathcal{L}(x, z)$$

$$\frac{\partial \mathcal{L}(s)}{\partial w} = \sum_{(x, z) \in S} \frac{\partial \mathcal{L}(x, z)}{\partial w}$$

$$\text{输出层: } \frac{\partial \mathcal{L}(x, z)}{\partial a_k} = \sum_{k=1}^K \frac{\partial \mathcal{L}(x, z)}{\partial y_k} \frac{\partial y_k}{\partial a_k}$$

$$\frac{\partial \mathcal{L}(x, z)}{\partial y_k} = \frac{\partial (- \sum_{i=1}^K z_i \ln y_i)}{\partial y_k} = - \frac{z_k}{y_k}$$

$$\frac{\partial y_k}{\partial a_k} = \frac{\partial (\frac{e^{a_k}}{\sum_{i=1}^K e^{a_i}})}{\partial a_k}$$

$$= y_k \delta_{kk} - y_k y_{k'} \quad (\text{if } k=k' \delta_{kk}=1 \text{ else } \delta_{kk}=0)$$

$$\text{因此: } \frac{\partial \mathcal{L}(x, z)}{\partial a_k} = - \sum_{k=1}^K \frac{z_k}{y_k} (y_k \delta_{kk} - y_k y_{k'})$$

$$= y_k - z_k \quad (\sum_{i=1}^K z_i = 1)$$

$$\text{这一层的梯度: } \delta w_{hk} = \frac{\partial \mathcal{L}(x, z)}{\partial a_k} \frac{\partial a_k}{\partial w_{hk}} = b_h (y_k - z_k)$$

定义  $\delta_j \stackrel{\text{def}}{=} \frac{\partial L(x, z)}{\partial a_j}$   $H_i$ : 第  $i$  层神经元数  
 $h_i$ : 第  $i$  层第  $h$  个神经元  
 第  $i$  层前隐层的梯度.  $i$ : 第  $i$  个输入, 即  $x_i$

$$\delta_{h2} = \frac{\partial L(x, z)}{\partial a_{h2}} = \sum_{k=1}^K \frac{\partial L(x, z)}{\partial a_k} \frac{\partial a_k}{\partial b_{h2}} \frac{\partial b_{h2}}{\partial a_{h2}}$$

$$= \theta'(a_{h2}) \sum_{k=1}^K \delta_k w_{h2k}$$

~~$$\Delta w_{h2k} = \frac{\partial L(x, z)}{\partial w_{h2k}} = \frac{\partial L(x, z)}{\partial a_{h2}} \frac{\partial a_{h2}}{\partial w_{h2k}}$$~~
~~$$= b_{h2} \delta_{h2}$$~~

$$\Delta w_{h2k} = \frac{\partial L(x, z)}{\partial w_{h2k}} = \frac{\partial L(x, z)}{\partial a_{h2}} \frac{\partial a_{h2}}{\partial w_{h2k}}$$

$$= \delta_{h2} b_{h1}$$

第  $i$  层隐层的梯度.

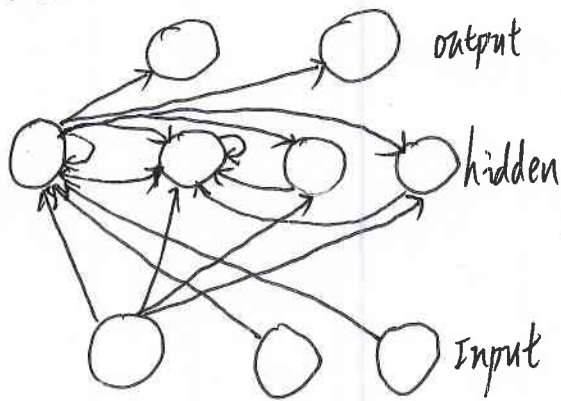
$$\delta_{h1} = \frac{\partial L(x, z)}{\partial a_{h1}} = \sum_{h2=1}^{H_2} \frac{\partial L(x, z)}{\partial a_{h2}} \frac{\partial a_{h2}}{\partial b_{h1}} \frac{\partial b_{h1}}{\partial a_{h1}}$$

$$= \theta'(a_{h1}) \sum_{h2=1}^{H_2} \delta_{h2} w_{h1h2}$$

$$w_{ih1} = \frac{\partial L(x, z)}{\partial w_{ih1}} = \frac{\partial L(x, z)}{\partial a_{h1}} \frac{\partial a_{h1}}{\partial w_{ih1}}$$

$$= \delta_{h1} x_i$$

RNN



反向传播:

输出层:

$$\delta_k^t = \frac{\partial L(x^t, z^t)}{\partial a_k^t} = y_k^t - z_k^t$$

$$\Delta w_{hk} = \frac{\partial L}{\partial w_{hk}} = \sum_{t=1}^T \frac{\partial L(x^t, z^t)}{\partial a_k^t} \frac{\partial a_k^t}{\partial w_{hk}} = \sum_{t=1}^T \delta_k^t b_h^t$$

前向传播 对于长度为T的序列X

H: 隐层神经元的个数, 在此为4个

K: 为类别个数

$$a_h^t = \sum_{i=1}^I w_{ih} x_i^t + \sum_{h'=1}^H w_{hh'} b_{h'}^t$$

$$b_h^t = \theta_h(a_h^t)$$

$$a_k^t = \sum_{h=1}^H w_{hk} b_h^t$$

$$y_k^t = \frac{e^{a_k^t}}{\sum_{k=1}^K e^{a_k^t}}$$

$$L(x^t, z^t) = - \sum_{k=1}^K z_k^t \ln y_k^t$$

损失函数  $L = \sum_{t=1}^T L(x^t, z^t)$

隐层:

$$\delta_h^t = \frac{\partial L}{\partial a_h^t} = \sum_{t=1}^T \frac{\partial L(x^t, z^t)}{\partial a_h^t} = \frac{\partial L(x^t, z^t)}{\partial a_h^t} + \frac{\partial L(x^{t+1}, z^{t+1})}{\partial a_h^t}$$

$$= \sum_{k=1}^K \frac{\partial L(x^t, z^t)}{\partial a_k^t} \frac{\partial a_k^t}{\partial b_h^t} \frac{\partial b_h^t}{\partial a_h^t} + \sum_{h'=1}^H \frac{\partial L(x^{t+1}, z^{t+1})}{\partial a_{h'}^{t+1}} \frac{\partial a_{h'}^{t+1}}{\partial b_h^t} \frac{\partial b_h^t}{\partial a_h^t}$$

$$= \theta'(a_h^t) \sum_{k=1}^K \delta_k^t w_{hk} + \theta'(a_h^t) \sum_{h'=1}^H \delta_{h'}^{t+1} w_{hh'}$$

$$= \theta'(a_h^t) \left( \sum_{k=1}^K \delta_k^t w_{hk} + \sum_{h'=1}^H \delta_{h'}^{t+1} w_{hh'} \right)$$

$$\Delta w_{ih} = \frac{\partial L}{\partial w_{ih}} = \sum_{t=1}^T \frac{\partial L(x^t, z^t)}{\partial a_h^t} \frac{\partial a_h^t}{\partial w_{ih}} = \sum_{t=1}^T \delta_h^t x_i^t$$

## 前向传播

Input gates

$$a_i^t = \sum_{i=1}^I w_{ii} x_i^t + \sum_{h=1}^H w_{hi} b_h^t + \sum_{c=1}^C w_{ci} s_c^t$$

$$b_i^t = f(a_i^t)$$

forget gates

$$a_{\phi}^t = \sum_{i=1}^I w_{i\phi} x_i^t + \sum_{h=1}^H w_{h\phi} b_h^t + \sum_{c=1}^C w_{c\phi} s_c^t$$

$$b_{\phi}^t = f(a_{\phi}^t)$$

cells

$$a_c^t = \sum_{i=1}^I w_{ic} x_i^t + \sum_{h=1}^H w_{hc} b_h^t$$

$$s_c^t = b_{\phi}^t s_c^{t-1} + b_c^t g(a_c^t)$$

output gates

$$a_w^t = \sum_{i=1}^I w_{iw} x_i^t + \sum_{h=1}^H w_{hw} b_h^t + \sum_{c=1}^C w_{cw} s_c^t$$

$$b_w^t = f(a_w^t)$$

cell outputs

$$b_c^t = b_w^t h(s_c^t)$$

反向传播

$$\epsilon_c^t = \frac{\partial \mathcal{L}}{\partial b_c^t} \quad \epsilon_s^t = \frac{\partial \mathcal{L}}{\partial s_c^t}$$

$$y_k^t = \frac{e^{a_k^t}}{\sum_{k=1}^K e^{a_k^t}}$$

$$\delta_j^t = \frac{\partial \mathcal{L}}{\partial a_j^t}$$

$$\mathcal{L} = \frac{1}{2} \sum_{t=1}^T \mathcal{L}(x^t, z^t)$$

隐层参数:

$$\mathcal{L} = - \sum_{k=1}^K z_k^t \ln y_k^t$$

$$a_k^t = \sum_{h=1}^H \sum_{c=1}^C b_c^t w_{ck}^h$$

$$\frac{\partial \mathcal{L}(x^t, z^t)}{\partial a_k^t} = y_k^t - z_k^t$$

$$\Delta w_{ck} = \frac{\partial \mathcal{L}}{\partial w_{ck}} = \frac{1}{2} \sum_{t=1}^T \frac{\partial \mathcal{L}(x^t, z^t)}{\partial w_{ck}}$$

$$= \frac{1}{2} \sum_{t=1}^T \frac{\partial \mathcal{L}(x^t, z^t)}{\partial a_k^t} \frac{\partial a_k^t}{\partial w_{ck}}$$

$$= \frac{1}{2} \sum_{t=1}^T \delta_k^t b_c^t$$

$$\begin{aligned} \epsilon_c^t = \frac{\partial \mathcal{L}}{\partial b_c^t} &= \sum_{k=1}^K \frac{\partial \mathcal{L}}{\partial a_k^t} \frac{\partial a_k^t}{\partial b_c^t} + \sum_{\omega=1}^H \frac{\partial \mathcal{L}}{\partial a_\omega^{t+1}} \frac{\partial a_\omega^{t+1}}{\partial b_c^t} \\ &+ \sum_{\phi=1}^H \frac{\partial \mathcal{L}}{\partial a_\phi^{t+1}} \frac{\partial a_\phi^{t+1}}{\partial b_c^t} + \sum_{l=1}^H \frac{\partial \mathcal{L}}{\partial a_l^{t+1}} \frac{\partial a_l^{t+1}}{\partial b_c^t} \\ &+ \sum_{m=1}^{HXC} \frac{\partial \mathcal{L}}{\partial a_m^{t+1}} \frac{\partial a_m^{t+1}}{\partial b_c^t} \\ &= \sum_{k=1}^K \delta_k^t w_{ck} + \sum_{\omega=1}^H \delta_\omega^{t+1} w_{c\omega} + \sum_{\phi=1}^H \delta_\phi^{t+1} w_{c\phi} \\ &+ \sum_{l=1}^H \delta_l^{t+1} w_{cl} + \sum_{m=1}^{HXC} \delta_m^{t+1} w_{cm} \\ &= \sum_{k=1}^K \delta_k^t w_{ck} + \sum_g \delta_g^{t+1} w_{cg} \end{aligned}$$

Output gates.

$$\delta_\omega^t = \sum_{c=1}^C \frac{\partial \mathcal{L}}{\partial b_c^t} \frac{\partial b_c^t}{\partial a_\omega^t} = f'(a_\omega^t) \sum_{c=1}^C \epsilon_c^t h(s_c^t)$$

$$\Delta w_{iw} = \frac{\partial \mathcal{L}}{\partial w_{iw}} = \frac{1}{2} \sum_{t=1}^T \frac{\partial \mathcal{L}}{\partial a_i^t} \frac{\partial a_i^t}{\partial w_{iw}} = \frac{1}{2} \sum_{t=1}^T \delta_i^t x_i$$

$$\Delta w_{iw} = \frac{\partial \mathcal{L}}{\partial w_{iw}} = \frac{1}{2} \sum_{t=1}^T \frac{\partial \mathcal{L}}{\partial a_i^t} \frac{\partial a_i^t}{\partial w_{iw}} = \frac{1}{2} \sum_{t=1}^T \delta_i^t x_i$$

$$\Delta w_{hw} = \frac{\partial \mathcal{L}}{\partial w_{hw}} = \frac{1}{2} \sum_{t=1}^T \frac{\partial \mathcal{L}}{\partial a_h^t} \frac{\partial a_h^t}{\partial w_{hw}} = \frac{1}{2} \sum_{t=1}^T \delta_h^t h_h$$

$$\Delta w_{cw} = \frac{1}{2} \sum_{t=1}^T \delta_\omega^t s_c^t$$

$$\epsilon_s^t = \frac{\partial L}{\partial s_c^t} \quad (\text{states})$$

$$\begin{aligned} \epsilon_s^t = \frac{\partial L}{\partial s_c^t} &= \frac{\partial L}{\partial b_0^t} \frac{\partial b_0^t}{\partial s_c^t} + \frac{\partial L}{\partial a_{w0}^t} \frac{\partial a_{w0}^t}{\partial s_c^t} + \frac{\partial L}{\partial a_{\phi}^{t+1}} \frac{\partial a_{\phi}^{t+1}}{\partial s_c^t} + \frac{\partial L}{\partial a_{w1}^{t+1}} \frac{\partial a_{w1}^{t+1}}{\partial s_c^t} + \frac{\partial L}{\partial a_{\phi}^{t+1}} \frac{\partial a_{\phi}^{t+1}}{\partial s_c^t} \\ &= b_{w0}^t \epsilon_c^t h'(s_c^t) + \delta_{w0}^t w_{w0} + \delta_{\phi}^{t+1} w_{\phi} + \delta_{w1}^{t+1} w_{w1} + b_{\phi}^{t+1} \epsilon_s^{t+1} \end{aligned}$$

cells:

$$\delta_c^t = \frac{\partial L}{\partial a_c^t} = \frac{\partial L}{\partial s_c^t} \frac{\partial s_c^t}{\partial a_c^t} = b_c^t \epsilon_s^t g'(a_c^t)$$

$$\Delta w_{ic} = \delta_c^t x_i \quad \Delta w_{hp} = \delta_c^t / b_h \quad \Delta w_{cp} = \delta_c^t s_c^{t+1}$$

$$\Delta w_{hc} = \delta_c^t / b_h$$

Forget Gates:

$$\begin{aligned} \delta_{\phi}^t &= \sum_{c=1}^C \frac{\partial L}{\partial s_c^t} \frac{\partial s_c^t}{\partial b_{\phi}^t} \frac{\partial b_{\phi}^t}{\partial a_{\phi}^t} \\ &= \sum_{c=1}^C \epsilon_c^t s_c^{t+1} f'(a_{\phi}^t) \end{aligned}$$

Input gates:

$$\begin{aligned} \delta_i^t &= \frac{\partial L}{\partial a_i^t} = \sum_{c=1}^C \frac{\partial L}{\partial s_c^t} \frac{\partial s_c^t}{\partial b_i^t} \frac{\partial b_i^t}{\partial a_i^t} \\ &= f'(a_i^t) \sum_{c=1}^C \epsilon_c^t g(a_i^t) \end{aligned}$$

$$\Delta w_{il} = \delta_i^t x_i^t$$

$$\Delta w_{hl} = \delta_l^t / b_h \quad \Delta w_{cl} = \delta_l^t s_c^{t+1}$$

$$\Delta w_{ip} = \delta_{\phi}^t x_i \quad \Delta w_{hp} = \delta_{\phi}^t / b_h \quad \Delta w_{cp} = \delta_{\phi}^t s_c^{t+1}$$