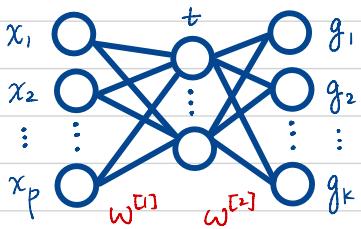


SPINN for Softmax-NN with gMP. (参考 Feng-Simon 的文章)



$$L = \prod_{i=1}^n \prod_{l=1}^k \left[\frac{\exp(g_l(x_i))}{\sum_{j=1}^k \exp(g_j(x_i))} \right]^{1\{y_i=l\}}$$

$$J = -\log L = -\sum_{i=1}^n \sum_{l=1}^k 1\{y_i=l\} \log \frac{\exp(g_l(x_i))}{\sum_{j=1}^k \exp(g_j(x_i))}$$

$$= -\sum_{i=1}^n \sum_{l=1}^k 1\{y_i=l\} g_l(x_i) + \sum_{i=1}^n \log \sum_{j=1}^k \exp(g_j(x_i))$$

Forward:

x_i	x^T	$p \times n$
\vec{h}_t	$= a(w^{[1]}x^T)$	$t \times n$
\vec{h}_t	$= a(w^{[1]}x^T)$	$t \times p$
\vec{g}_k	$= w^{[k]}\vec{h}_t$	$K \times n$
		$K \times t$
		$t \times n$

$$\nabla_{g_k} J = -\sum_{i=1}^n 1\{y_i=k\} + \sum_{i=1}^n \frac{\exp(g_k(x_i))}{\sum_{j=1}^k \exp(g_j(x_i))}$$

Backprop:

$$\frac{\partial J}{\partial w^{[1]}} = \sum_l \frac{\partial J}{\partial g_l} \cdot \frac{\partial g_l}{\partial w^{[1]}} = \sum_l \frac{\partial J}{\partial g_l} \cdot \begin{bmatrix} - & 0 & - \\ - & \frac{\partial g_l}{\partial w^{[1]}} & - \\ - & 0 & - \end{bmatrix} = \text{diag}\left(\frac{\partial J}{\partial g_l}\right) \cdot \mathbf{1}_k \cdot \vec{h}_t^T$$

$$\frac{\partial J}{\partial w^{[1]}} = \sum_l \frac{\partial J}{\partial g_l} \cdot \frac{\partial g_l}{\partial \vec{h}_t} \cdot \frac{\partial \vec{h}_t}{\partial (w^{[1]}x^T)} \cdot \frac{\partial (w^{[1]}x^T)}{\partial w^{[1]}}$$

$$\begin{aligned} \vec{h}_t^T &= \sum_l \frac{\partial J}{\partial g_l} \cdot w_l^{[1]T} * \begin{bmatrix} 1 \\ \vdots \\ 0 \end{bmatrix} \cdot x_i^T \\ &= \underbrace{\begin{bmatrix} | & | & | \\ w_1^{[1]T} & \dots & w_k^{[1]T} \\ | & | & | \end{bmatrix}}_{t \times k} \cdot \underbrace{\begin{bmatrix} \nabla g_1 \\ \vdots \\ \nabla g_k \end{bmatrix}}_{K \times 1} * \underbrace{\begin{bmatrix} 1 \\ \vdots \\ 0 \end{bmatrix}}_{w.r.t. w^{[1]}x^T} \cdot \underbrace{x_i^T}_{1 \times p} \end{aligned}$$

$$= \text{diag}\left(\begin{bmatrix} 1 \\ \vdots \\ 0 \end{bmatrix}\right) \cdot w^{[1]T} \cdot \nabla g_j \cdot x_i^T$$

Objective:

$$\min_w L(y, w) + \frac{\lambda_0}{2} \|w^{(0)}\|_2^2 + \sum_{j=1}^P \Pr\left(\sum_{t=1}^T \Pr(|w_{jt}^{(t)}|; \lambda); \lambda\right).$$

Solution:

$$(i) \text{ Smooth part: } L(y, w) + \frac{\lambda_0}{2} \|w^{(0)}\|_2^2$$

$$\text{Therefore, } \frac{\partial L}{\partial w^{(0)}} = \text{diag}\left(\frac{\partial \Pr}{\partial w^{(0)}}\right) \cdot \mathbf{1}_K \cdot \vec{u}_t + \lambda_0 \cdot w^{(0)}.$$

(ii) Thresholding part:

$$R = \frac{1}{2} \sum_{j=1}^P \|w_j^{(0)} - u_t\|_2^2 + \sum_{j=1}^P \Pr\left(\sum_{t=1}^T \Pr(|w_{jt}^{(t)}|; \lambda); \lambda\right).$$

$$\frac{\partial R}{\partial w_{jt}^{(0)}} = w_{jt}^{(0)} - u_t + \underbrace{\partial \Pr\left[\sum_{t=1}^T \Pr(|\tilde{w}_{jt}^{(t)}|; \lambda); \lambda\right] \cdot \partial \Pr(|\tilde{w}_{jt}^{(t)}|; \lambda) \cdot \partial(|w_{jt}^{(t)}|)}_{\alpha}.$$

Let $\frac{\partial R}{\partial w_{jt}^{(0)}} = 0$, we have

$$0 = \begin{cases} w_{jt}^{(0)} - u_t + \alpha, & \text{if } w_{jt}^{(0)} > 0, \\ -u_t + \alpha \cdot [-1, 1], & \text{if } w_{jt}^{(0)} = 0, \\ w_{jt}^{(0)} - u_t - \alpha, & \text{if } w_{jt}^{(0)} < 0. \end{cases}$$

so

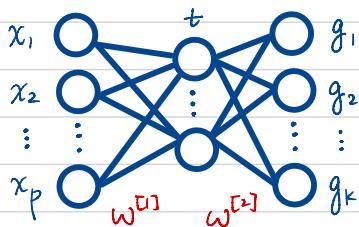
$$w_{jt}^{(0)} = \begin{cases} u_t - \alpha, & u_t > \alpha, \\ 0, & -\alpha \leq u_t \leq \alpha, \\ u_t + \alpha, & u_t < -\alpha. \end{cases}$$

which is $w_{jt}^{(0)} \leftarrow \text{sign}(u_t) \cdot (|u_t| - \alpha)_+$.

Part (i) is embedded in the back prop.

Part (ii) is the proximal operator.

SPINN + Integrative analysis



$$\text{Forward: } \begin{aligned} \vec{h}_t &= \alpha(w^{[1]}x^T) \\ &= \frac{tx \sum_m}{K} \cdot \frac{pm}{K} \cdot \frac{\sum_n}{m} \\ \vec{g} &= w^{[2]} \vec{h}_t \\ &= \frac{Kx}{m} \cdot \frac{Kt}{pm} \cdot \frac{\sum_n}{m} \end{aligned}$$

$$L = \prod_{i=1}^n \prod_{l=1}^k \left[\frac{\exp(g_l(x_i))}{\sum_{j=1}^k \exp(g_j(x_i))} \right]^{1\{y_i=l\}}$$

$$J = -\log L = -\sum_{i=1}^n \sum_{l=1}^k 1\{y_i=l\} \log \frac{\exp(g_l(x_i))}{\sum_{j=1}^k \exp(g_j(x_i))}$$

$$= -\sum_{i=1}^n \sum_{l=1}^k 1\{y_i=l\} g_l(x_i) + \sum_{i=1}^n \log \sum_{j=1}^k \exp(g_j(x_i))$$

$$\nabla_{g_k} J = -\sum_{i=1}^n 1\{y_i=k\} + \sum_{i=1}^n \frac{\exp(g_k(x_i))}{\sum_{j=1}^k \exp(g_j(x_i))}$$

Backprop:

$$\frac{\partial J}{\partial w^{[1]}} = \sum_l \frac{\partial J}{\partial g_l} \cdot \frac{\partial g_l}{\partial w^{[1]}} = \sum_l \frac{\partial J}{\partial g_l} \cdot \begin{bmatrix} 0 & \dots & 0 \\ -\frac{\partial g_l}{\partial w^{[1]}} & \dots & 0 \\ 0 & \dots & 0 \end{bmatrix} = \text{diag}\left(\frac{\partial J}{\partial g_l}\right) \cdot \mathbf{1}_k \cdot \vec{h}_t^T$$

$$\begin{aligned} \frac{\partial J}{\partial w^{[1]}} &= \sum_l \frac{\partial J}{\partial g_l} \cdot \frac{\partial g_l}{\partial \vec{h}_t} \cdot \frac{\partial \vec{h}_t}{\partial (w^{[1]}x^T)} \cdot \frac{\partial (w^{[1]}x^T)}{\partial w^{[1]}} \\ t \times pm &= \sum_l \frac{\partial J}{\partial g_l} \cdot w_l^{[2]T} * \begin{bmatrix} 1 \\ \vdots \\ 0 \end{bmatrix} \cdot x_i^T \\ &= \underbrace{\begin{bmatrix} 1 & \dots & 1 \\ w_1^{[2]T} & \dots & w_K^{[2]T} \\ \vdots & & \vdots \\ 1 & \dots & 1 \end{bmatrix}}_{t \times k} \cdot \underbrace{\begin{bmatrix} \nabla g_1 J \\ \vdots \\ \nabla g_K J \end{bmatrix}}_{K \times 1} * \underbrace{\begin{bmatrix} 1 \\ \vdots \\ 0 \end{bmatrix}}_{t \times 1} \cdot \underbrace{x_i^T}_{1 \times p} \\ &= \text{diag}\left(\begin{bmatrix} 1 \\ \vdots \\ 0 \end{bmatrix}\right) \cdot w^{[2]T} \cdot \nabla g_j \cdot x_i^T \cdot \underbrace{1}_{1 \times pm} \end{aligned}$$

Objective:

$$\min_{\omega} \mathcal{L}(y, \omega) + \frac{\lambda_0}{2} \|\omega^{[0]}\|_2^2 + \sum_{j=1}^P P_j \left[\sum_{m=1}^M \Omega_\alpha(\omega_{jm}^{[0]}; \lambda) \right]$$

$$\text{where } \Omega_\alpha(\theta) = (1-\alpha) \cdot \|\theta\|_1 + \alpha \cdot \|\theta\|_2.$$

Solution:

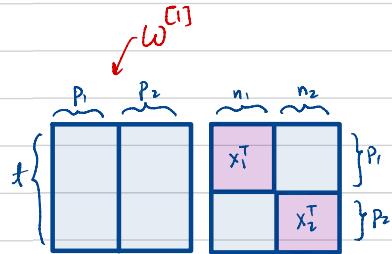
(i) Smooth part: $\mathcal{L}(y, \omega) + \frac{\lambda_0}{2} \|\omega^{[0]}\|_2^2$

Therefore, $\frac{\partial \mathcal{L}}{\partial \omega^{[0]}} = \text{diag}\left(\frac{\partial \mathcal{L}}{\partial \omega^{[0]}}\right) \cdot \mathbf{1}_k \cdot \tilde{\mathbf{h}} + \lambda_0 \cdot \omega^{[0]}$.

(ii) Thresholding part:

$$R = \frac{1}{2} \sum_{j=1}^P \left[\sum_{m=1}^M \|\omega_{jm}^{[0]} - \ell\|_2^2 + \sum_{m=1}^M \Omega_\alpha(\omega_{jm}^{[0]}; \lambda) \right]$$

$$\begin{aligned} \frac{\partial R}{\partial \omega_{jm}^{[0]}} &= \omega_{jm}^{[0]} - \ell + \partial P_j \left[\sum_{m=1}^M \Omega_\alpha(\tilde{\omega}_{jm}^{[0]}; \lambda) \right] \cdot \partial (\Omega_\alpha(\omega_{jm}^{[0]})) \\ &= \omega_{jm}^{[0]} - \ell + \text{const.} \cdot \left[(1-\alpha) \cdot \partial \|\omega_{jm}\|_1 + \alpha \cdot \partial \|\omega_{jm}\|_2 \right]. \end{aligned}$$

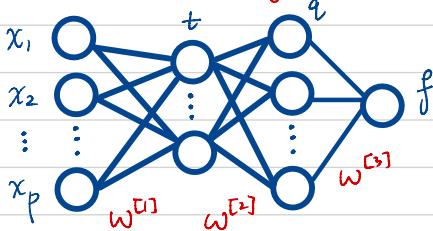


Therefore, $\omega_{jm}^{[0]} \leftarrow S(\omega_{jm}^{[0]}, \text{const.} \cdot (1-\alpha) \cdot \gamma_k)$,
 $\omega_{jm}^{[0]} \leftarrow \left(1 - \frac{\text{const.} \cdot \alpha \cdot \gamma_k}{\|\omega_{jm}\|_2} \right)_+ \omega_{jm}^{[0]}$ } learning rate for $j=1, \dots, P$.

Part (i) is embedded in the back prop.

Part (ii) is the proximal operator.

SPINN + Integrative analysis (Regression, MSE).



$$\mathcal{L} = \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\nabla_{f_i} \mathcal{L} = -2(y_i - f(x_i)) = 2f(x_i) - 2y_i$$

$$\nabla_f \mathcal{L} = (\nabla_{f_1} \mathcal{L}, \dots, \nabla_{f_n} \mathcal{L})^T \quad n \times 1$$

Forward: $x^T \underset{pm \times \sum_m}{\overbrace{\text{pm}}}$

 $\vec{h}^t = \alpha(w^{[1]} x^T)$

$\underset{tx \sum_m}{\overbrace{t \times pm}} \underset{pm \times \sum_m}{\overbrace{\text{pm}}} \underset{\sum_m}{\overbrace{t}}$

$\vec{h}^q = \alpha(w^{[2]} \vec{h}^t)$

$\underset{q \times \sum_m}{\overbrace{q \times t}} \underset{t \times \sum_m}{\overbrace{\sum_m}}$

$f = \alpha(w^{[3]} \vec{h}^q)$

$\underset{1 \times \sum_m}{\overbrace{1 \times q}} \underset{q \times \sum_m}{\overbrace{\sum_m}}$

Integrative

Single

Forward: $x^T \underset{pxn}{\overbrace{p \times n}}$

$\vec{h}^t = \alpha(w^{[1]} x^T)$

$\underset{txn}{\overbrace{t \times n}} \underset{txp}{\overbrace{t \times p}} \underset{pxn}{\overbrace{p \times n}}$

$\vec{h}^q = \alpha(w^{[2]} \vec{h}^t)$

$\underset{q \times n}{\overbrace{q \times t}} \underset{t \times n}{\overbrace{t \times n}}$

$f = w^{[3]} \vec{h}^q$

$\underset{1 \times n}{\overbrace{1 \times q}} \underset{q \times n}{\overbrace{q \times n}}$

Backprop:

$$\frac{\partial \mathcal{L}}{\partial w^{[3]}} = \sum_i \frac{\partial \mathcal{L}}{\partial f_i} \cdot \frac{\partial f_i}{\partial (w^{[3]} \vec{h}^q)} \cdot \frac{\partial (w^{[3]} \vec{h}^q)}{\partial w^{[3]}} = \frac{\partial \mathcal{L}}{\partial f_i} \cdot [0/1] \cdot \vec{h}^q$$

$$\frac{\partial \mathcal{L}}{\partial w^{[3]}} = \sum_i \frac{\partial \mathcal{L}}{\partial f_i} \cdot \frac{\partial f_i}{\partial (w^{[3]} \vec{h}^q)} \cdot \frac{\partial (w^{[3]} \vec{h}^q)}{\partial \vec{h}^q} \cdot \frac{\partial \vec{h}^q}{\partial (w^{[2]} \vec{h}^t)} \cdot \frac{\partial (w^{[2]} \vec{h}^t)}{\partial w^{[2]}}$$

$$q \times t = \sum_i \frac{\partial \mathcal{L}}{\partial f_i} \cdot [0/1] \cdot w^{[3]T} \cdot \begin{bmatrix} 0 \\ \vdots \\ i \end{bmatrix} \cdot \vec{h}^t = \text{diag}(\begin{bmatrix} 0 \\ \vdots \\ i \end{bmatrix}) \cdot w^{[3]T} \cdot \frac{\partial \mathcal{L}}{\partial f_i} \cdot [0/1] \cdot \vec{h}^t$$

w.r.t. $w^{[3]} \vec{h}^q$ w.r.t. $w^{[2]} \vec{h}^t$

$$\frac{\partial \mathcal{L}}{\partial w^{[1]}} = \sum_i \frac{\partial \mathcal{L}}{\partial f_i} \cdot \frac{\partial f_i}{\partial (w^{[3]} \vec{h}^q)} \cdot \frac{\partial (w^{[3]} \vec{h}^q)}{\partial \vec{h}^q} \cdot \frac{\partial \vec{h}^q}{\partial (w^{[2]} \vec{h}^t)} \cdot \frac{\partial (w^{[2]} \vec{h}^t)}{\partial \vec{h}^t} \cdot \frac{\partial \vec{h}^t}{\partial (w^{[1]} x_i)} \cdot \frac{\partial (w^{[1]} x_i)}{\partial w^{[1]}}$$

$$t \times p = \sum_i \frac{\partial \mathcal{L}}{\partial f_i} \cdot [0/1] \cdot w^{[3]T} \cdot \begin{bmatrix} 0 \\ \vdots \\ i \end{bmatrix} \cdot w^{[2]T} \cdot \begin{bmatrix} 0 \\ \vdots \\ i \end{bmatrix} \cdot x_i^T$$

w.r.t. $w^{[3]} \vec{h}^q$ w.r.t. $w^{[2]} \vec{h}^t$ w.r.t. $w^{[1]} x_i$

$$= \sum_i \text{diag}(\begin{bmatrix} 0 \\ \vdots \\ i \end{bmatrix}) \cdot w^{[3]T} \cdot \text{diag}(\begin{bmatrix} 0 \\ \vdots \\ i \end{bmatrix}) \cdot w^{[2]T} \cdot \frac{\partial \mathcal{L}}{\partial f_i} \cdot [0/1] \cdot x_i^T$$