From The Neural Network Side, since there is one hidden layer.

$$f_k(x) = \beta_0 + \beta_k^T \theta(\alpha_{on} + \alpha_n^T x)$$

From the PPR side we know:

$$f(x) = \sum_{m=1}^{N} g_n(w_n^T x)$$

Also we know:

$$g_n(w_m^T x) = \beta_n \theta(\alpha_{on} + S_n(w_n^T x))$$

Then since is a PPR. we know that $w = \dfrac{\hat{w}}{\|\hat{w}\|}$,

Then, we need to find some $S_n$ and $\hat{w}$, such will allow us to equiparate both weights; this $S = \|\alpha_n\|$ and

$$w_n^T = \frac{\alpha_n}{\|\alpha_n\|}$$

Then $g_n(w^T x) = F_n(x) = \beta_0 + \beta_n \theta(\alpha_{on} + \alpha_n^T x)$ and

$$f(x) = \sum_{k=1}^{K} f_m(x).$$

The same argument could be done if we use a logistic regression for the $g_n(w^T x)$ of PPR and a classification network from the neural network where $g(T) = \dfrac{e^{T_k}}{\sum_{l=1}^{\cdot} e^{T_l}}$.

IF the function is linear we know that $f'' = 0$.

Then if $y = \frac{1}{1+e^{-v}}$

$$\frac{\partial y}{\partial v} = \frac{e^{-v}}{(1+e^{-v})^2}$$

$$\frac{\partial^2 y}{\partial v^2} = e^{-v}(1+e^{-v})^2\left(2\cdot(1+e^{-v})^{-1} - 1\right)$$

Then if $v = 0$.

$$\left.\frac{\partial^2 y}{\partial v^2}\right|_{v=0} = \frac{1}{4}\cdot\left(\frac{2}{2} - 1\right)$$

$$= 0 \;//\; \Rightarrow \text{The we could argue that is linear at } v=0.$$

# ESL 11.4

$$R(\theta) = -\sum_{i=1}^{N} \sum_{k=1}^{K} y_{ik} \log f_k(x_i)$$

$$f_k(x) = \frac{e^{T_k}}{\sum_{l=1}^{K} e^{T_l}}$$

Since there is no hidden layer

$$T_k = \beta_{0k} + \beta_k^T x.$$

And now we assume that the $k=1$ is base category. Then

$$T_k - T_1 = (\beta_{0k} - \beta_{01}) + (\beta_k - \beta_1)^T x$$

$$\therefore \quad Pr(k \cdot k) = f_k(x) = \frac{e^{T_k - T_1}}{1 + e^{T_k - T_1}}.$$

Wich transform for multiple $k$ and using cross Entropy.

$$-\sum y_{ik} \log(f_k(x))$$

$$\Rightarrow \text{wich is the multilogistic objective function.}$$

# 5.1 Bishop

First we need to found the conversion or relation between $\tanh(a)$ and $\theta(a)$.

$$\frac{e^a - e^{-a}}{e^a + e^{-a}} = \frac{e^a}{e^a + e^{-a}} * - \frac{e^{-a}}{e^a + e^{-a}} = \frac{e^a}{e^a + e^{-a}} - \frac{e^{-a} + e^a - e^a}{e^a + e^{-a}}$$

$$= \frac{e^a}{e^a + e^{-a}} - \left(1 - \frac{e^a}{e^a + e^{-a}}\right)$$

$$= \frac{2e^a}{e^a + e^{-a}} \cdot \frac{e^{-a}}{e^{-a}} - 1$$

$$= \frac{2}{1 + e^{-2a}} - 1$$

$$\frac{e^a - e^{-a}}{e^a + e^{-a}} = 2\,\theta(2a) - 1$$

Then if we replace this into

$$a_k = \sum_{j=1}^{n} w_{kj}^{(k)} \cdot (2\theta(2a) - 1) + w_{k_0}^{(k)}$$

$$a_k = \sum_{j=1}^{k} 2w^{(k)}\theta(2a) \quad 2w_{kj}^{k} + w_{k_0}^{(k)}$$

$$a_k = \sum_{j} 2w_{kj}^{(k)}\theta(2a) + \sum_{j}\left[\cdot 2w_{kj}^{k} + w_{k_0}^{k}\right]$$

Then if for a regular NEURAL network we will have to convert into :

$$\omega_p^{(2p)} = 2\omega_k^{(2)}, \quad a_p^{(p)} = 2a_j^{(2)}, \quad \omega_{ko}^{(2p)} = -\sum_{j=1}^{n} \omega_{k_j}^{(2)} + \omega_{ko}^2$$

Using this equality, we transform $a_k$ into:

$$a_p = \sum_{j=1}^{n} \omega_p^{(2p)} \theta(a) + \omega_{po}^{(2p)}$$

# Bishop 5.5

$$P(T/W_1, \ldots, W_K) = \prod_{K=1}^{K} y_K^{T_K} \Rightarrow \text{Network output for } K \text{ class}$$
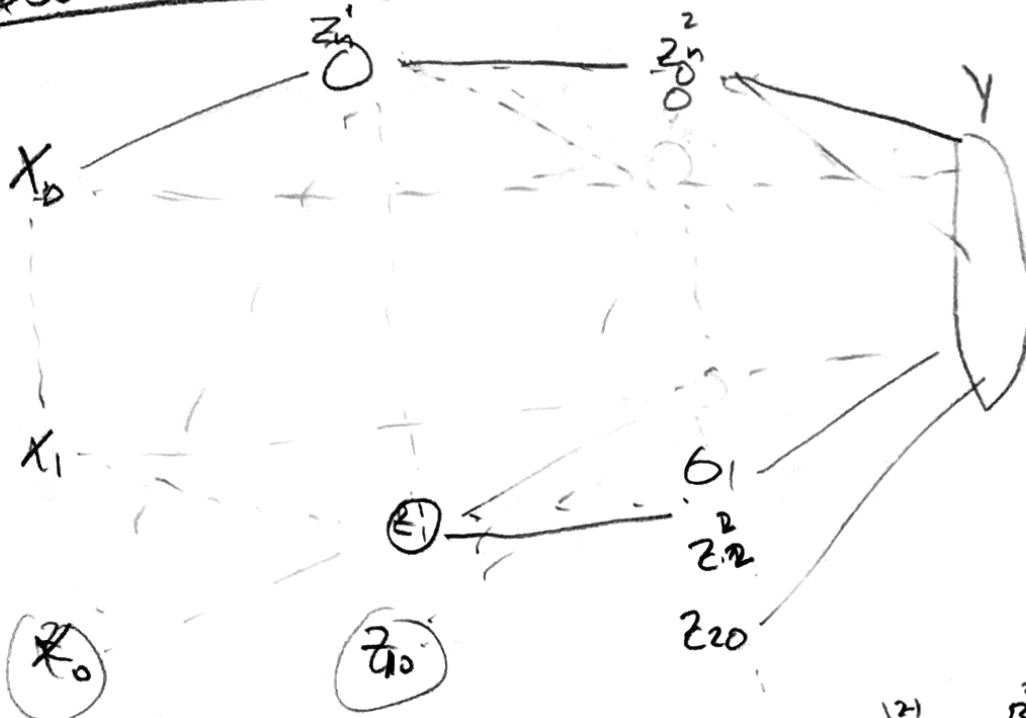
Then for $N$ data point we have that

$$P(T/W_1, \ldots, W_K) = \prod_{K=1}^{K} \prod_{n=1}^{N} y_{nk}^{T_{nk}}$$

Then if we compare for the log function.

$$\ell(T/W_1, \ldots, W_K) = - \sum_{K=1}^{K} \sum_{n=1}^{N} T_{nk} \cdot \log(y_{nk})$$

Which is the log-likelihood function of the 5.24.

# Question 6



**Equations:**

$$Z_n^1 = \Theta(\alpha_{0n} + \alpha_n^T X)$$

$$\underbrace{}_{(1)^T P}$$

$$Z_n^2 = \Theta(\beta_{0K}^{(2)} + \beta^{\pi^{(2)}} Z_n^1)$$

$$\hat{Y}_K = \Theta(\beta_{0K}^1 + \beta_K Z^2 + \delta_K^T X + \delta_{0K})$$

$$R = \frac{1}{n} \sum_{i=1}^{k} (\hat{Y}_K - y)^2$$

$$\frac{\partial R}{\partial \beta_K} = \delta_{Ki} Z_{ni}, \qquad S_{ni} = \Theta'(\alpha_n^T X_i) \sum_{K=1}^{k} \beta_{Kn} \delta_{Ki}$$

$S_{Ki} \Rightarrow$ are errors in output layer

$$\frac{\partial R}{\partial \alpha_n} = S_{ni} X_{il}$$

$$\frac{\partial R}{\partial \delta_K} = 2(\hat{Y}_n - y) \cdot \hat{Y}_K (1 - \hat{y}_n) \cdot X$$

\* The main change is that we have a another set of parameters that go from the input layer to the output layer, creating a new partial derivate.