Notations: $\odot$ 表示 Hadamard product，$\otimes$ 表示 kronecker product，$Sgn(x)=\begin{cases} 1 & x\geq 0 \\ 0 & else \end{cases}$，$f(\Xi)=\begin{bmatrix} f(x_{1,1}) & \cdots & f(x_{1,n}) \\ \vdots & & \vdots \\ f(x_{m,1}) & \cdots & f(x_{m,n}) \end{bmatrix}\in\mathbb{R}^{m\times n}$

1. (1) 由 $y_i=BN_{\gamma,\beta}(\hat{x_i})=\gamma\hat{x_i}+\beta$ 可知

$\dfrac{\partial y_i}{\partial\gamma}=\hat{x_i}$，$\dfrac{\partial y_i}{\partial\beta}=1$，其中 $\hat{x_i}$ 的计算即按照 BN 的处理流程，而：$\begin{cases} \mu_B=\frac{1}{m}\sum\limits_{i=1}^{m}x_i \\ \sigma_B^2=\frac{1}{m}\sum\limits_{i=1}^{m}(x_i-\mu_B)^2 \\ \hat{x_i}=\frac{x_i-\mu_B}{\sqrt{\sigma_B^2+\varepsilon}} \end{cases}$

(2) 考虑 softmax 函数 $f:\mathbb{R}^N\to\mathbb{R}^C$（可微），对于 $x\in\mathbb{R}^N$，$y=f(x)\in\mathbb{R}^C$，其中

$y_k=P(类别k|x,W)=\dfrac{\exp(W_k^Tx)}{\sum\limits_{n=1}^{C}\exp(W_n^Tx)}$，其中 $W_i\in\mathbb{R}^N$（$1\leq i\leq C$ 表示权重）

考虑输出 $y$ 对输入 $x$ 的偏导：$\dfrac{\partial y_k}{\partial x}=\begin{bmatrix} \frac{\partial y_k}{\partial x_1} \\ \vdots \\ \frac{\partial y_k}{\partial x_N} \end{bmatrix}$

$\dfrac{\partial y_k}{\partial x_i}=\dfrac{1}{(\sum\limits_{n=1}^{C}\exp(W_n^Tx))^2}\left(\left(\sum\limits_{n=1}^{C}\exp(W_n^Tx)\right)\exp(W_k^Tx)W_{k,i}-\exp(W_k^Tx)\left(\sum\limits_{n=1}^{C}W_{n,i}\exp(W_n^Tx)\right)\right)$

$=W_{k,i}y_k-y_k\cdot\dfrac{\sum\limits_{n=1}^{C}W_{n,i}\exp(W_n^Tx)}{\sum\limits_{n=1}^{C}\exp(W_n^Tx)}=y_k\left(W_{k,i}-\dfrac{\sum\limits_{n=1}^{C}W_{n,i}\exp(W_n^Tx)}{\sum\limits_{n=1}^{C}\exp(W_n^Tx)}\right)$

$\therefore$ 输出 $y$ 相对于输入 $x$ 的 Jacobi 矩阵 $\left[\dfrac{\partial y}{\partial x^T}\right]_{k,i}=\dfrac{\partial y_k}{\partial x_i}$ 如上所示

(3) 对于前向传播过程，首先计算单个样本对应的预测例 $\hat{y}_A^i\in\mathbb{R}^a$ 与 $\hat{y}_B^i\in\mathbb{R}^b$

令 $FC_{1A}$，$DP$，$FC_{1B}$，$BN$ 的输出分别为 $x_{1A}^i\in\mathbb{R}^s$，$x_{DP}^i\in\mathbb{R}^s$，$x_{1B}^i\in\mathbb{R}^t$，$x_{BN}^i\in\mathbb{R}^t$，

则有 $x_{1A}^i=\sin(\theta_{1A}x^i+b_{1A})$，

$\begin{cases} x_{DP}^i=M\odot x_{1A}^i & \text{对应 Task A 的前向传播结果} \Rightarrow \hat{y}_A^i=\theta_{2A}(M\odot\sin(\theta_{1A}x^i+b_{1A}))+b_{2A} \\ \hat{y}_A^i=\theta_{2A}x_{DP}^i+b_{2A} \end{cases}$

对于 Task B，则有 $x_{1B}^i=\theta_{1B}\Sigma^i$

在计算 BN 输出时，$\mu=\frac{1}{m}\sum\limits_{i=1}^{m}x_{1B}^i$

放 $x_{BN}^i=ReLu(x_{1B}^i-\mu+b_{1B})$.

$=ReLu(\theta_{1B}\Sigma^i-\frac{1}{m}\theta_{1B}\sum\limits_{i=1}^{m}\Sigma^i+b_{1B})$

$\therefore$ $\hat{y}_B^i=softmax(\theta_{2B}(\hat{y}_A^i+x_{BN}^i)+b_{2B})$

更进一步，$\hat{y}_B^i=softmax(\theta_{2B}(\hat{y}_A^i+x_{BN}^i)+b_{2B})=\dfrac{\exp(\theta_{2B}(\hat{y}_A^i+x_{BN}^i)+b_{2B})}{1_b^T\exp(\theta_{2B}(\hat{y}_A^i+x_{BN}^i)+b_{2B})}$

其中 $\hat{y}_{B,k}^i=\dfrac{\exp(\theta_{2B,k}(\hat{y}_A^i+x_{BN}^i)+b_{2B})}{\sum\limits_{p=1}^{b}\exp(\theta_{2B,p}(\hat{y}_A^i+x_{BN}^i)+b_{2B})}$ 其中 $\theta_{2B,k}$ 表示矩阵 $\theta_{2B}$ 的第 $k$ 行

(4) 将损失函数 $\mathcal{L}$ 写成向量的形式：$\mathcal{L}=\frac{1}{m}\sum\limits_{i=1}^{m}\left[\frac{1}{2}\|\hat{y}_A^i-y_A^i\|_2-y_B^{i\,T}\log\hat{y}_B^i\right]$

$\therefore$ 可得 $d\mathcal{L}=\frac{1}{m}\sum\limits_{i=1}^{m}\left(\frac{\partial\mathcal{L}}{\partial\hat{y}_A^i}\right)^Td\hat{y}_A^i+\frac{1}{m}\sum\limits_{i=1}^{m}\left(\frac{\partial\mathcal{L}}{\partial\hat{y}_B^i}\right)^Td\hat{y}_B^i$

$=\frac{1}{m}\sum\limits_{i=1}^{m}(\hat{y}_A^i-y_A^i)^Td\hat{y}_A^i+\frac{1}{m}\sum\limits_{i=1}^{m}\left(y_B^{i\,T}(\frac{1}{\hat{y}_B^i}\odot d\hat{y}_B^i)\right)$

$=\frac{1}{m}\sum\limits_{i=1}^{m}(\hat{y}_A^i-y_A^i)^Td\hat{y}_A^i+\frac{1}{m}\sum\limits_{i=1}^{m}\left(y_B^{i\,T}\odot\frac{1}{\hat{y}_B^i})^Td\hat{y}_B^i\right)$

下面讲述将 $d\hat{y_A^i}$ 与 $d\hat{y_B^i}$ 使用 $d\theta_{1A}, d\theta_{2A}, d\theta_{1B}, d\theta_{2B}$ 表示

$$d\hat{y_B^i} = \hat{y_B^i} \odot \left(d\theta_{2B}(\hat{y_A^i} + \Sigma_{BN}^i) + db_{2B}\right) - \hat{y_B^i} \cdot \hat{y_B^i}^T \left(d\theta_{2B}(\hat{y_A^i} + \Sigma_{BN}^i) + db_{2B}\right)$$

$$\exists\, dL = \frac{1}{m}\sum_{i=1}^{m}(\hat{y_A^i} - y_A^i)^T d\hat{y_A^i} + \frac{1}{m}\sum_{i=1}^{m}(\hat{y_B^i}^T \odot \hat{y_B^i}^T)\left(\hat{y_B^i} \odot (d\theta_{2B}(\hat{y_A^i} + \Sigma_{BN}^i) + db_{2B}) - \hat{y_B^i}\hat{y_B^i}^T(d\theta_{2B}(\hat{y_A^i} + \Sigma_{BN}^i) + db_{2B})\right)$$

$$= \frac{1}{m}\sum_{i=1}^{m}(\hat{y_A^i} - y_A^i)^T d\hat{y_A^i} + \frac{1}{m}\sum_{i=1}^{m}\left(y_B^i{}^T d\theta_{2B}(\hat{y_A^i} + \Sigma_{BN}^i) + db_{2B}) - \hat{y_B^i}^T(d\theta_{2B}(\hat{y_A^i} + \Sigma_{BN}^i) + db_{2B})\right)$$

由于 $\hat{y_A^i}$ 与 $\theta_{2B}$ 无关，故 $dL = \frac{1}{m}\sum_{i=1}^{m}(\hat{y_A^i} - y_A^i)^T d\hat{y_A^i} + tr\left(\frac{1}{m}\sum_{i=1}^{m}(\hat{y_A^i} + \Sigma_{BN}^i)(y_B^i - \hat{y_B^i})^T\right)d\theta_{2B}\right) + tr\left(\frac{1}{m}\sum_{i=1}^{m}(y_B^i - \hat{y_B^i})^T db_{2B}\right)$

故 $\dfrac{\partial L}{\partial \theta_{2B}} = \dfrac{1}{m}\sum_{i=1}^{m}(y_B^i - \hat{y_B^i})(\hat{y_A^i} + \Sigma_{BN}^i)^T$ , $\dfrac{\partial L}{\partial b_{2B}} = \dfrac{1}{m}\sum_{i=1}^{m}(y_B^i - \hat{y_B^i})$

这后，求解 $\dfrac{\partial L}{\partial \theta_{1B}}$ ，先进一步将 $d\hat{y_B^i}$ 表示为 $d\hat{y_A^i}$ 与 $d\Sigma_{BN}^i$ 的组合形式

$$d\hat{y_B^i} = d\left(\frac{\exp(\theta_{2B}(\hat{y_A^i} + \Sigma_{BN}^i) + b_{2B})}{\mathbb{1}_b^T \exp(\theta_{2B}(\hat{y_A^i} + \Sigma_{BN}^i) + b_{2B})}\right) = \hat{y_B^i} \odot (\theta_{2B}(d\hat{y_A^i} + d\Sigma_{BN}^i)) - \hat{y_B^i}\cdot\hat{y_B^i}^T \theta_{2B}(d\hat{y_A^i} + d\Sigma_{BN}^i)$$

$$= \hat{y_B^i} \odot (\theta_{2B}d\hat{y_A^i}) + \hat{y_B^i} \odot (\theta_{2B}d\Sigma_{BN}^i) - \hat{y_B^i}\hat{y_B^i}^T \theta_{2B}d\hat{y_A^i} - \hat{y_B^i}\hat{y_B^i}^T\theta_{2B}d\Sigma_{BN}^i$$

由于 $\hat{y_A^i}$ 不是 $\theta_{1B}$ 与 $b_{1B}$ 的函数，故只用表示 $d\Sigma_{BN}^i$

一步表示 $\Sigma_{BN}^i = ReLu\left(\Sigma_{1B}^i - \frac{1}{m}\sum_{s=1}^{m}\Sigma_{1B}^s + b_{1B}\right)$

$$= ReLu\left(\theta_{1B}\left[\Sigma^i - \frac{1}{m}\sum_{s=1}^{m}\Sigma^s\right] + b_{1B}\right)$$

故 $d\Sigma_{BN}^i = sgn\left(\theta_{1B}[\Sigma^i - \bar{\Sigma}] + b_{1B}\right) \odot (d\theta_{1B}[\Sigma^i - \bar{\Sigma}] + db_{1B})$   $\left(\bar{\Sigma} = \frac{1}{m}\sum_{s=1}^{m}\Sigma^s\right)$

代入可得: $dL = c + \frac{1}{m}\sum_{i=1}^{m}\left(y_B^i{}^T \odot \hat{y_B^i}^T\right)\left(\hat{y_B^i} \odot (\theta_{2B}d\Sigma_{BN}^i) - \hat{y_B^i}\hat{y_B^i}^T \theta_{2B}d\Sigma_{BN}^i\right)$

$$= c + \frac{1}{m}\sum_{i=1}^{m}(y_B^i - \hat{y_B^i})^T \theta_{2B}d\Sigma_{BN}^i$$

$$= c + \frac{1}{m}\sum_{i=1}^{m}(y_B^i - \hat{y_B^i})^T \theta_{2B}\left(sgn(\theta_{1B}[\Sigma^i - \bar{\Sigma}] + b_{1B}) \odot (d\theta_{1B}[\Sigma^i - \bar{\Sigma}] + db_{1B})\right)$$

$$= c + \frac{1}{m}\sum_{i=1}^{m}tr\left(\left[(y_B^i - \hat{y_B^i})^T \theta_{2B} \odot sgn^T(\theta_{1B}[\Sigma^i - \bar{\Sigma}] + b_{1B})\right](d\theta_{1B}[\Sigma^i - \bar{\Sigma}] + db_{1B})\right)$$

$$= c + \frac{1}{m}\sum_{i=1}^{m}tr\left([\Sigma^i - \bar{\Sigma}]\left[(y_B^i - \hat{y_B^i})^T \theta_{2B} \odot sgn^T(\theta_{1B}[\Sigma^i - \bar{\Sigma}] + b_{1B})\right]d\theta_{1B}\right)$$

$$+ \frac{1}{m}\sum_{i=1}^{m}tr\left([(y_B^i - \hat{y_B^i})^T \theta_{2B} \odot sgn^T(\theta_{1B}[\Sigma^i - \bar{\Sigma}] + b_{1B})]db_{1B}\right) \quad \text{其中 } c \text{ 为无关项.}$$

故 $\dfrac{\partial L}{\partial \theta_{1B}} = \dfrac{1}{m}\sum_{i=1}^{m}\left(\theta_{2B}^T(y_B^i - \hat{y_B^i}) \odot sgn(\theta_{1B}[\Sigma^i - \bar{\Sigma}] + b_{1B})\right)[\Sigma^i - \bar{\Sigma}]^T$

$\dfrac{\partial L}{\partial b_{1B}} = \dfrac{1}{m}\sum_{i=1}^{m}\left(\theta_{2B}^T(y_B^i - \hat{y_B^i}) \odot sgn(\theta_{1B}[\Sigma^i - \bar{\Sigma}] + b_{1B})\right)$

之后求解 $\frac{\partial L}{\partial \theta_{2A}}$, $\frac{\partial L}{\partial b_{2A}}$, $\frac{\partial L}{\partial \theta_{1A}}$, $\frac{\partial L}{\partial b_{1A}}$

由于 $\hat{y_A^{\flat}}$ 与 $\hat{y_B^{\flat}}$ 与 $\theta_{2A}$, $\theta_{2B}$, $b_{2A}$, $b_{2B}$ 均有关, 故需同时展开

$$d\hat{y_A^{\flat}} = d(\theta_{2A} \Sigma_{DP}^{\flat} + b_{2A}) = d\theta_{2A} \Sigma_{DP}^{\flat} + db_{2A}$$

$$d\hat{y_B^{\flat}} = \hat{y_B^{\flat}} \odot (\theta_{2B} d\hat{y_A^{\flat}}) + \hat{y_B^{\flat}} \odot (\theta_{2B} d\Sigma_{BN}^{\flat}) - \hat{y_B^{\flat}} \hat{y_B^{\flat T}} \theta_{2B} d\hat{y_A^{\flat}} - \hat{y_B^{\flat}} \hat{y_B^{\flat T}} \theta_{2B} d\Sigma_{BN}^{\flat}$$

首先求解 $\frac{\partial L}{\partial \hat{y_A^{\flat}}}$,

$$dL = \frac{1}{m} \sum_{i=1}^{m} (\hat{y_A^{\flat}} - y_A^{\flat})^T d\hat{y_A^{\flat}} + \frac{1}{m} \sum_{i=1}^{m} [(y_B^{\flat T} \odot \frac{1}{\hat{y_B^{\flat}}}) d\hat{y_B^{\flat}}],$$

$$= \frac{1}{m} \sum_{i=1}^{m} (\hat{y_A^{\flat}} - y_A^{\flat})^T d\hat{y_A^{\flat}} + \frac{1}{m} \sum_{i=1}^{m} [(y_B^{\flat T} \odot \frac{1}{\hat{y_B^{\flat}}})(\hat{y_B^{\flat}} \odot (\theta_{2B} d\hat{y_A^{\flat}}) - \hat{y_B^{\flat}} \hat{y_B^{\flat T}} \theta_{2B} d\hat{y_A^{\flat}}] + C \quad (C \text{ 并标与 } d\hat{y_A^{\flat}} \text{ 无关的项})$$

$$= \frac{1}{m} \sum_{i=1}^{m} (\hat{y_A^{\flat}} - y_A^{\flat})^T d\hat{y_A^{\flat}} + \frac{1}{m} \sum_{i=1}^{m} (y_B^{\flat} - \hat{y_B^{\flat}})^T \theta_{2B} d\hat{y_A^{\flat}} + C$$

$$= \frac{1}{m} \sum_{i=1}^{m} ((\hat{y_A^{\flat}} - y_A^{\flat}) + \theta_{2B}^T (y_B^{\flat} - \hat{y_B^{\flat}}))^T d\hat{y_A^{\flat}} + C$$

又 $d\hat{y_A^{\flat}} = d(\theta_{2A} \Sigma_{DP}^{\flat} + b_{2A}) = d\theta_{2A} \Sigma_{DP}^{\flat} + db_{2A}$

$$\Rightarrow dL = \frac{1}{m} \sum_{i=1}^{m} [((\hat{y_A^{\flat}} - y_A^{\flat}) + \theta_{2B}^T (y_B^{\flat} - \hat{y_B^{\flat}}))^T (d\theta_{2A} \Sigma_{DP}^{\flat} + db_{2A})$$

$$= tr(\frac{1}{m} \sum_{i=1}^{m} ((\hat{y_A^{\flat}} - y_A^{\flat}) + \theta_{2B}^T (y_B^{\flat} - \hat{y_B^{\flat}}) \Sigma_{DP}^{\flat})^T d\theta_{2A})$$

$$+ tr(\frac{1}{m} \sum_{i=1}^{m} ((\hat{y_A^{\flat}} - y_A^{\flat}) + \theta_{2B}^T (y_B^{\flat} - \hat{y_B^{\flat}}))^T db_{2A})$$

故 $\frac{\partial L}{\partial \theta_{2A}} = \frac{1}{m} \sum_{i=1}^{m} [(\hat{y_A^{\flat}} - y_A^{\flat}) + \theta_{2B}^T (y_B^{\flat} - \hat{y_B^{\flat}})] \cdot \Sigma_{DP}^{\flat T}$

$\quad \frac{\partial L}{\partial b_{2A}} = \frac{1}{m} \sum_{i=1}^{m} ((\hat{y_A^{\flat}} - y_A^{\flat}) + \theta_{2B}^T (y_B^{\flat} - \hat{y_B^{\flat}}))$

同理, 继续求解残差

$y_A^{\flat} = \theta_{2A} \Sigma_{DP}^{\flat} + b_{2A} = \theta_{2A} (M \odot \Sigma_{1A}^{\flat}) + b_{2A}$

故 $d\hat{y_A^{\flat}} = \theta_{2A} (M \odot d\Sigma_{1A}^{\flat})$

又 $d\Sigma_{1A}^{\flat} = d[\sin(\theta_{1A} \Sigma^{\flat} + b_{1A})] = \cos(\theta_{1A} \Sigma^{\flat} + b_{1A}) \odot (d\theta_{1A} \Sigma^{\flat} + db_{1A})$

故 $d\hat{y_A^{\flat}} = \theta_{2A} (M \odot (\cos(\theta_{1A} \Sigma^{\flat} + b_{1A}) \odot (d\theta_{1A} \Sigma^{\flat} + db_{1A}))$

$\quad = \theta_{2A} ((M \odot \cos(\theta_{1A} \Sigma^{\flat} + b_{1A})) \odot (d\theta_{1A} \Sigma^{\flat} + db_{1A}))$

代入得 $dL = \frac{1}{m} \sum_{i=1}^{m} [(\hat{y_A^{\flat}} - y_A^{\flat}) + \theta_{2B}^T (y_B^{\flat} - \hat{y_B^{\flat}})]^T \theta_{2A} ((M \odot \cos(\theta_{1A} \Sigma^{\flat} + b_{1A})) \odot (d\theta_{1A} \Sigma^{\flat} + db_{1A}))$

$\quad = tr(\frac{1}{m} \sum_{i=1}^{m} (M \odot \cos(\theta_{1A} \Sigma^{\flat} + b_{1A}))^T \odot ([(\hat{y_A^{\flat}} - y_A^{\flat}) + \theta_{2B}^T (y_B^{\flat} - \hat{y_B^{\flat}})]^T \theta_{2A}) (d\theta_{1A} \Sigma^{\flat} + db_{1A}))$

$\quad = tr(\frac{1}{m} \sum_{i=1}^{m} \Sigma_i^{\flat} (M \odot \cos(\theta_{1A} \Sigma^{\flat} + b_{1A}))^T \odot ([(\hat{y_A^{\flat}} - y_A^{\flat}) + \theta_{2B}^T (y_B^{\flat} - \hat{y_B^{\flat}})]^T \theta_{2A}) d\theta_{1A})$

$\quad + tr((M \odot \cos(\theta_{1A} \Sigma^{\flat} + b_{1A}))^T \odot ([(\hat{y_A^{\flat}} - y_A^{\flat}) + \theta_{2B}^T (y_B^{\flat} - \hat{y_B^{\flat}})]^T \theta_{2A}) db_{1A})$

$$\frac{\partial \mathcal{L}}{\partial \Theta_{1A}} = \frac{1}{m} \sum_{i=1}^{m} \left( \left( m \odot \cos(\Theta_{1A} \Sigma^i + b_{1A}) \right) \odot \left[ \Theta_{2A}^T (l \hat{y}_A^i - y_A^i) + \Theta_{2B}^T (y_B^i - \hat{y}_B^i) \right] \right) \Sigma^{i \, T}$$

$$\frac{\partial \mathcal{L}}{\partial b_A} = \frac{1}{m} \sum_{i=1}^{m} \left( \left( m \odot \cos(\Theta_{1A} \Sigma^i + b_{1A}) \right) \odot \left( \Theta_{2A}^T (l \hat{y}_A^i - y_A^i) + \Theta_{2B}^T (y_B^i - \hat{y}_B^i) \right) \right)$$