

# Report of Homework 1

Hou Lintao 2018310076

## 1. Part One: Code Implementation

### 1.1 Multilayer Perceptron (MLP)

The algorithms were completed in *mlp.py*. The hyperparameters of hidden dim, initial learning rate, batch size and epochs were set as 50, 0.1, 16 and 60. The learning rate  $lr$  will be 0.1 times the before value after 20 epochs each time. The loss curve was shown as Figure 1.1, and the test accuracy was 95.18% and exceeds 90%.

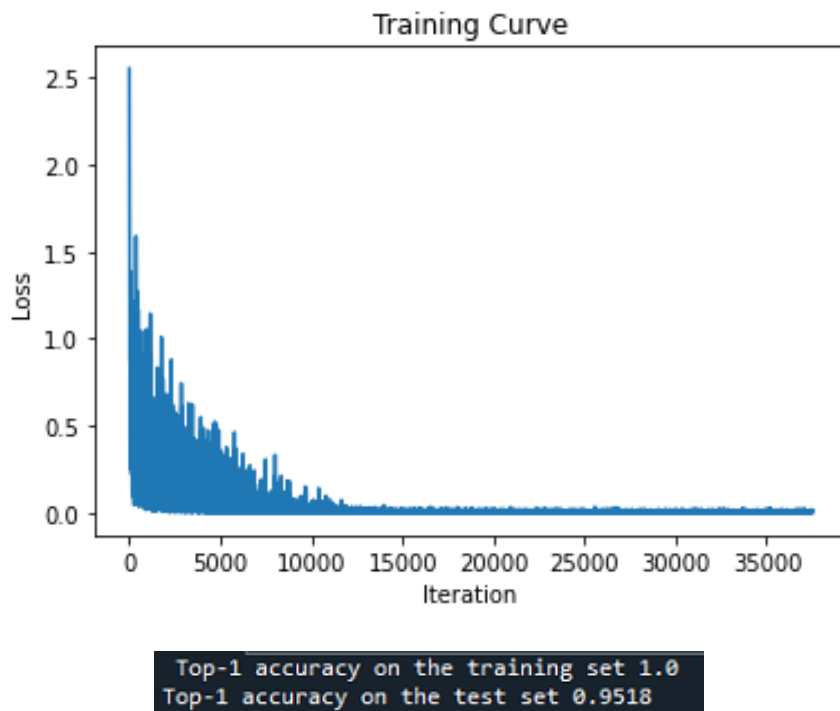


Figure 1.1 Loss curve and the test accuracy.

### 1.2 Variational Autoencoder (VAE)

The algorithms were completed in *vae.py*. and the training results were shown in file folder visualization. The sample results with Decoder Noise were shown in Figure 1.2, and the data distribution is almost the same as the given data distribution.

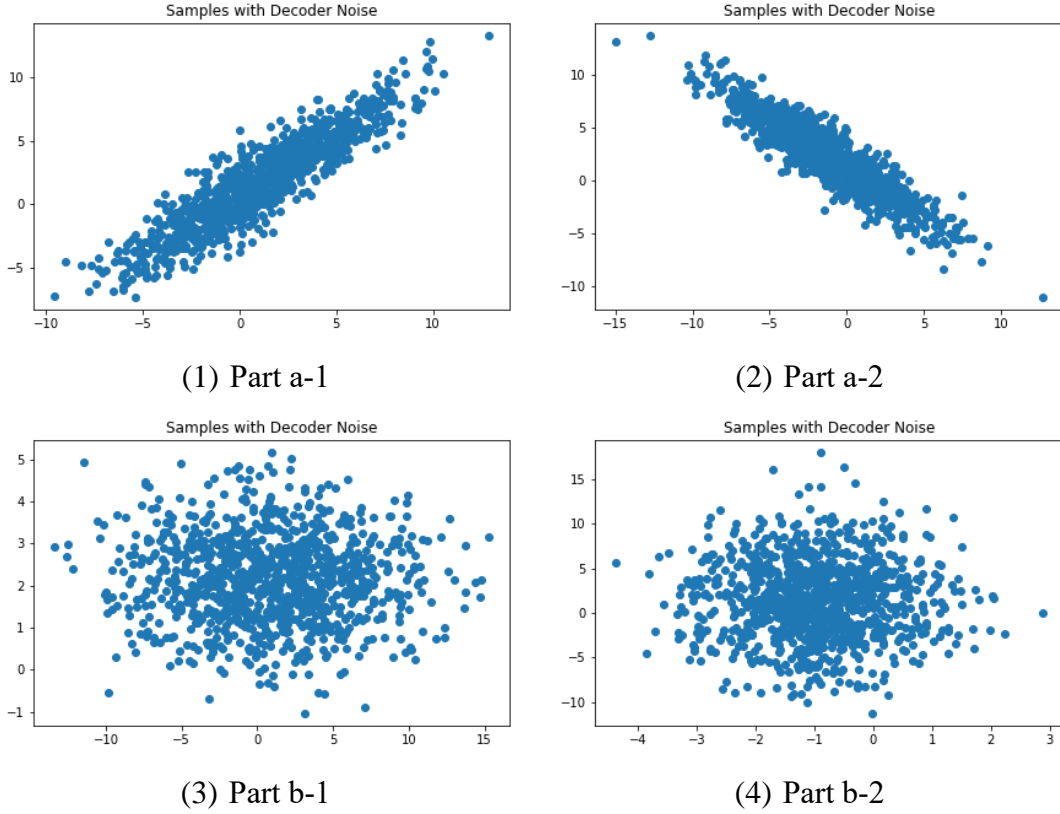


Figure 1.2 Samples with Decoder Noise

## 2. Part Two: Back-propagation

### 2.1 Block One: gradients of some basic layer

- (i) The gradients of the output  $y_i = \mathbf{BN}_{\gamma, \beta}(x_i)$  with respect to the parameters of  $\gamma, \beta$ :

$$\frac{\partial y_i}{\partial \gamma} = \frac{\partial(\gamma \hat{x}_i + \beta)}{\partial \gamma} = \hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$\frac{\partial y_i}{\partial \beta} = \frac{\partial(\gamma \hat{x}_i + \beta)}{\partial \beta} = 1$$

- (ii) The forward computations of the dropout layer is:

$$y_{DP,j} = x_j \cdot \mathbf{M}_j$$

$$\mathbf{M}_j = \begin{cases} 0, & r_j < p \\ 1/(1-p), & r_j \geq p \end{cases}$$

$$= \frac{1}{1-p} \text{sgn}(r_j - p)$$

The function  $\text{sgn}(x) = (x + |x|)/2$ .

Therefore, the gradients of the output of a dropout layer with respect to its input are:

$$\frac{\partial y_j}{\partial x_j} = M_j$$

(iii) The forward computations of the softmax function is:

$$y_{soft,i} = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} = 1 - \frac{\sum_{j=1, j \neq i}^n e^{x_j}}{\sum_{j=1}^n e^{x_j}}$$

Therefore, the gradients of the output of a softmax function with respect to its input are calculated as following.

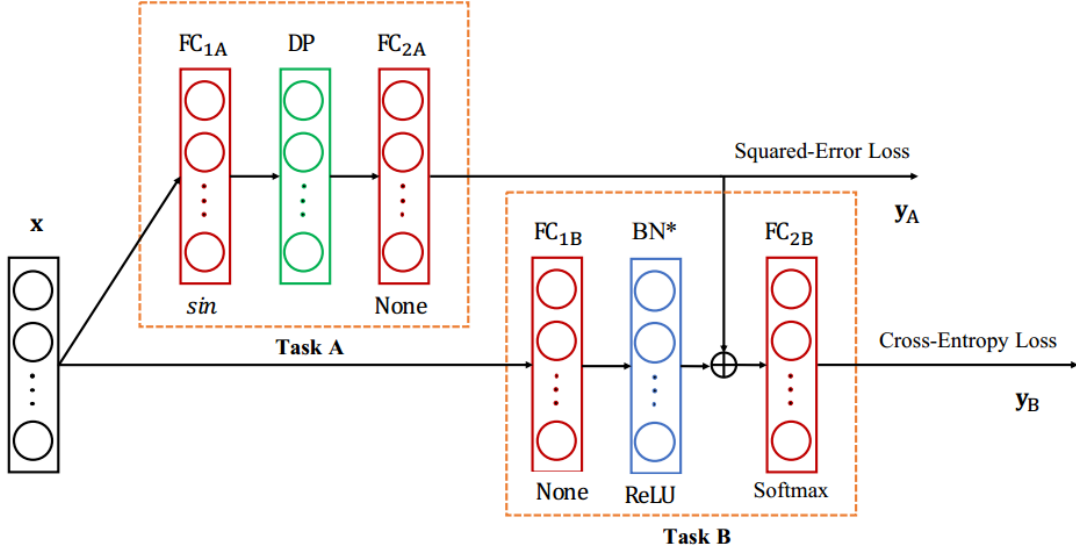
When  $i = k$  :

$$\begin{aligned} \frac{\partial y_{soft,i}}{\partial x_i} &= \left( 1 - \frac{\sum_{j=1, j \neq i}^n e^{x_j}}{\sum_{j=1}^n e^{x_j}} \right) \bigg/ \frac{\partial \sum_{j=1}^n e^{x_j}}{\partial x_i} = \frac{\sum_{j=1, j \neq i}^n e^{x_j}}{\left( \sum_{j=1}^n e^{x_j} \right)^2} \frac{\partial \sum_{j=1}^n e^{x_j}}{\partial x_i} \\ &= \frac{\sum_{j=1, j \neq i}^n e^{x_j}}{\sum_{j=1}^n e^{x_j}} \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} = (1 - y_{soft,i}) y_{soft,i} \end{aligned}$$

When  $i \neq k$ :

$$\begin{aligned} \frac{\partial y_{soft,i}}{\partial x_k} &= \left( \frac{\partial e^{x_i}}{\sum_{j=1}^n e^{x_j}} \right) \bigg/ \frac{\partial \sum_{j=1}^n e^{x_j}}{\partial x_k} = - \frac{e^{x_i}}{\left( \sum_{j=1}^n e^{x_j} \right)^2} \frac{\partial \sum_{j=1}^n e^{x_j}}{\partial x_k} \\ &= - \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} \frac{e^{x_k}}{\sum_{j=1}^n e^{x_j}} = -y_{soft,i} y_{soft,k} \end{aligned}$$

## 2.2 Block Two: feed-forward and backpropagation of the multi-task network



### (i) Task A:

The result of layer  $FC_{1A}$  is:

$$\mathbf{z}_{1A} = \theta_{1A} \mathbf{x} + \mathbf{b}_{1A}$$

$$\mathbf{x}_{1A} = \sin(\mathbf{z}_{1A})$$

The result of layer  $FC_{DP}$  is:

$$\mathbf{r} = \text{random}(\mathbf{x}_{1A}) \in U(0,1)$$

$$\mathbf{M} = \frac{1}{1-p} \text{sgn}(\mathbf{r} - p \cdot \mathbf{1})$$

$$\mathbf{x}_{DP} = \mathbf{x}_{1A} \odot \mathbf{M}$$

The vector  $\mathbf{r}$  is a random vector with uniform distribution between 0 and 1 for each neuron of  $\mathbf{x}_{1A}$ , and  $\mathbf{1}$  is the vector that all entries are one. The definition of mask vector  $\mathbf{M}$  can be seen in 2.1(ii).

The result of layer  $FC_{2A}$  is:

$$\hat{\mathbf{y}}_A = \mathbf{x}_{2A} = \theta_{2A} \mathbf{x}_{DP} + \mathbf{b}_{2A}$$

### Task B:

The result of layer  $FC_{1B}$  is:

$$\mathbf{x}_{1B} = \theta_{1B} \mathbf{x}$$

The result of layer  $BN^*$  is:

$$\mu = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_{1B}^i$$

$$\mathbf{x}_{BN} = \mathbf{x}_{1B} - \mu + \mathbf{b}_{1B}$$

$$\mathbf{x}_{BN^*} = \text{ReLU}(\mathbf{x}_{BN}) = \max(\mathbf{x}_{BN}, 0)$$

The result of layer  $\oplus$  is:

$$\mathbf{x}_{\oplus} = \mathbf{x}_{2A} + \mathbf{x}_{BN*}$$

The result of layer FC<sub>2B</sub> is:

$$\begin{aligned}\mathbf{z}_{2B} &= \theta_{2B} \mathbf{x}_{\oplus} + \mathbf{b} \\ \hat{\mathbf{y}}_B &= \mathbf{x}_{2B} = \text{Softmax}(\mathbf{z}_{2B})\end{aligned}$$

(ii) The loss function of Task A and Task B is:

$$L = \frac{1}{m} \sum_{i=1}^m \left[ \frac{1}{2} \|\hat{\mathbf{y}}_A^i - \mathbf{y}_A^i\|_2^2 - \sum_{k=1}^b y_{B,k}^i \log \hat{y}_{B,k}^i \right]$$

We split the loss function into two parts:

$$\begin{aligned}L &= L_A + L_B \\ L_A &= \frac{1}{m} \sum_{i=1}^m \frac{1}{2} \|\hat{\mathbf{y}}_A^i - \mathbf{y}_A^i\|_2^2 \\ L_B &= -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^b y_{B,k}^i \log \hat{y}_{B,k}^i\end{aligned}$$

For Task A and Task B, we can first calculate the residual of each layer.

**Task B:**

The residual of layer FC<sub>2B</sub>:

$$\begin{aligned}\frac{\partial L}{\partial \hat{y}_{B,j}^i} &= \frac{\partial L_B}{\partial \hat{y}_{B,j}^i} = -\frac{1}{m} \sum_{k=1}^b \frac{\partial (y_{B,k}^i \log \hat{y}_{B,k}^i)}{\partial \hat{y}_{B,j}^i} = -\frac{1}{m} \frac{y_{B,j}^i}{\hat{y}_{B,j}^i} \\ \frac{\partial L}{\partial z_{2B,j}^i} &= \sum_{k=1}^b \frac{\partial L}{\partial \hat{y}_{B,k}^i} \frac{\partial \hat{y}_{B,k}^i}{\partial z_{2B,j}^i} = \sum_{k=1, k \neq j}^b \frac{\partial L}{\partial \hat{y}_{B,k}^i} \frac{\partial \hat{y}_{B,k}^i}{\partial z_{2B,j}^i} + \frac{\partial L}{\partial \hat{y}_{B,j}^i} \frac{\partial \hat{y}_{B,j}^i}{\partial z_{2B,j}^i} \\ &= \frac{1}{m} \left[ \sum_{k=1, k \neq j}^b \left( \frac{y_{B,k}^i}{\hat{y}_{B,k}^i} \hat{y}_{B,k}^i \hat{y}_{B,j}^i \right) + \frac{y_{B,j}^i}{\hat{y}_{B,j}^i} (\hat{y}_{B,j}^i - 1) \hat{y}_{B,j}^i \right] \\ &= \frac{1}{m} \left[ \sum_{k=1, k \neq j}^b y_{B,k}^i \hat{y}_{B,j}^i + y_{B,j}^i (\hat{y}_{B,j}^i - 1) \right] \\ &= \frac{1}{m} \left( \hat{y}_{B,j}^i \sum_{k=1}^b y_{B,k}^i - y_{B,j}^i \right) = \frac{1}{m} (\hat{y}_{B,j}^i - y_{B,j}^i) \\ &\quad \left( \sum_{k=1}^b y_{B,k}^i = 1 \text{ for one-hot question} \right)\end{aligned}$$

Therefore the residual of  $\mathbf{z}_{2B}^i$  is:

$$\frac{\partial L}{\partial \mathbf{z}_{2B}^i} = \frac{1}{m} (\hat{\mathbf{y}}_B^i - \mathbf{y}_B^i)$$

The residual of layer  $\oplus$ :

$$\frac{\partial L}{\partial \mathbf{x}_{\oplus}^i} = \frac{\partial L}{\partial \mathbf{z}_{2B}^i} \frac{\partial \mathbf{z}_{2B}^i}{\partial \mathbf{x}_{\oplus}^i} = \frac{1}{m} \theta_{2B}^T (\hat{\mathbf{y}}_B^i - \mathbf{y}_B^i)$$

The residual of layer BN\*:

$$\frac{\partial L}{\partial \mathbf{x}_{BN*}^i} = \frac{\partial L}{\partial \mathbf{x}_{\oplus}^i} \frac{\partial \mathbf{x}_{\oplus}^i}{\partial \mathbf{x}_{BN*}^i} = \frac{1}{m} \theta_{2B}^T (\hat{\mathbf{y}}_B^i - \mathbf{y}_B^i)$$

$$\frac{\partial L}{\partial \mathbf{x}_{\text{BN}}^i} = \frac{\partial L}{\partial \mathbf{x}_{\text{BN}^*}^i} \frac{\partial \mathbf{x}_{\text{BN}^*}^i}{\partial \mathbf{x}_{\text{BN}}^i} = \frac{1}{m} \theta_{2\text{B}}^T (\hat{\mathbf{y}}_{\text{B}}^i - \mathbf{y}_{\text{B}}^i) \circ \text{sgn}(\mathbf{x}_{\text{BN}}^i)$$

The residual of layer FC<sub>1B</sub>:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{x}_{1\text{B}}^i} &= \sum_{j=1}^m \frac{\partial L}{\partial \mathbf{x}_{\text{BN}}^j} \frac{\partial \mathbf{x}_{\text{BN}}^j}{\partial \mathbf{x}_{1\text{B}}^i} = \sum_{j=1, j \neq i}^m \frac{\partial L}{\partial \mathbf{x}_{\text{BN}}^j} \frac{\partial (\mathbf{x}_{1\text{B}}^j - \mu + \mathbf{b}_{1\text{B}})}{\partial \mathbf{x}_{1\text{B}}^i} \\ &= \sum_{j=1}^m \frac{\partial L}{\partial \mathbf{x}_{\text{BN}}^j} \left( \frac{\partial \mathbf{x}_{1\text{B}}^j}{\partial \mathbf{x}_{1\text{B}}^i} - \frac{\partial \mu}{\partial \mathbf{x}_{1\text{B}}^i} + 0 \right) \\ &= \frac{\partial L}{\partial \mathbf{x}_{\text{BN}}^i} - \frac{1}{m} \sum_{j=1}^m \frac{\partial L}{\partial \mathbf{x}_{\text{BN}}^j} \\ &= \frac{1}{m} \theta_{2\text{B}}^T (\hat{\mathbf{y}}_{\text{B}}^i - \mathbf{y}_{\text{B}}^i) \circ \text{sgn}(\mathbf{x}_{\text{BN}}^i) - \frac{1}{m^2} \sum_{j=1}^m \left[ \theta_{2\text{B}}^T (\hat{\mathbf{y}}_{\text{B}}^j - \mathbf{y}_{\text{B}}^j) \circ \text{sgn}(\mathbf{x}_{\text{BN}}^j) \right] \end{aligned}$$

### Task A:

The residual of layer FC<sub>2A</sub>:

$$\begin{aligned} \frac{\partial L}{\partial \hat{\mathbf{y}}_{\text{A}}^i} &= \frac{\partial L_{\text{A}}}{\partial \hat{\mathbf{y}}_{\text{A}}^i} + \frac{\partial L_{\text{B}}}{\partial \hat{\mathbf{y}}_{\text{A}}^i} \\ &= \frac{1}{m} (\hat{\mathbf{y}}_{\text{A}}^i - \mathbf{y}_{\text{A}}^i) + \frac{\partial L_{\text{B}}}{\partial \mathbf{x}_{\oplus}^i} \frac{\partial \mathbf{x}_{\oplus}^i}{\partial \hat{\mathbf{y}}_{\text{A}}^i} \\ &= \frac{1}{m} \left[ (\hat{\mathbf{y}}_{\text{A}}^i - \mathbf{y}_{\text{A}}^i) + \theta_{2\text{B}}^T (\hat{\mathbf{y}}_{\text{B}}^i - \mathbf{y}_{\text{B}}^i) \right] \end{aligned}$$

The residual of layer DP:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{x}_{\text{DP}}^i} &= \frac{\partial L}{\partial \hat{\mathbf{y}}_{\text{A}}^i} \frac{\partial \hat{\mathbf{y}}_{\text{A}}^i}{\partial \mathbf{x}_{\text{DP}}^i} \\ &= \frac{1}{m} \theta_{2\text{A}}^T \left[ (\hat{\mathbf{y}}_{\text{A}}^i - \mathbf{y}_{\text{A}}^i) + \theta_{2\text{B}}^T (\hat{\mathbf{y}}_{\text{B}}^i - \mathbf{y}_{\text{B}}^i) \right] \end{aligned}$$

The residual of layer FC<sub>1A</sub>:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{x}_{1\text{A}}^i} &= \frac{\partial L}{\partial \mathbf{x}_{\text{DP}}^i} \frac{\partial \mathbf{x}_{\text{DP}}^i}{\partial \mathbf{x}_{1\text{A}}^i} \\ &= \frac{1}{m} \theta_{2\text{A}}^T \left[ (\hat{\mathbf{y}}_{\text{A}}^i - \mathbf{y}_{\text{A}}^i) + \theta_{2\text{B}}^T (\hat{\mathbf{y}}_{\text{B}}^i - \mathbf{y}_{\text{B}}^i) \right] \circ \mathbf{M}^i \\ &\quad \left( \mathbf{M}^i = \frac{1}{1-p} \text{sgn}(\mathbf{r}^i(\mathbf{x}_{1\text{A}}^i) - p \cdot \mathbf{1}) \right) \\ \frac{\partial L}{\partial \mathbf{z}_{1\text{A}}^i} &= \frac{\partial L}{\partial \mathbf{x}_{1\text{A}}^i} \frac{\partial \mathbf{x}_{1\text{A}}^i}{\partial \mathbf{z}_{1\text{A}}^i} = \frac{\partial L}{\partial \mathbf{x}_{1\text{A}}^i} \circ \frac{\partial \sin(\mathbf{z}_{1\text{A}}^i)}{\partial \mathbf{z}_{1\text{A}}^i} \\ &= \frac{1}{m} \theta_{2\text{A}}^T \left[ (\hat{\mathbf{y}}_{\text{A}}^i - \mathbf{y}_{\text{A}}^i) + \theta_{2\text{B}}^T (\hat{\mathbf{y}}_{\text{B}}^i - \mathbf{y}_{\text{B}}^i) \right] \circ \mathbf{M}^i \circ \cos(\mathbf{z}_{1\text{A}}^i) \end{aligned}$$

### Backward Propagation:

Therefore, we can get the gradients of the overall loss in a mini-batch with respect to the parameters at each layer.

### Task A:

The gradients of layer FC<sub>2A</sub>:

$$\begin{aligned}\frac{\partial L}{\partial \theta_{2A}} &= \sum_{i=1}^m \frac{\partial L}{\partial \hat{\mathbf{y}}_A^i} \frac{\partial \hat{\mathbf{y}}_A^i}{\partial \theta_{2A}} = \frac{1}{m} \sum_{i=1}^m \left[ (\hat{\mathbf{y}}_A^i - \mathbf{y}_A^i) + \theta_{2B}^T (\hat{\mathbf{y}}_B^i - \mathbf{y}_B^i) \right] (\mathbf{x}_{\text{DP}}^i)^T \\ \frac{\partial L}{\partial \mathbf{b}_{2A}} &= \sum_{i=1}^m \frac{\partial L}{\partial \hat{\mathbf{y}}_A^i} \frac{\partial \hat{\mathbf{y}}_A^i}{\partial \mathbf{b}_{2A}} = \frac{1}{m} \sum_{i=1}^m \left[ (\hat{\mathbf{y}}_A^i - \mathbf{y}_A^i) + \theta_{2B}^T (\hat{\mathbf{y}}_B^i - \mathbf{y}_B^i) \right]\end{aligned}$$

The gradients of layer FC<sub>1A</sub>:

$$\begin{aligned}\frac{\partial L}{\partial \theta_{1A}} &= \sum_{i=1}^m \frac{\partial L}{\partial \mathbf{z}_{1A}^i} \frac{\partial \mathbf{z}_{1A}^i}{\partial \theta_{1A}} = \frac{1}{m} \sum_{i=1}^m \left\{ \theta_{2A}^T \left[ (\hat{\mathbf{y}}_A^i - \mathbf{y}_A^i) + \theta_{2B}^T (\hat{\mathbf{y}}_B^i - \mathbf{y}_B^i) \right] \odot \mathbf{M}^i \odot \cos(\mathbf{z}_{1A}^i) (\mathbf{x}^i)^T \right\} \\ \frac{\partial L}{\partial \mathbf{b}_{1A}} &= \sum_{i=1}^m \frac{\partial L}{\partial \mathbf{z}_{1A}^i} \frac{\partial \mathbf{z}_{1A}^i}{\partial \mathbf{b}_{1A}} = \frac{1}{m} \sum_{i=1}^m \left\{ \theta_{2A}^T \left[ (\hat{\mathbf{y}}_A^i - \mathbf{y}_A^i) + \theta_{2B}^T (\hat{\mathbf{y}}_B^i - \mathbf{y}_B^i) \right] \odot \mathbf{M}^i \odot \cos(\mathbf{z}_{1A}^i) \right\}\end{aligned}$$

## Task B:

The gradients of layer FC<sub>2B</sub>:

$$\begin{aligned}\frac{\partial L}{\partial \theta_{2B}} &= \sum_{i=1}^m \frac{\partial L}{\partial \mathbf{z}_{2B}^i} \frac{\partial \mathbf{z}_{2B}^i}{\partial \theta_{2B}} = \frac{1}{m} \sum_{i=1}^m \left[ (\hat{\mathbf{y}}_B^i - \mathbf{y}_B^i) (\mathbf{x}_{\oplus}^i)^T \right] \\ \frac{\partial L}{\partial \mathbf{b}_{2B}} &= \sum_{i=1}^m \frac{\partial L}{\partial \mathbf{z}_{2B}^i} \frac{\partial \mathbf{z}_{2B}^i}{\partial \mathbf{b}_{2B}} = \frac{1}{m} \sum_{i=1}^m (\hat{\mathbf{y}}_B^i - \mathbf{y}_B^i)\end{aligned}$$

The gradients of layer FC<sub>1B</sub> and BN\*:

$$\begin{aligned}\frac{\partial L}{\partial \theta_{1B}} &= \sum_{i=1}^m \frac{\partial L}{\partial \mathbf{x}_{1B}^i} \frac{\partial \mathbf{x}_{1B}^i}{\partial \theta_{1B}} \\ &= \frac{1}{m} \sum_{i=1}^m \left\{ \left\{ \theta_{2B}^T (\hat{\mathbf{y}}_B^i - \mathbf{y}_B^i) \odot \text{sgn}(\mathbf{x}_{\text{BN}}^i) - \frac{1}{m} \sum_{j=1}^m \left[ \theta_{2B}^T (\hat{\mathbf{y}}_B^j - \mathbf{y}_B^j) \odot \text{sgn}(\mathbf{x}_{\text{BN}}^j) \right] \right\} (\mathbf{x}^i)^T \right\} \\ \frac{\partial L}{\partial \mathbf{b}_{1B}} &= \sum_{i=1}^m \frac{\partial L}{\partial \mathbf{x}_{1B}^i} \frac{\partial \mathbf{x}_{1B}^i}{\partial \mathbf{b}_{1B}} = \frac{1}{m} \sum_{i=1}^m \left[ \theta_{2B}^T (\hat{\mathbf{y}}_B^i - \mathbf{y}_B^i) \odot \text{sgn}(\mathbf{x}_{\text{BN}}^i) \right]\end{aligned}$$