

# MathDNN Homework 6

Department of Computer Science and Engineering  
2021-16988 Jaewan Park

## Problem 1

Let  $A_i$ ,  $b_i$  the parameters for the linear relationship between layers  $y_i$  and  $y_{i+1}$ . ( $A_i \in \mathbb{R}^{m \times n}$ ,  $b_i \in \mathbb{R}^m$ ,  $y_i \in \mathbb{R}^n$ ,  $y_{i+1} \in \mathbb{R}^m$ ) Then applying dropout before the activation function results in

$$y_{i+1} = \sigma(A_i y'_i + b_i), \quad [y'_i]_j = \begin{cases} 0 & \text{(with probability } p) \\ \frac{[y_i]_j}{1-p} & \text{(otherwise)} \end{cases}$$

while the opposite gives

$$[y_{i+1}]_j = \begin{cases} 0 & \text{(with probability } p) \\ \frac{[\sigma(A_i y_i + b_i)]_j}{1-p} & \text{(otherwise)} \end{cases}.$$

For notational convenience, denote the resulting  $y_{i+1}$  from each cases  $y_{i+1}^1$  and  $y_{i+1}^2$ . We should find the appropriate  $\sigma$  where  $\mathbf{E}[y_{i+1}^1]_j = \mathbf{E}[y_{i+1}^2]_j$  for  $j = 1, 2, \dots, m$ . We can generally obtain the following.

$$\begin{aligned} \mathbf{E}[y_{i+1}^1]_j &= \mathbf{E}[\sigma(A_i y'_i + b_i)]_j = \mathbf{E}[\sigma([A_i y'_i + b_i]_j)] \\ &= \mathbf{E}[\sigma([A_i]_{j,:} \cdot y'_i + [b_i]_j)] \\ &= \mathbf{E}[\sigma([A_i]_{j,1}[y'_i]_1 + \dots + [A_i]_{j,n}[y'_i]_n + [b_i]_j)] \\ &= \mathbf{E}[\sigma([A_i]_{j,1}[y'_i]_1)] + \dots + \mathbf{E}[\sigma([A_i]_{j,n}[y'_i]_n)] + \mathbf{E}[\sigma([b_i]_j)] \\ &= \sigma\left([A_i]_{j,1} \frac{[y_i]_1}{1-p}\right) \times (1-p) + \dots + \sigma\left([A_i]_{j,n} \frac{[y_i]_n}{1-p}\right) \times (1-p) + \sigma([b_i]_j) \\ \mathbf{E}[y_{i+1}^2]_j &= \frac{[\sigma(A_i y_i + b_i)]_j}{1-p} \times (1-p) \\ &= \sigma([A_i]_{j,1}[y_i]_1) + \dots + \sigma([A_i]_{j,n}[y_i]_n) + \sigma([b_i]_j) \end{aligned}$$

For the three cases, results differ as the following.

(1) **nn.ReLU()**

For  $a > 0$ ,  $\text{ReLU}(ax) = a\text{ReLU}(x)$  from the definition of ReLU. Therefore the two results are equivalent.

$$\begin{aligned} \mathbf{E}[y_{i+1}^1]_j &= \sigma\left([A_i]_{j,1} \frac{[y_i]_1}{1-p}\right) \times (1-p) + \dots + \sigma\left([A_i]_{j,n} \frac{[y_i]_n}{1-p}\right) \times (1-p) + \sigma([b_i]_j) \\ &= \sigma([A_i]_{j,1}[y_i]_1) + \dots + \sigma([A_i]_{j,n}[y_i]_n) + \sigma([b_i]_j) \\ &= \mathbf{E}[y_{i+1}^2]_j \end{aligned}$$

(2) `nn.Sigmoid()`

Generally,  $\sigma\left([A_i]_{j,k} \frac{[y_i]_k}{1-p}\right)(1-p) = \sigma\left([A_i]_{j,1}[y_i]_1\right)$  does not stand for the sigmoid function. Therefore changing the order gives nonequivalent results.

(3) `nn.LeakyReLU()`

For  $a > 0$ ,  $\text{LeakyReLU}(ax) = a\text{LeakyReLU}(x)$  from the definition of leaky ReLU. Therefore similar to the ReLU function, changing the order gives equivalent results.

Therefore in cases of ReLU and leaky ReLU, the order of dropout and activation functions does not matter.

## Problem 2

The weight initializing code of PyTorch's `nn.Linear` model is given as the following.

```
def reset_parameters(self) -> None:
    # Setting a=sqrt(5) in kaiming_uniform is the same as initializing with
    # uniform(-1/sqrt(in_features), 1/sqrt(in_features)). For details, see
    # https://github.com/pytorch/pytorch/issues/57109
    init.kaiming_uniform_(self.weight, a=math.sqrt(5))
    if self.bias is not None:
        fan_in, _ = init._calculate_fan_in_and_fan_out(self.weight)
        bound = 1 / math.sqrt(fan_in) if fan_in > 0 else 0
        init.uniform_(self.bias, -bound, bound)
```

The weight and bias elements are all initialized from  $\mathcal{U}(-\text{bound}, \text{bound})$  where  $\text{bound} = \frac{1}{\sqrt{\text{fan\_in}}}$ . In this case,  $\text{fan\_in} = n_\ell$  for each step. Therefore all elements of  $[A_k]$  and  $[b_k]$  follow  $\mathcal{U}\left(-\frac{1}{\sqrt{n_{k-1}}}, \frac{1}{\sqrt{n_{k-1}}}\right)$ , and have mean 0, variance  $\frac{1}{3n_{k-1}}$ .

For  $j = 1, 2, \dots, n_0$ ,  $\mathbf{E}[x_j] = 0$  and  $\mathbf{Var}[x_j] = 1$ . Since all variables are uncorrelated, for  $j = 1, 2, \dots, n_1$ ,

$$\begin{aligned}\mathbf{E}[y_1]_j &= \mathbf{E}[A_1]_{j,1}[x]_1 + \dots + [A_1]_{j,n_0}[x]_{n_0} + [b_1]_j \\ &= \mathbf{E}[A_1]_{j,1}\mathbf{E}[x]_1 + \dots + \mathbf{E}[A_1]_{j,n_0}\mathbf{E}[x]_{n_0} + \mathbf{E}[b_1]_j \\ &= 0 \\ \mathbf{Var}[y_1]_j &= \mathbf{Var}[A_1]_{j,1}[x]_1 + \dots + [A_1]_{j,n_0}[x]_{n_0} + [b_1]_j \\ &= \mathbf{Var}[A_1]_{j,1}\mathbf{Var}[x]_1 + \dots + \mathbf{Var}[A_1]_{j,n_0}\mathbf{Var}[x]_{n_0} + \mathbf{Var}[b_1]_j \\ &= \frac{n_0 + 1}{3n_0}.\end{aligned}$$

Let  $y_0 = x$ . Then generally, we can inductively show that for  $l = 1, 2, \dots, L$ , all elements of  $y_\ell$  have the

same mean and variance as the following.

$$\begin{aligned}
\mathbf{E}[y_\ell]_j &= \mathbf{E}[A_\ell]_{j,1} \mathbf{E}[y_{\ell-1}]_1 + \cdots + \mathbf{E}[A_\ell]_{j,n_{\ell-1}} \mathbf{E}[y_{\ell-1}]_{n_{\ell-1}} + \mathbf{E}[b_\ell]_j \\
&= 0 \\
\mathbf{Var}[y_\ell]_j &= \mathbf{Var}[A_\ell]_{j,1} \mathbf{Var}[y_{\ell-1}]_1 + \cdots + \mathbf{Var}[A_\ell]_{j,n_{\ell-1}} \mathbf{Var}[y_{\ell-1}]_{n_{\ell-1}} + \mathbf{Var}[b_\ell]_j \\
&= \frac{\mathbf{Var}[y_{\ell-1}]_1 \cdot n_{\ell-1} + 1}{3n_{\ell-1}}
\end{aligned}$$

The variance cannot be further simplified, but if the bias is 0, we can simply say  $\mathbf{Var}[y_\ell]_j = \frac{1}{3^\ell}$ .

### Problem 3

*Notation* For a matrix or vector  $X$ , the notation  $[X]_{i,j}$  or  $[X]_i$  refers to the element of  $X$  at that index, and  $\{f(i,j)\}_{i,j}$  refers to a matrix of which element at  $(i,j)$  is  $f(i,j)$ .

(1) For  $\ell = L$ ,

$$\frac{\partial y_\ell}{\partial y_{\ell-1}} = \frac{\partial y_L}{\partial y_{L-1}} = A_L.$$

For  $\ell = 2, \dots, L-1$ ,

$$\begin{aligned}
\frac{\partial y_\ell}{\partial y_{\ell-1}} &= \frac{\partial}{\partial y_{\ell-1}} (\sigma(A_\ell y_{\ell-1} + b_\ell) + y_{\ell-1}) \\
&= \left\{ \frac{\partial}{\partial [y_{\ell-1}]_j} [\sigma(A_\ell y_{\ell-1} + b_\ell) + y_{\ell-1}]_i \right\}_{i,j} \\
&= \left\{ \frac{\partial}{\partial [y_{\ell-1}]_j} (\sigma([A_\ell y_{\ell-1}]_i + [b_\ell]_i) + [y_{\ell-1}]_i) \right\}_{i,j} \\
&= \left\{ \sigma'([A_\ell y_{\ell-1}]_i + [b_\ell]_i) \cdot \frac{\partial}{\partial [y_{\ell-1}]_j} ([A_\ell y_{\ell-1}]_i + [b_\ell]_i) + 1 \right\}_{i,j} \\
&= \left\{ \sigma'([A_\ell y_{\ell-1}]_i + [b_\ell]_i) \cdot \frac{\partial}{\partial [y_{\ell-1}]_j} \left( \sum_{k=1}^{n_{\ell-1}} [A_\ell]_{i,k} [y_{\ell-1}]_k + [b_\ell]_i \right) \right\}_{i,j} \\
&= \left\{ \sigma'([A_\ell y_{\ell-1}]_i + [b_\ell]_i) \cdot [A_\ell]_{i,j} + 1 \right\}_{i,j} \\
&= \text{diag} \left( \sigma'(A_\ell y_{\ell-1} + b_\ell) \right) A_\ell + \{1\}_{i,j}.
\end{aligned}$$

(2) For  $\ell = L$ ,

$$\frac{\partial y_L}{\partial b_\ell} = \frac{\partial y_L}{\partial b_L} = \frac{\partial}{\partial b_L} (A_L y_{L-1} + b_L) = 1.$$

For  $\ell = 1, \dots, L-1$ , we can obtain the following.

$$\begin{aligned}
\frac{\partial y_\ell}{\partial b_\ell} &= \frac{\partial}{\partial b_\ell} (\sigma(A_\ell y_{\ell-1} + b_\ell) + y_{\ell-1}) = \frac{\partial}{\partial b_\ell} \sigma(A_\ell y_{\ell-1} + b_\ell) \\
&= \left\{ \frac{\partial}{\partial [b_\ell]_j} [\sigma(A_\ell y_{\ell-1} + b_\ell)]_i \right\}_{i,j} = \left\{ \frac{\partial}{\partial [b_\ell]_j} \sigma([A_\ell y_{\ell-1} + b_\ell]_i) \right\}_{i,j} = \left\{ \frac{\partial}{\partial [b_\ell]_j} \sigma([A_\ell y_{\ell-1}]_i + [b_\ell]_i) \right\}_{i,j} \\
&= \left\{ \sigma'([A_\ell y_{\ell-1}]_i + [b_\ell]_i) \cdot \frac{\partial}{\partial [b_\ell]_j} ([A_\ell y_{\ell-1}]_i + [b_\ell]_i) \right\}_{i,j} \\
&= \{\sigma'([A_\ell y_{\ell-1}]_i + [b_\ell]_i) \cdot \delta_{ij}\}_{i,j} \quad (\delta_{ij}: \text{Kronecker Delta}) \\
&= \text{diag} \left( \sigma'(A_\ell y_{\ell-1} + b_\ell) \right)
\end{aligned}$$

Therefore

$$\begin{aligned}
\frac{\partial y_L}{\partial b_\ell} &= \frac{\partial y_L}{\partial y_{L-1}} \frac{\partial y_{L-1}}{\partial y_{L-2}} \dots \frac{\partial y_{\ell+1}}{\partial y_\ell} \frac{\partial y_\ell}{\partial b_\ell} \\
&= A_L \left( \text{diag} \left( \sigma'(A_{L-1} y_{L-2} + b_{L-1}) \right) A_{L-1} + \{1\}_{i,j} \right) \\
&\quad \dots \left( \text{diag} \left( \sigma'(A_{\ell+1} y_\ell + b_{\ell+1}) \right) A_{\ell+1} + \{1\}_{i,j} \right) \text{diag}(\sigma'(A_\ell y_{\ell-1} + b_\ell)).
\end{aligned}$$

For  $\ell = L$ ,

$$\frac{\partial y_L}{\partial A_\ell} = \frac{\partial y_L}{\partial A_L} = \frac{\partial}{\partial A_L} (A_L y_{L-1} + b_L) = y_{L-1}.$$

For  $\ell = 1, \dots, L-1$ , we can obtain the following.

$$\begin{aligned}
\left[ \frac{\partial y_L}{\partial A_\ell} \right]_{i,j} &= \frac{\partial y_L}{\partial [A_\ell]_{i,j}} = \frac{\partial y_L}{\partial y_\ell} \frac{\partial y_\ell}{\partial [A_\ell]_{i,j}} = \frac{\partial y_L}{\partial y_\ell} \left\{ \frac{\partial [y_\ell]_k}{\partial [A_\ell]_{i,j}} \right\}_k \\
&= \frac{\partial y_L}{\partial y_\ell} \left\{ \frac{\partial}{\partial [A_\ell]_{i,j}} \left( \sigma([A_\ell y_{\ell-1}]_k + [b_\ell]_k) + [y_{\ell-1}]_k \right) \right\}_k = \frac{\partial y_L}{\partial y_\ell} \left\{ \frac{\partial}{\partial [A_\ell]_{i,j}} \sigma([A_\ell y_{\ell-1}]_k + [b_\ell]_k) \right\}_k \\
&= \frac{\partial y_L}{\partial y_\ell} \left\{ \sigma'([A_\ell y_{\ell-1}]_k + [b_\ell]_k) \left( \frac{\partial}{\partial [A_\ell]_{i,j}} ([A_\ell y_{\ell-1}]_k + [b_\ell]_k) \right) \right\}_k \\
&= \frac{\partial y_L}{\partial y_\ell} \begin{bmatrix} 0 \\ \vdots \\ \sigma'([A_\ell y_{\ell-1}]_i + [b_\ell]_i) [y_{\ell-1}]_j \\ \vdots \\ 0 \end{bmatrix} \quad (\text{All elements except the } i\text{-th element are 0.}) \\
&= \left[ \frac{\partial y_L}{\partial y_\ell} \right]_i \sigma'([A_\ell y_{\ell-1}]_i + [b_\ell]_i) [y_{\ell-1}]_j \\
&= [\sigma'(A_\ell y_{\ell-1} + b_\ell)]_i \left[ \frac{\partial y_L}{\partial y_\ell} \right]_i [y_{\ell-1}]_j
\end{aligned}$$

Therefore

$$\begin{aligned}
\frac{\partial y_L}{\partial A_\ell} &= \text{diag}(\sigma'(A_\ell y_{\ell-1} + b_\ell)) \left( \frac{\partial y_L}{\partial y_\ell} \right)^\top y_{\ell-1}^\top \\
&= \text{diag}(\sigma'(A_\ell y_{\ell-1} + b_\ell)) \left( \frac{\partial y_L}{\partial y_{L-1}} \right)^\top \left( \frac{\partial y_{L-1}}{\partial y_{L-2}} \right)^\top \cdots \left( \frac{\partial y_{\ell+1}}{\partial y_\ell} \right)^\top y_{\ell-1}^\top \\
&= \text{diag}(\sigma'(A_\ell y_{\ell-1} + b_\ell)) A_L^\top \left( \text{diag}(\sigma'(A_{L-1} y_{L-2} + b_{L-1})) A_{L-1} + \{1\}_{i,j} \right)^\top \\
&\quad \cdots \left( \text{diag}(\sigma'(A_{\ell+1} y_\ell + b_{\ell+1})) A_{\ell+1} + \{1\}_{i,j} \right)^\top y_{\ell-1}^\top.
\end{aligned}$$

(3) For  $i \leq \ell$ ,

$$\begin{aligned}
\frac{\partial y_L}{\partial b_i} &= A_L \left( \text{diag}(\sigma'(A_{L-1} y_{L-2} + b_{L-1})) A_{L-1} + \{1\}_{i,j} \right) \\
&\quad \cdots \left( \text{diag}(\sigma'(A_{i+1} y_i + b_{i+1})) A_{i+1} + \{1\}_{i,j} \right) \text{diag}(\sigma'(A_i y_{i-1} + b_i)) \\
\frac{\partial y_L}{\partial A_i} &= \text{diag}(\sigma'(A_i y_{i-1} + b_i)) A_L^\top \left( \text{diag}(\sigma'(A_{L-1} y_{L-2} + b_{L-1})) A_{L-1} + \{1\}_{i,j} \right)^\top \\
&\quad \cdots \left( \text{diag}(\sigma'(A_{i+1} y_i + b_{i+1})) A_{i+1} + \{1\}_{i,j} \right)^\top y_{i-1}^\top.
\end{aligned}$$

The given conditions guarantee that for any  $\ell$  and  $i \leq \ell$ , at least one of  $A_{L-1}, \dots, A_{i+1}$  and one of  $\text{diag}(\sigma'(A_{L-1} y_{L-2} + b_{L-1})), \dots, \text{diag}(\sigma'(A_{i+1} y_i + b_{i+1}))$  is 0. However, due to the  $\{1\}_{i,j}$  term in each calculation, the gradients do not vanish, rather remain in large values.

## Problem 4

- (1) The number of trainable parameters between two layers can be calculated as  $C_{\text{out}} \times (C_{\text{in}} \times F \times F + 1)$ , where  $F$  is the size of the convolution filter and  $C_{\text{in}}$  and  $C_{\text{out}}$  are each the number of channels for input and output. The addition of 1 is made due to the bias. In the original construction, the total number of parameters is

$$128 \times (256 \times 1 \times 1 + 1) + 128 \times (128 \times 3 \times 3 + 1) + 256 \times (128 \times 1 \times 1 + 1) = 213504.$$

In the modified construction, the total number of parameters is

$$(4 \times (256 \times 1 \times 1 + 1)) \times 32 + (4 \times (4 \times 3 \times 3 + 1)) \times 32 + (256 \times (4 \times 1 \times 1 + 1)) \times 32 = 78592.$$

Therefore the modified version reduces the number of trainable parameters considerably.

- (2) The convolution model can be implemented via PyTorch as the following. (The final summation is done for 32 objects, thus for convenience the repetition is abbreviated as ‘...’ in the code.)

```

class STMConvLayer(nn.Module):
    def __init__(self):
        super(STMConvLayer, self).__init__()
        self.conv1layers = [nn.Conv2d(256, 4, 1) for _ in range(32)]
        self.conv2layers = [nn.Conv2d(4, 4, 3) for _ in range(32)]
        self.conv3layers = [nn.Conv2d(4, 256, 1) for _ in range(32)]
    def forward(self, x):
        # [apply 1x1conv with 4 output channels
        #  apply 3x3conv with 4 output channels (with padding=1)
        #  apply 1x1conv with 256 output channels] X 32
        # Add all 32 outputs
        out_list = [nn.ReLU(self.conv1layers[i](x)) for i in range(32)]
        out_list = [nn.ReLU(self.conv2layers[i](out_list[i])) for i in range(32)]
        out_list = [nn.ReLU(self.conv3layers[i](out_list[i])) for i in range(32)]
        out = out_list[0] + out_list[1] + ... + out_list[31]
        return out

```