

Backpropagation
Lab Report # 1

By
312581020
許瀚丰

Deep Learning
Spring 2024
Date Submitted: March 26, 2024

Contents

1	Introduction	4
1.1	Lab Description	4
1.2	Objective	4
2	Experiment setups	4
2.1	Sigmoid	5
2.1.1	Method	5
2.1.2	Implementation	6
2.2	Neural Network	6
2.2.1	Linear Layer	6
2.2.2	Loss Function	6
2.3	Backpropagation	7
2.3.1	Method	7
2.3.2	Implementation	8
3	Results of your testing	10
3.1	Hyperparameter Settings	10
3.2	Screenshot and comparison figure	10
3.3	Show the accuracy of your prediction	10
3.4	Learning curve (loss, epoch curve)	12
4	Discussion	12
4.1	Try different learning rates	12
4.1.1	Linear	12
4.1.2	XOR	13
4.2	Try different numbers of hidden units	13
4.2.1	Linear	13
4.2.2	XOR	14

4.3	Try without activation functions	14
4.3.1	Linear	14
4.3.2	XOR	15
5	Extra	15
5.1	Implement different optimizers	15
5.1.1	Linear	16
5.1.2	XOR	16
5.2	Implement different activation functions	17
5.2.1	Linear	17
5.2.2	XOR	17
5.3	Implement convolutional layers (Not Finished Yet)	18
5.4	Multiclass Architecture	18
A	Optimizer	21
A.1	SGDM	21
A.2	RMSPProp	21
A.3	Adam	22
B	Activation Function	24
B.1	ReLU	24
B.2	LeakyReLU	24
B.3	Tanh	24

Listings

1	Activation Class	5
2	Layer Class	5
3	Sigmoid	6
4	Linear Layer	6

5	Binary Cross Entropy Loss	7
6	Backpropagation	9
7	SGD	9
8	Convolution Layer	18
9	Softmax	19
10	Cross Entropy Loss	19
11	Neural Network with multiclass	19
12	SGDM	21
13	RMSPProp	21
14	Adam	22
15	ReLU	24
16	LeakyReLU	24
17	Tanh	24

1 Introduction

1.1 Lab Description

本次實驗是實作一個簡單的Neural Network，並包含了利用Forward來預測答案，再透過Backward來更新網路的參數。而本實驗的目標就是希望讓預測的答案與實際的答案越像越好，並透過兩個任務來驗證模型的好壞，分別為一個簡單的線性可分的Linear問題與一個線性不可分的XOR問題。

1.2 Objective

在此實驗中，我希望能夠透過嘗試撰寫自己的Package的實作方式，希望能夠做出類似於Pytorch的工具，在未來遇到簡單的分類問題就能夠直接套用。而在實作上，我分別將一些常用的功能分開撰寫，分別為以下幾種：

nn.py

包含常見的Activation Function和NN的初始化與Forward、Backward等。

loss.py

包含在分類模型中常見的loss function，如BCELoss, CrossEntropyLoss等。

optim.py

包含一些常見的optimizer與一個簡單的learning rate scheduler。

utils.py

包含save model、load model 與一些常用的功能。

2 Experiment setups

在本實驗中，我共有兩種不同的Abstract Class，分別為給Activation Function繼承的Activation Class與給不同Network繼承的Layer Class，分別都需要實作Forward與Backward，並紀錄從前一層來的Net Input，以利之後需要更新時能夠直接使用，如程式碼1與2，分別為Activation Class 與Layer Class。

```

1 class activation(ABC):
2     def __init__(self):
3         self.a = None
4     @abstractmethod
5     def forward(self):
6         pass
7     @abstractmethod
8     def backward(self):
9         pass

```

Listing 1: Activation Class

```

1 class layer(ABC):
2     def __init__(self):
3         self.w = None
4         self.b = None
5         self.a = None
6     @abstractmethod
7     def forward(self):
8         pass
9     @abstractmethod
10    def backward(self):
11        pass

```

Listing 2: Layer Class

2.1 Sigmoid

2.1.1 Method

sigmoid function 能夠將輸入從 $(-\infty, \infty)$ mapping 成 $(-1, 1)$ ，且處處可導。其公式為 $\sigma(x) = \frac{1}{1+e^{-x}}$ ，而在計算時，我們需要對於其求導，化簡後可以得到一個簡單的結果 $\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$ ，其證明如下：

$$\begin{aligned}
 \frac{d\sigma(x)}{dx} &= \left(\frac{1}{1+e^{-x}}\right)'(e^{-x})' = \frac{-1}{(1+e^{-x})^{-2}} * -e^{-x} = \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1}{(1+e^{-x})} \frac{e^{-x}}{(1+e^{-x})} \\
 &= \sigma(x) * \frac{(1+e^{-x}) - 1}{(1+e^{-x})} = \sigma(x) * \left(1 - \frac{1}{(1+e^{-x})}\right) = \sigma(x)(1 - \sigma(x))
 \end{aligned}$$

2.1.2 Implementation

我的實作方式參考了助教給的程式碼，在forward時，就直接回傳input經過sigmoid的結果，而在backward時，由於是對於已經經過sigmoid的結果做偏微分，因此直接對其計算即可，如程式碼3所示：

```
1 class sigmoid(activation):
2     def forward(self, x):
3         self.a = x
4         return 1.0 / (1.0 + np.exp(-x))
5     def backward(self, x):
6         return np.multiply(x, 1 - x)
```

Listing 3: Sigmoid

2.2 Neural Network

2.2.1 Linear Layer

對於Neural Network的部分，首先我先實作出Linear Layer的架構，一開始會先初始化每層的weights與bias，而Forward就是簡單的將weights乘上x再加上bias，而内部的a則是記錄每層輸入進來的input，用以在對權重做更新時能夠使用，如程式碼4所示：

```
1 class Linear(layer):
2     def __init__(self, in_features, out_features):
3         self.w = np.random.randn(in_features, out_features).T
4         self.b = np.random.randn(1, out_features).T
5         self.a = None
6     def forward(self, x):
7         self.a = x
8         return self.w @ x + self.b
9     def backward(self, delta):
10        return self.w.T @ delta
```

Listing 4: Linear Layer

2.2.2 Loss Function

在本次實驗中Loss Function我選擇的是Binary Cross Entropy，原因為其相比與Mean

Square Error更適合處理這種分類的問題，使用時就將預測結果 \hat{y} 與真實的Label y 帶入公式 $\mathcal{L}_{BCE}(\hat{y}, y) = -(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y}))$ ，並回傳其對於 \hat{y} 的偏微分 $\frac{\partial L}{\partial \hat{y}} = -(\frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}})$ ，如程式碼5所示。

```

1 class BCELoss():
2     def __init__(self):
3         pass
4     def __call__(self, y_pred, y_target, eps = 1e-6):
5         y_pred = np.clip(y_pred, eps, 1 - eps)
6         loss = -np.mean(y_target * np.log(y_pred) + (1 - y_target) * np.log(1 - y_pred))
7         delta = -(y_target / y_pred - (1 - y_target) / (1 - y_pred))
8         return loss, delta

```

Listing 5: Binary Cross Entropy Loss

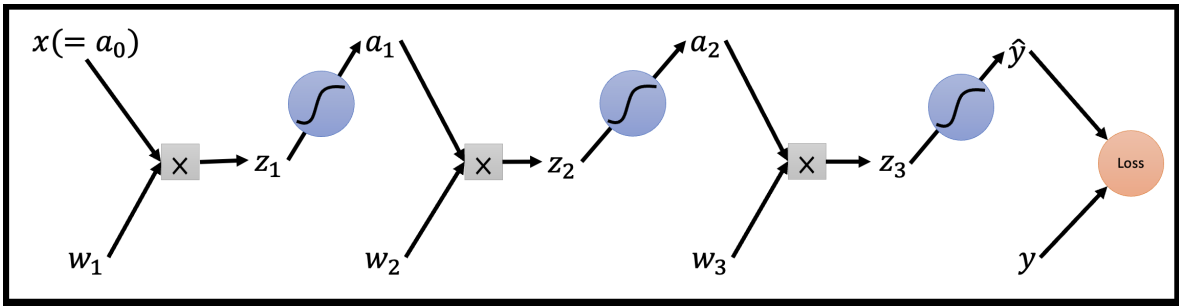


Figure 1: 此圖為一個簡單的Computational Graph範例。

2.3 Backpropagation

2.3.1 Method

在Backpropagation的部分，為了減少一些重複的計算，我使用Delta Rule的方式來更新網路。假設在模型中Loss Function為Binary Cross Entropy(BCELoss)，而最後一層則透過Sigmoid Function來讓結果介於(-1,1)之間，通過簡化我們可以得到最後一層L的 $\Delta_L = \hat{y} - y$ ，而對於其他層 $L - 1$ 到1，其 $\Delta_l = ((\mathbf{W}_{l+1})^T \Delta_{l+1}) \odot [\mathbf{a}_l \odot (1 - \mathbf{a}_l)]$ ，其中 a_l 為Foward時輸入此Layer的值，並令 $a_0 = x$ ，而 \odot 則是Elementwise product，而

更新則是透過 $w_l = w_l - \Delta_l * a_{l-1}$ 來進行。以下為透過圖1為例子作為說明：

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_3} \frac{\partial z_3}{\partial w_3} = -\left(\frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}}\right)(\hat{y}(1-\hat{y}))a_2 = (\hat{y} - y)a_2$$

對於最後一層的 w_3 對 Loss Function 的偏微分計算如上所示，由於 BCE Loss 微分後的分母與 Sigmoid 微分後相乘剛好可以化簡，因此我們可以得到 $\Delta_3 = (\hat{y} - y)$ 。但若不是用這個組合，也只需計算個別計算 $\frac{\partial L}{\partial \hat{y}}$ 與 $\frac{\partial \hat{y}}{\partial z_3}$ 再相乘即可。

$$\begin{aligned} \frac{\partial L}{\partial w_2} &= \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_3} \frac{\partial z_3}{\partial a_2} \frac{\partial a_2}{\partial z_2} \frac{\partial z_2}{\partial w_2} = (\hat{y} - y)w_3 \odot (a^2 \odot (1 - a^2))a^1 \\ &= ((w_3)^T \Delta_2) \odot [a_2 \odot (1 - a_2)]a_1 \end{aligned}$$

倒數第二層的 w_2 對於 Loss Function 的偏微分為 $((w_3)^T \Delta_2) \odot [a_2 \odot (1 - a_2)]a_1$ ，我們可以發現得出來的結論中在 a_1 前面那項剛好就是 $((w_3)^T \Delta_2) \odot [a_2 \odot (1 - a_2)]$ ，我們將其定為 Δ_2 。

$$\begin{aligned} \frac{\partial L}{\partial w_1} &= \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_3} \frac{\partial z_3}{\partial a_2} \frac{\partial a_2}{\partial z_2} \frac{\partial z_2}{\partial a_1} \frac{\partial a_1}{\partial z_1} \frac{\partial z_1}{\partial x_1} \\ &= (\hat{y} - y)w_3 \odot (a^2 \odot (1 - a^2))w_2 \odot (a^1 \odot (1 - a^1))a_0 \\ &= ((w_2)^T \Delta_2) \odot [a_1 \odot (1 - a_1)]a_0 \end{aligned}$$

倒數第三層也就是 w_1 對於 Loss Function 的偏微分為 $((w_2)^T \Delta_2) \odot [a_1 \odot (1 - a_1)]a_0$ ，因此 Δ_1 也與上面推導的結果一樣，就是 $((w_2)^T \Delta_2) \odot [a_1 \odot (1 - a_1)]$

由此方法我們在計算時就能減少許多重複計算的部分，在程式的撰寫上也更有彈性，讓使用者可以自由的設定模型的參數與想要層數上，不需再額外推導繁瑣的更新公式。

2.3.2 Implementation

在實作上，建立模型需要先定義此模型的 Forward 與 Backward 的方式，如程式碼6所示，Backward 就如在 Method 所提的方式，一開始的 Δ_L 是由 Loss Function 所回傳

的 $\frac{\partial L}{\partial \hat{y}}$ 值，進入迴圈後會將 Δ_L 乘上 $\frac{\partial \hat{y}}{\partial z_3}$ ，並依據規則依序計算出其他的 Δ_l ，最後再由程式碼7中的方式，將每一層linear的wieght與bias更新，以此就能夠訓練整個模型。

```

1 class Model(nn.Module):
2     def __init__(self, input_size, output_size, activation, hidden_size):
3         super().__init__()
4         if activation in nn.activation_functions_list:
5             self.act = eval(f"nn.{activation}")
6         else:
7             raise ValueError(f"Activation function {activation} is not supported")
8         self.net = [
9             nn.Linear(input_size, hidden_size),
10            self.act(),
11            nn.Linear(hidden_size, hidden_size),
12            self.act(),
13            nn.Linear(hidden_size, output_size),
14            nn.sigmoid(),
15        ]
16        self.delta = np.empty(len(self.net) + 1, dtype = object)
17    def forward(self, x):
18        for layer in self.net:
19            x = layer.forward(x)
20        return x
21    def backward(self, delta, y_pred):
22        self.delta[-1] = delta
23        for i in range(len(self.net) - 1, -1, -1):
24            if isinstance(self.net[i], nn.Layer):
25                self.delta[i] = self.net[i].backward(self.delta[i + 1])
26            elif isinstance(self.net[i], nn.activation):
27                if i + 1 >= len(self.net):
28                    self.delta[i] = self.delta[i + 1] * self.net[i].backward(y_pred)
29                else:
30                    self.delta[i] = self.delta[i + 1] * self.net[i].backward(self.net[
31            i + 1].a)
32            else:
33                assert False, "Unknown layer type"

```

Listing 6: Backpropagation

```

1 class SGD():
2     def __init__(self, model, lr = 0.001):
3         self.lr = lr
4         self.model = model
5     def step(self):

```

```

6         for i, layer in enumerate(self.model.net):
7             if isinstance(layer, nn.Layer):
8                 layer.w -= self.lr * (self.model.delta[i + 1] @ layer.a.T)
9                 layer.b -= self.lr * self.model.delta[i + 1]

```

Listing 7: SGD

3 Results of your testing

3.1 Hyperparameter Settings

以下為本次實驗的超參數設定:

- **Activation Function:** Sigmoid
- **Loss Function:** Binary Cross Entropy
- **Optimizer:** SGD
- **NN Architecture:** $x \rightarrow \text{linear}(2, 32) \xrightarrow{\sigma} \text{linear}(32, 32) \xrightarrow{\sigma} \text{linear}(32, 1) \xrightarrow{\sigma} \hat{y}$
- **Epoch:** 5000
- **Learning Rate:** $5 * 10^{-2}$

3.2 Screenshot and comparison figure

實驗結果如表1所示，可見模型確實成功訓練並完成了此分類的任務。

3.3 Show the accuracy of your prediction

實驗結果如表2所示，第一個row為兩個任務各自的Epoch與Loss的關係，第二個row則是部分擷取最後九個點的預測結果與對於整個資料集的預測準確率，我們可以觀察到Linear的結果相較於XOR的結果更靠近0或1，可見XOR的任務較為困難。

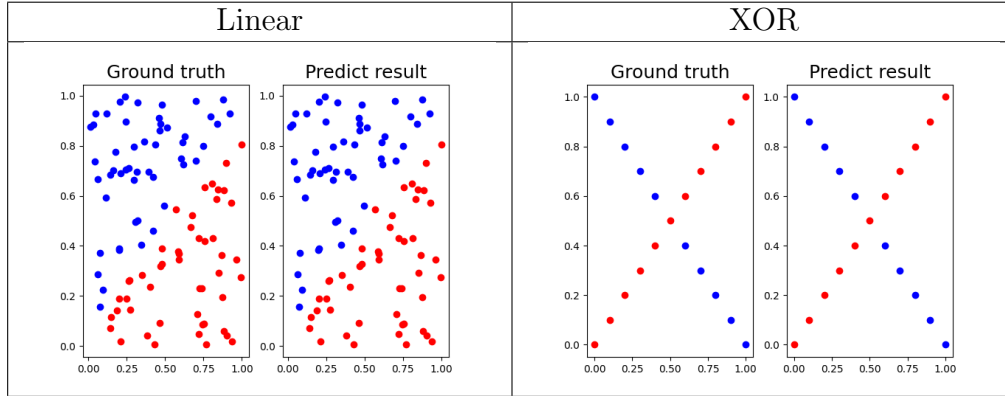


Table 1: Result

Linear	XOR
Epoch 500/5000 Loss: 0.05663244	Epoch 500/5000 Loss: 0.31611481
Epoch 1000/5000 Loss: 0.03983622	Epoch 1000/5000 Loss: 0.23784658
Epoch 1500/5000 Loss: 0.03402205	Epoch 1500/5000 Loss: 0.20701075
Epoch 2000/5000 Loss: 0.03133604	Epoch 2000/5000 Loss: 0.19210126
Epoch 2500/5000 Loss: 0.02993964	Epoch 2500/5000 Loss: 0.18418954
Epoch 3000/5000 Loss: 0.02917074	Epoch 3000/5000 Loss: 0.17978816
Epoch 3500/5000 Loss: 0.02873435	Epoch 3500/5000 Loss: 0.17727625
Epoch 4000/5000 Loss: 0.02848248	Epoch 4000/5000 Loss: 0.17582199
Epoch 4500/5000 Loss: 0.02833573	Epoch 4500/5000 Loss: 0.17497311
Epoch 5000/5000 Loss: 0.02824974	Epoch 5000/5000 Loss: 0.17447523
[6.13203018e-07]	[0.53532701]
[3.63784915e-01]	[0.15402676]
[8.20798810e-01]	[0.79981006]
[1.06650747e-06]	[0.07087889]
[1.02071072e-02]	[0.93510711]
[2.62242979e-06]	[0.03745177]
[1.39522461e-07]	[0.97180459]
[9.99971459e-01]	[0.02357602]
[1.62103941e-07]]	[0.98207442]]
Final Accuracy: 1.00	Final Accuracy: 1.00

Table 2: Prediction

3.4 Learning curve (loss, epoch curve)

實驗結果如表3所示，x軸方向為Epoch，y軸方向則是Loss的大小，由於兩者所使用的超參數皆相同，因此我們同樣可以從Loss下降的幅度與收斂值觀察XOR的任務相比與Linear確實是較為具有挑戰性的。

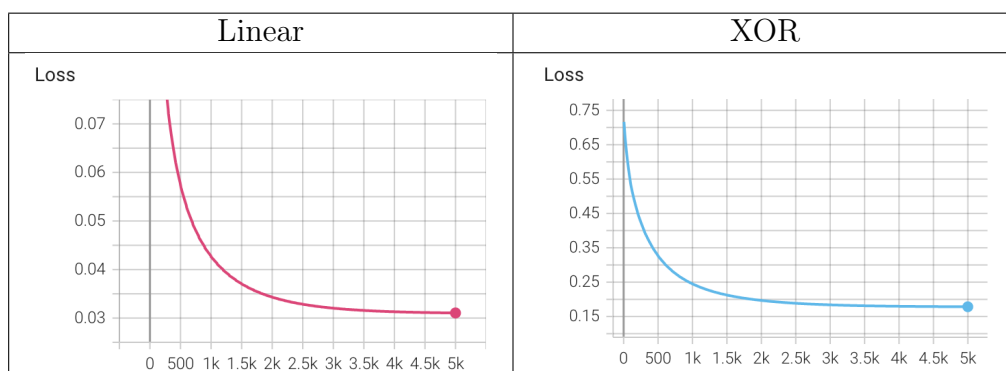


Table 3: Learning Curve

4 Discussion

4.1 Try different learning rates

以下實驗分別測試五種不同的Learning Rate，分別為：0.0005, 0.005, 0.05, 0.5, 1

4.1.1 Linear

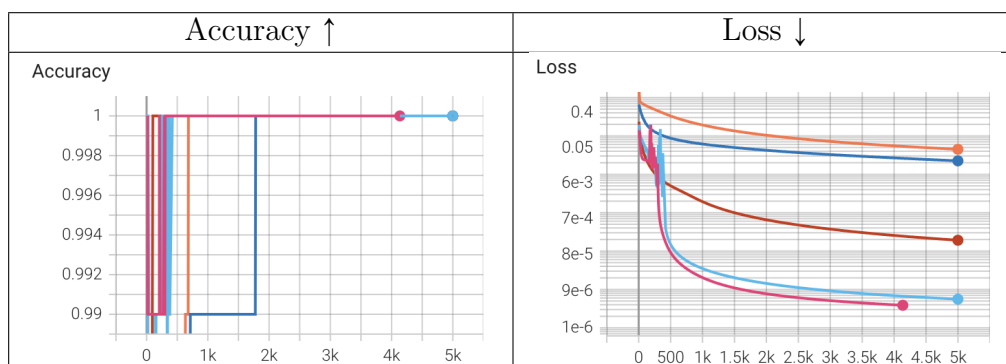


Table 4: Different Learning Rates on Linear

實驗結果如表4所示，其顏色與LR的對照為0.0005(Orange)，0.005(Blue)，0.05(Red)，0.5(Cyan)，1(Magenta)。我們可以觀察到在Accuracy的部分，不同的LR對Acc的差異不大，都為0.99至1.0，但在LR大於 $5 * 10^{-1}$ 的結果在Loss的前期可以看到很明顯的震盪。

4.1.2 XOR

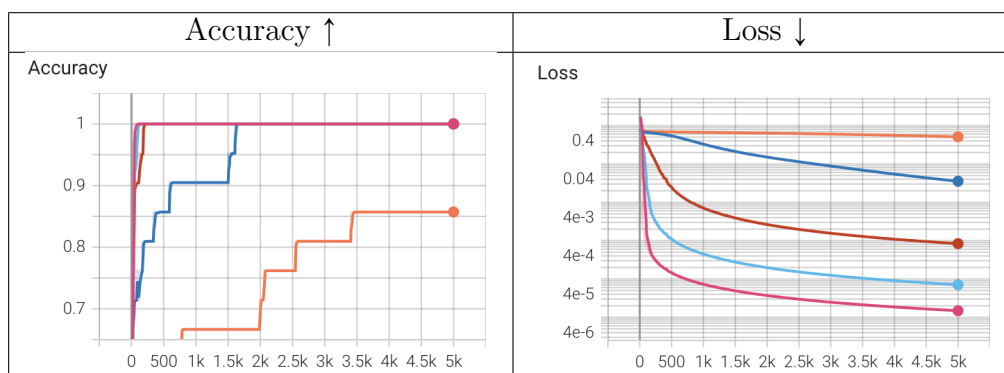


Table 5: Different Learning Rates on XOR

實驗結果如表5所示，其顏色與LR的對照為0.0005(Orange)，0.005(Blue)，0.05(Red)，0.5(Cyan)，1(Magenta)。我們可以觀察到除了LR為0.0005的情況，其餘不同的settings到最後的Acc皆為1，我認為是因為XOR這個任務的資料非常少，因此LR太小反而讓模型困在Local minimum進而提早收斂，才導致此結果。

4.2 Try different numbers of hidden units

以下實驗分別測試五種不同的hidden units，分別為：64，32，16，8，4，2，1

4.2.1 Linear

實驗結果如表6所示，其顏色與hidden units的對照為64(Orange)，32(Blue)，16(Red)，8(Cyan)，4(Magenta)，2(Green)，1(Gray)，從實驗中可以發現由於Linear是一個非常簡單且線性可分的問題，只要一個hidden layer就能夠完成分類的任務，因此不管在Acc與Loss上選擇不同hidden unit的差異不大。

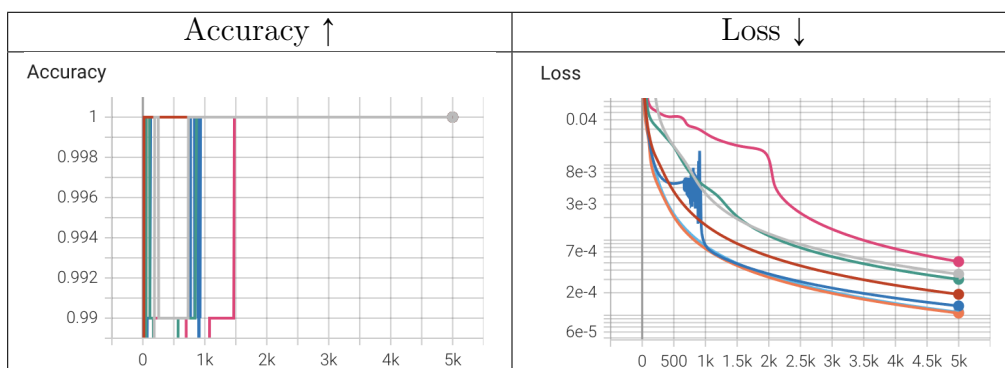


Table 6: Different hidden units on Linear

4.2.2 XOR

實驗結果如表7所示，其顏色與hidden units的對照為64(Orange)，32(Blue)，16(Red)，8(Cyan)，4(Magenta)，2(Green)，1(Gray)，從實驗中可以發現XOR的問題除了hidden unit為1的情況，其餘不管是在Acc與Loss皆差異不大。

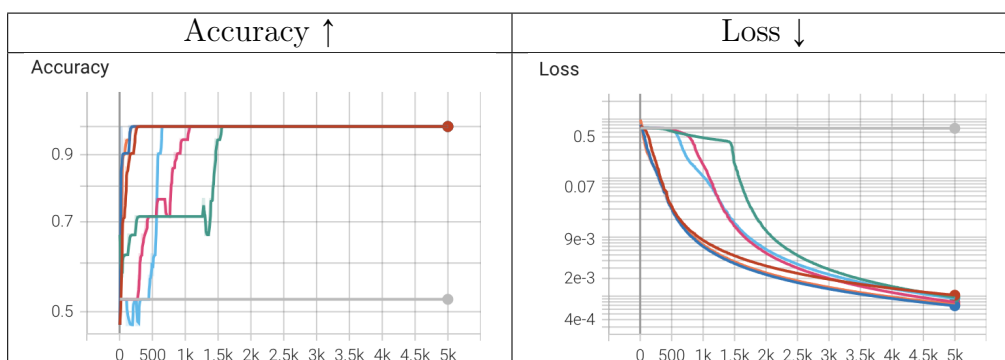


Table 7: Different hidden units on XOR

4.3 Try without activation functions

4.3.1 Linear

實驗結果如表8所示，有activation function(Orange)，沒有activation function(Blue)，從實驗中結果我發現對於Linear這種線性可分的任務有無使用activation function結果差異不大，但不使用activation function所產生的結果可能會較不穩定。

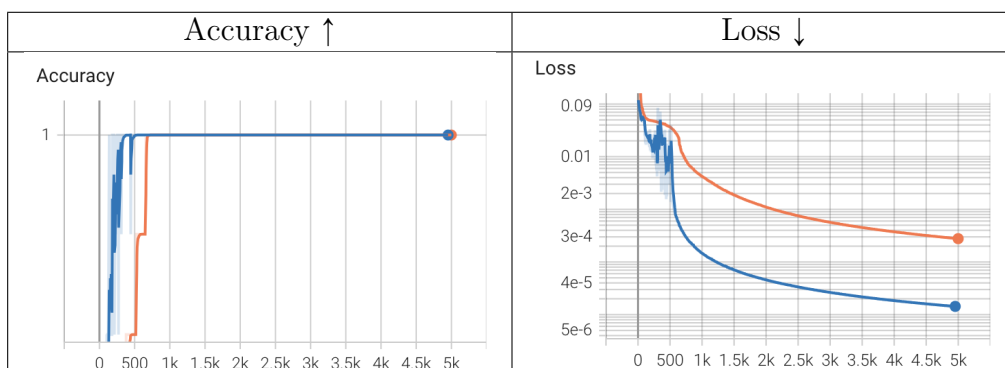


Table 8: Comparison of Using Activation Function on Linear

4.3.2 XOR

實驗結果如表9所示，有activation function(Orange)，沒有activation function(Blue)，從實驗中結果我發現對於XOR這種線性不可分的問題，若不使用activation function，不管如何運算其輸出都還是輸入的線性轉換，因此無法解決這種問題。

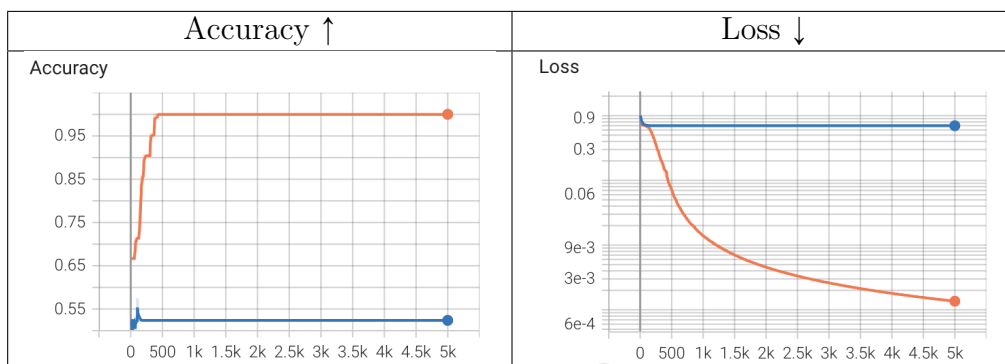


Table 9: Comparison of Using Activation Function on XOR

5 Extra

5.1 Implement different optimizers

在本次實驗中，除了原始的SGD外，我還另外實作了SGDM(SGD with Momentum)、RMSProp與Adam，請見A。

5.1.1 Linear

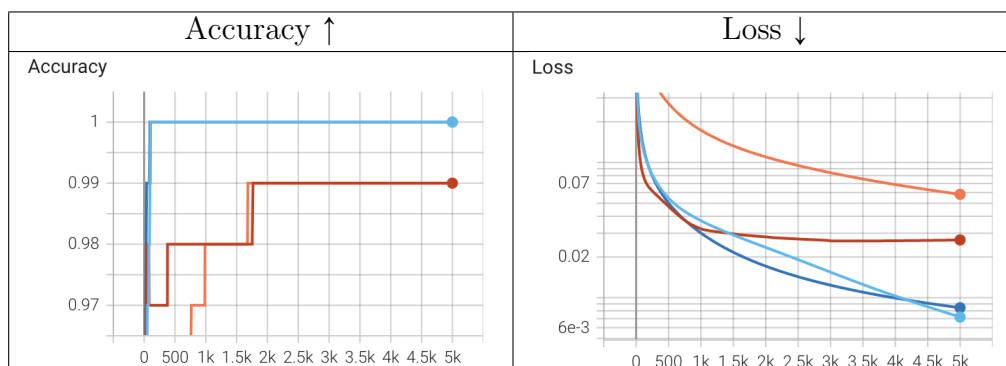


Table 10: Different Optimizers on Linear

實驗結果如表10分別為SGD(Orange)、SGDM(Blue)、RMSProp(Red)、Adam(Cyan)，可以觀察到從SGD更換成其他Optimizer皆能有效的提升模型更新的效率，在Linear這個任務中，使用Adam與SGDM的效果是最好的。

5.1.2 XOR

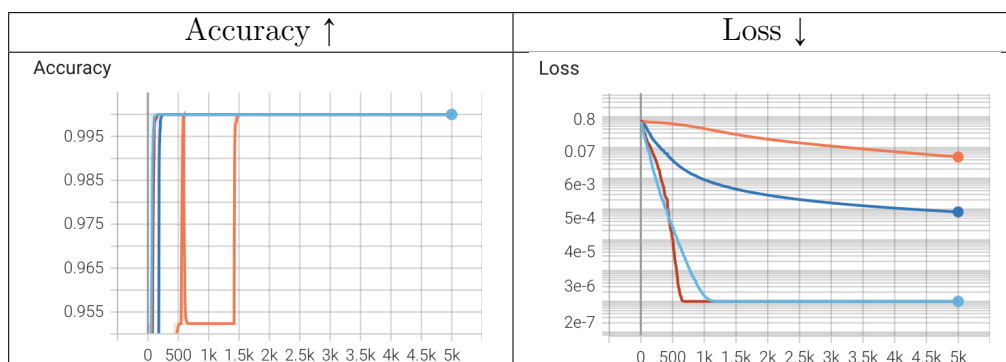


Table 11: Different Optimizers on Linear

實驗結果如表11分別為SGD(Orange)、SGDM(Blue)、RMSProp(Red)、Adam(Cyan)，可以觀察到從SGD更換成其他Optimizer皆能有效的提升模型更新的效率，在XOR這個任務中，使用Adam與RMSProp的效果最好的。

5.2 Implement different activation functions

在本次實驗中，除了原始的sigmoid外，我還另外實作了ReLU、LeakyReLU與Tanh，實作部分請參見B。

5.2.1 Linear

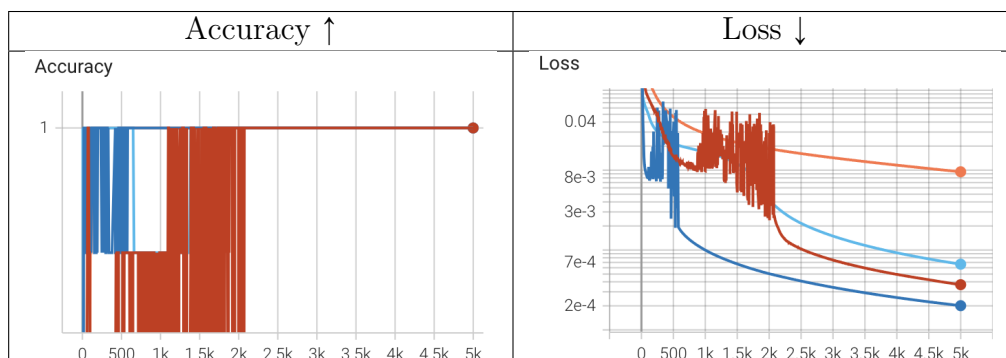


Table 12: Different Activation Function on Linear

實驗結果如表12分別為sigmoid(Orange)、ReLU(Blue)、LeakyReLU(Red)、Tanh(Cyan)，可以觀察到從sigmoid更換成其他activation function皆能有效的使loss下降，但也可以發現使用ReLU與LeakyReLU都會在前期產生不小的震盪，後期則會趨於穩定，最後則是ReLU的Loss最低。

5.2.2 XOR

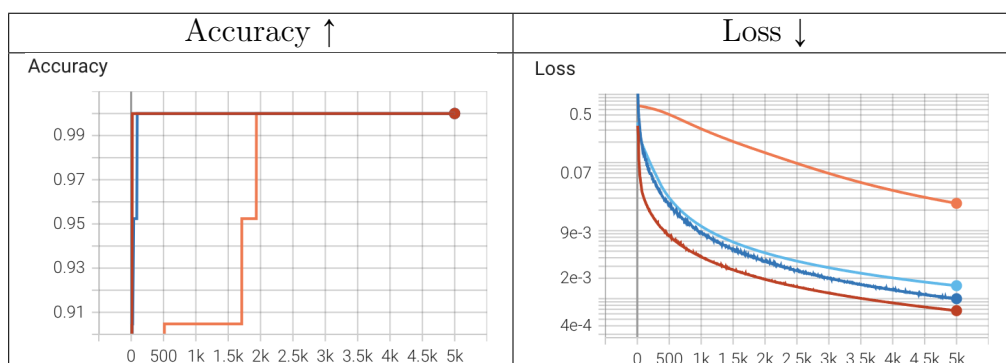


Table 13: Different Activation Function on XOR

實驗結果如表13分別為sigmoid(Orange)、ReLU(Blue)、LeakyReLU(Red)、Tanh(Cyan)，可以觀察到從sigmoid更換成其他activation function皆能有效的使loss下降，而最後效果最好的則是LeakyReLU。

5.3 Implement convolutional layers (Not Finished Yet)

在本實驗中，我嘗試實作了Conv1d的Forward，透過對於kernel對輸入的點進行運算並加總，但對於Backward的實作上遇到一些困難因此無法完成。

```
1 class Simple_Conv1d(layer):
2     def __init__(self, kernel_size = 3):
3         self.w = np.random.randn((kernel_size, 1))
4         self.b = np.random.randn(1)
5         self.a = None
6     def forward(self, x):
7         self.a = x
8         h = np.zeros(x.shape[-2] - self.w.shape[0] + 1)
9         for i in range(h.shape[0]):
10             h[i] = np.sum(x[i:i + self.w.shape[0]] * self.w)
11         return h + self.b
12     def backward(self, delta):
13         pass
```

Listing 8: Convolution Layer

5.4 Multiclass Architecture

在實作上，由於本次實驗的兩個任務皆為一個二分類問題，因此在結果上只需要輸出為一個神經元 \hat{y} ，並以是否大於0.5作為依據來判斷是某個class即可。然而若遇到其他多分類的任務此架構便不再適用。因此為了讓模型能夠加彈性，我另外實作了multiclass的版本，將輸出的結果變為一個vector，並透過Softmax將其轉換此輸入為某個class的機率，並透過Cross Entropy來計算Loss來更新神經網路，Softmax與Cross Entropy的實作如程式碼9與10。

在實作時，由於Softmax的計算為單調性的，因此預測出label實際上就是值最大的那個index，因此不需要真的經過Softmax，而實際上Softmax會直接放在CrossEntropy Loss中才會使用，會先將預測的結果做完Softmax後在與實

際label的one-hot encoding的結果計算Cross Entropy來計算Loss，而其對於前一層輸出的偏微分為 $\frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} = (\hat{y} - y)$ ，剛好與過去使用的Sigmoid加上BCELoss的結果相同，其Neural Network的架構如程式碼11。

```

1 class Softmax(activation):
2     def forward(self, x):
3         self.a = x
4         exp_x = np.exp(x - np.max(x))
5         return exp_x / np.sum(exp_x, axis = 0)
6
7     def backward(self, x):
8         return x * (1 - x)

```

Listing 9: Softmax

```

1 class CrossEntropyLoss():
2     def __init__(self):
3         pass
4     def __call__(self, y_pred, y_target, eps = 1e-8):
5         y_target = one_hot_encoding(y_target, y_pred.shape[0]).T
6         y_pred = np.exp(y_pred - np.max(y_pred)) / np.sum(np.exp(y_pred - np.max(
7         y_pred)), axis = 0)
8         y_pred = np.clip(y_pred, eps, 1 - eps)
9         loss = -np.mean(y_target * np.log(y_pred))
10        delta = (y_pred - y_target)
11        return loss, delta

```

Listing 10: Cross Entropy Loss

```

1 class Model(nn.Module):
2     def __init__(self, input_size, output_size, activation, hidden_size):
3         super().__init__()
4         if activation in nn.activation_functions_list:
5             self.act = eval(f"nn.{activation}")
6         else:
7             raise ValueError(f"Activation function {activation} is not supported")
8         self.net = [
9             nn.Linear(input_size, hidden_size),
10            self.act(),
11            nn.Linear(hidden_size, hidden_size),
12            self.act(),
13            nn.Linear(hidden_size, output_size),
14        ]

```

```

15     self.delta = np.empty(len(self.net) + 1, dtype = object)
16
17     def forward(self, x):
18         for layer in self.net:
19             x = layer.forward(x)
20         return x
21
22     def backward(self, delta, y_pred):
23         self.delta[-1] = delta
24         for i in range(len(self.net) - 1, -1, -1):
25             if isinstance(self.net[i], nn.layer):
26                 self.delta[i] = self.net[i].backward(self.delta[i + 1])
27             elif isinstance(self.net[i], nn.activation):
28                 if i + 1 >= len(self.net):
29                     self.delta[i] = self.delta[i + 1] * self.net[i].backward(y_pred)
30                 else:
31                     self.delta[i] = self.delta[i + 1] * self.net[i].backward(self.net[
32 i + 1].a)
33             else:
34                 assert False, "Unknown layer type"

```

Listing 11: Neural Network with multiclass

A Optimizer

A.1 SGDM

```
1 class SGDM():
2     def __init__(self, model, lr = 0.001, momentum = 0.9):
3         assert isinstance(model, nn.module), "model should be an instance of module"
4         assert 0 < momentum < 1, "momentum should be in (0, 1)"
5         self.lr = lr
6         self.momentum = momentum
7         self.model = model
8         self.v_w = []
9         self.v_b = []
10        for layer in self.model.net:
11            if isinstance(layer, nn.layer):
12                self.v_w.append(np.zeros_like(layer.w))
13                self.v_b.append(np.zeros_like(layer.b))
14            else:
15                self.v_w.append(None)
16                self.v_b.append(None)
17
18        def step(self):
19            for i, layer in enumerate(self.model.net):
20                if isinstance(layer, nn.layer):
21                    self.v_w[i] = self.momentum * self.v_w[i] - self.lr * self.model.delta
22                    [i + 1] @ layer.a.T
23                    self.v_b[i] = self.momentum * self.v_b[i] - self.lr * self.model.delta
24                    [i + 1]
25
26                    layer.w += self.v_w[i]
27                    layer.b += self.v_b[i]
```

Listing 12: SGDM

A.2 RMSProp

```
1 class RMSprop():
2     def __init__(self, model, lr = 0.001, alpha = 0.9, epsilon = 1e-8):
3         assert isinstance(model, nn.module), "model should be an instance of module"
4         assert 0 < alpha < 1, "alpha should be in (0, 1)"
5         assert 0 < epsilon, "epsilon should be positive"
6         self.lr = lr
7         self.alpha = alpha
```

```

8     self.epsilon = epsilon
9     self.model = model
10    self.v_w = []
11    self.v_b = []
12    for layer in self.model.net:
13        if isinstance(layer, nn.Layer):
14            self.v_w.append(np.zeros_like(layer.w))
15            self.v_b.append(np.zeros_like(layer.b))
16        else:
17            self.v_w.append(None)
18            self.v_b.append(None)
19
20    def step(self):
21        for i, layer in enumerate(self.model.net):
22            if isinstance(layer, nn.Layer):
23                self.v_w[i] = self.alpha * self.v_w[i] + (1 - self.alpha) * (self.
model.delta[i + 1] @ layer.a.T) ** 2
24                self.v_b[i] = self.alpha * self.v_b[i] + (1 - self.alpha) * (self.
model.delta[i + 1]) ** 2
25                layer.w -= self.lr * self.model.delta[i + 1] @ layer.a.T / (np.sqrt(
self.v_w[i]) + self.epsilon)
26                layer.b -= self.lr * self.model.delta[i + 1] / (np.sqrt(self.v_b[i]) +
self.epsilon)

```

Listing 13: RMSProp

A.3 Adam

```

1 class Adam():
2     def __init__(self, model, lr = 0.001, beta1 = 0.9, beta2 = 0.999, epsilon = 1e-8):
3         assert isinstance(model, nn.Module), "model should be an instance of module"
4         assert 0 < beta1 < 1 and 0 < beta2 < 1, "beta1 and beta2 should be in (0, 1)"
5         assert 0 < epsilon, "epsilon should be positive"
6         self.lr = lr
7         self.beta1 = beta1
8         self.beta2 = beta2
9         self.epsilon = epsilon
10        self.model = model
11        self.m_w, self.v_w, self.m_b, self.v_b = [], [], [], []
12        for layer in self.model.net:
13            if isinstance(layer, nn.Layer):
14                self.m_w.append(np.zeros_like(layer.w))

```

```

15         self.v_w.append(np.zeros_like(layer.w))
16         self.m_b.append(np.zeros_like(layer.b))
17         self.v_b.append(np.zeros_like(layer.b))
18     else:
19         self.m_w.append(None)
20         self.v_w.append(None)
21         self.m_b.append(None)
22         self.v_b.append(None)
23     self.t = 0
24
25     def step(self):
26         self.t += 1
27         for i, layer in enumerate(self.model.net):
28             if isinstance(layer, nn.Layer):
29                 self.m_w[i] = self.beta1 * self.m_w[i] + (1 - self.beta1) * (self.
model.delta[i + 1] @ layer.a.T)
30                 self.v_w[i] = self.beta2 * self.v_w[i] + (1 - self.beta2) * (self.
model.delta[i + 1] @ layer.a.T) ** 2
31                 m_hat_w = self.m_w[i] / (1 - self.beta1 ** self.t)
32                 v_hat_w = self.v_w[i] / (1 - self.beta2 ** self.t)
33                 layer.w -= self.lr * m_hat_w / (np.sqrt(v_hat_w) + self.epsilon)
34
35                 self.m_b[i] = self.beta1 * self.m_b[i] + (1 - self.beta1) * self.model
.delta[i + 1]
36                 self.v_b[i] = self.beta2 * self.v_b[i] + (1 - self.beta2) * (self.
model.delta[i + 1]) ** 2
37                 m_hat_b = self.m_b[i] / (1 - self.beta1 ** self.t)
38                 v_hat_b = self.v_b[i] / (1 - self.beta2 ** self.t)
39                 layer.b -= self.lr * m_hat_b / (np.sqrt(v_hat_b) + self.epsilon)

```

Listing 14: Adam

B Activation Function

B.1 ReLU

```
1 class ReLU(activation):
2     def forward(self, x):
3         self.a = x
4         return np.maximum(0, x)
5     def backward(self, x):
6         return (x > 0).astype(int)
```

Listing 15: ReLU

B.2 LeakyReLU

```
1 class LeakyReLU(activation):
2     def __init__(self, alpha = 0.02):
3         self.alpha = alpha
4     def forward(self, x):
5         self.a = x
6         return np.maximum(self.alpha * x, x)
7     def backward(self, x):
8         return (x > 0).astype(int) + self.alpha * (x <= 0).astype(int)
```

Listing 16: LeakyReLU

B.3 Tanh

```
1 class Tanh(activation):
2     def forward(self, x):
3         self.a = x
4         return np.tanh(x)
5     def backward(self, x):
6         return 1 - x ** 2
```

Listing 17: Tanh