



# DL Class Lab1

## Lab1 : Backpropagation

### ▼ 1. Introduction



Lab1 是僅以 numpy library 去實現 **Backpropagation** 的算法，並建立 regression model 訓練 **linear** 與 **XOR** 分佈資料。透過這次業加強了對於 gradient-based training 的 neuron network 之實作與訓練技能。藉由這次機會練習建立 python 之模組化的物件，增加整體程式的有效擴充性，進而達到易維護與易測試的目的，下方為個物件的介：

#### ▼ A. Activation:

- Attributes:
  - forward\_func → dict, for calling the forward pass easier.
  - backward\_func → dict, for calling the backward pass easier.
  - func\_name → str, to determine which activation method should be called.
  - pre\_input → numpy array, to record the previous input.
  - pre\_output → numpy array, to record the previous forward output.
- Methods:
  - For forward pass:
    - forward
    - sigmoid
    - relu
    - softplus
  - For backward pass:
    - backward
    - derivative\_sigmoid
    - derivative\_relu
    - derivative\_softplus
  - For update some trainable activation functions:
    - update

#### ▼ B. Linear\_Layer:

- Attributes:
  - input\_dim → int, the dimension of the input.

- output\_dim → int, the dimension of the output.
- w → numpy array, the weights of linear layer with the shape of (output\_dim, input\_dim)
- delta\_w → numpy array, to record the gradient of the weights.
- momentum\_w → numpy array, to record the momentum of the weights.
- beta → float, hyper parameter for updating the momentum\_w.
- pre\_input → numpy array, to record the previous input of forward pass.
- Methods:
  - forward
  - backward
  - update
  - L2\_Norm → for normalize the weights with the method of l2 norm.

### ▼ C. Model:

- Attributes:
  - history → dict, to record the information of training and testing.
  - layers → list, the list contains the whole layers of model
- Methods
  - Build
  - fit → for training the model
  - forward
  - backward
  - update
  - Record\_History
  - Show\_Training\_Stage\_Info
  - predict\_2\_binary
  - Compute\_Accuracy
  - evaluate

## ▼ 2. Experiment Setups

- Device & Software:
  - Device: Azure Server
  - OS: ubuntu 18.04
  - python version: 3.8
  - numpy version: 1.22.4
- Standar Model Setups:
  - Architecture:

- 2-3-Sigmoid-3-Sigmoid-1, which means the code below:

```
model.Build([
    Linear_Layer(2, 3),
    Activation('sigmoid'),
    Linear_Layer(3, 3),
    Activation('sigmoid'),
    Linear_Layer(3, 1),
])
```

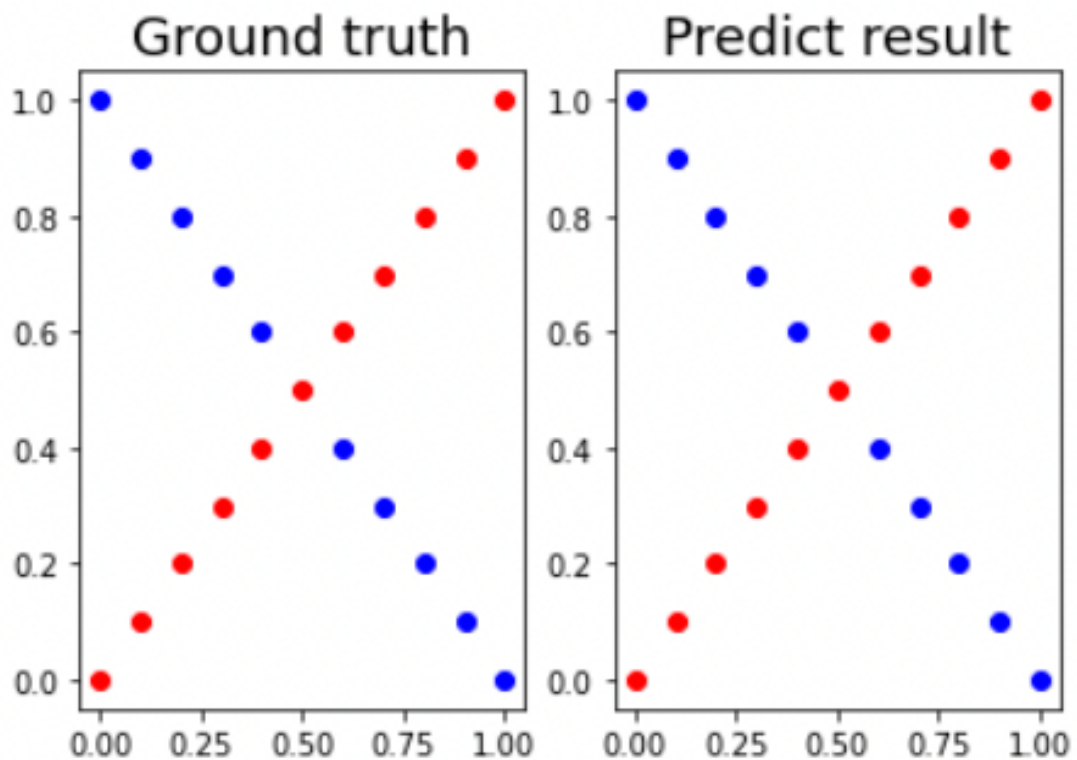
- Epochs: 30000
- Learning Rate: 0.1
- Optimizer: sgdm
- Loss Function: Mean Square Error

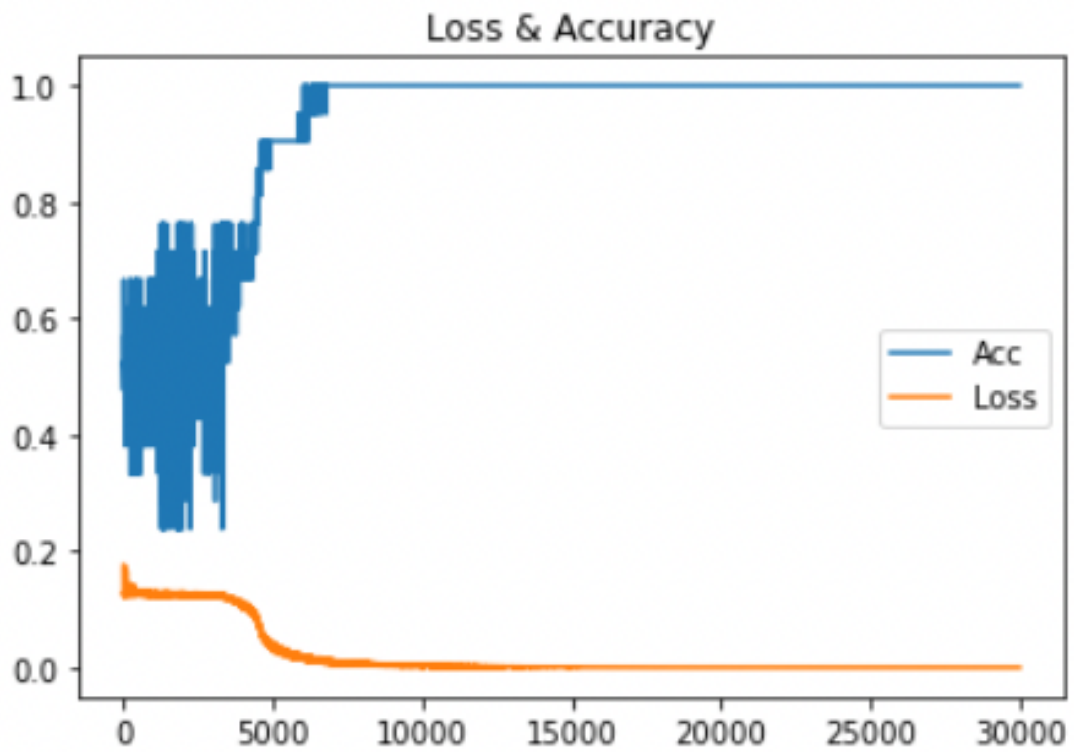
### ▼ 3. Results of Your Testing

#### ▼ A. XOR Dataset

Standar Model Result-XOR

Datasets	Learning Rate	Optimizer	Architectures	Epochs	Accuracy
XOR	0.1	sgdm	2-3-Sigmoid-3-Sigmoid-1	30000	100%

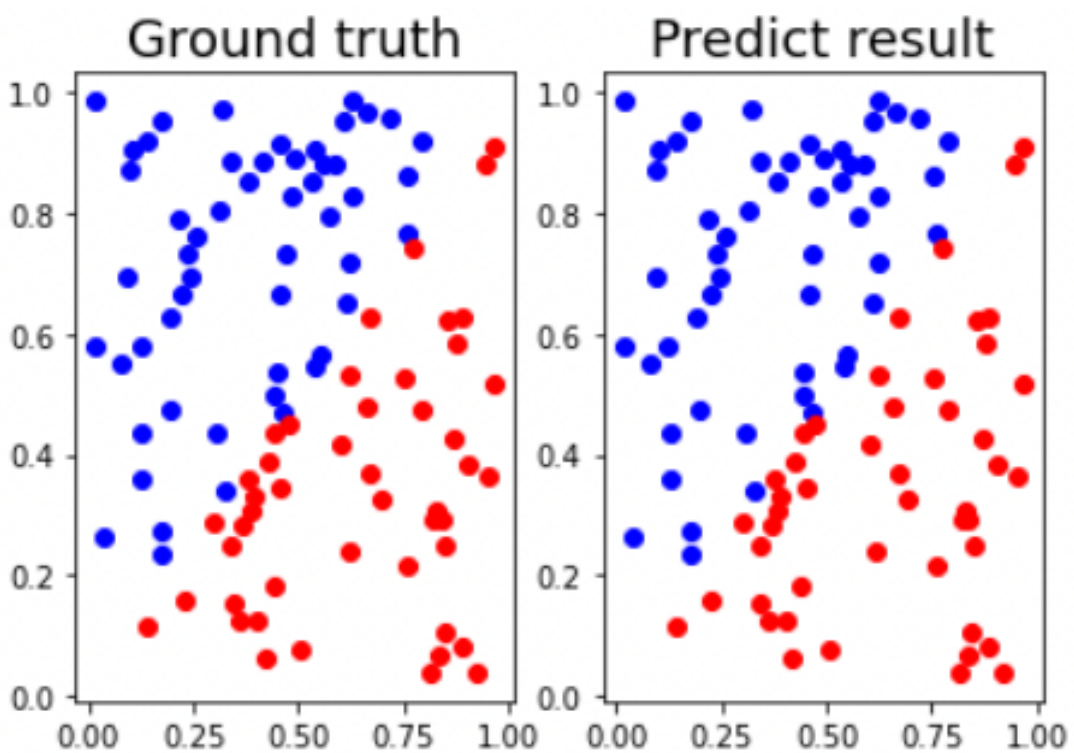


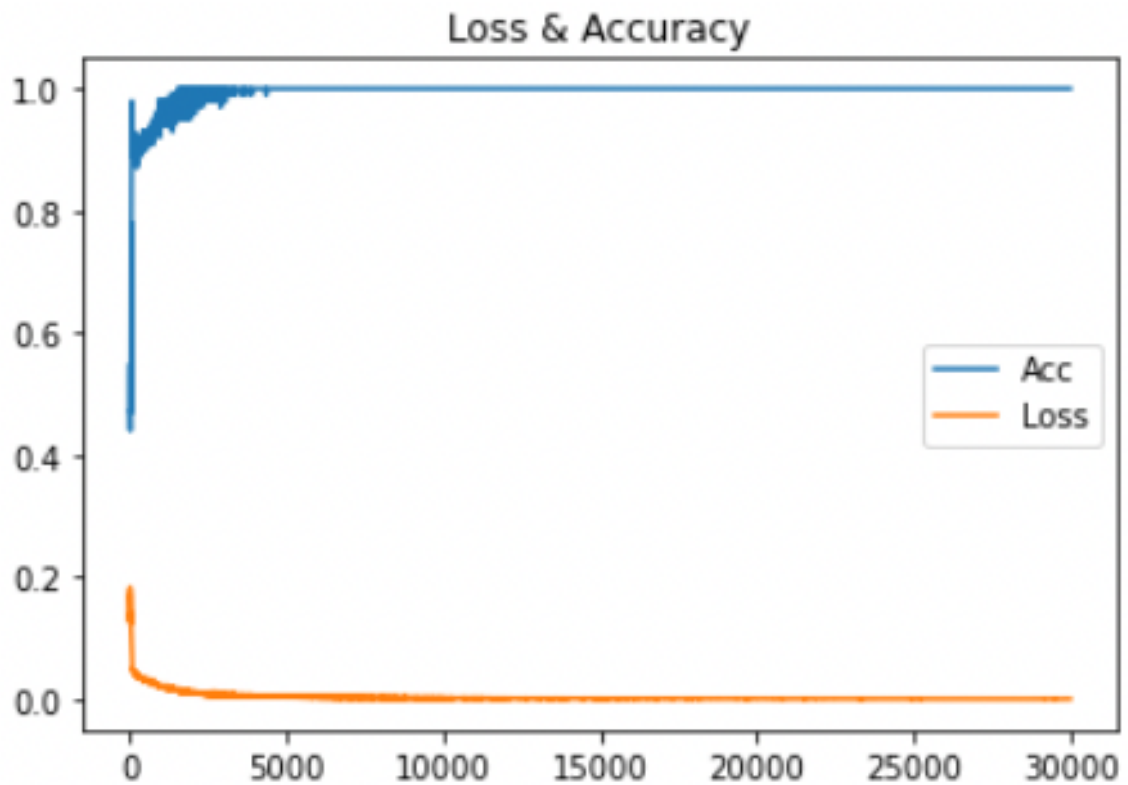


#### ▼ B. Linear Dataset

Standar Model Result-Linear

Datasets	Learning Rate	Optimizer	Architectures	Epochs	Accuracy
Linear	0.1	sgdm	2-3-Sigmoid-3-Sigmoid-1	30000	100%





## ▼ 4. Discussion

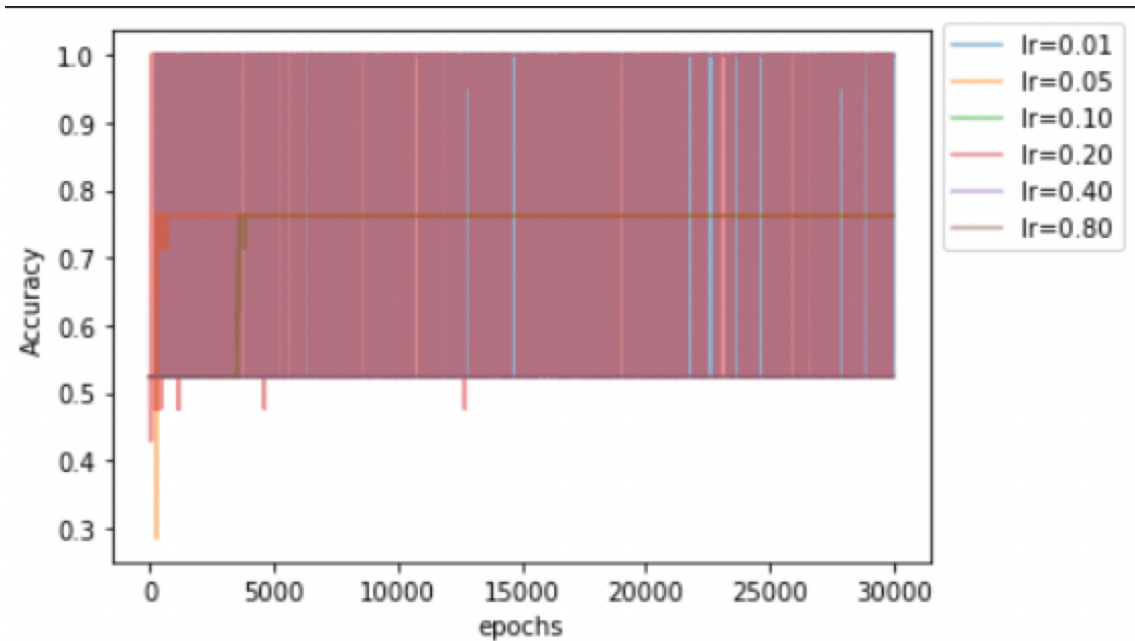
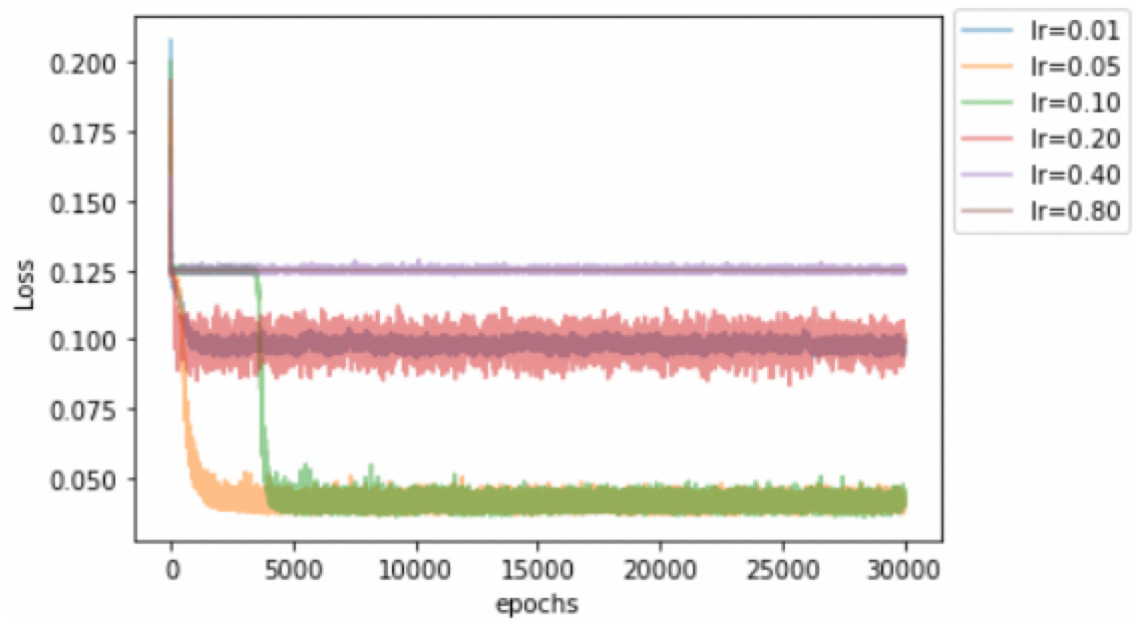
### ▼ A. Try different learning rates



以下方的表可以顯示出 Learning Rate 對訓練模型影響是非常顯著的，所以如何調整 Learning Rate 並更新參數這也是個很大的坑。

Differenet Learning Rate

Datasets	Learning Rate	Optimizer	Architectures	Epochs	Accuracy
XOR	0.01	sgdm	2-3-ReLU-3-ReLU-1-Sigmoid	30000	52%
XOR	0.05	sgdm	2-3-ReLU-3-ReLU-1-Sigmoid	30000	76%
XOR	0.1	sgdm	2-3-ReLU-3-ReLU-1-Sigmoid	30000	100%
XOR	0.2	sgdm	2-3-ReLU-3-ReLU-1-Sigmoid	30000	52%
XOR	0.4	sgdm	2-3-ReLU-3-ReLU-1-Sigmoid	30000	52%
XOR	0.8	sgdm	2-3-ReLU-3-ReLU-1-Sigmoid	30000	52%



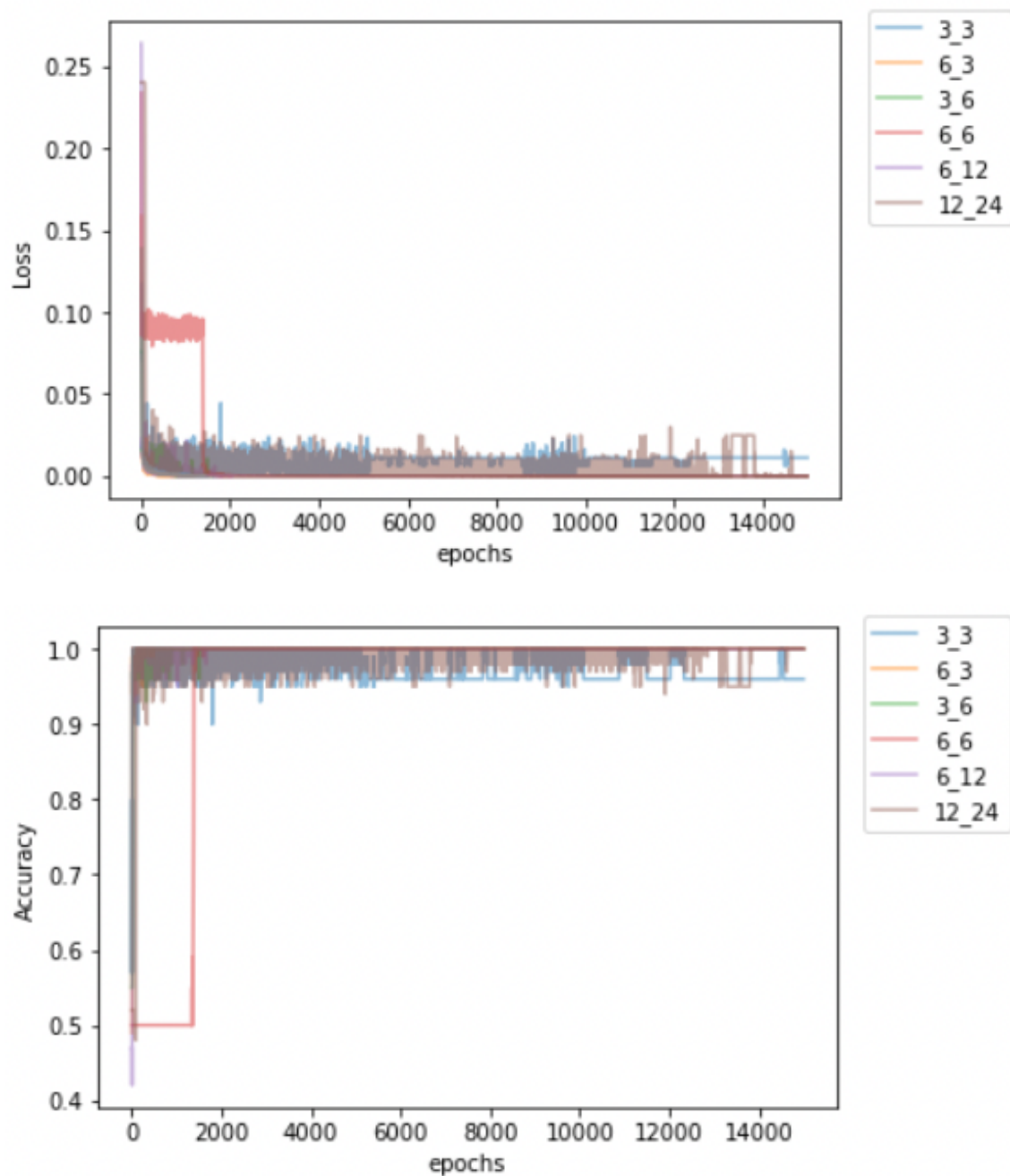
## ▼ B. Try different numbers of hidden units



以下方的圖表得知模型的參數量足夠多的情況下，就能達到不錯的效果。

Differenet Number of Hidden Units

Datasets	Learning Rate	Optimizer	Architectures	Epochs	Accuracy
Linear	0.1	sgdm	2-3-ReLU-3-ReLU-1-Sigmoid	15000	96%
Linear	0.1	sgdm	2-6-ReLU-3-ReLU-1-Sigmoid	15000	100%
Linear	0.1	sgdm	2-3-ReLU-6-ReLU-1-Sigmoid	15000	100%
Linear	0.1	sgdm	2-6-ReLU-6-ReLU-1-Sigmoid	15000	100%
Linear	0.1	sgdm	2-6-ReLU-12-ReLU-1-Sigmoid	15000	100%
Linear	0.1	sgdm	2-12-ReLU-24-ReLU-1-Sigmoid	15000	100%



### ▼ C. Try without activation functions

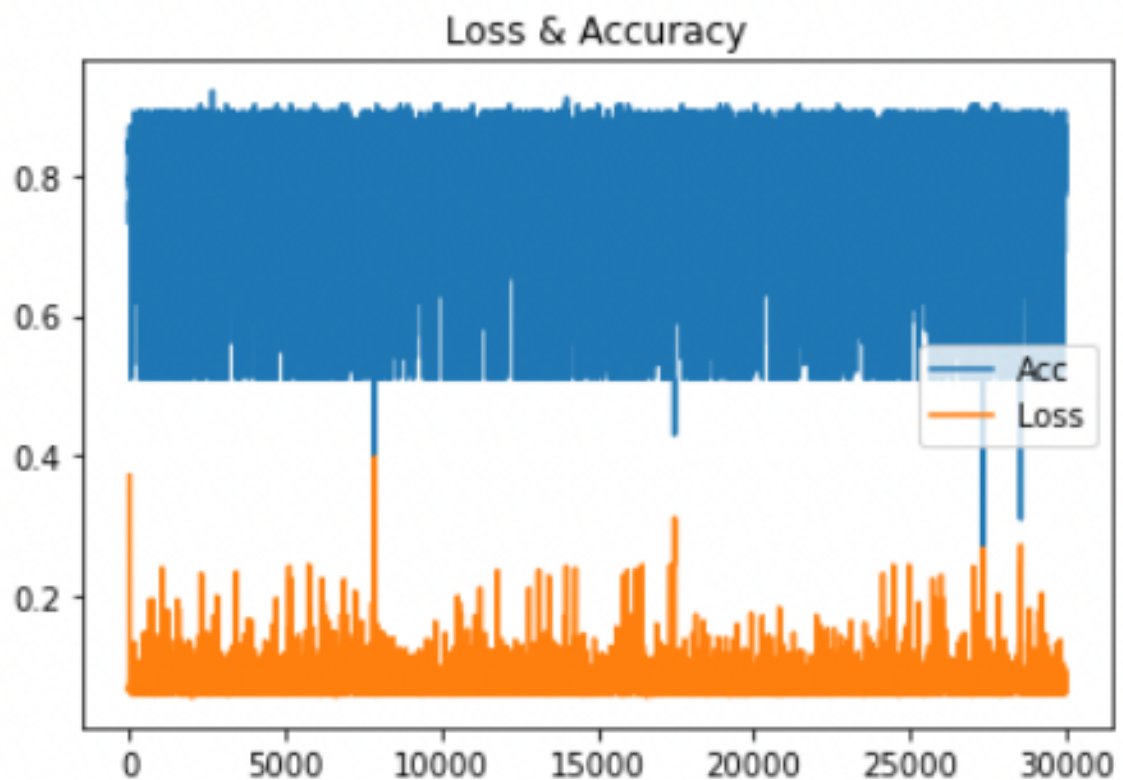
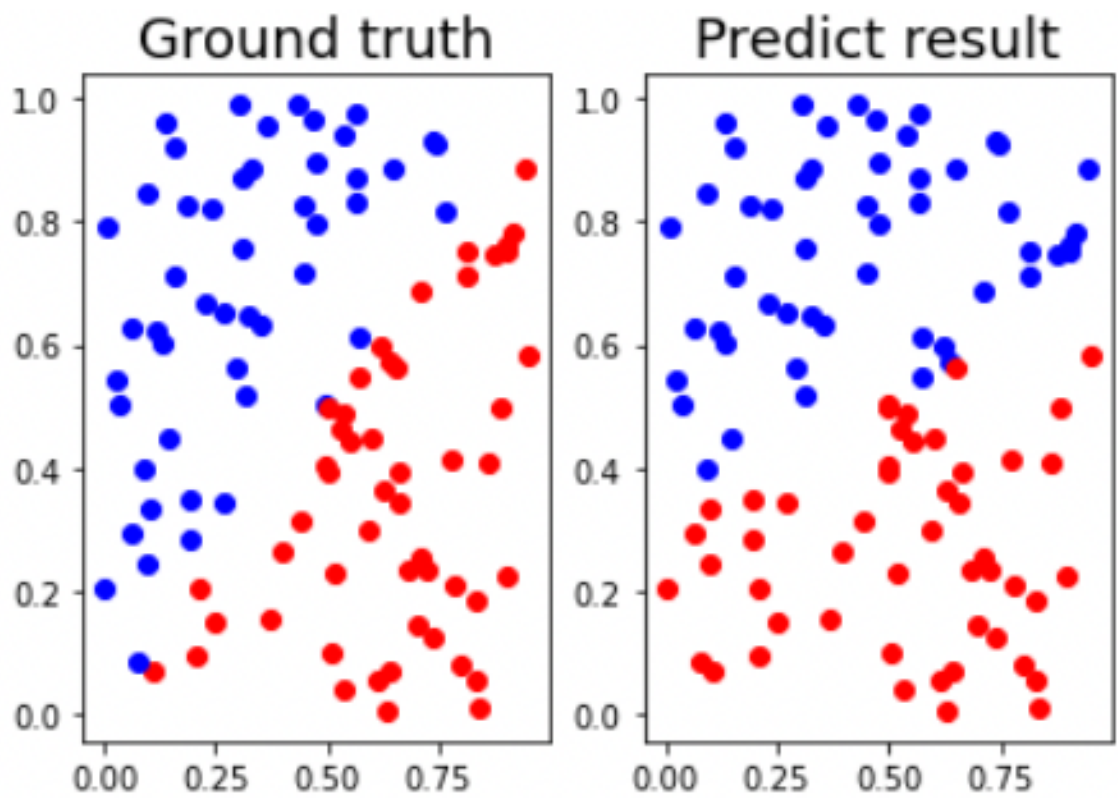


下面圖表表示出缺少 activation functions 會導致訓練困難以及準確度下降。

#### Linear Data:

Without Activation Functions Comparision

Datasets	Learning Rate	Optimizer	Architectures	Epochs	Accuracy
Linear	0.1	sgdm	2-3-3-1	30000	80%

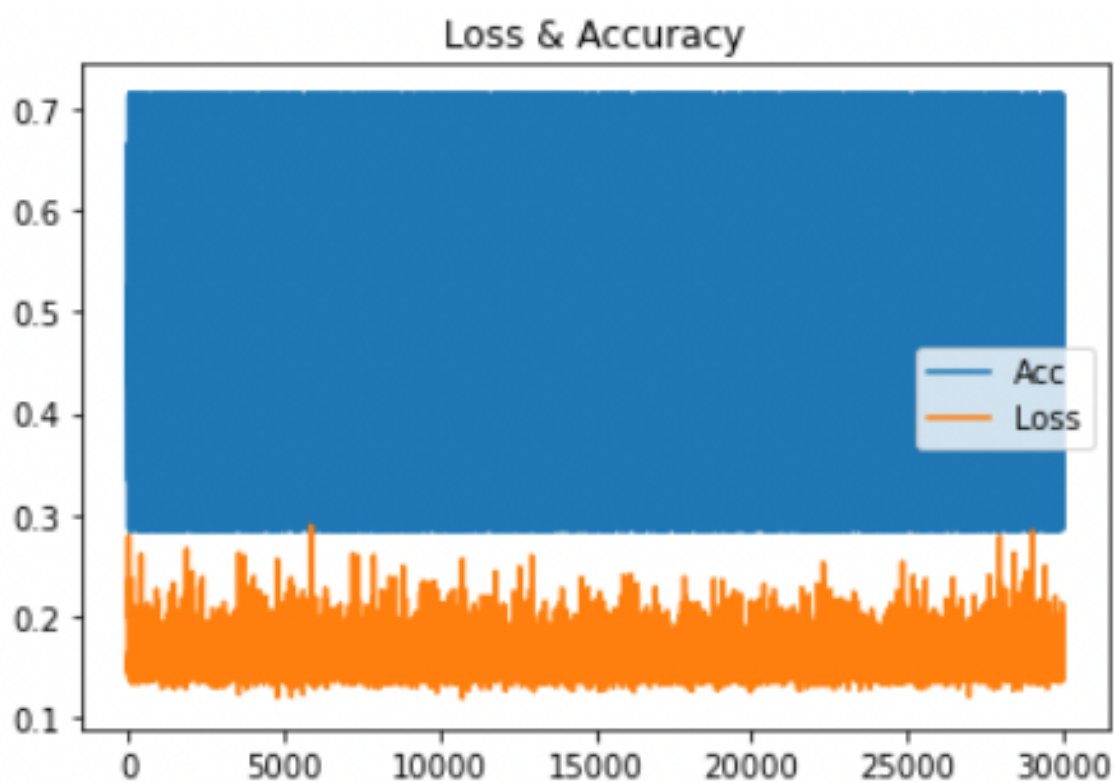
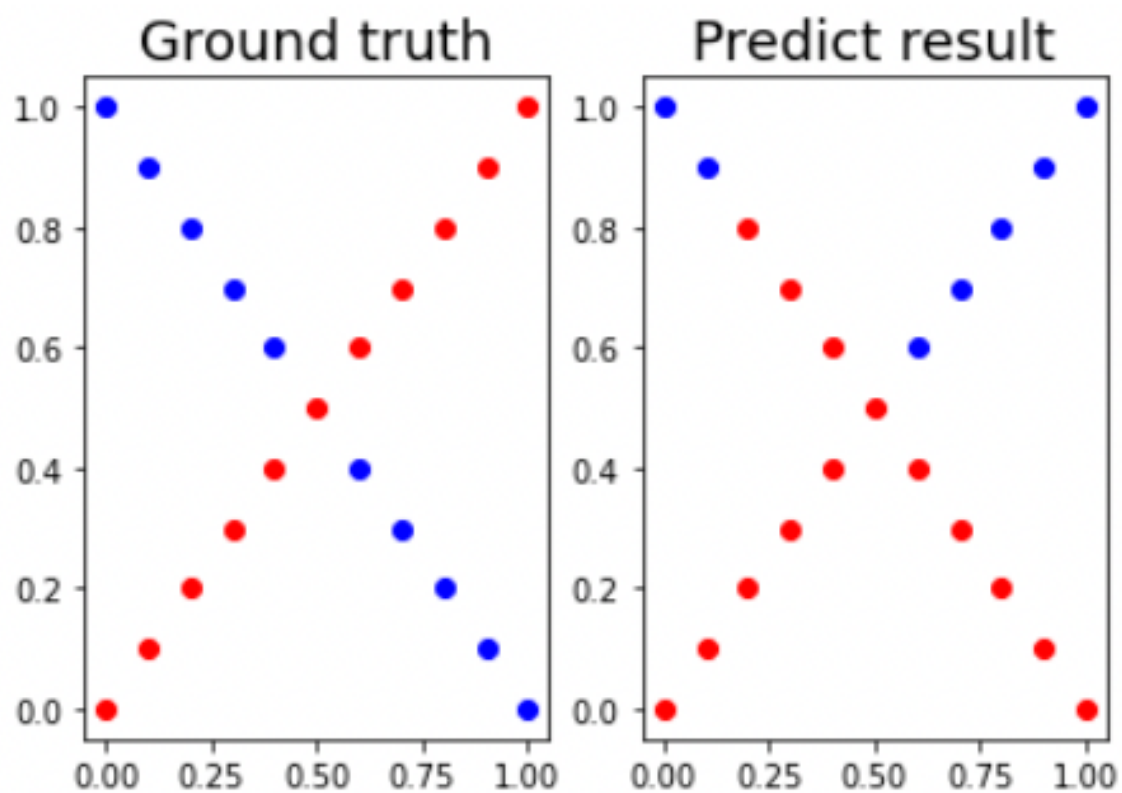


### XOR Data:

Without Activation Functions Comparision

Datasets	Learning Rate	Optimizer	Architectures	Epochs	Accuracy
XOR	0.1	sgdm	2-3-3-1	30000	38%





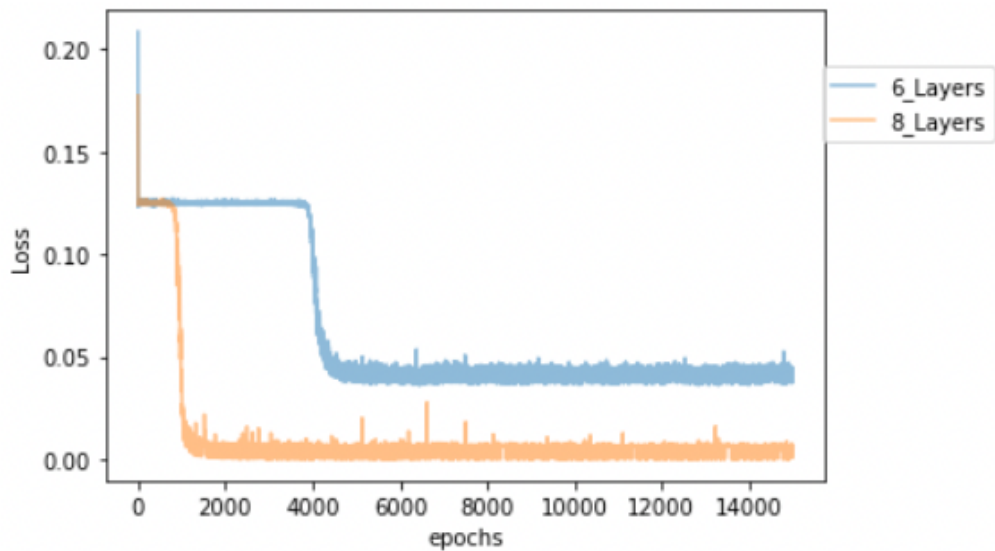
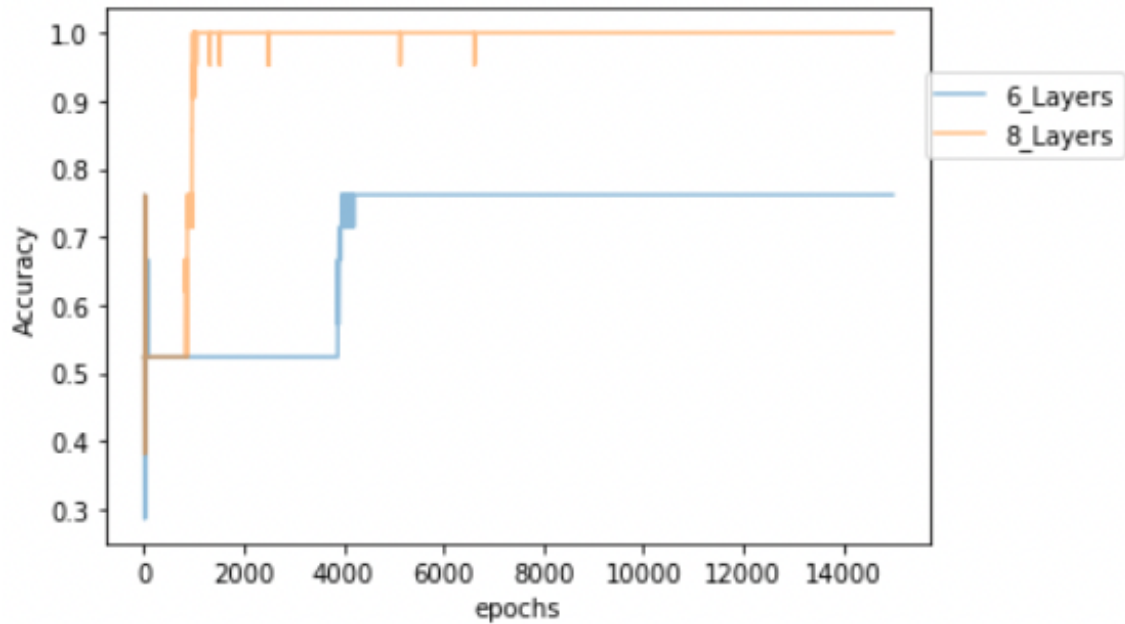
#### ▼ D. Try different number of layers



下方圖表表示，層數較多的網路可以有較佳的效果。

Different Number of Layers

Datasets	Learning Rate	Optimizer	Architectures	Epochs	Accuracy
XOR	0.1	sgdm	2-3-ReLU-3-ReLU-1-Sigmoid	15000	76%
XOR	0.1	sgdm	2-3-ReLU-3-ReLU-3-ReLU-1-Sigmoid	15000	100%



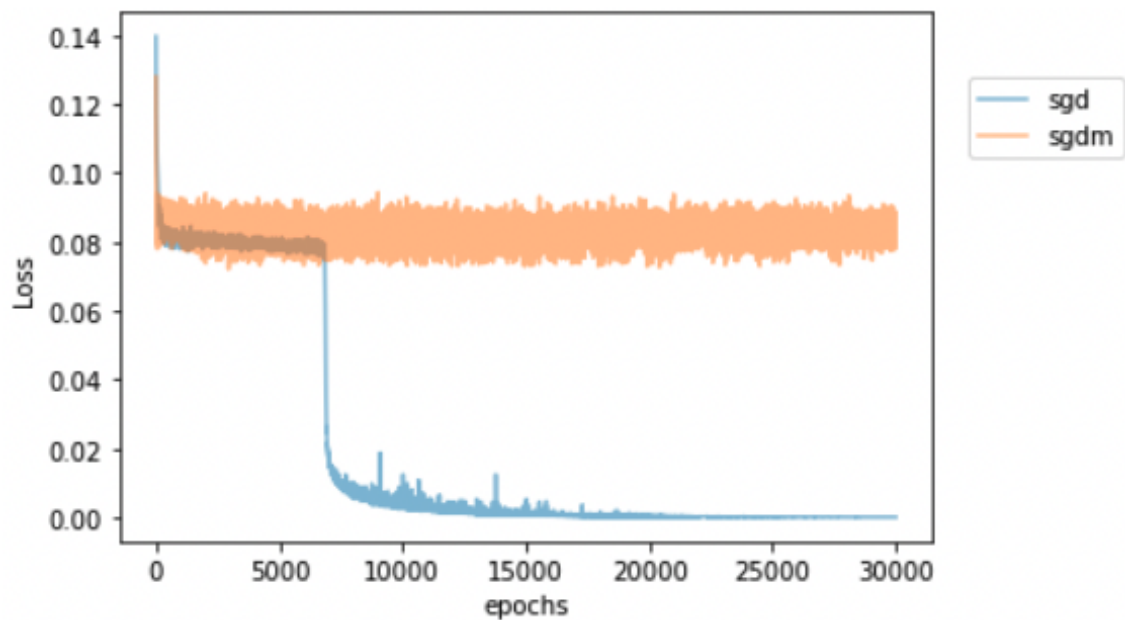
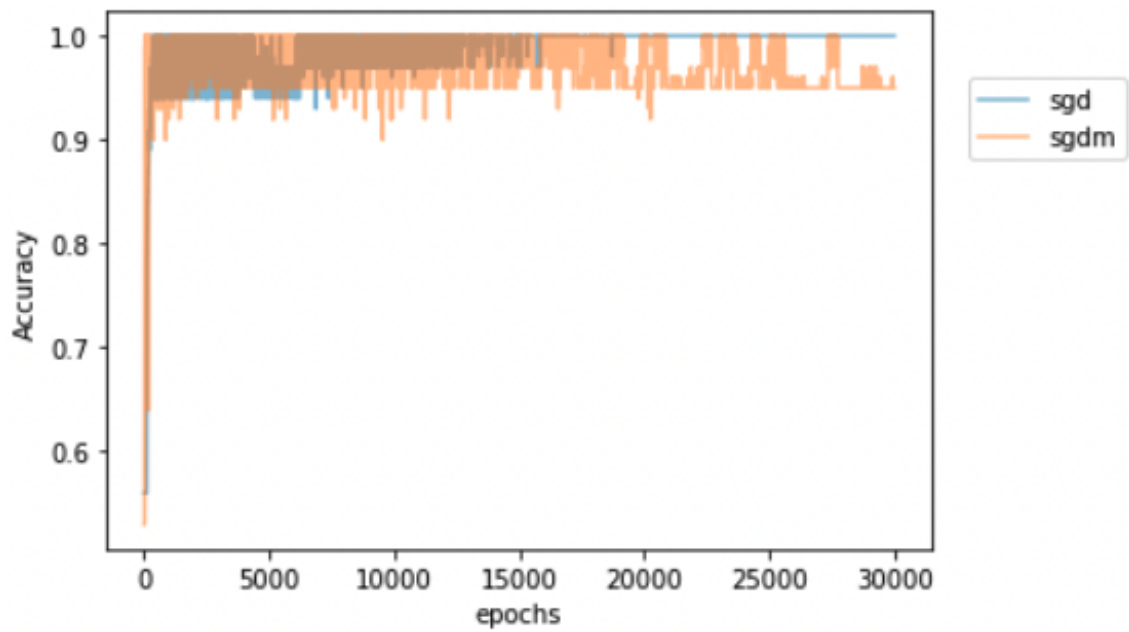
### ▼ E. Implement different optimizers



圖顯示 sgdm 的震盪幅度比 sgd 高。

Different Optimizers

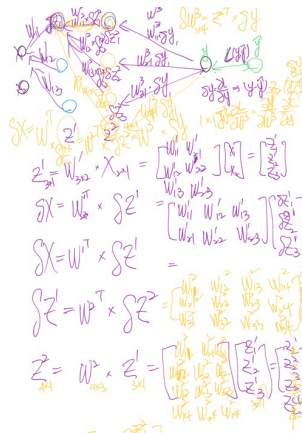
Datasets	Learning Rate	Optimizer	Architectures	Epochs	Accuracy
Linear	0.1	sgd	2-3-ReLU-3-ReLU-1-Sigmoid	30000	100%
Linear	0.1	sgdm	2-3-ReLU-3-ReLU-1-Sigmoid	30000	95%



### ▼ F. Manuscript



在實作 backpropagation 時，手算證明的過程：



$$Z'_1 = W'_{31} \times X_{11} = \begin{bmatrix} W'_{31} & W'_{32} & W'_{33} \end{bmatrix} \begin{bmatrix} X_{11} \\ X_{12} \\ X_{13} \end{bmatrix} = \begin{bmatrix} Z'_{11} \\ Z'_{12} \\ Z'_{13} \end{bmatrix}$$

$$\delta X = W'_{31} \times \delta Z'_1 = \begin{bmatrix} W'_{31} & W'_{32} & W'_{33} \end{bmatrix} \begin{bmatrix} \delta Z'_{11} \\ \delta Z'_{12} \\ \delta Z'_{13} \end{bmatrix}$$

$$\delta X = W'^T \times \delta Z'_1 = \begin{bmatrix} W'_{31} & W'_{32} & W'_{33} \end{bmatrix} \begin{bmatrix} \delta Z'_{11} \\ \delta Z'_{12} \\ \delta Z'_{13} \end{bmatrix}$$

$$\delta Z'_1 = W'^T \times \delta Z'_2 = \begin{bmatrix} W'_{31} & W'_{32} & W'_{33} \end{bmatrix} \begin{bmatrix} \delta Z'_{21} \\ \delta Z'_{22} \\ \delta Z'_{23} \end{bmatrix}$$

$$Z'_2 = W'_{31} \times Z'_1 = \begin{bmatrix} W'_{31} & W'_{32} & W'_{33} \end{bmatrix} \begin{bmatrix} Z'_{11} \\ Z'_{12} \\ Z'_{13} \end{bmatrix} = \begin{bmatrix} Z'_{21} \\ Z'_{22} \\ Z'_{23} \end{bmatrix}$$

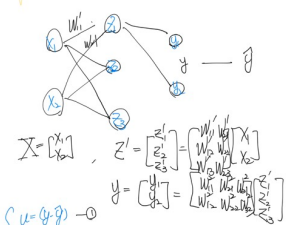
$$\delta Z'_{21} \times Z'^T_{13} = \begin{bmatrix} \delta Z'_{21} \\ \delta Z'_{22} \\ \delta Z'_{23} \end{bmatrix} \begin{bmatrix} Z'_{11} & Z'_{12} & Z'_{13} \end{bmatrix}$$

$$= \begin{bmatrix} W'_{31} & W'_{32} & W'_{33} \end{bmatrix}$$

$$\delta Z'_2 = W'^T \times \delta Z'_1$$

$$\delta W = \delta \times Z'^T$$

$$W^{t+1} = W^t + \eta \times \delta W^{t+1}$$



$$X = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix}, Z' = \begin{bmatrix} Z'_1 \\ Z'_2 \\ Z'_3 \end{bmatrix} = \begin{bmatrix} W'_{31} & W'_{32} & W'_{33} \\ W'_{21} & W'_{22} & W'_{23} \\ W'_{11} & W'_{12} & W'_{13} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix}$$

$$Y = \begin{bmatrix} Y_1 \end{bmatrix} = \begin{bmatrix} W_{31} & W_{32} & W_{33} \end{bmatrix} \begin{bmatrix} Z'_1 \\ Z'_2 \\ Z'_3 \end{bmatrix}$$

$$u = (y, y) - 0$$

$$\mathcal{L}(y, y) = \frac{1}{2} (y, y) = \frac{1}{2} u^2 - 0$$

$$\frac{\partial \mathcal{L}}{\partial u} = u = (y, y) - 0$$

$$\frac{\partial \mathcal{L}}{\partial y} = \frac{\partial \mathcal{L}}{\partial u} \frac{\partial u}{\partial y} = (y, y) - 0$$

$$\frac{\partial \mathcal{L}}{\partial Z'_2} = \frac{\partial \mathcal{L}}{\partial Z'_1} \frac{\partial Z'_1}{\partial Z'_2} = \frac{\partial \mathcal{L}}{\partial Z'_1} (y, y) = W'^T (y, y)$$

$$\delta Z'_1 = \begin{bmatrix} W'_{31} & W'_{32} & W'_{33} \\ W'_{21} & W'_{22} & W'_{23} \\ W'_{11} & W'_{12} & W'_{13} \end{bmatrix} \begin{bmatrix} \delta Y_1 \\ \delta Y_2 \\ \delta Y_3 \end{bmatrix}$$

$$\delta Z'_2 = \begin{bmatrix} W'_{31} & W'_{32} & W'_{33} \\ W'_{21} & W'_{22} & W'_{23} \\ W'_{11} & W'_{12} & W'_{13} \end{bmatrix} \begin{bmatrix} \delta Z'_1 \\ \delta Z'_2 \\ \delta Z'_3 \end{bmatrix}$$

$$\delta W = \frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}}{\partial Z'_1} \frac{\partial Z'_1}{\partial W} = \frac{\partial \mathcal{L}}{\partial Z'_1} (y, y)$$

$$\frac{\partial \mathcal{L}}{\partial W} = \begin{bmatrix} W'_{31} & W'_{32} & W'_{33} \\ W'_{21} & W'_{22} & W'_{23} \\ W'_{11} & W'_{12} & W'_{13} \end{bmatrix} \begin{bmatrix} \delta Z'_1 \\ \delta Z'_2 \\ \delta Z'_3 \end{bmatrix}$$

$$\delta W = \begin{bmatrix} \delta W_{31} & \delta W_{32} & \delta W_{33} \\ \delta W_{21} & \delta W_{22} & \delta W_{23} \\ \delta W_{11} & \delta W_{12} & \delta W_{13} \end{bmatrix}$$

$$= \begin{bmatrix} \delta Z'_1 & \delta Z'_2 & \delta Z'_3 \end{bmatrix}$$

$$u = \begin{bmatrix} \delta Z'_1 \\ \delta Z'_2 \\ \delta Z'_3 \end{bmatrix} \begin{bmatrix} Z'_1 & Z'_2 & Z'_3 \end{bmatrix} = \delta Y \cdot Z'^T$$

$$= \begin{bmatrix} \delta Z'_1 & \delta Z'_2 & \delta Z'_3 \end{bmatrix} \begin{bmatrix} Z'_1 & Z'_2 & Z'_3 \end{bmatrix}$$

$$\delta Z = W'^T \cdot \delta Y$$

$$\delta W = \delta Y \cdot Z'^T$$

$$y = \frac{1}{1 + e^{-x}}$$

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

