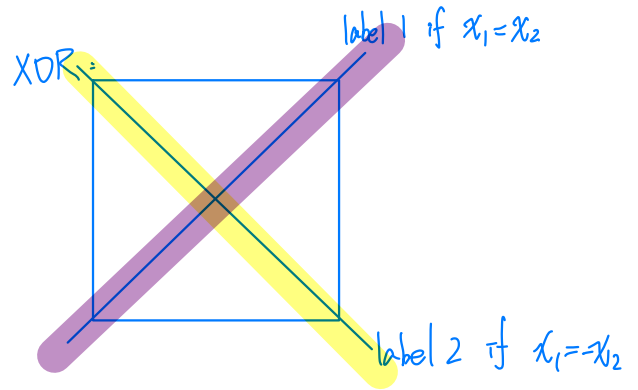
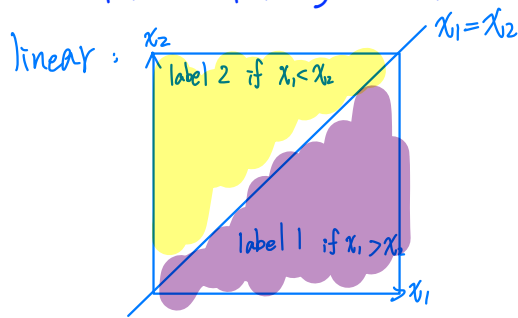


1. Introduction (20%)

Ground true data generation



Construct neural network model

suppose the amount of hidden units for first and second layer are h_1 and h_2 respectively
original version only deals with one data for each iteration, but in fact we can deal with multiple (batch size) data for each iteration

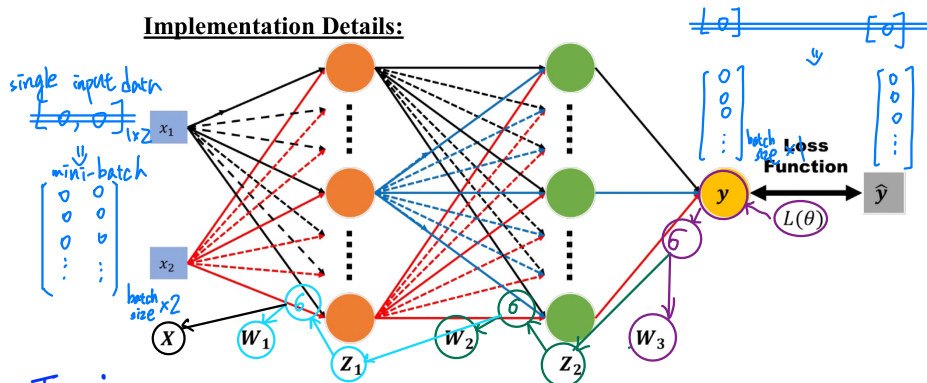


Figure 2. Forward pass

Train

randomly initialize the network weights (no bias)

$$W_1, W_2, W_3$$

while not converge (for each epoch check whether loss is smaller than epsilon or not)

the design of loss function is related to the problem and network design. it will be explained precisely in 2.B.

$$\text{Loss} = \frac{-1}{\text{batchsize}} \sum_{i \in \text{mini-batch}} (\hat{y}_i \log y_i + (1 - \hat{y}_i) \log (1 - y_i))$$

forward

$$\begin{aligned} b(xW_1) &= z_1 \\ b(z_1W_2) &= z_2 \\ b(z_2W_3) &= y \end{aligned} \quad , \text{ where } b(x) = \frac{1}{1 + e^{-x}}, \quad b'(x) = b(x)(1 - b(x))$$

backward

$$\begin{aligned} \frac{\partial L}{\partial W_3} &= \frac{\partial L}{\partial y} \otimes \frac{\partial y}{\partial (z_2 W_3)} \otimes \frac{\partial (z_2 W_3)}{\partial z_2} \otimes \frac{\partial z_2}{\partial (z_1 W_2)} \otimes \frac{\partial (z_1 W_2)}{\partial z_1} \otimes \frac{\partial z_1}{\partial (x W_1)} \otimes \frac{\partial (x W_1)}{\partial x} \\ \frac{\partial L}{\partial W_2} &= \frac{\partial L}{\partial y} \otimes \frac{\partial y}{\partial (z_2 W_3)} \otimes \frac{\partial (z_2 W_3)}{\partial z_2} \otimes \frac{\partial z_2}{\partial (z_1 W_2)} \otimes \frac{\partial (z_1 W_2)}{\partial z_1} \otimes \frac{\partial z_1}{\partial (x W_1)} \otimes \frac{\partial (x W_1)}{\partial x} \\ \frac{\partial L}{\partial W_1} &= \frac{\partial L}{\partial y} \otimes \frac{\partial y}{\partial (z_2 W_3)} \otimes \frac{\partial (z_2 W_3)}{\partial z_2} \otimes \frac{\partial z_2}{\partial (z_1 W_2)} \otimes \frac{\partial (z_1 W_2)}{\partial z_1} \otimes \frac{\partial z_1}{\partial (x W_1)} \otimes \frac{\partial (x W_1)}{\partial x} \end{aligned}$$

update network weights

$$\begin{aligned} W_1 &= W_1 - \text{learning rate} \cdot \frac{\partial L}{\partial W_1} \\ W_2 &= W_2 - \text{learning rate} \cdot \frac{\partial L}{\partial W_2} \\ W_3 &= W_3 - \text{learning rate} \cdot \frac{\partial L}{\partial W_3} \end{aligned}$$

we want to adjust model weights W_1, W_2 , and W_3 to lower down loss function, so we compute the gradient which represents the steepest direction to update them the computation details will be discussed in 2.C.

* W_1, W_2, W_3 个臭相關性, 只有 z_1, z_2, y 才臭相關性!!!

Test

forward the trained network and output the predicted y
print the accuracy

2. Experiment setups (30%):

A. Sigmoid functions

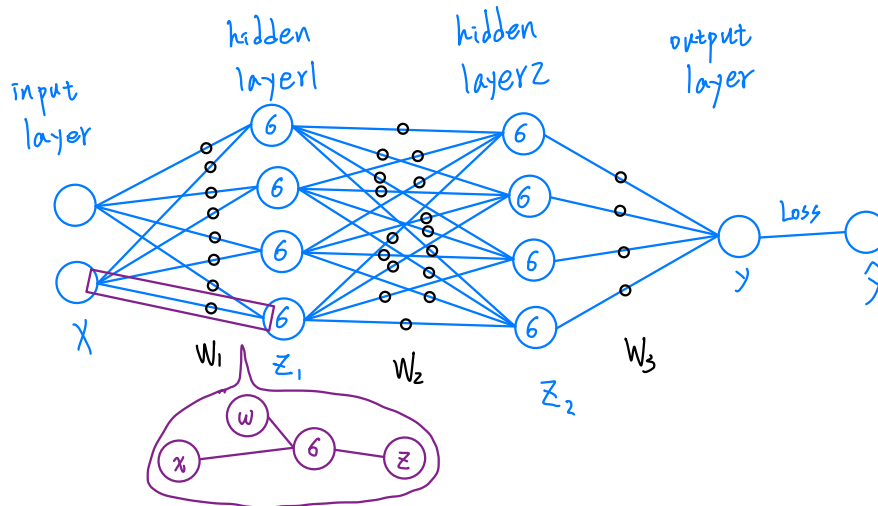
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

```
def sigmoid(M):
    return 1.0/(1.0+np.exp(-M))

def derivative_sigmoid(M):
    return sigmoid(M)*(1-sigmoid(M))
```

B. Neural network



amount of hidden units for first and second layer are $h_1=10$ and $h_2=10$ respectively

learning rate is 0.3

epsilon is 0.01

model weights W_1 , W_2 , W_3 are randomly initialized with size $(2, h_1)$, (h_1, h_2) , $(h_2, 1)$ respectively

```
#initialize model parameter
nHiddenUnits=(10,10) #amount of hidden units for each layer
learningRate = 0.3 #learning rate
epsilon = 0.01 #to judge converge or not
def __init__(self, nHiddenUnits, learningRate, epsilon):
    (h1, h2) = nHiddenUnits
    self.lr = learningRate
    self.eps = epsilon
    self.W1 = np.random.randn(2, h1)
    self.W2 = np.random.randn(h1, h2)
    self.W3 = np.random.randn(h2, 1)
```

the network forward parameters like this:

```
def forward(self, bx):
    self.inputs = bx
    self.Z1 = sigmoid(self.inputs@self.W1)
    self.Z2 = sigmoid(self.Z1@self.W2)
    pred_y = sigmoid(self.Z2@self.W3)
    return pred_y
```

$\sigma(xW_1) = Z_1$
 $\sigma(Z_1W_2) = Z_2$
 $\sigma(Z_2W_3) = y$

the loss function is defined like this:

since it's a binary classification problem and we embed the label using one hot encoding method, we can view it as logistic regression problem which use sigmoid function as activation function and use binary cross entropy as loss function

$$\text{Loss} = \frac{1}{\text{batchsize}} \sum_{i \in \text{batch}} (\hat{y}_i \log y_i + (1 - \hat{y}_i) \log (1 - y_i))$$

y has been activated by sigmoid function

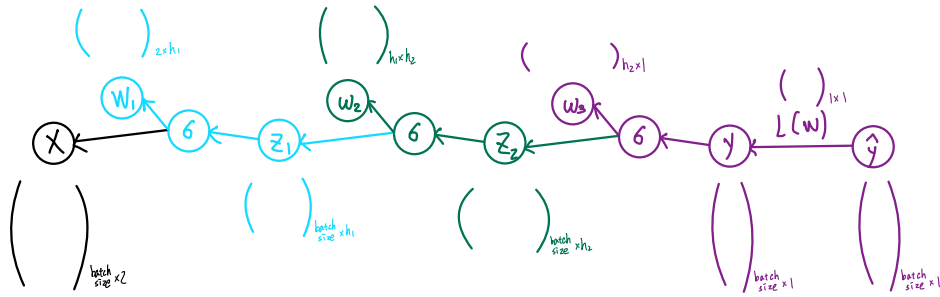
```
def loss(self, gt_y, pred_y):
    batchsize = gt_y.shape[0]
    return (-1/batchsize) * np.sum(
        gt_y*np.log(gt_y+self.eps)
        +(1-pred_y)*np.log(1-pred_y+self.eps))
```

C. Backpropagation

here shows the detailed backward computation using chain rule:

⊙: Hardamard product

@: standard matrix product



it's actually a diagonal matrix (only the diagonal has value, and the other elements are all zero), so the standard matrix product of two matrices are equal to the Hardamard product (element-wise multiplication) of two matrices

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial y} @ \frac{\partial y}{\partial (z_2 w_3)} @ \frac{\partial (z_2 w_3)}{\partial w_3}$$

$$= \left[\begin{array}{c} -1 \\ \text{batchsize} \end{array} \begin{array}{c} \frac{\partial y}{\partial y} \cdot \frac{\partial y}{\partial (z_2 w_3)} \cdot \frac{\partial (z_2 w_3)}{\partial w_3} \\ \frac{\partial y}{\partial y} \cdot \frac{\partial y}{\partial (z_2 w_3)} \cdot \frac{\partial (z_2 w_3)}{\partial w_3} \end{array} \right]^T @ \left[\begin{array}{c} \frac{\partial (z_2 w_3)}{\partial w_3} \\ \frac{\partial (z_2 w_3)}{\partial w_3} \end{array} \right]^T @ z_2$$

$$\frac{\partial L}{\partial w_2} = \left[\begin{array}{c} -1 \\ \text{batchsize} \end{array} \begin{array}{c} \frac{\partial L}{\partial y} @ \frac{\partial y}{\partial (z_2 w_3)} @ \frac{\partial (z_2 w_3)}{\partial w_3} \end{array} \right]^T @ \left[\begin{array}{c} \frac{\partial (z_2 w_3)}{\partial w_3} \\ \frac{\partial (z_2 w_3)}{\partial w_3} \end{array} \right]^T @ z_2 @ \frac{\partial (z_1 w_2)}{\partial w_2}$$

$$\frac{\partial L}{\partial w_1} = \left[\begin{array}{c} -1 \\ \text{batchsize} \end{array} \begin{array}{c} \frac{\partial L}{\partial y} @ \frac{\partial y}{\partial (z_2 w_3)} @ \frac{\partial (z_2 w_3)}{\partial w_3} @ \frac{\partial (z_1 w_2)}{\partial w_2} @ \frac{\partial (z_1 w_2)}{\partial w_1} \end{array} \right]^T @ \left[\begin{array}{c} \frac{\partial (z_1 w_2)}{\partial w_1} \\ \frac{\partial (z_1 w_2)}{\partial w_1} \end{array} \right]^T @ x$$

$$w_1 = w_1 - \text{learning rate} \cdot \frac{\partial L}{\partial w_1}$$

$$w_2 = w_2 - \text{learning rate} \cdot \frac{\partial L}{\partial w_2}$$

$$w_3 = w_3 - \text{learning rate} \cdot \frac{\partial L}{\partial w_3}$$

```
def backward(self, gt_y, pred_y):
    #backward propagation
    #dL/d{w3} = dL/dy * dy/d{Z2W3} * d{Z2W3}/d{w3}
    batchsize = gt_y.shape[0]

    grad_L_y = (-1/batchsize) * (gt_y/pred_y - (1-gt_y)/(1-pred_y))
    grad_L_Z2W3 = grad_L_y * derivative_sigmoid(self.Z2@self.W3)
    grad_L_W3 = (grad_L_Z2W3.T @ self.Z2).T

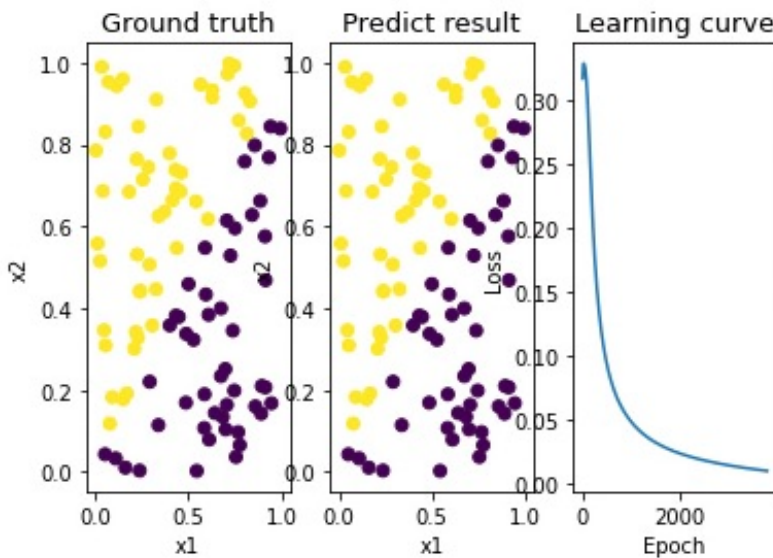
    #dL/d{w2} = dL/d{Z2W3} * d{Z2W3}/d{Z2} * d{Z2}/d{Z1W2} * d{Z1W2}/d{w2}
    grad_L_Z2 = grad_L_Z2W3 @ self.W3.T
    grad_L_Z1W2 = grad_L_Z2 * derivative_sigmoid(self.Z1@self.W2)
    grad_L_W2 = (grad_L_Z1W2.T @ self.Z1).T

    #dL/d{w1} = dL/d{Z1W2} * d{Z1W2}/d{Z1} * d{Z1}/d{XW1} * d{XW1}/d{w1}
    grad_L_Z1 = grad_L_Z1W2 @ self.W2.T
    grad_L_XW1 = grad_L_Z1 * derivative_sigmoid(self.inputs@self.W1)
    grad_L_W1 = (grad_L_XW1.T @ self.inputs).T

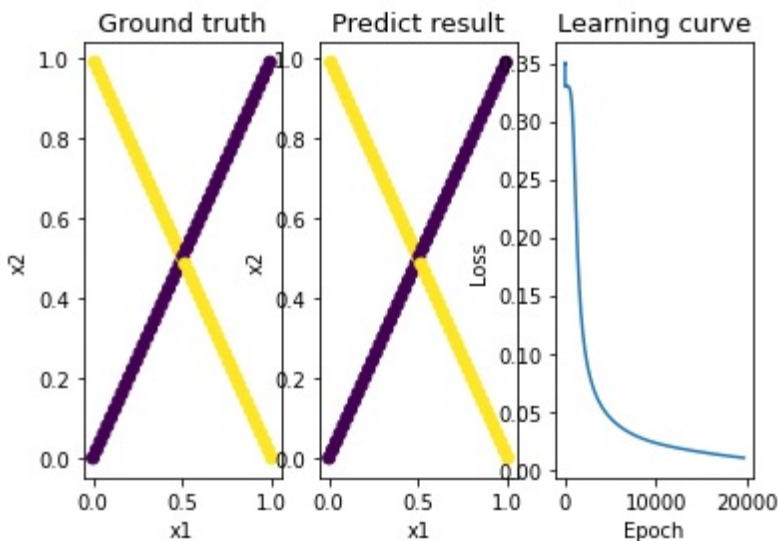
    #update model weights
    self.W1 = self.W1 - self.lr*grad_L_W1
    self.W2 = self.W2 - self.lr*grad_L_W2
    self.W3 = self.W3 - self.lr*grad_L_W3
```

3. Result of your testing (20%)

- A. Screenshot and comparison figure
- B. Show the accuracy of your prediction
- C. Learning curve (loss, epoch curve)



```
=====data type : linear=====
training ... epoch:500, loss:0.09334, acc:1.00
training ... epoch:1000, loss:0.04846, acc:1.00
training ... epoch:1500, loss:0.03231, acc:1.00
training ... epoch:2000, loss:0.02372, acc:1.00
training ... epoch:2500, loss:0.01828, acc:1.00
training ... epoch:3000, loss:0.01443, acc:1.00
training ... epoch:3500, loss:0.01148, acc:1.00
testing ... accuracy:1.0
```



```
=====data type : XOR=====
training ... epoch:500, loss:0.32837, acc:0.80
training ... epoch:1000, loss:0.27320, acc:0.87
training ... epoch:1500, loss:0.16714, acc:0.92
training ... epoch:2000, loss:0.11478, acc:0.93
training ... epoch:2500, loss:0.08860, acc:0.95
training ... epoch:3000, loss:0.07297, acc:0.95
training ... epoch:3500, loss:0.06246, acc:0.96
training ... epoch:4000, loss:0.05483, acc:0.96
training ... epoch:4500, loss:0.04899, acc:0.97
training ... epoch:5000, loss:0.04435, acc:0.97
training ... epoch:5500, loss:0.04054, acc:0.97
training ... epoch:6000, loss:0.03734, acc:0.97
training ... epoch:6500, loss:0.03460, acc:0.97
training ... epoch:7000, loss:0.03223, acc:0.97
training ... epoch:7500, loss:0.03015, acc:0.97
training ... epoch:8000, loss:0.02833, acc:0.97
training ... epoch:8500, loss:0.02671, acc:0.97
training ... epoch:9000, loss:0.02527, acc:0.98
training ... epoch:9500, loss:0.02398, acc:0.98
:
training ... epoch:16000, loss:0.01376, acc:0.99
training ... epoch:16500, loss:0.01319, acc:0.99
training ... epoch:17000, loss:0.01264, acc:0.99
training ... epoch:17500, loss:0.01211, acc:0.99
training ... epoch:18000, loss:0.01160, acc:0.99
training ... epoch:18500, loss:0.01111, acc:0.99
training ... epoch:19000, loss:0.01064, acc:0.99
training ... epoch:19500, loss:0.01020, acc:0.99
testing ... accuracy:0.99
```

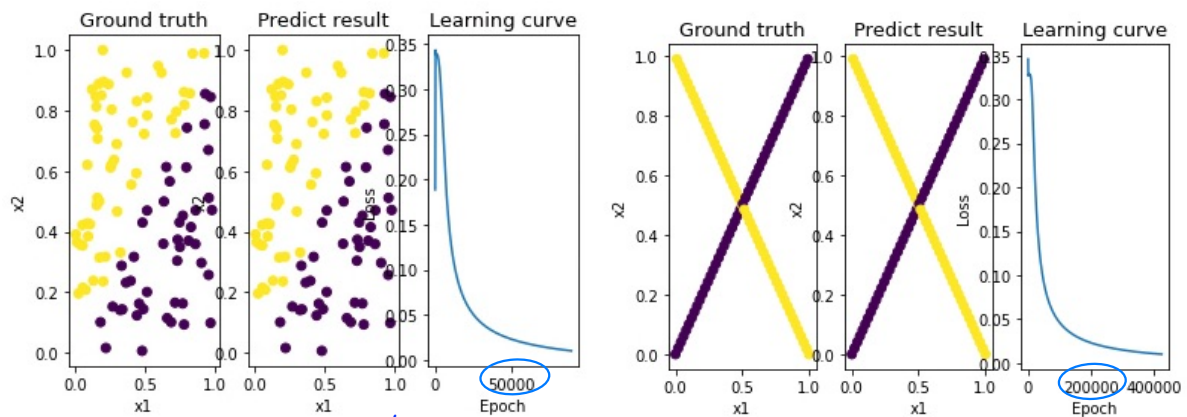
D. anything you want to present

- in the backward propagation, because some matrices are actually diagonal matrices, the result of multiplying the diagonal matrices with other matrix is actually the same as directly multiplying two matrix's elements one by one. hence, in the implementation, I just use * instead of @.
- notice that the model weight of many neural network graphs are represented by lines, but the weights are actually also nodes. hence when we calculate backward propagation by hand, we need to draw the weights with nodes instead of lines, otherwise we can't know exactly where the chain rule does.

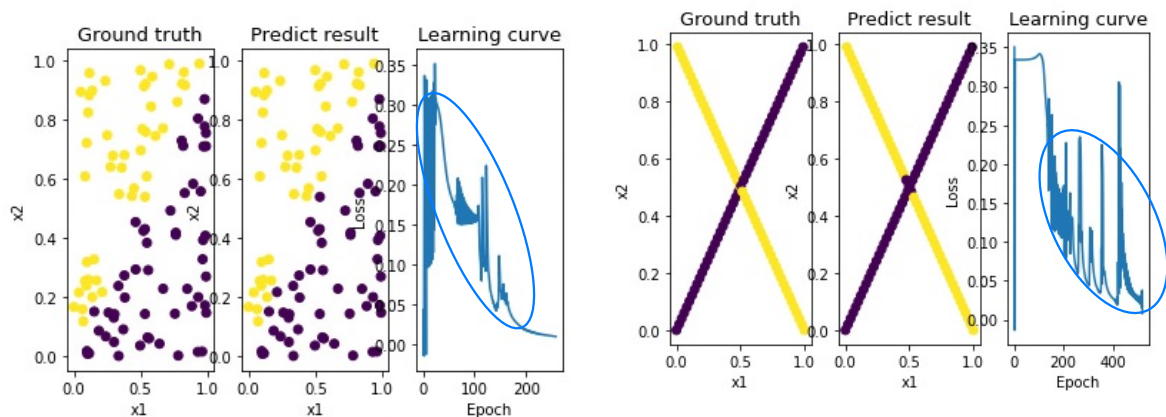
4. Discussion (30%)

A. Try different learning rates

learning rate = 0.01, the network requires more epochs to converge

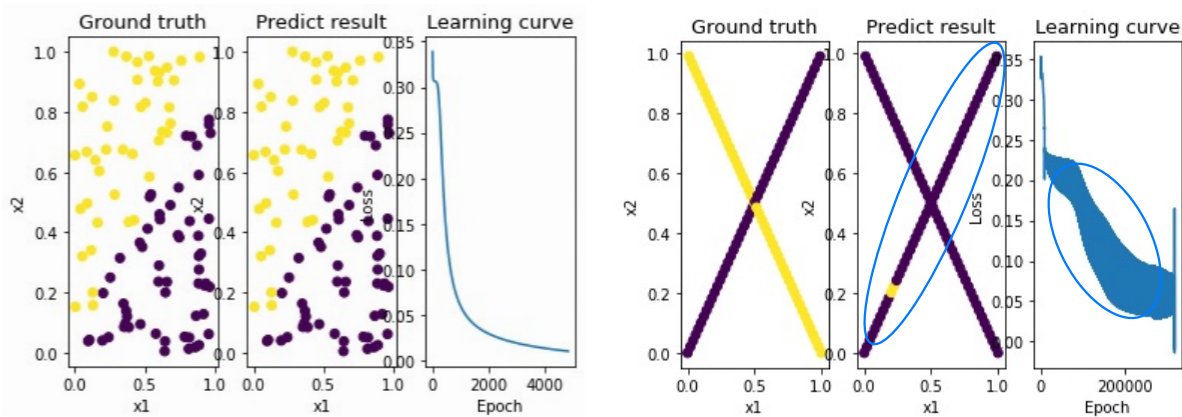


learning rate = 10, the network's convergence process is not stable

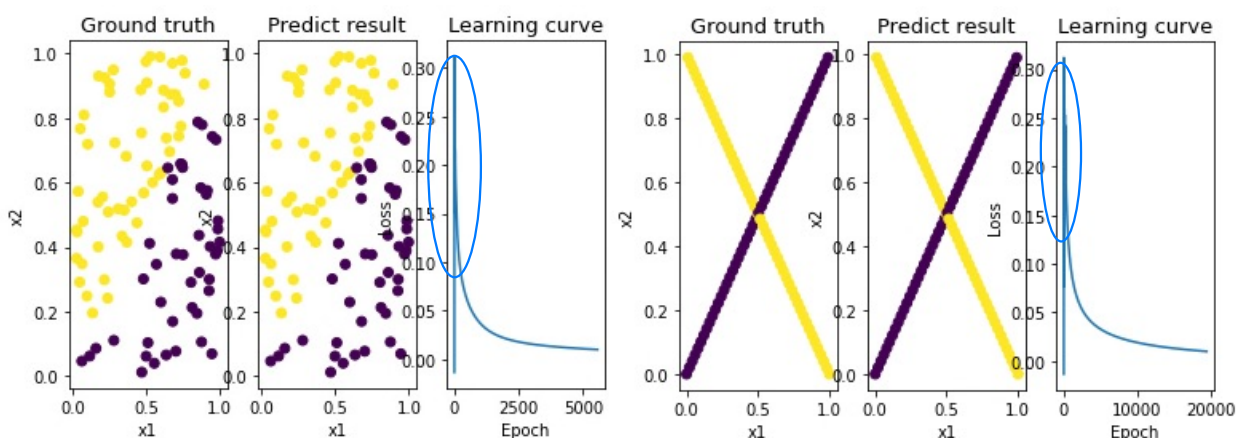


B. Try different numbers of hidden units

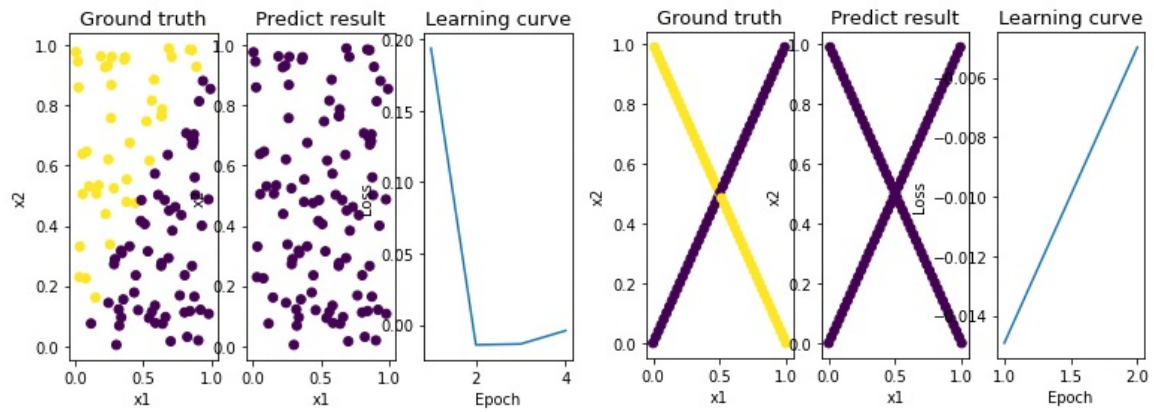
of hidden units for first and second hidden layer = (3,3), the network is not stable and may not predict well



of hidden units for first and second hidden layer = (70,70), the loss decrease fast in the beginning



of hidden units for first and second hidden layer = (1000,1000) , the network predict nothing



C. Try without activation function

the training process won't converge

```
training ... epoch:500, loss:-34363.45243, acc:0.50
training ... epoch:1000, loss:-104405.06171, acc:0.50
training ... epoch:1500, loss:-201080.83092, acc:0.50
training ... epoch:2000, loss:-320438.23154, acc:0.50
training ... epoch:2500, loss:-460040.43090, acc:0.50
training ... epoch:3000, loss:-618173.92873, acc:0.50
training ... epoch:3500, loss:-793542.07872, acc:0.50
training ... epoch:4000, loss:-985115.95048, acc:0.50
training ... epoch:4500, loss:-1192051.30267, acc:0.50
training ... epoch:5000, loss:-1413637.99221, acc:0.50
training ... epoch:5500, loss:-1649267.06546, acc:0.50
```

