

DLP\_LAB2\_310605001\_王盈驊

tags: deep learning

# Introduction

---

This work will implement a two hidden layers neural network with forwarding pass and back propagation only use Numpy and other standard libraries, the deep learning framework is not allowed to use in this homework. In the part of forward pass, we use sigmoid function to be our active function, In part of back propagation, we use chain rule and gradient descent to complete the work. The details will be introduced below.

## Experiment setups

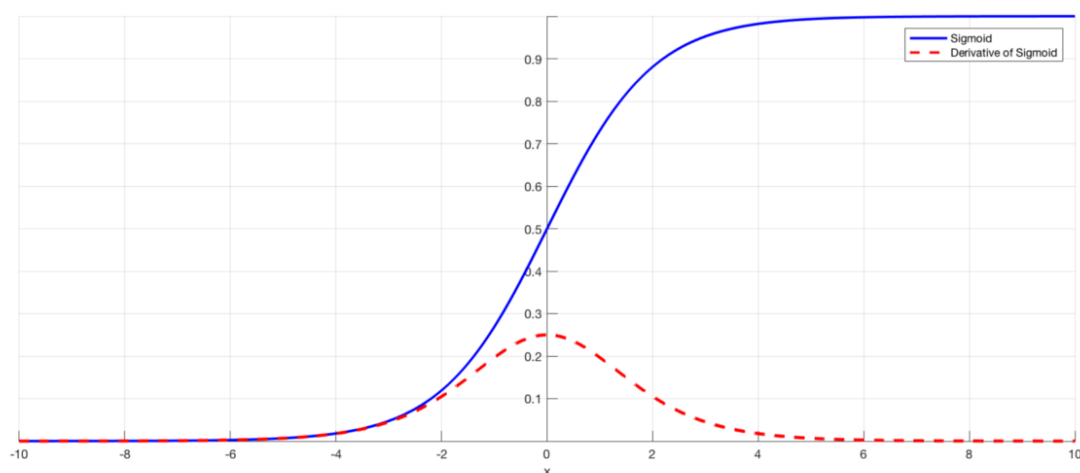
---

### A. Sigmoid functions

The sigmoid functions as activate function on neural network, this function can conquer nonlinear problems, likes XOR. The derivation formula and graph of sigmoid are shown below.

Since all equations needed for the derivative of the Sigmoid function are already found during the feedforward step it saves us a ton of computations, and this is one of the benefits of using the Sigmoid function as an activation function in the layers of a neural network.





### B. Neural network

## Implementation Details:

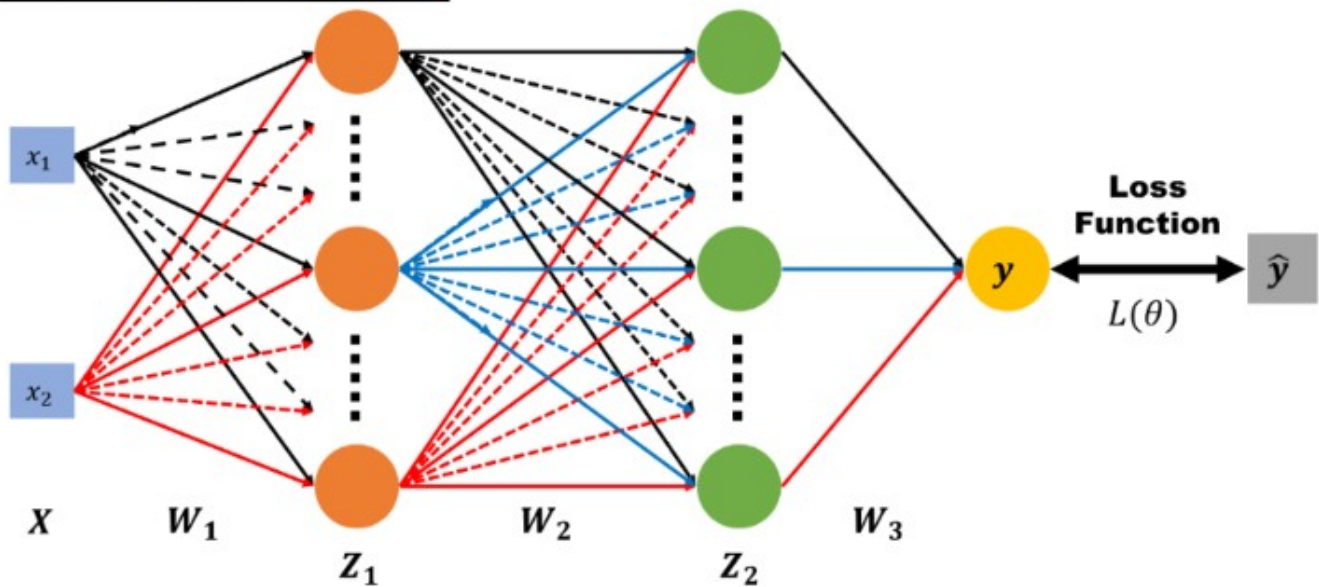


Figure 2. Forward pass

- In the figure 2, we use the following definitions for the notations:
  - $x_1, x_2$  : *nerual network inputs*
  - $X : [x_1, x_2]$
  - $y$  : *nerual network outputs*
  - $\hat{y}$  : *ground truth*
  - $L(\theta)$  : *loss function*
  - $W_1, W_2, W_3$  : *weight matrix of network layers*
- Here are the computations represented:
 
$$Z_1 = \sigma(XW_1) \quad Z_2 = \sigma(Z_1W_2) \quad y = \sigma(Z_2W_3)$$
- In the equations, the  $\sigma$  is sigmoid function that refers to the special case of the **logistic** function and defined by the formula:
 
$$\sigma(x) = \frac{1}{1+e^{-x}}$$

Neural networks are set of algorithms inspired by the functioning of human brian.

- ==Input units==** : The activity of the input units represents the raw information that is fed into the network. this also called input layer.
- ==Hidden units==** : The activity of each hidden unit is determined by the activities of the input units and the weights on the connections between the input and the hidden units. this also called hidden layer.
- ==Output units==** : The behaviour of the output units depends on the activity of the hidden units and the weights between the hidden and output units. this also called output layer.
- ==Active function==** : The purpose of the activation function is to introduce non-linearity into the output of a neuron. This is important because most real world data is nonlinear. Every activation

function (or non-linearity) takes a single number and performs a certain fixed mathematical operation on it. There are several activation functions you may encounter in practice: Sigmoid, tanh, ReLU...

## C. Backpropagation

- Loss function: In this work, we use Mean-Square Error as loss function.
- gradient descent: There's an important parameter  $\eta$  which scales the gradient and thus controls the step size. In machine learning, it is called learning rate and have a strong influence on performance

$$p_{n+1} = p_n - \eta \nabla f(p_n)$$

- chain rule: We use chain rule to get gradient of every weight and scale with learning rate to modify weight.

$$L(\hat{\theta}) = \sum_{n=1}^N (\theta - \hat{\theta})^2$$

$$z_1 = \sigma(xw_1)$$

$$z_2 = \sigma(z_1w_2)$$

$$y = \sigma(z_2w_3)$$

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial y} \frac{\partial y}{\partial w_3} = z_2^T \cdot \sigma'(y) \circ 2(y - \hat{y})$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial y} \frac{\partial y}{\partial z_2} \frac{\partial z_2}{\partial w_2} = z_1^T \cdot \sigma'(z_2) \circ (\sigma'(y) \circ 2(y - \hat{y}) \cdot w_3^T)$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial y} \frac{\partial y}{\partial z_2} \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1} = x^T \cdot \sigma'(z_1) \circ ((\sigma'(z_2) \circ (\sigma'(y) \circ 2(y - \hat{y}) \cdot w_3^T)) \cdot w_2^T)$$

operator:

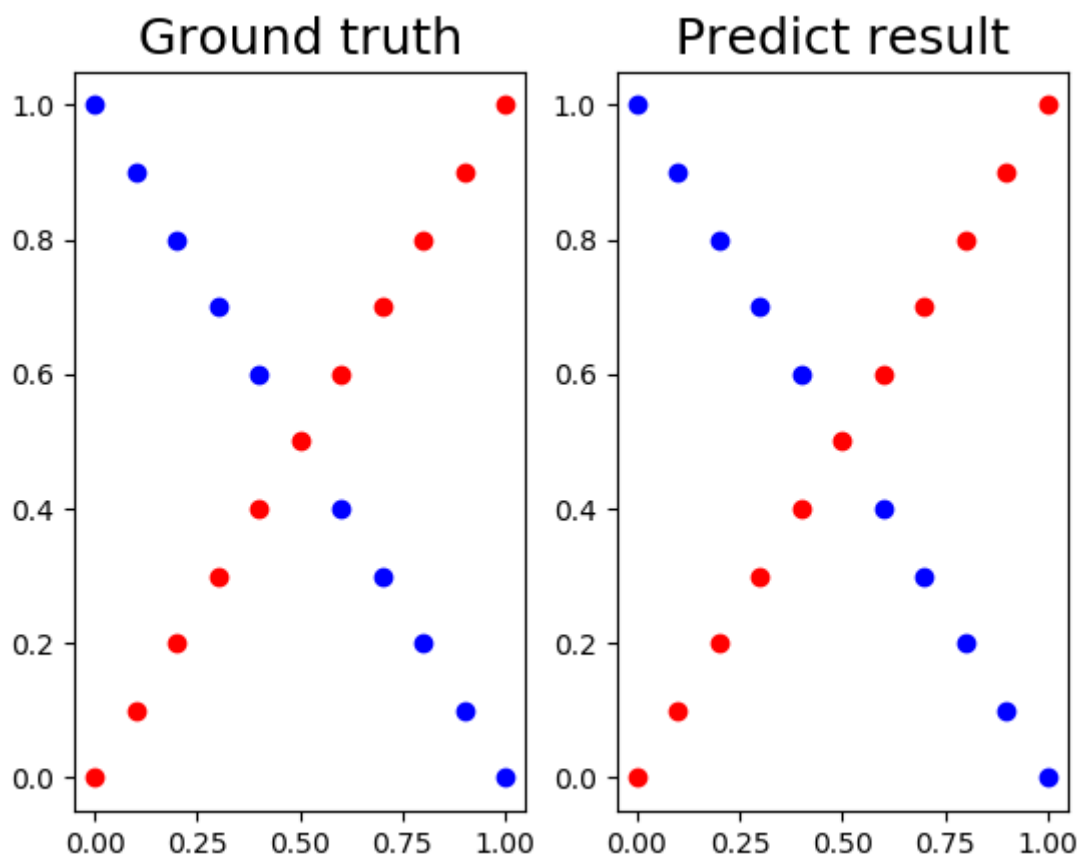
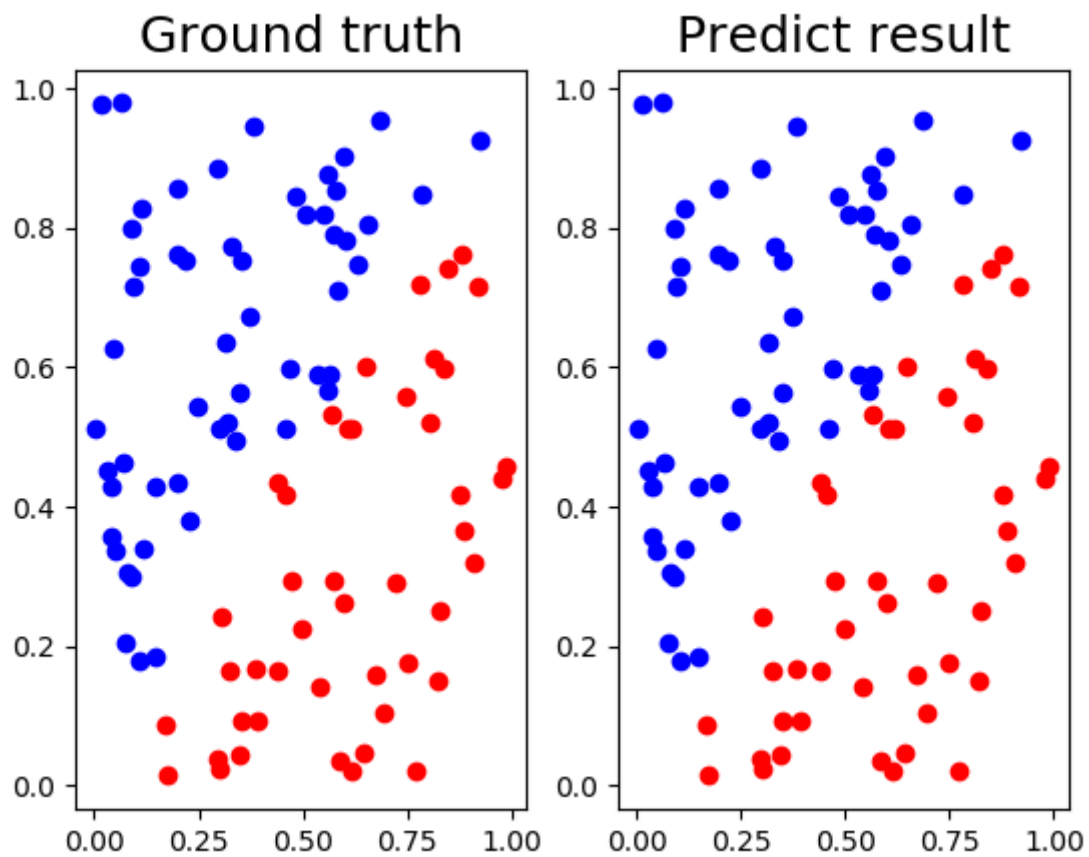
①  $\cdot$  : matrix multiplication

②  $\circ$  : Hadamard product

## Results of your testing

### A. Screenshot and comparison figure

## 1. linear



## 2. XOR

B. Show the accuracy of your prediction

1. ==linear data== : accuracy 100%

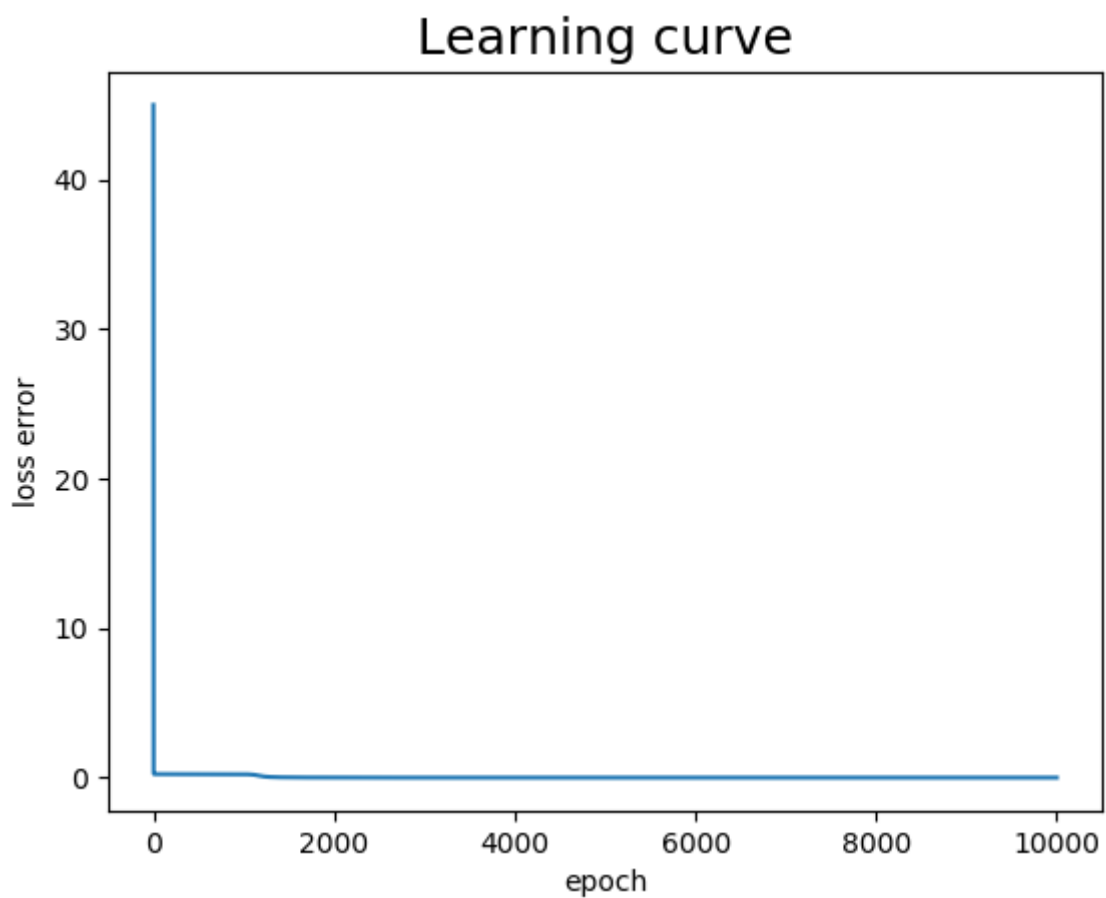
```
[4.20966184e-04] [8.09919930e-04]
[9.68027706e-01] [8.08096737e-02]
[2.71107417e-03] [1.58173267e-03]
[9.99193583e-01] [2.46851349e-03]
[6.74276417e-02] [1.27538829e-01]
[9.96636033e-01] [6.18574088e-01]
[9.99075584e-01] [9.76597203e-03]
[9.21743490e-03] [9.99492281e-01]
[9.08906174e-04] [9.82225973e-01]
[9.49627918e-03] [6.26698438e-04]
[1.02651554e-03] [9.91694389e-01]
[9.97336149e-01] [9.96804513e-01]
[4.63382765e-01] [9.98547391e-01]
[5.99086182e-02] [2.00262732e-01]
[2.42264826e-01] [5.91691956e-04]
[2.78378499e-03] [9.93007831e-01]
[5.44382982e-04] [9.92860572e-01]
[9.91851783e-01] [6.54351815e-04]
[9.99384817e-01] [9.96776675e-01]
[9.99230519e-01] [8.56857263e-01]
[9.96661184e-01] [6.52332705e-04]
[9.99328233e-01] [9.99133291e-01]
[9.98611341e-01] [9.99310557e-01]
[9.99232142e-01] [9.72525453e-01]
[1.30798653e-02] [6.96193499e-01]
[2.02684046e-03] [5.31698600e-04]
[4.67519450e-04] [9.97103430e-01]
[9.99178956e-01] [9.80551291e-01]
[6.36651411e-04] [6.90777928e-01]
[5.60921235e-04] [3.17915646e-01]
[9.97515432e-01] [9.99509869e-01]
[8.63628907e-01] [9.94084646e-01]
[9.97544232e-01] [9.97691154e-01]
[9.93456934e-01] [9.72713329e-01]
[9.98834180e-01] [9.96619955e-01]
[9.98489086e-01] [5.40358760e-04]
[9.97720497e-01] [2.65636187e-03]
[1.92523571e-03] [9.97611173e-01]
[5.94928797e-02] [9.99382274e-01]
[8.09919930e-04] [7.27285492e-01]
[9.94686492e-01]
[2.95974074e-03]
[9.98720760e-01]
[1.09587843e-02]
[5.24824491e-03]
[4.99923156e-03]
[9.15298167e-01]
[9.89625874e-01]
[5.33469057e-04]
[5.59577461e-04]
[8.49607064e-01]
[9.98266226e-01]
[1.17407897e-02]
[7.47336196e-02]
[2.48895429e-01]
[9.93207328e-01]
[9.99290572e-01]
[9.62705043e-01]
[2.61466962e-03]
[9.83313013e-01]
[3.04702040e-03]
accuracy : 100.0%, loss : 0.009996
```

2. ==XOR data== : accuracy 100%

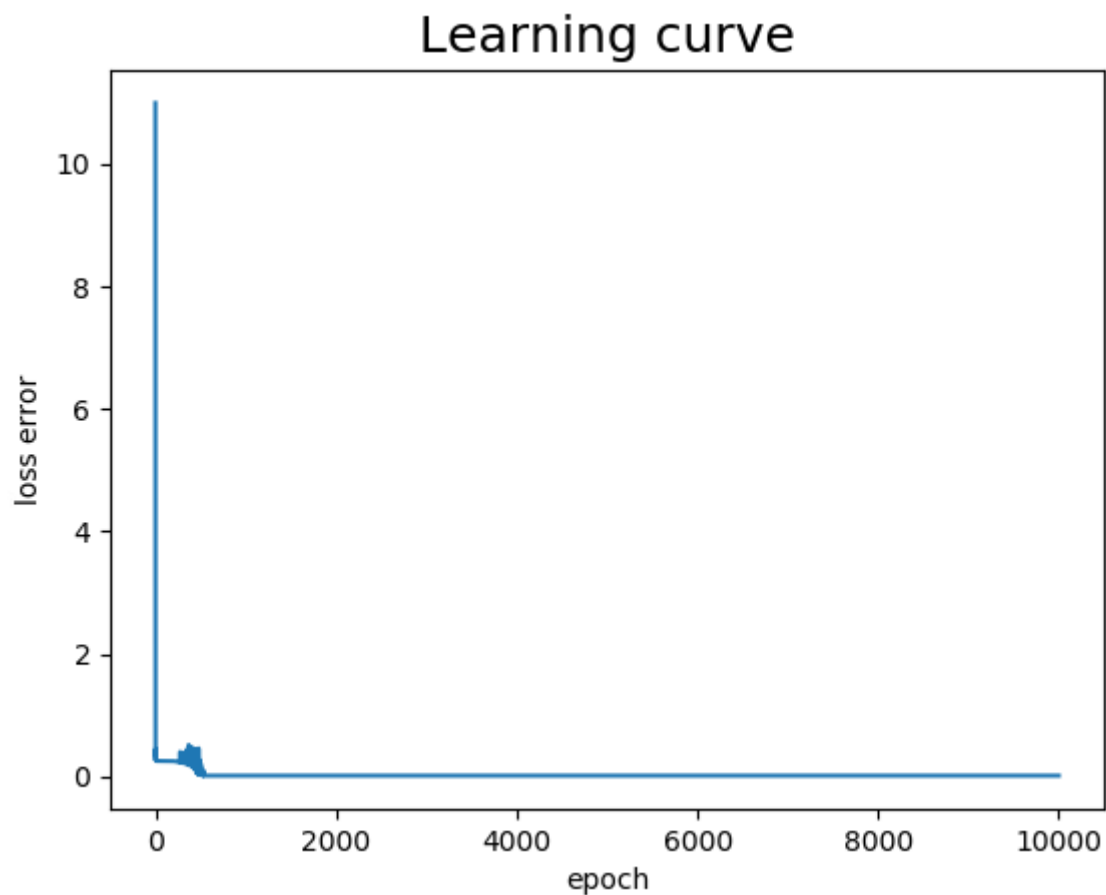
```
[ [0.00436719]  
[0.97280763]  
[0.02090975]  
[0.97232855]  
[0.06099143]  
[0.97065677]  
[0.09449886]  
[0.95885126]  
[0.09270912]  
[0.68874875]  
[0.06899492]  
[0.04457432]  
[0.74821348]  
[0.02768852]  
[0.98544072]  
[0.01769212]  
[0.99127726]  
[0.01202818]  
[0.99174142]  
[0.00878233]  
[0.99162215]]  
accuracy : 100.0%, loss : 0.009262
```

## C. Learning curve (loss, epoch curve)

### 1. linear data



## 2. XOR data



## D. anything you want to present

## 1. parameter of linear data

- hidden\_layer1\_size = 5
- hidden\_layer2\_size = 5
- epochs = 10000
- learning\_rate = 0.01

## 2. parameter of XOR data

- hidden\_layer1\_size = 5
- hidden\_layer2\_size = 5
- epochs = 10000
- learning\_rate = 0.8

## Discussion

---

## A. Try different learning rates

## 1. linear data

learning rate	0.1	0.01	0.003
accuracy	99%	100%	100%

<b>learning rate</b>	<b>0.1</b>	<b>0.01</b>	<b>0.003</b>
convergence rate	fast	fast	slow

2. XOR data

<b>learning rate</b>	<b>0.8</b>	<b>0.01</b>
accuracy	99%	52.4%

## B. Try different numbers of hidden units

1. linear data

<b>hidden units</b>	<b>(2,2)</b>	<b>(5,5)</b>
accuracy	66.7%	100%
convergence rate	slow	fast

2. XOR data

<b>hidden units</b>	<b>(2,2)</b>	<b>(5,5)</b>
accuracy	66.7%	100%

## C. Try without activation functions

1. linear data

<b>activation functions</b>	<b>no</b>	<b>yes</b>
accuracy	45.0%	100%

2. XOR data

<b>activation functions</b>	<b>no</b>	<b>yes</b>
accuracy	52.0%	100%

## D. Anything you want to share

## Extra

---

Implement different optimizers.

Implement different activation functions.

Implement convolutional layers.