

```
import torch

import torch.nn as nn

import torch.optim as optim

import numpy as np
```

#0

```
# Generate some toy data for the neural network
# We'll use a simple XOR problem for demonstration
X = torch.tensor(data: [[0, 0], [0, 1], [1, 0], [1, 1]], dtype=torch.float32)
y = torch.tensor(data: [[0], [1], [1], [0]], dtype=torch.float32)
```

#1

```
# Define the neural network architecture
2 usages
class XORNeuralNetwork(nn.Module): # Defining the class, format: ClassName(ParentClassName)
    def __init__(self, input_dim, hidden_dim, output_dim): # Create a class instance
        super(XORNeuralNetwork, self).__init__() # Initializing parent class inheritance, same as class ClassName : public ParentClass in C++
        self.hidden_layer = nn.Linear(input_dim, hidden_dim)
        self.output_layer = nn.Linear(hidden_dim, output_dim)

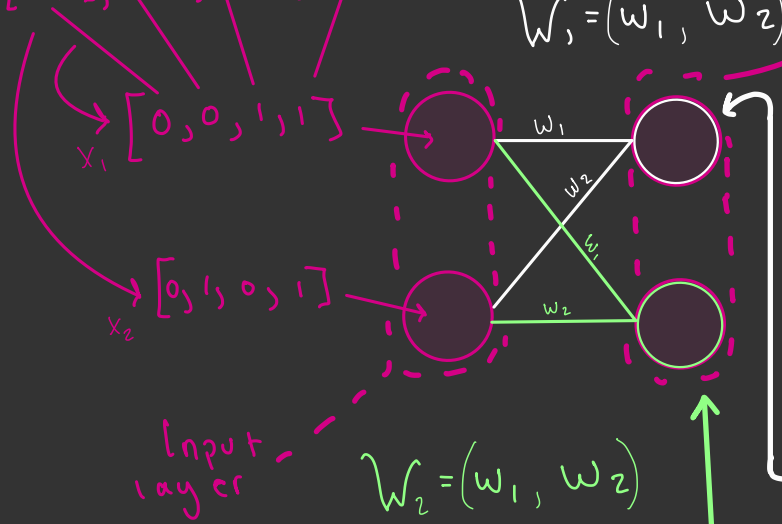
    def forward(self, x):
        hidden_output = torch.sigmoid(self.hidden_layer(x))
        output = torch.sigmoid(self.output_layer(hidden_output))
        return output
```

#2

#2.2

$$\text{self.hidden\_layer} = \text{nn.linear}(\text{input\_dim}, \text{hidden\_dim}) \equiv y = x W^T + b$$

Input  $X$   
 $= [0, 0], [0, 1], [1, 0], [1, 1]$



Hidden layer #1

$Z = \text{weighted sum of inputs} + \text{bias}$   
 $\hookrightarrow$  "raw score"

$$Z = (x_1 \cdot w_1) + (x_2 \cdot w_2) + b$$

$$z[0,0] = (0 \cdot w_1) + (0 \cdot w_2) + b = b$$

$$z[0,1] = (0 \cdot w_1) + (1 \cdot w_2) + b = w_2 + b$$

$$z[1,0] = (1 \cdot w_1) + (0 \cdot w_2) + b = w_1 + b$$

$$z[1,1] = (1 \cdot w_1) + (1 \cdot w_2) + b = w_1 + w_2 + b$$

$$\text{Output} = \text{sigmoid}(z[n]) = \frac{1}{1 + e^{-z}} \quad \left. \vphantom{\frac{1}{1 + e^{-z}}} \right\} \begin{array}{l} \text{Activation} \\ \text{Function} \end{array}$$

$$Z = (x_1 \cdot w_1) + (x_2 \cdot w_2) + b$$

$$z[0,0] = (0 \cdot w_1) + (0 \cdot w_2) + b = b$$

$$z[0,1] = (0 \cdot w_1) + (1 \cdot w_2) + b = w_2 + b$$

$$z[1,0] = (1 \cdot w_1) + (0 \cdot w_2) + b = w_1 + b$$

$$z[1,1] = (1 \cdot w_1) + (1 \cdot w_2) + b = w_1 + w_2 + b$$

$$\text{Output} = \text{sigmoid}(z[n]) = \frac{1}{1 + e^{-z}}$$

$$z[0,0] \quad z[0,1] \quad z[1,0] \quad z[1,1]$$

$$\text{Raw output} = [b, w_2 + b, w_1 + b, w_1 + w_2 + b]$$

$\hookrightarrow$  Activation function applied to all

$$\hookrightarrow \text{Final neuron output} = [0[0,0], 0[0,1], 0[1,0], 0[1,1]]$$

\*Although representations are the same, actual values of weights and biases differ.

$$z[0,0] \quad z[0,1] \quad z[1,0] \quad z[1,1]$$

$$\text{Raw output} = [b, w_2 + b, w_1 + b, w_1 + w_2 + b]$$

$\hookrightarrow$  Activation function applied to all

$$\hookrightarrow \text{Final neuron output} = [0[0,0], 0[0,1], 0[1,0], 0[1,1]]$$

Forward Pass

Step #1:

Input  $\rightarrow$  1<sup>st</sup> Hidden layer

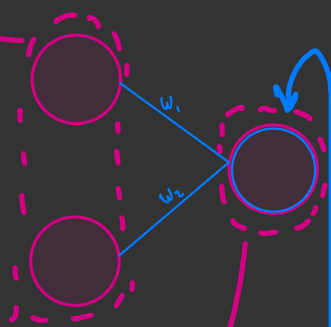
Calculations for all neurons within 1<sup>st</sup> layer.

Layer - to - layer process is the same

Hidden layer #1

$$W_1 = (w_1, w_2)$$

$z$  = weighted sum of inputs + bias  
↳ "raw score"



$$z = (x_1 \cdot w_1) + (x_2 \cdot w_2) + b$$

$$z[0,0] = (0 \cdot w_1) + (0 \cdot w_2) + b = b$$

$$z[0,1] = (0 \cdot w_1) + (1 \cdot w_2) + b = w_2 + b$$

$$z[1,0] = (1 \cdot w_1) + (0 \cdot w_2) + b = w_1 + b$$

$$z[1,1] = (1 \cdot w_1) + (1 \cdot w_2) + b = w_1 + w_2 + b$$

$$\text{Output} = \text{sigmoid}(z[n]) = \frac{1}{1 + e^{-z}} \quad \left. \vphantom{\frac{1}{1 + e^{-z}}} \right\} \begin{array}{l} \text{Activation} \\ \text{Function} \end{array}$$

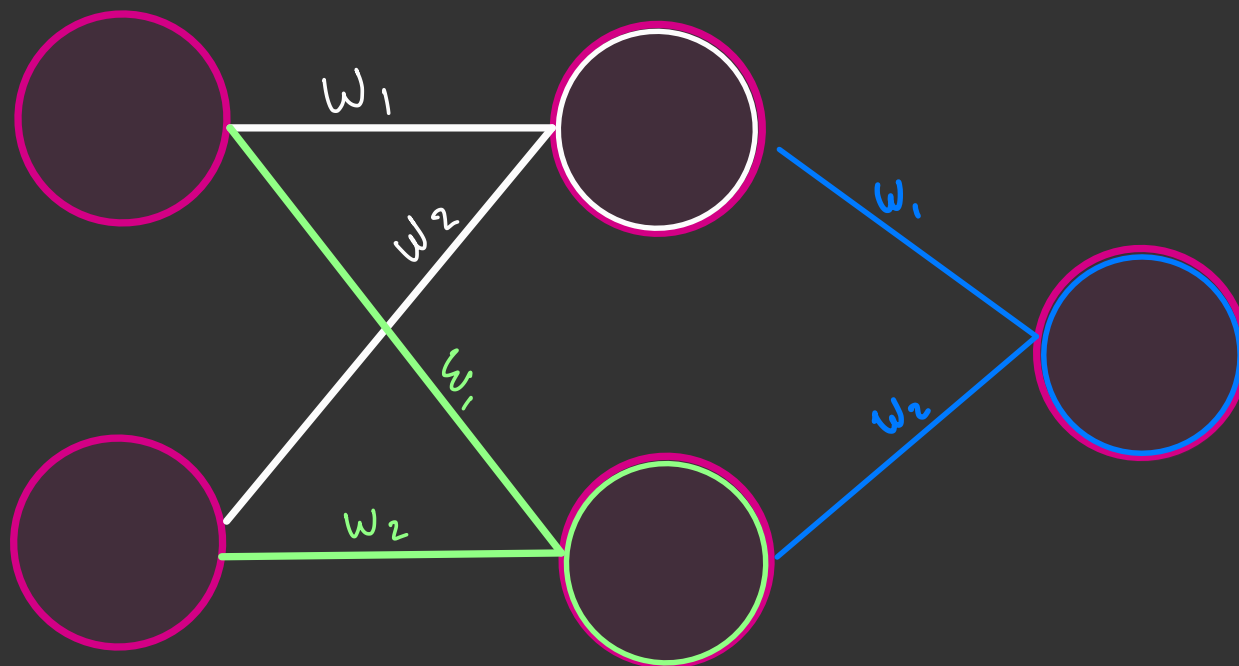
Output layer

$$z[0,0] \quad z[0,1] \quad z[1,0] \quad z[1,1]$$

$$\text{Raw output} = [b, w_2 + b, w_1 + b, w_1 + w_2 + b]$$

↳ Activation function applied to all

$$\text{Final neuron output} = [0[0,0], 0[0,1], 0[1,0], 0[1,1]]$$



```
# Initialize the neural network
```

```
input_dim = 2
```

```
hidden_dim = 2
```

```
output_dim = 1
```

```
model = XORNeuralNetwork(input_dim, hidden_dim, output_dim) → create model
```

#3

```
# Define the loss function (Mean Squared Error) and the optimizer (Stochastic Gradient Descent)
```

```
criterion = nn.MSELoss() → class | we're creating a class instance
```

```
optimizer = optim.SGD(model.parameters(), lr=0.1)
```

#4

```
# Training loop
```

```
epochs = 100000
```

#5

```
for epoch in range(epochs): 5.1
```

```
optimizer.zero_grad() # reset the gradients of the params
```

```
# Forward pass 5.2
```

```
predictions = model(X) → model.forward(x) → (step# 2.2)
```

```
# Calculate the loss 5.3
```

```
loss = criterion(predictions, y) → criterion.forward(predictions, y) → nn.MSELoss.forward(criterion, predictions, y)  
↘ distance from predicted to target
```

```
# Backpropagation
```

```
loss.backward() 5.4 → compute gradients
```

```
# Update weights
```

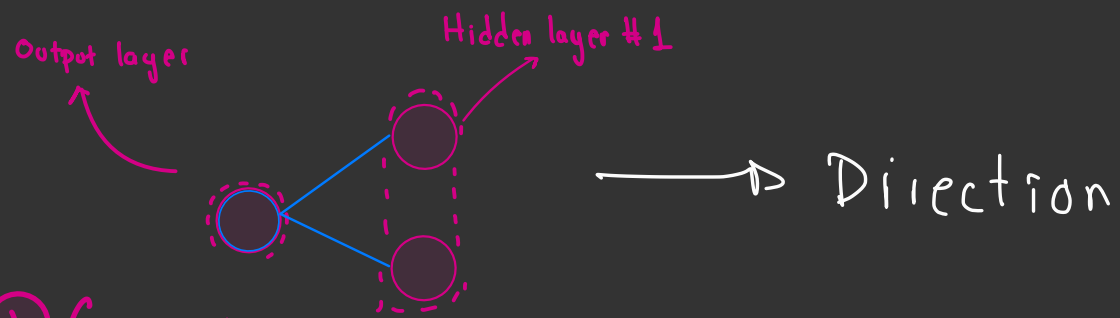
```
optimizer.step() 5.5 → update params based on gradients to minimize loss
```

```
if epoch % 1000 == 0:
```

```
    print(f"Epoch {epoch}, Loss: {loss.item()}")
```



## 5.4 Backpropagation: calculate neurons' contribution to error



① Compute gradient (Partial derivative) of loss w/ respect to output  $\hat{y}$  (within output neuron)

$$\frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)$$

② Chain rule: Propagate error through activation function | Basically calculate gradient w/ respect to pre-activation func output ( $z$ ).  
Also taking sigmoided  $z$  in mind, or rather, its derivative:  $\hat{y}(1 - \hat{y})$

$$\frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y) \cdot \hat{y}(1 - \hat{y})$$

③ Gradient w/ respect to weights and biases connected to hidden layer.

Weights gradient:  $\frac{\partial L}{\partial W_{\text{output}}} = \frac{\partial L}{\partial z} \cdot h$  Hidden layer output

Bias gradient:  $\frac{\partial L}{\partial b_{\text{output}}} = \frac{\partial L}{\partial z}$

How much the weights of the output layer should change to reduce loss.

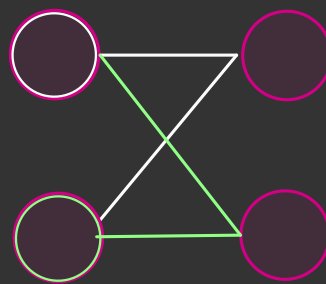
How much bias adjustment in output neuron to reduce loss.

## Backprop to hidden layer

Hidden layer #1

Input layer

Within neurons  
↓



→ Direction

4 Gradient w/respect to hidden layer outputs, pre-activation values

For each neuron, the gradient of loss w/respect to its output  $h$  is calculated by summing the contributions of neurons in the output layer.

$$\frac{\partial L}{\partial h} = \sum_{i=1}^{n_{\text{output}}} \left( \frac{\partial L}{\partial z_{\text{output}}} \cdot W_{\text{output}, i} \right)$$

$W_{\text{output}} =$  weight connecting hidden neuron to the  $i$ -th output  $n$ .

5 Gradient w/respect to hidden layer weights and biases

Weights gradient:  $\frac{\partial L}{\partial W_{\text{hidden}}} = \frac{\partial L}{\partial h} \cdot x$

Bias gradient:  $\frac{\partial L}{\partial b_{\text{hidden}}} = \frac{\partial L}{\partial h}$

- After backprop, the gradients for each parameter (weights and biases) have been calculated and stored in their `.grad` attribute. In the next step, the optimizer's job is to use these gradients to update the parameters in order to minimize loss.

