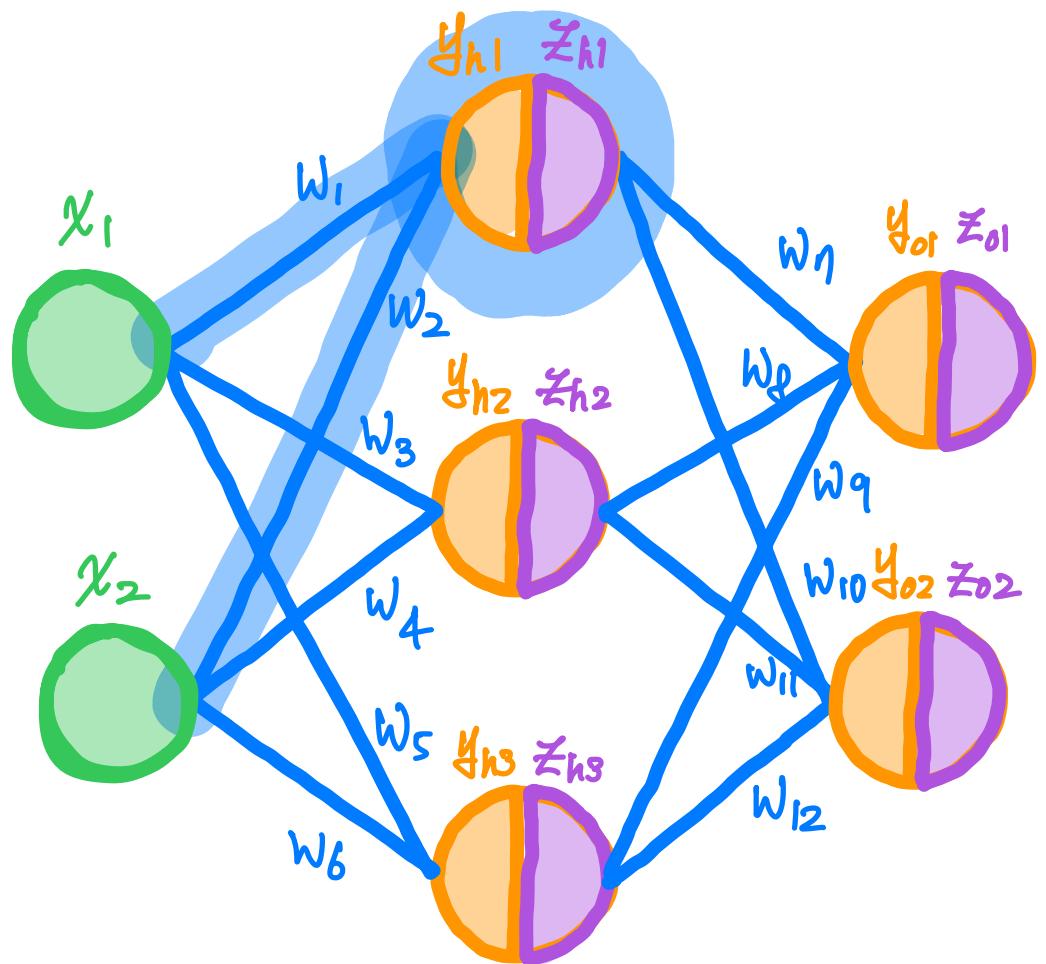


1 star-bits

# Forward Propagation

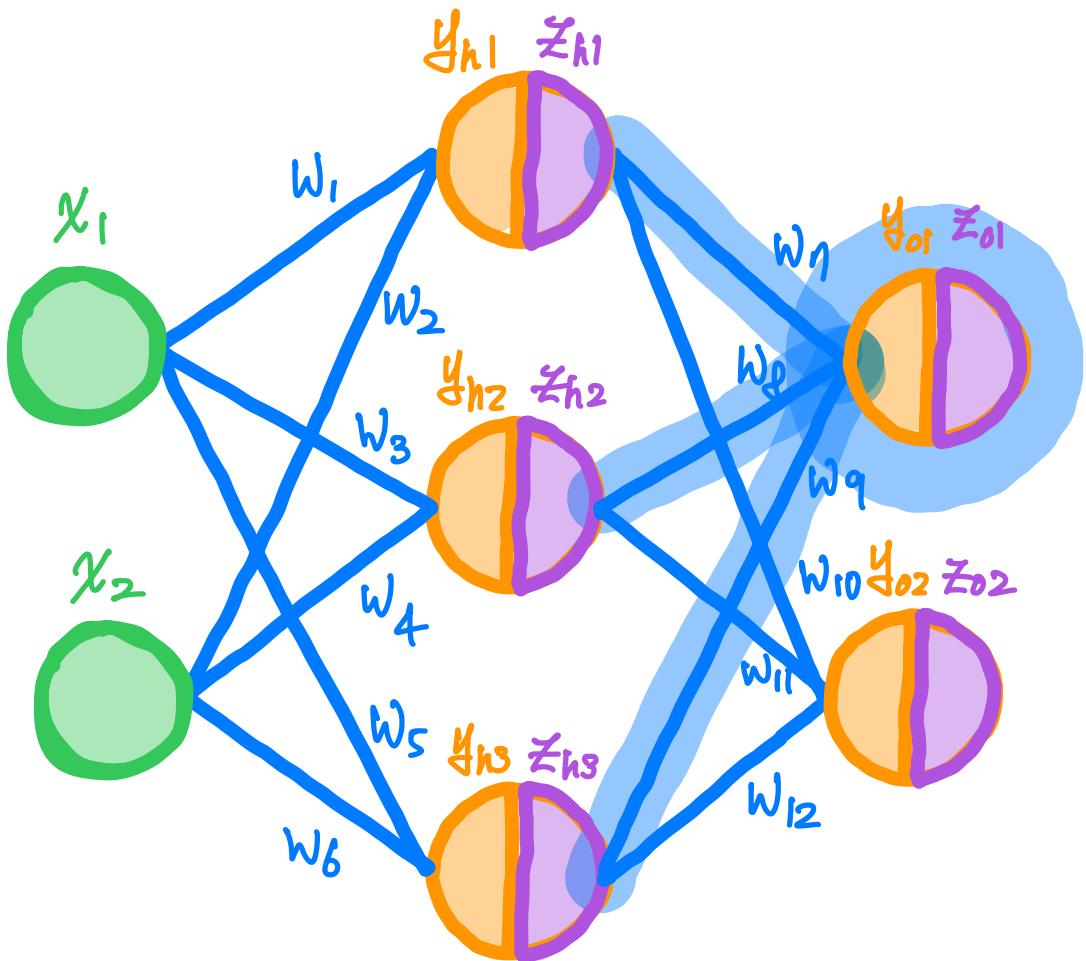


$$y_{h1} = w_1 \cdot x_1 + w_2 \cdot x_2$$

$$z_{h1} = \text{sigmoid}(y_{h1})$$

$$= \frac{1}{1 + e^{-y_{h1}}}$$

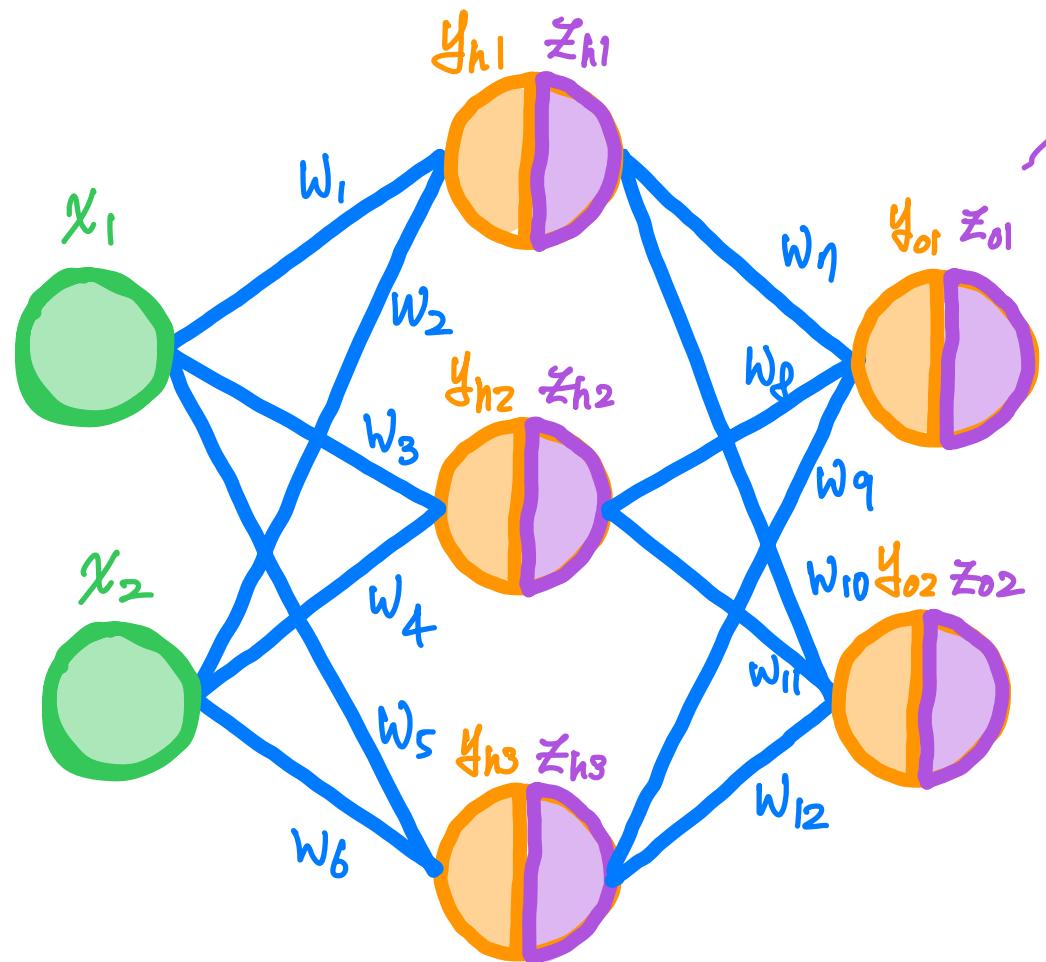
# Forward Propagation



weight은 어딨지 알? trainable elements

$$\begin{aligned}
 y_{01} &= w_1 \cdot z_{h1} + w_2 \cdot z_{h2} + w_3 \cdot z_{h3} \\
 z_{01} &= \text{sigmoid}(y_{01}) \\
 &= \frac{1}{1 + e^{-y_{01}}} \\
 &= \boxed{\text{01 } \rightarrow \text{[고2]에 대한 예측값}}
 \end{aligned}$$

02 | 3D | 일반화 훈련하기



$$L = - \sum_i t_i \ln(z_i)$$

$$= - [t_1 \ln(z_{01}) + t_2 \ln(z_{02})]$$

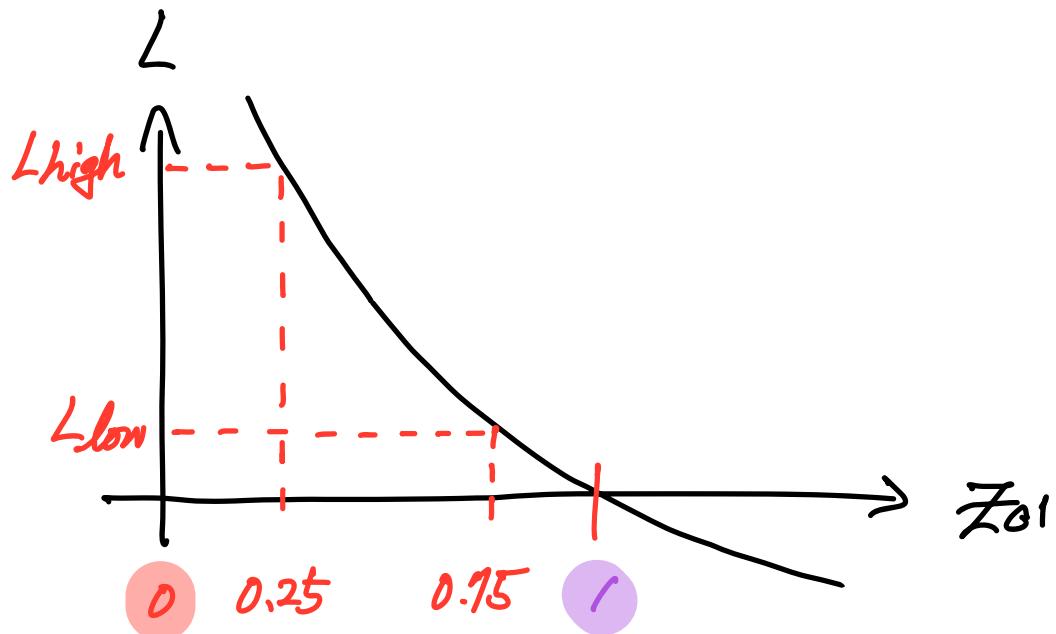
$$= - [t_1 \ln(z_{01}) + (1-t_1) \ln(1-z_{01})]$$

$$= L_1 + L_2 = L_{\text{total}}$$

QH cross entropy loss  $\frac{L}{\lambda}$ ?

$$L = -[t_1 \ln(\lambda_{01}) + (1-t_1) \ln(1-\lambda_{01})]$$

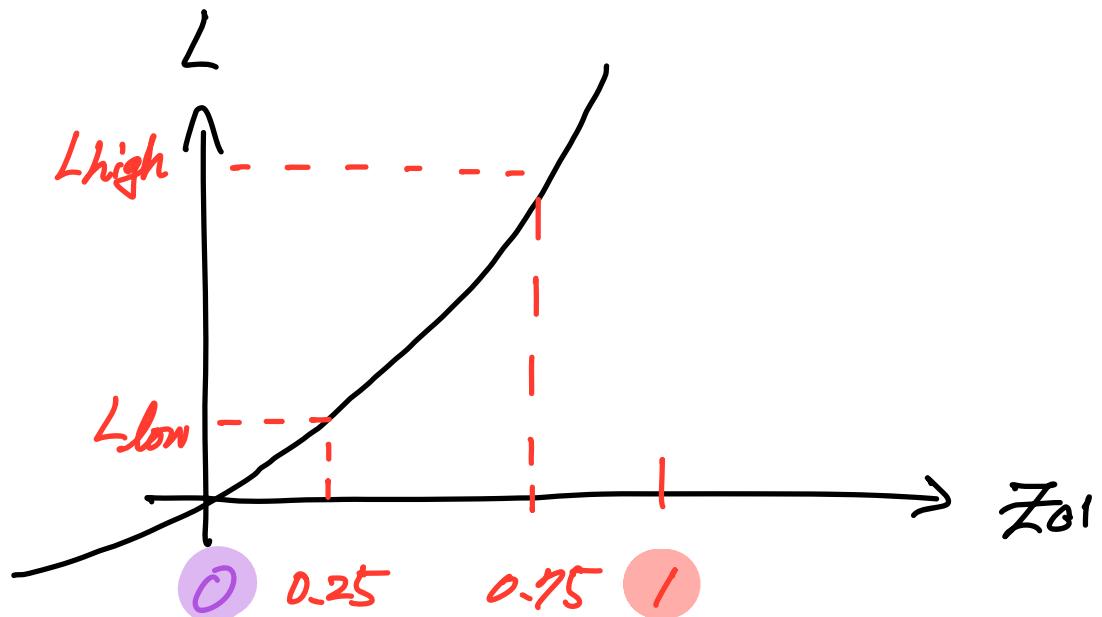
e.g.  $t_1 = 1$ :  $L = -[\ln \lambda_{01}]$



QH cross entropy loss  $\frac{L}{\lambda}$ ?

$$L = -[t_1 \ln(z_{01}) + (1-t_1) \ln(1-z_{01})]$$

e.g.  $t_1 = 0$ :  $L = -[\ln(1-z_{01})]$



What 'cross' entropy diff?

$$S = k_B \ln \Omega. \quad \text{As } S_{\text{total}} = S_1 + S_2, \quad S_{\text{total}} = \Omega_1 \times \Omega_2$$

$$S = -k_B \sum_i p_i \ln p_i \quad \text{Accounts for probabilities per microstate.}$$

$$\text{If } P = \frac{1}{\Omega}, \quad S = -k_B \Omega \frac{1}{\Omega} \ln \left( \frac{1}{\Omega} \right)$$

$$H = - \sum_i p_i \ln p_i \quad \text{Average amount of information content in a message}$$



$$C = - \sum_i t_i p_i \ln p_i$$

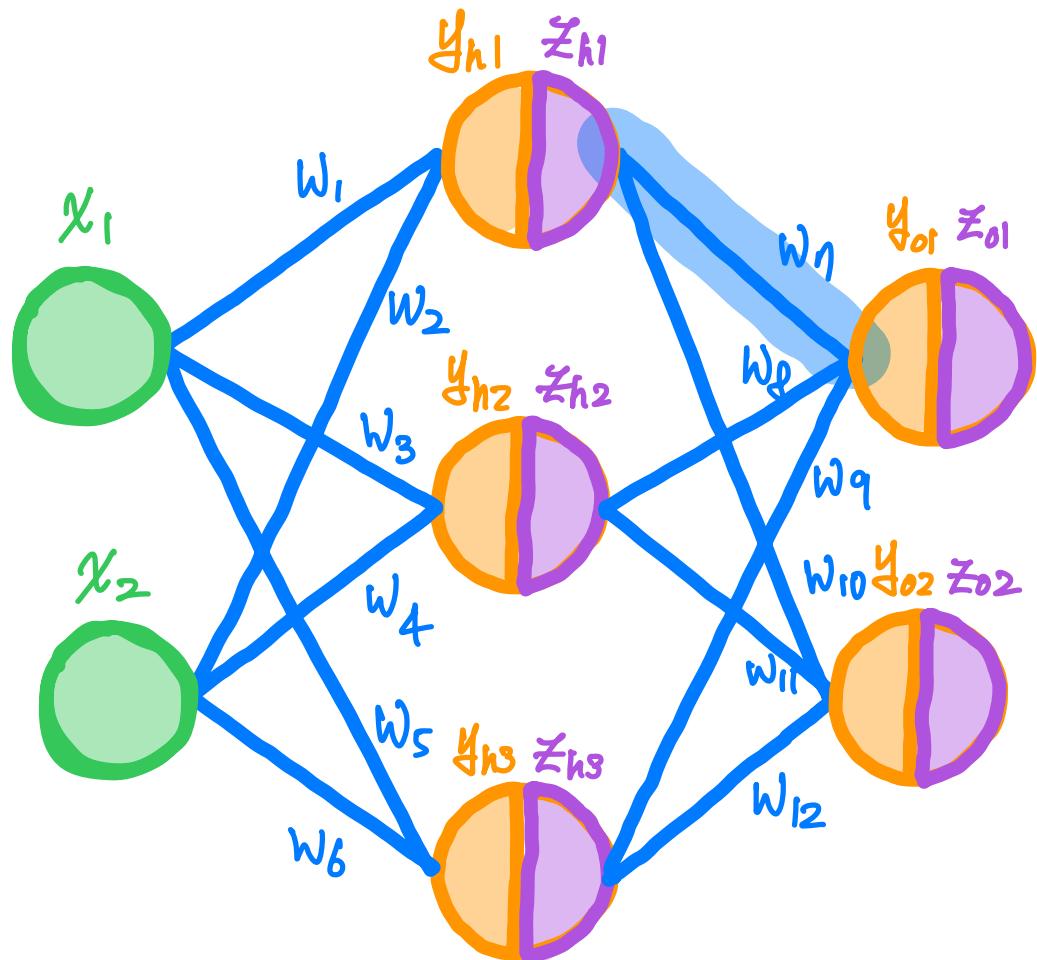
↑                   ↑  
'cross'

$-\ln p_i$  : information content of an event  $i$ .

As  $p_i \rightarrow 1$ , surprisal  $\rightarrow 0$ .

As  $p_i \rightarrow 0$ , surprisal  $\rightarrow \infty$ .

어디까지나 틀린 정도에 따라 weight를 조정하자

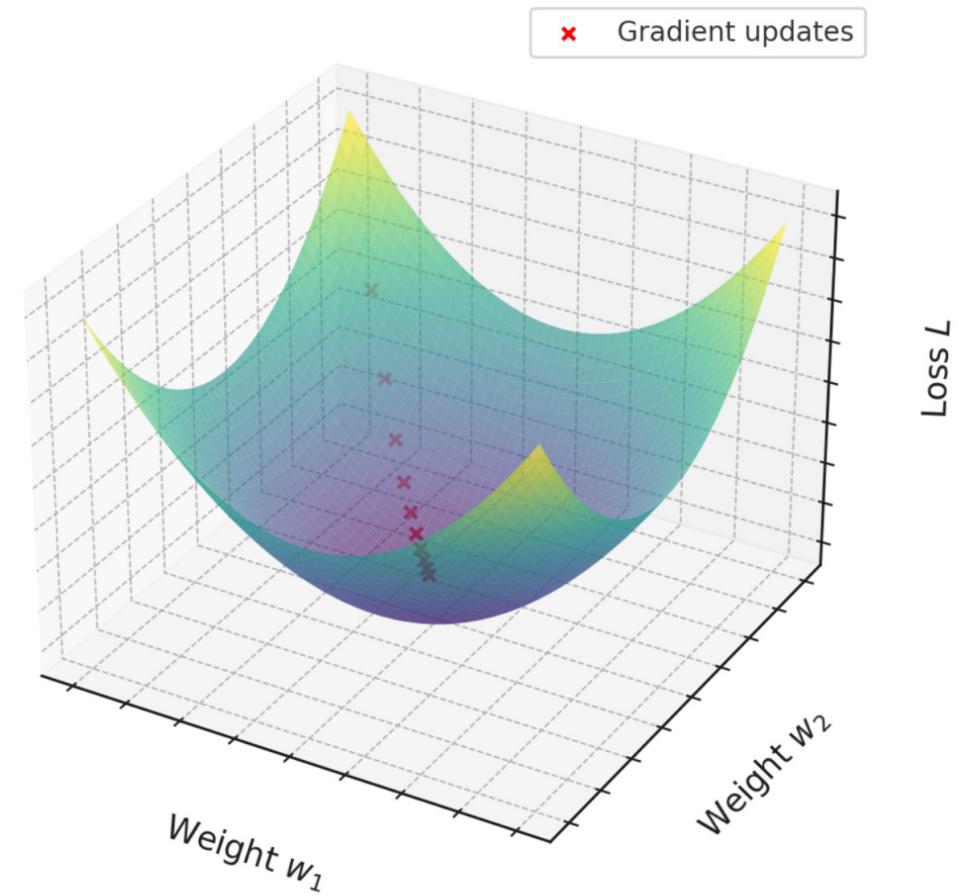
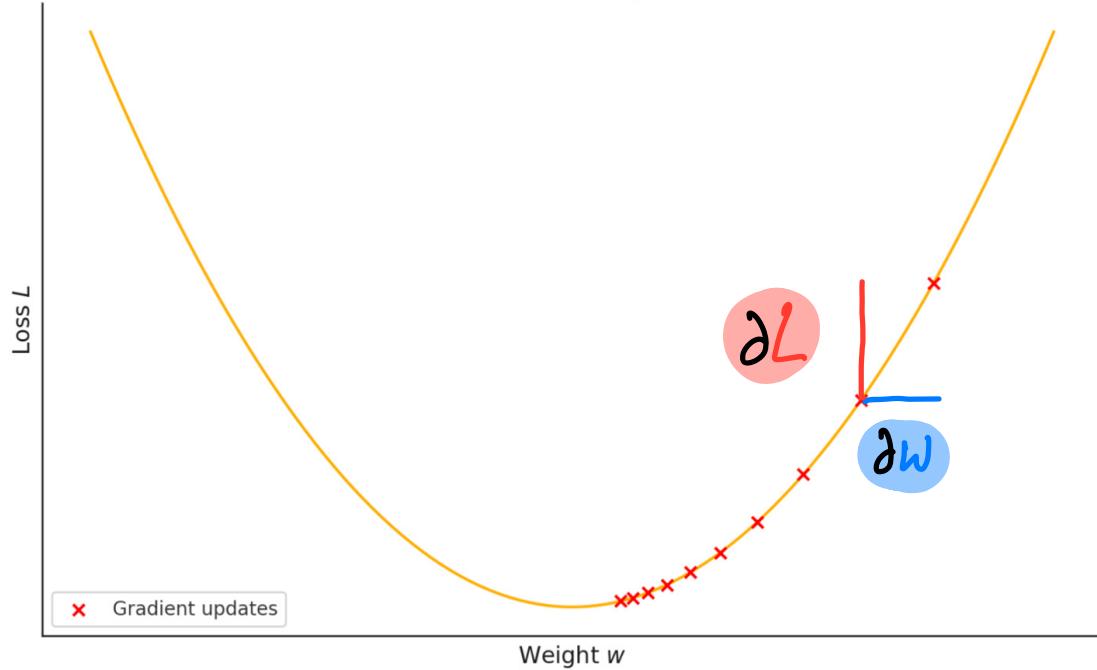


numerical gradient

$$w_{\text{new}} = w_{\text{old}} - LR \cdot \boxed{\frac{\partial L_{\text{total}}}{\partial w_{\text{old}}}}$$

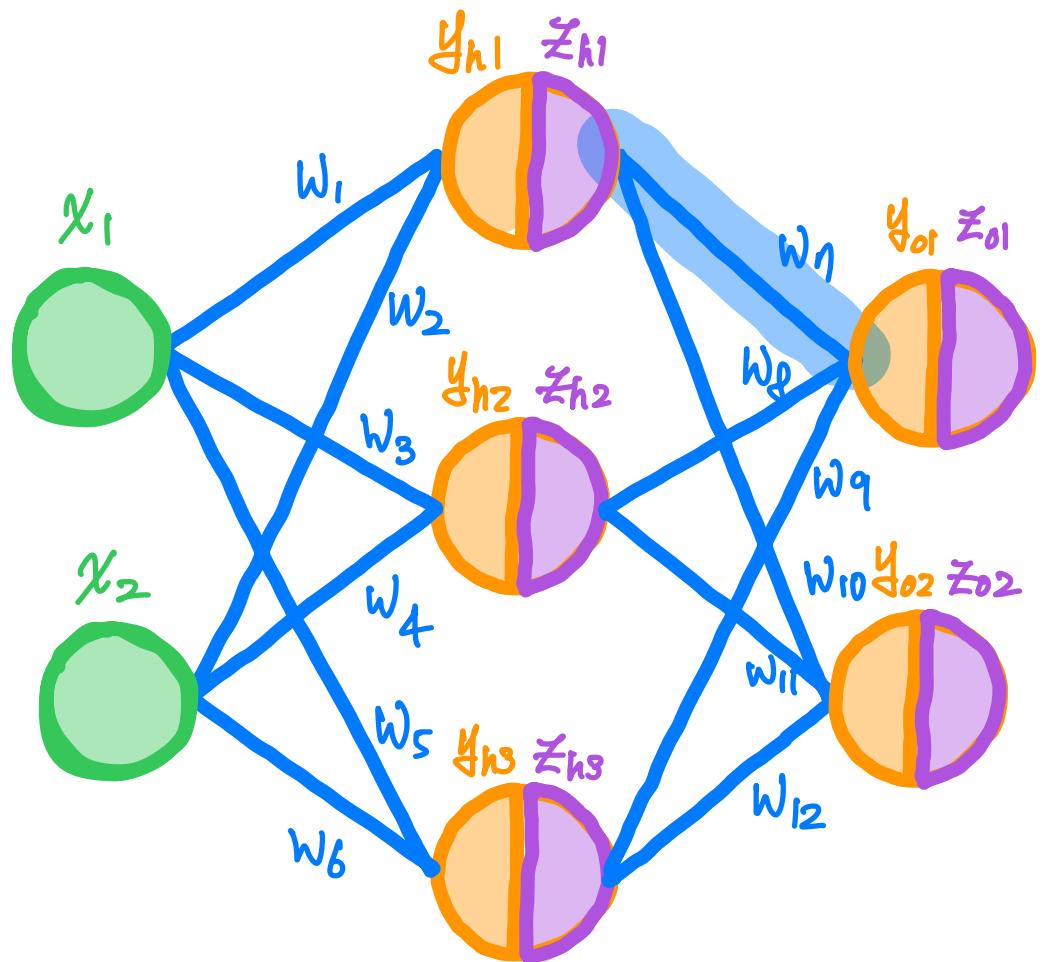
$w$ 을 조절하면서  $L_{\text{total}}$ 이 변화하는 정도.

$\partial L / \partial w$ ?  $Loss \frac{N}{2} (A - \hat{A})^2 + \epsilon$



LR

等式： $\frac{\partial L_{\text{total}}}{\partial w_1}$  ⇒ 3f7



$$\frac{\partial L_{\text{total}}}{\partial w_1} = \frac{\partial L_1}{\partial w_1}$$

$$L_1 = \text{Loss}(\text{sigmoid}(z_{o1}))$$

$$\frac{\partial L_1}{\partial w_1} = \boxed{\frac{\partial L_1}{\partial z_{o1}}} \cdot \boxed{\frac{\partial z_{o1}}{\partial y_{o1}}} \cdot \boxed{\frac{\partial y_{o1}}{\partial w_1}}$$

(1)                    (2)                    (3)

①

$$L_1 = -[x_1 \ln(z_{01}) + (1-x_1) \ln(1-z_{01})]$$

$$\frac{\partial L_1}{\partial z_{01}} = - \left[ x_1 \cdot \frac{1}{z_{01}} + (1-x_1) \frac{1}{1-z_{01}} \cdot (-1) \right] = \boxed{\frac{z_{01} - x_1}{z_{01}(1-z_{01})}}$$

①

②

$$\cdot z_{01} = \text{sigmoid}(y_{01}) = \frac{1}{1 + e^{-y_{01}}}$$

$$\frac{\partial z_{01}}{\partial y_{01}} = \frac{(1)' \cdot (1 + e^{-y_{01}}) - 1 \cdot (1 + e^{-y_{01}})'}{(1 + e^{-y_{01}})^2} = \frac{0 - (1 \cdot (e^{-y_{01}}) \cdot (-1))}{(1 + e^{-y_{01}})^2}$$

$$= \boxed{\frac{e^{-y_{01}}}{(1 + e^{-y_{01}})^2}}$$

$\frac{\partial Z_{01}}{\partial y_{01}}$   $\stackrel{N}{=} \text{in terms of } Z_{01} \geq \text{적어보자.}$

$$\frac{\partial Z_{01}}{\partial y_{01}} = -\frac{e^{-y_{01}}}{(1 + e^{-y_{01}})^2} = e^{-y_{01}} \cdot (Z_{01})^2$$

As  $Z_{01} = \frac{1}{1 + e^{-y_{01}}}$ ,  $e^{-y_{01}} = \frac{1}{Z_{01}} - 1 = \frac{1 - Z_{01}}{Z_{01}}$

$$\frac{\partial Z_{01}}{\partial y_{01}} = \frac{1 - Z_{01}}{Z_{01}} \cdot (Z_{01})^2 = \boxed{(1 - Z_{01}) \cdot Z_{01}}$$

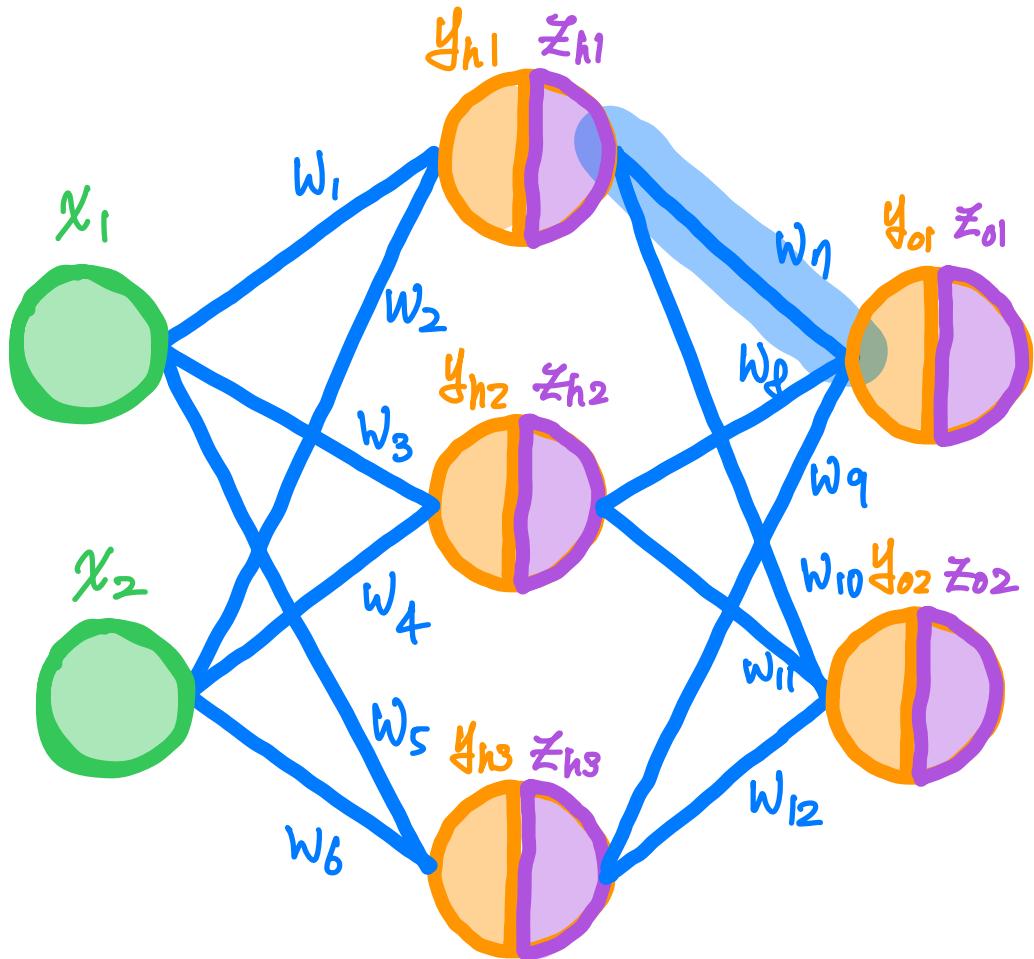
②

③

$$\cdot y_{01} = w_1 \cdot z_{h1} + w_2 \cdot z_{h2} + w_3 \cdot z_{h3}$$

$$\frac{\partial y_{01}}{\partial w_1} = [z_{h1}]$$

③



$$\frac{\partial L_{\text{total}}}{\partial w_\eta}$$

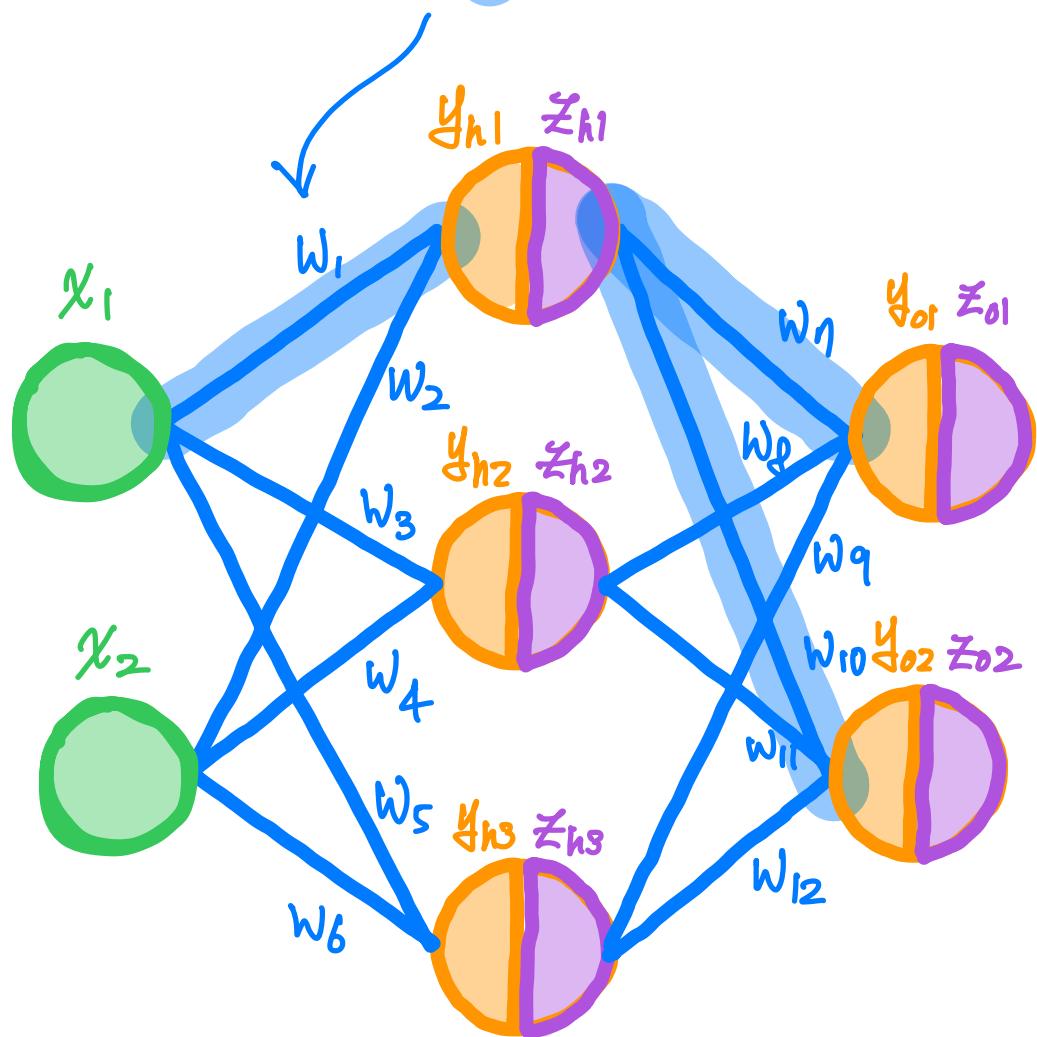
$$= \frac{\partial L_1}{\partial w_\eta}$$

$$= \frac{\partial L_1}{\partial z_{01}} \cdot \frac{\partial z_{01}}{\partial y_{01}} \cdot \frac{\partial y_{01}}{\partial w_\eta} = ① \cdot ② \cdot ③$$

$$= \frac{z_{01} - t_1}{z_{01}(1 - z_{01})} \cdot z_{01}(1 - z_{01}) \cdot z_{h1}$$

$$= [(z_{01} - t_1) \cdot z_{h1}]$$

$$\frac{\partial L}{\partial w_1} : \frac{\partial L_{total}}{\partial w_1} \rightarrow \text{Backprop}$$



$$\frac{\partial L_{total}}{\partial w_1} = \frac{\partial L_1}{\partial w_1} + \frac{\partial L_2}{\partial w_1}$$

$$= \frac{\partial L_1}{\partial z_{o1}} \cdot \frac{\partial z_{o1}}{\partial y_{o1}} \cdot \frac{\partial y_{o1}}{\partial z_{h1}} \cdot \frac{\partial z_{h1}}{\partial y_{h1}} \cdot \frac{\partial y_{h1}}{\partial w_1}$$

$$+ \frac{\partial L_2}{\partial z_{o2}} \cdot \frac{\partial z_{o2}}{\partial y_{o2}} \cdot \frac{\partial y_{o2}}{\partial z_{h2}} \cdot \frac{\partial z_{h2}}{\partial y_{h1}} \cdot \frac{\partial y_{h1}}{\partial w_1}$$

$$= \frac{z_{o1} - t_1}{z_{o1}(1-z_{o1})} \cdot \cancel{z_{o1}(1-z_{o1})} \cdot w_7 \cdot z_{h1}(1-z_{h1}) \cdot x_1$$

$$+ \frac{z_{o2} - t_2}{z_{o2}(1-z_{o2})} \cdot \cancel{z_{o2}(1-z_{o2})} \cdot w_{10} \cdot z_{h1}(1-z_{h1}) \cdot x_1$$

$$= \boxed{(z_{o1} - t_1) \cdot w_7 \cdot z_{h1}(1-z_{h1}) \cdot x_1 \\ + (z_{o2} - t_2) \cdot w_{10} \cdot z_{h1}(1-z_{h1}) \cdot x_1}$$

# Backward Propagation

$$\text{e.g. } \frac{\partial L_{\text{total}}}{\partial w_1} = (z_{o1} - t_1) \cdot z_{h1}$$

$$\frac{\partial L_{\text{total}}}{\partial w_1} = (z_{o1} - t_1) \cdot w_1 \cdot z_{h1}(1-z_{h1}) \cdot x_1 + (z_{o2} - t_2) \cdot w_{10} \cdot z_{h1}(1-z_{h1}) \cdot x_1$$

지금까지 한 것.

Numerical gradient 없이 analytic하게 미분 계산 가능

Sigmoid + cross entropy loss의 backprop

$$(z - t) \frac{\partial}{\partial z}$$

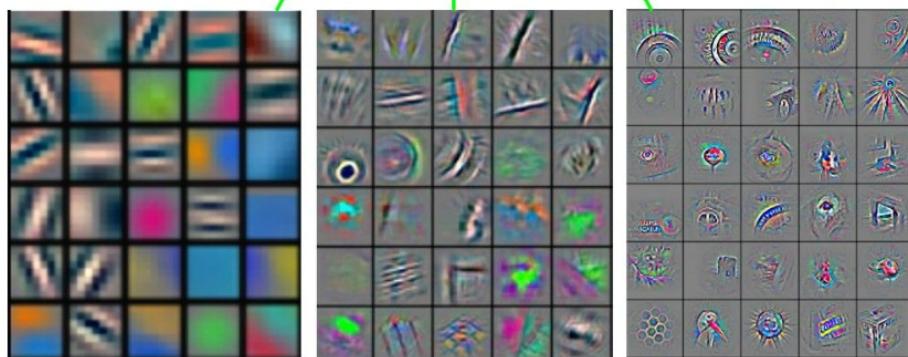
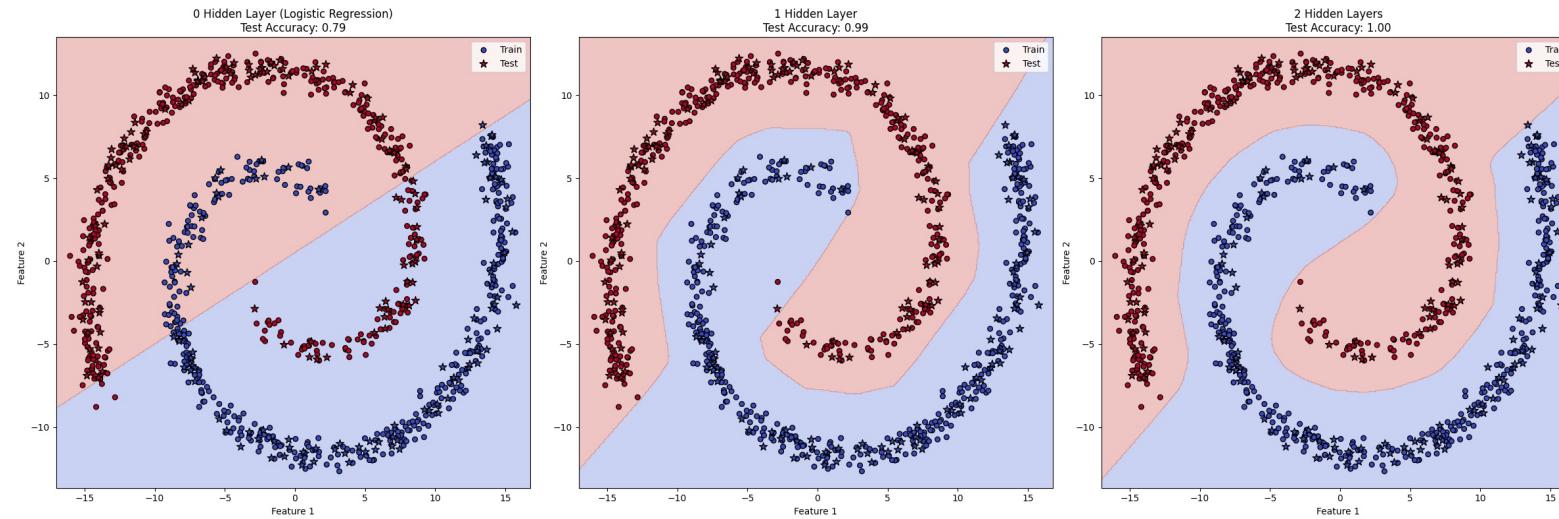
$$L = \text{Loss}(\text{sigmoid}(y))$$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial y} = \textcircled{1} \cdot \textcircled{2} = \frac{z-t}{z(1-z)} \cdot (1-z)z = z - t$$

이 loss function은 Keras supervised learning  
0이면 예상 (pred - truth)이 정답인 LR을  
기반으로 gradient update 할 수 있다? **아니다!**

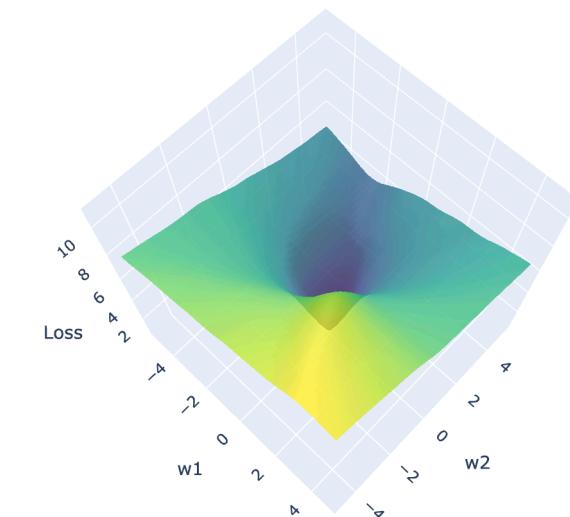
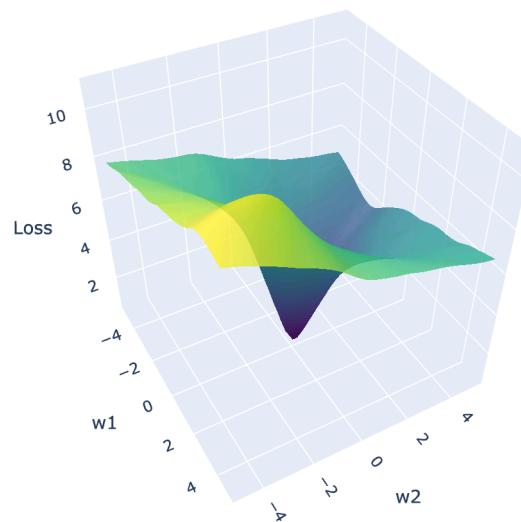
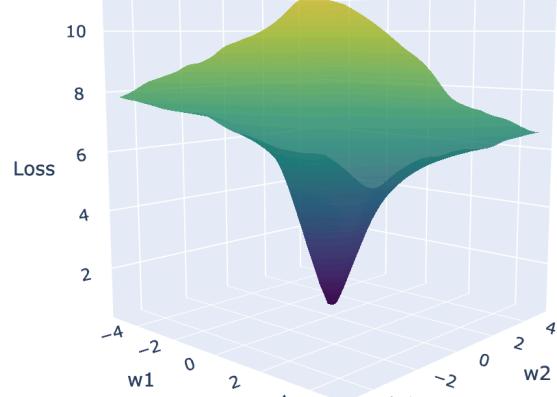
Loss is whatever function we've decided to use to optimize the parameters.

深层 layer が deep learning の特徴



복잡한 이미지의 표현이 가능해짐

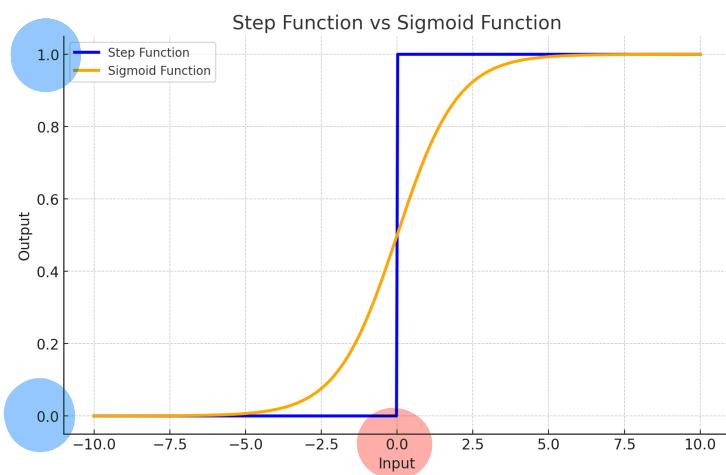
이전 slide의 spiral dataset의 경우 weight 2개 ( $w_1$ ,  $w_2$ )가  
각각 hidden layer를通過시킬 때 출력 loss surface



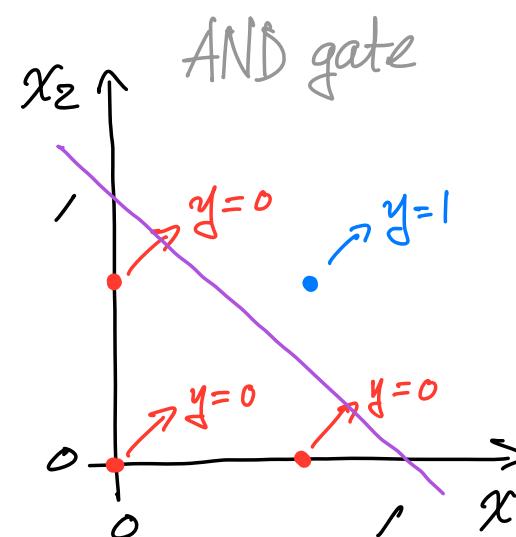
이 경우 dimension이 2이지만 weight의 경우는 dim이  
생각해보면  $\mathbb{R}^2$ , dim이 높아지면 local minima가 많아질 것이다

Q1 activation function이 뭘까?

Q2 sigmoid function은 알고자.



가장 가능성 있는 step function은

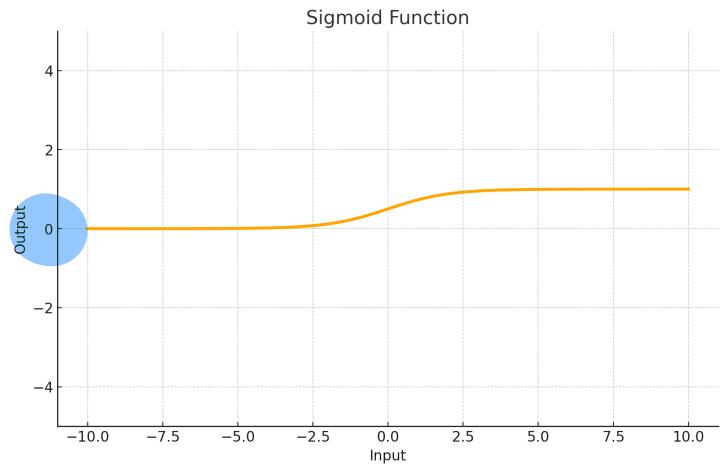


$x_1$	$x_2$	$y$
0	0	0
0	1	0
1	0	0
1	1	1

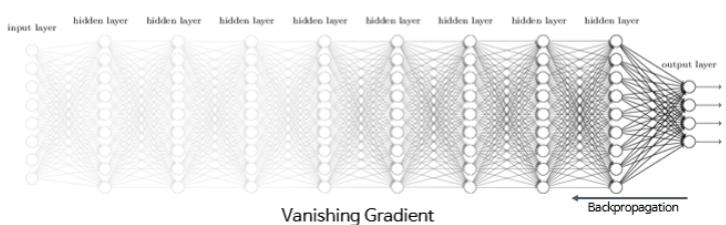
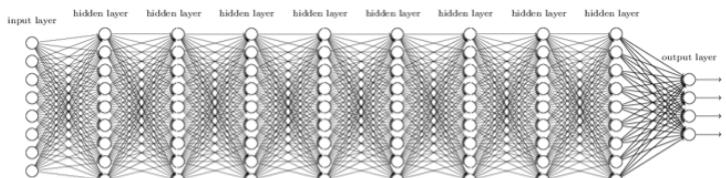
$$y = \text{step}(w_1 \cdot x_1 + w_2 \cdot x_2 + b)$$

without them, the entire network reduces to a linear relationship

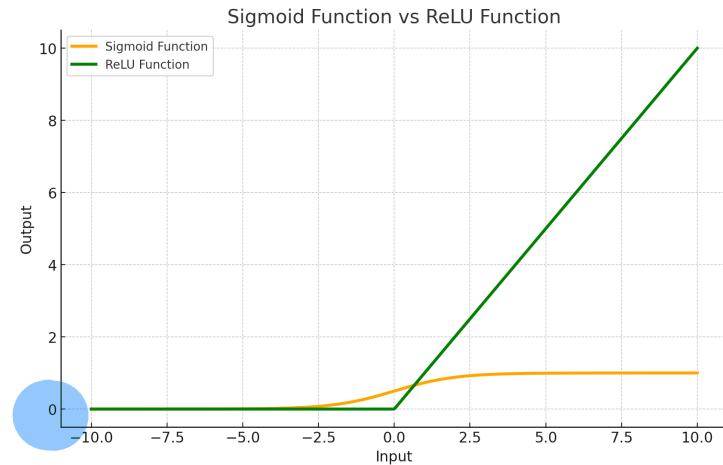
# Sigmoidal $\sigma(x)$ - Vanishing gradient.



↪ input  $x$   $\rightarrow$   $\sigma(x)$   $\rightarrow$   $\sigma^2(x)$



Nonlinearity 를 구현하기위한 활성화 함수 - ReLU



$$z = \max(0, y)$$

Sigmoid 를 사용해 학습률 + weight 를  
자꾸 유지하는 심리 - Weight decay

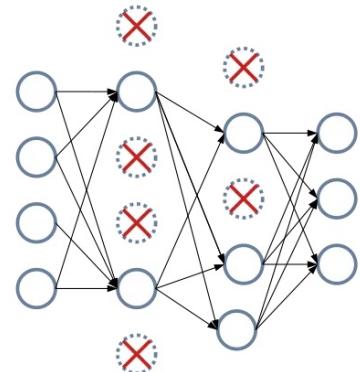
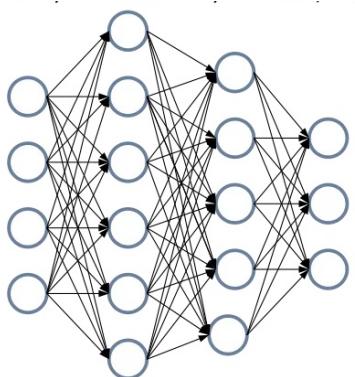
$$\text{Loss}_{\text{after}} = \text{Loss}_{\text{before}} + \lambda \sum_j w_j^2 \quad (\text{L2 regularization})$$

Again, the loss is whatever function we've decided to use to optimize the parameters.

그냥 previous layer output의 distribution을  
직접 조정해주는 안전? - Batch normalization

Inference 때 train data의 같은 mean & variance 정보를 갖고 있어야 한.  
Bias가 필요없어진다.

Training 때 누린 막대식 한계니까 더 robust 해지려면 -  
Dropout



레이어의 구성요소들을 순서대로 쌓아보자

linear layer

BatchNorm

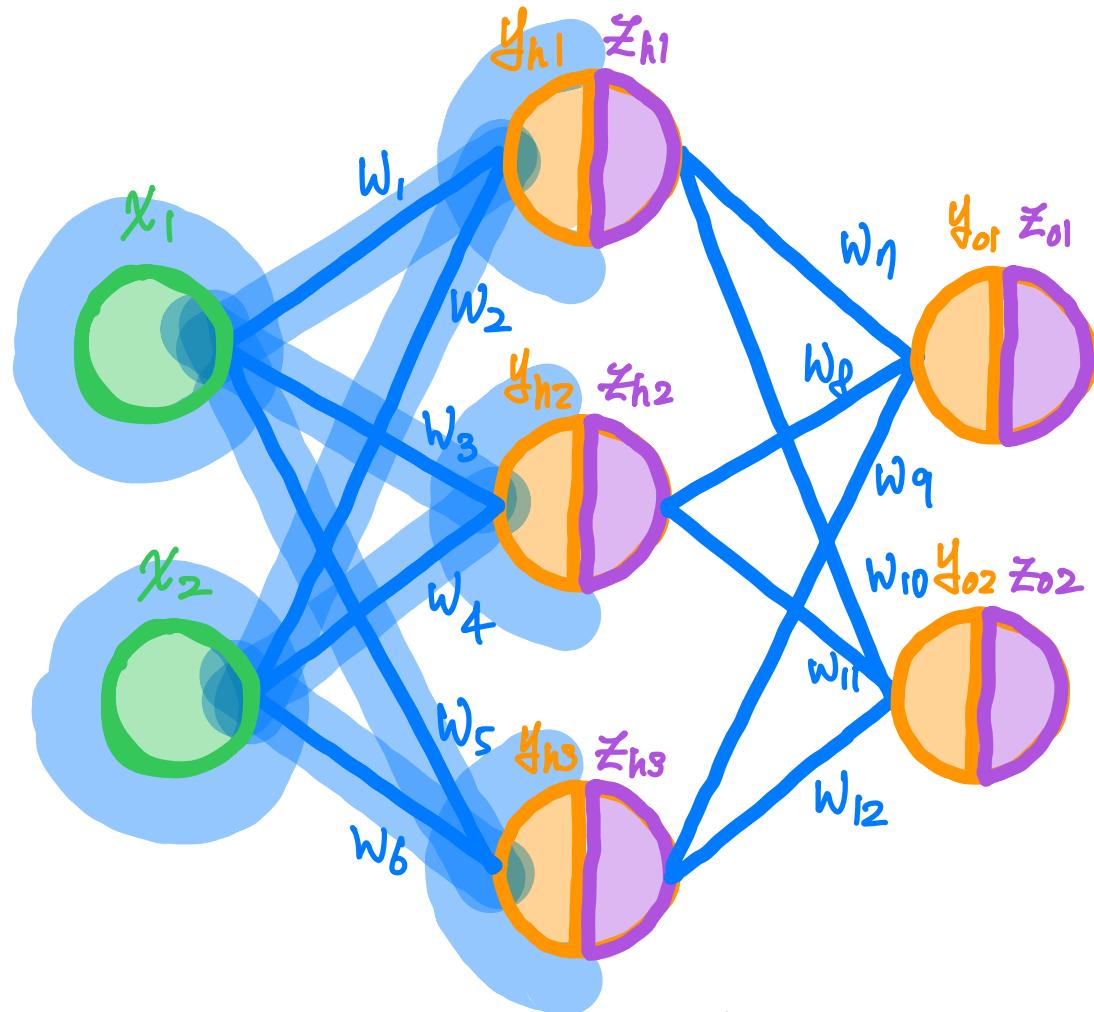
ReLU

Dropout

```
layer1 = nn.Sequential(  
    nn.Linear(input_dim, output_dim),  
    nn.BatchNorm(output_dim),  
    nn.ReLU(),  
    nn.Dropout(dropout_ratio),  
)
```

왜 이 순서인지 생각해보자

`nn.Linear(input_dim, output_dim)`



`input_dim x output_dim` matrix

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \times \begin{bmatrix} w_1 & w_3 & w_5 \\ w_2 & w_4 & w_6 \end{bmatrix} = \begin{bmatrix} y_{h1} \\ y_{h2} \\ y_{h3} \end{bmatrix}$$

$m$  by  $n$ . the first number ( $m$  here) represents the outermost dimension.  
`matrix.shape`  $\rightarrow (2, 3)$   
`matrix`  $\rightarrow [[w_1, w_3, w_5], [w_2, w_4, w_6]]$

Training은 코드상에서 어떤 형식으로 진행된다

num\_epochs 만큼 loop을 들면서,

각 카트리고리에 대한 pred를 생성한다,

pred에 따른 loss를 계산하고

weight를 업데이트 할 gradient를 초기화하고

backprop으로 gradient를 계산하고

gradient를 이용해 weight를 업데이트 할.

(hidden layer 없이 solely weight 학습만 있는 경우의 예)

```
# 여기가 데이터입니다
X = torch.FloatTensor(X_np) # 200개의 데이터포인트에 대해 Feature 1과 Feature 2의 2차원 정보가 있겠죠? 그럼 행렬의 사이즈는? --> (200, 2)
y = torch.FloatTensor(y_np) # 200개의 데이터포인트에 대해 0 또는 1의 레이블이 있겠죠? 그럼 행렬의 사이즈는? --> (200,)

# 여기가 가중치입니다
W = torch.zeros(2, requires_grad=True) # (z = w1 * x1 + w2 * x2)이니까 W에는 w1과 w2가 필요하겠죠? 그럼 행렬의 사이즈는? --> (2,)

learning_rate = 0.01
optimizer = optim.SGD([W], lr=learning_rate)

epochs = 100

# 100번 반복합니다
for epoch in range(epochs + 1):

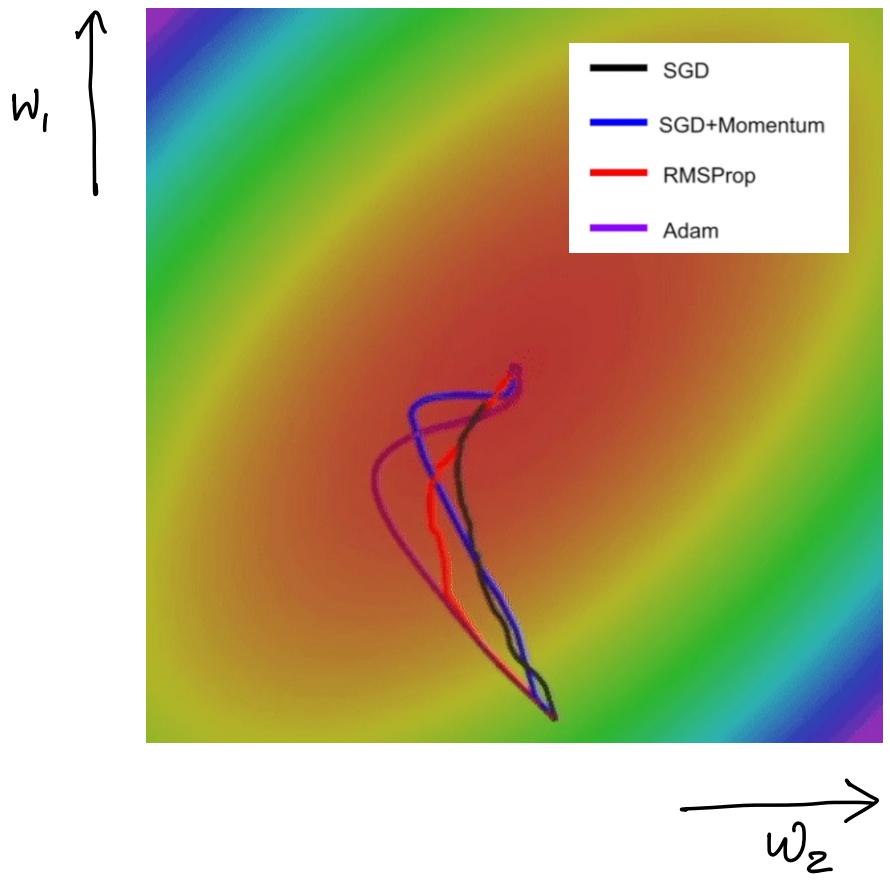
    z = X @ W # 200개의 데이터포인트에 대해 행렬연산을 해줍니다. 아웃풋 행렬의 사이즈는? --> (200,)

    # 여기가 pred입니다. 그러면 forward pass를 한 번 한거죠.
    pred = torch.sigmoid(z) # 200개의 데이터포인트에 대해 시그모이드 함수를 통과시켜줍니다

    # 여기가 loss function입니다
    loss = -(y * torch.log(pred + 1e-7) + (1 - y) * torch.log(1 - pred + 1e-7)).mean() # 200개의 prediction에 대한 loss의 평균값

    # 여기가 backprop입니다. 파이토치에서는,
    optimizer.zero_grad() # 이전 epoch의 gradient(즉, dL/dw)를 리셋시켜 주는 과정이 필요합니다
    loss.backward() # 이렇게 gradient를 계산하고,
    optimizer.step() # weight들을 gradient를 이용해 업데이트 시켜줍니다
```

# 경사 풀기 방법과 빙판效应 - Optimizer



GD: 전체 train dataset 기반 업데이트

SGD: minibatch 기반 업데이트

SGD + momentum: 과거의 gradient가 현재의 업데이트에 영향.

Adagrad: 각 dimension 기반 업데이트가 동일 배율로 일어나지 않음.

RMSProp: Adagrad의 속도 조절기가 시간이 지남수록 decay 함.

Adam: Adagrad + Momentum

이제 카운터를 제거하고, parameter 개수 확인

```
1 class Model(nn.Module):
2     def __init__(self):
3         super().__init__()
4
5         self.layer1 = nn.Sequential(
6             nn.Linear(28*28, 50, bias=False),
7             nn.BatchNorm1d(50),
8             nn.ReLU(),
9             nn.Dropout(0.5),
10        )
11        self.layer2 = nn.Sequential(
12            nn.Linear(50, 10, bias=False),
13        )
14
15    def forward(self, x):
16        x = x.view(-1, 28*28)
17        x = self.layer1(x)
18        x = self.layer2(x)
19
20        return x
21
22 model = Model()
23 print(model)
```

Model(  
  (layer1): Sequential(  
    (0): Linear(in\_features=784, out\_features=50, bias=False)  
    (1): BatchNorm1d(50, eps=1e-05, momentum=0.1, affine=True, track\_running\_stats=True)  
    (2): ReLU()  
    (3): Dropout(p=0.5, inplace=False)  
  )  
  (layer2): Sequential(  
    (0): Linear(in\_features=50, out\_features=10, bias=False)  
  )  
)

```
1 summary(model, input_size=(100, 784))
-----
Layer (type)           Output Shape      Param #
Linear-1              [-1, 50]           39,200
BatchNorm1d-2          [-1, 50]           100
ReLU-3                [-1, 50]           0
Dropout-4              [-1, 50]           0
Linear-5              [-1, 10]           500
=====
Total params: 39,800
Trainable params: 39,800
Non-trainable params: 0
-----
Input size (MB): 0.30
Forward/backward pass size (MB): 0.00
Params size (MB): 0.15
Estimated Total Size (MB): 0.45
```

# Training loop의 징후

```

1 criterion = nn.CrossEntropyLoss().to(device)
2 optimizer = torch.optim.SGD(model.parameters(), lr=0.1)
3
4 epochs = 20
5 for epoch in range(epochs + 1):
6
7     avg_loss = 0
8     total_batch = len(dataloader)
9
10    for X, y in dataloader:
11        X = X.view(-1, 28*28).to(device)
12        y = y.to(device)
13
14        pred = model(X)
15
16        loss = criterion(pred, y)
17
18        optimizer.zero_grad()
19        loss.backward()
20        optimizer.step()
21
22        avg_loss += loss / total_batch
23
24    print(f'Epoch: {epoch:4d} Loss: {avg_loss:.6f}')

```

Epoch:	Loss:
0	0.576404
1	0.416413
2	0.378474
3	0.352650
4	0.335536
5	0.324588
6	0.316307
7	0.304487
8	0.299579
9	0.293268
10	0.288154
11	0.282893
12	0.280186
13	0.275586
14	0.277942
15	0.270528
16	0.267779
17	0.267178
18	0.265811
19	0.258809
20	0.260383

이전 시점의 loss 보다↓

$$L = -\sum_{i=1}^{\text{num\_classes}} t_i \ln(p_i)$$

Highly DHL 경우  $t_i = 0$

그리고  $\frac{1}{\text{num\_classes}}$

$$L = -\ln\left(\frac{1}{\text{num\_classes}}\right)$$

$$L = \ln(\text{num\_classes})$$

$$\ln(10) = 2.30$$

코드 + 고지 파일의 위치 :

[github.com / star-bits / sogangparrot](https://github.com/star-bits/sogangparrot)