layer $0$ (Input)　　　$1$　　　$2$　　　$3$ (Output)

$W_{01}^{(2)} \Rightarrow$ 2번째 layer에서 1번째 유닛으로.

$g(S_1^{(2)})$

$g(S_2^{(2)})$

경1

$$9|_1^{(3)} = g(S_1^{(2)}), \quad S_1^{(2)} = W_{01}^{(2)} + W_{11}^{(2)} 9|_1^{(2)} + W_{21}^{(2)} 9|_1^{(2)}$$

$$9|_2^{(3)} = g(S_2^{(2)}), \quad S_2^{(2)} = W_{02}^{(2)} + W_{12}^{(2)} 9|_1^{(2)} + W_{22}^{(2)} 9|_2^{(2)}$$

$\times W_{22}^{(2)}$

Error function

$$f = (9_1 - 9|_1^{(3)})^2 + (9_2 - 9|_2^{(3)})^2$$

$\llcorner$ Output1　　$\llcorner$ Output2

$$\frac{\partial f}{\partial W_{21}^{(2)}} = \boxed{\frac{\partial f}{\partial 9|_1^{(3)}} \cdot \frac{\partial 9|_1^{(3)}}{\partial S_1^{(2)}}} \cdot \boxed{\frac{\partial S_1^{(2)}}{\partial W_{21}^{(2)}}}$$

$\delta_1^{(2)} \triangleq \dfrac{\partial f}{\partial S_1^{(2)}}$

$9|_1^{(2)}$

$$\frac{\partial f}{\partial W_{22}^{(2)}} = \boxed{\frac{\partial f}{\partial 9|_2^{(3)}} \cdot \frac{\partial 9|_2^{(3)}}{\partial S_2^{(2)}}} \cdot \frac{\partial S_2^{(2)}}{\partial W_{22}^{(2)}}$$

$\delta_2^{(2)} \triangleq \dfrac{\partial f}{\partial S_2^{(2)}}$

$9|_2^{(2)}$

$W_{12}$가 2개의 output에 영향을 준다.

경1　　　　　　　　　경 2

$$\frac{\partial f}{\partial W_{12}^{(1)}} = \boxed{\frac{\partial f}{\partial 9|_1^{(3)}} \cdot \frac{\partial 9|_1^{(3)}}{\partial S_1^{(2)}}} \cdot \boxed{\frac{\partial S_1^{(2)}}{\partial 9|_2^{(2)}}} \cdot \boxed{\frac{\partial 9|_2^{(2)}}{\partial S_2^{(1)}} \cdot \frac{\partial S_2^{(1)}}{\partial W_{12}^{(1)}}} \oplus \boxed{\frac{\partial f}{\partial 9|_2^{(3)}} \cdot \frac{\partial 9|_2^{(3)}}{\partial S_2^{(2)}}} \cdot \boxed{\frac{\partial S_2^{(2)}}{\partial 9|_2^{(1)}}} \cdot \boxed{\frac{\partial 9|_2^{(2)}}{\partial S_2^{(1)}} \cdot \frac{\partial S_2^{(1)}}{\partial W_{12}^{(1)}}}$$

$\delta_1^{(2)}$　　　$W_{21}^{(1)}$　　　$9|_1^{(1)}$　　　$\delta_2^{(2)}$　　$W_{22}^{(2)}$　　　$9|_1^{(1)}$

$$= \begin{bmatrix} \delta_1^{(2)} & \delta_2^{(2)} \end{bmatrix} \begin{bmatrix} W_{21}^{(2)} \\ W_{22}^{(2)} \end{bmatrix} \frac{\partial 9|_2^{(2)}}{\partial S_2^{(1)}} 9|_1^{(1)} = S^{(2)T} \cdot W_2^{(2)} \cdot \frac{\partial 9|_2^{(2)}}{\partial S_2^{(1)}} \cdot 9|_1^{(1)}$$

$\delta^{(2)}$　　　　$W_2^{(2)}$

다 씨에요.

$$= \frac{\partial f}{\partial S_2^{(1)}} = \delta_2^{(1)}$$

이 값은 forward 계산시 저장해 둘 수 있음.

$$\frac{\partial f}{\partial W_{ij}^{(l)}} = \delta_j^{(l)} 9|_i^{(l)}$$

$$\delta_j^{(l)} = (\delta^{(l+1)})^T W_j^{(l+1)} \cdot \frac{\partial 9|_j^{(l+1)}}{\partial S_j^{(l)}} \qquad (l \le L-2)$$

$$W_{ij}^{(l)} \longleftarrow W_{ij}^{(l)} - \alpha \frac{\partial T}{\partial W_{ij}^{(l)}}$$

- ## Initial Weight

  초기값 ⇒ 0 이면 안됨.

- ## Stochastic Gradient Descent.

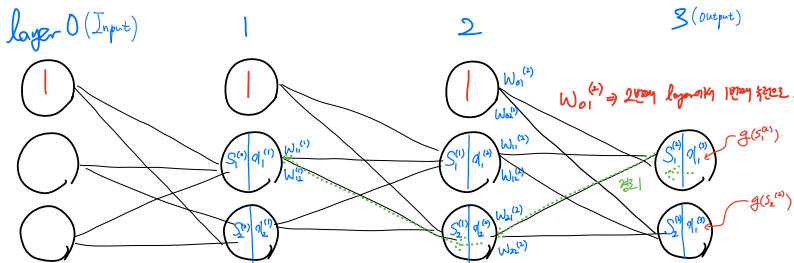  목적함수 $f_n$ 은 $f_1, f_2, \cdots f_n$ 로 여러 데이터들 Input 값들 넣어야 하는 것이야.

  괜히 $W_{ij}^{(l)} = W_{ij}^{(l)} - \alpha \sum_n \frac{\partial f_n}{\partial W_{ij}}$ 로 계산한다. 하면 아주 계산량이 너무 多.

  2개씩 Sum 하지 않고 1개로 가중치를 update.
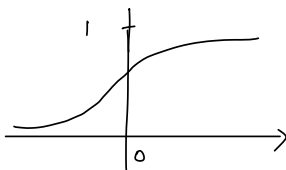
- ## Minibatch SGD

  $m$개씩 조금씩 사용해서 update.

## ⊙ Gradient Vanishing.



Sigmoid

0 에서의 미분값은 0.25 이다.
가장큰 기울기.

동과하기 전에는
sigmoid 편미분으로
있음.      ∴ 최대한 기울기여야 0.25

$W_{01}^{(2)}$ ⇒ 2번째 layer에서 1번째 뉴런으로.

$$\frac{\partial f}{\partial W_{ij}^{(l)}} = \delta_j^{(l)} q_i^{(l)}$$

$$\delta_j^{(l)} = (\delta^{(l+1)})^T W_j^{(l+1)} \cdot \boxed{\frac{\partial q_j^{(l+1)}}{\partial S_j^{(l)}}}$$

$$W_{ij}^{(l)} \longleftarrow W_{ij}^{(l)} - \alpha \frac{\partial f}{\partial W_{ij}^{(l)}}$$

$(l \leq L-2)$

· DNN에서 첫번째 layer의 weight를 경정할때 $S_1^{(1)}$을 구하려면 뒤쪽의 layer에서 제선 $\delta$들이 계속 곱해진 것이다.

2계번 기울기 $\dfrac{\partial g_J^{(l+1)}}{\partial S_J^{(l)}}$가 계속 곱해지게 돼는데 deep한 Net에서는 최대 기울기 0.25가 곱해진다고 해도 $S_1$이 이러서는 거의 0이 되기 때문에 gradient가

Gradient Vanishing 문제가 발생하는 것이다.

## ○ ReLU

· 2개서 기울기가 더 큰 Activation으로 개선된 함수이다.



가울기 1

※반약음수는 배끼지 않고 기울기가 1인 linear func를 쓰면

linear Regression이 된다.

ReLU는 linear Regression 이다.

$$a \xrightarrow{W_1} h_1 \xrightarrow{W_2} h_2 \xrightarrow{W_3} h_3 \xrightarrow{W_4} y \qquad Cost = f = \frac{1}{2}(y-t)^2$$

$$h = g(a \cdot w_1) \qquad y = g(h \cdot w_4)$$

$$W_4 = W_4 - \lambda \frac{\partial Cost(W_4)}{\partial W_4} \qquad W_4 = W_4 - h_3^T \lambda \delta_4$$

$$\frac{\partial Cost}{\partial W_4} = \boxed{\frac{\partial(ost}{\partial y} \cdot \frac{\partial y}{\partial(h_3 \cdot w_4)}}_{\delta_4} \cdot \frac{\partial(h_3 \cdot w_4)}{\partial W_4} = (y-t)g(h \cdot w_4)(1-g(h \cdot w_4)) \cdot h_3 = h_3^T \underline{(y-t)y(1-y)}_{\delta_4}$$

$$W_3 = W_3 - \lambda \frac{\partial Cost}{\partial W_3} \qquad W_3 = W_3 - h_2^T \lambda \delta_3$$

$$\frac{\partial Cost}{\partial W_3} = \boxed{\frac{\partial Cost}{\partial y} \cdot \frac{\partial y}{\partial(h_3 \cdot w_4)} \cdot \frac{\partial(h_3 \cdot w_4)}{\partial h_3} \cdot \frac{\partial h_3}{\partial(h_2 \cdot w_3)}}_{\delta_3} \cdot \frac{\partial(h_2 \cdot w_3)}{\partial W_3} = h_2^T \overset{\delta_3}{\underline{\underline{(y-t)y(1-y)}_{\delta_4} W_4^T h_3(1-h_3)}}$$

$$W_2 = W_2 - \lambda \frac{\partial Cost}{\partial W_2} \qquad W_2 = W_2 - h_1^T \lambda \delta_2$$

$$\frac{\partial Cost}{\partial W_2} = \boxed{\frac{\partial Cost}{\partial y} \cdot \frac{\partial y}{\partial(h_3 \cdot w_4)} \cdot \frac{\partial(h_3 \cdot w_4)}{\partial h_3} \cdot \frac{\partial h_3}{\partial(h_2 \cdot w_3)} \cdot \frac{\partial(h_2 \cdot w_3)}{\partial h_2} \cdot \frac{\partial h_2}{\partial(h_1 \cdot w_2)}}_{\delta_2} \cdot \frac{\partial(h_1 \cdot w_2)}{\partial W_2} = h_1^T \delta_3 \cdot W_3^T h_2(1-h_2)$$

- 일반화 식

$$W_t = W_t - h_{t-1}^T \lambda \delta_{t+1}$$
$$\delta_t = \delta_{t+1} \cdot W_{t+1}^T \boxed{h_t(1-h_t)}$$
$$\hookrightarrow Sigmoid \text{ 미분값}.$$

- Stochastic Gradient Descent.

  Initial Weight. $\Rightarrow$ Gaussian Dist.