

Back propagation

1. Optimization


2. Differentiation

3. Fully-connected

4. Chain Rule

5. Matrix form

6. Convolution



1. Optimization

Consider linear regression model such that $Y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \varepsilon$ where, $\varepsilon \sim N(0, \sigma^2)$

Data given \Rightarrow we can write model with data as $Y = X\beta + \varepsilon$ where $Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$ $X = \begin{pmatrix} 1 & x_{11} & \dots & x_{1k} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{nk} \end{pmatrix}$

Then, we can find out $\hat{\beta}$ by minimizing ε^2 .

$$\left\{ \begin{array}{l} \beta = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_k \end{pmatrix} \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix} \end{array} \right.$$

Because it's obvious that ε^2 is convex, what we need to do is finding $\hat{\beta}$ s.t. $\frac{\partial \varepsilon^2}{\partial \beta} \Big|_{\beta=\hat{\beta}} = 0$

$$\therefore \hat{\beta} = (X'X)^{-1} X'Y$$

However its computational complexity is $O(k^3)$ it means too high.

Thus, we use gradient descent.

Although, first method which figures out β as minimizer of \mathcal{E} , is closed form, gradient descent is more efficient.

iteration $\beta^{(t+1)} = \beta^{(t)} - \eta \frac{\partial \mathcal{E}}{\partial \beta}$ 이리 하여 순차적으로 update 함

오래 output = $f(\text{input}, \text{parameter})$, loss = $h(\text{output}, \text{target})$

오래 input: constant, parameter: variable 이므로

Optimization : loss를 최소화 만들기 위한 parameter를 찾는 것

2. Differentiation (Gradient)

흔히 계수가라 하는 것은 $\frac{\Delta y}{\Delta x}$, 즉, $\frac{y_2 - y_1}{x_2 - x_1}$ 일이다.

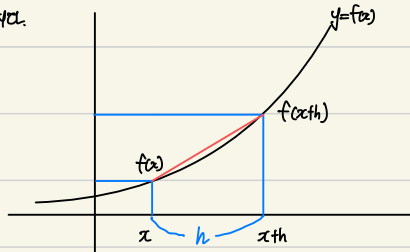
특히 이를 $\frac{dy}{dx}$ 와 혼용하여 사용되므로 하는데 엄밀 하에는 Δx 는 $x_2 - x_1$ 이고 dx 는 $\forall \epsilon > 0, \|x_2 - x_1\| < \epsilon$ 정수로 작은것을 의미함이다.

그래서 보통 $\lim_{\Delta x \rightarrow 0} \frac{\Delta y}{\Delta x} = \frac{dy}{dx}$ 라고 하며 이를 "미분" or "순간 기울기" 라고 함이다.

이때 $y = f(x)$ 가 다변수 함수이면 사실 같은 의미일게 어쨌든 느낄수 있습니다.

다시 일변수 일때 미분을 봅시다

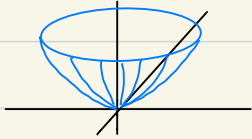
$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{(x+h) - (x)} \text{ 이때 이를 graph 상에서 나타낸다면}$$



변화량의 극한의 개념의 계수가 $h \rightarrow 0$ 이므로 $x+h \rightarrow x$ 이고 이는 x 에서의 순간 기울기라 볼수 있습니다.

따라서 $f'(x)$ 는 $f(x)$ 의 접선 기울기라 도할수있고 함이다.

이제 이번엔 $z = f(x, y) = x^2 + y^2$ 를 생각해봅시다.



와 같은 모양입니다.

그럼 이제 미분은 어떻게 정의할까요?

예를 들어 위의 모양에서 $f'(a)$ 는 $x=a$ 에서의 $f(x)$ 의 기울기 를 의미합니다.

그럼 이 변수에서 $f'(a, b)$ 는 어떤 값을 가져다 줘나요?

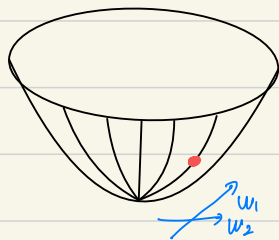
일 변수는 방향이 2만 있으므로 scalar 만, 다변수는 방향이 여러개이므로 각 방향에 대한 기울기를 주어져야 합니다.

즉, $f'(a, b) = \left(\lim_{h \rightarrow 0} \frac{f(a+h, b) - f(a, b)}{h}, \lim_{h \rightarrow 0} \frac{f(a, b+h) - f(a, b)}{h} \right)$ 가 되고 이는 각 방향에 대한 기울기가 됩니다.

Thus, we denote $\nabla f = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)$

이제 실재 어떻게 update 되려 보일수야.

Remind : 우리는 loss value 를 최소 만들어 주는 w 를 찾는 것이다.

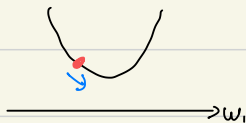


변경량이 최소의 loss 라면 $\nabla loss | \cdot$ 를 작고 여를 더함하여

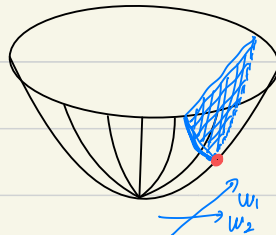
증가된 양의 방향의 용해가 되면 된다.

이제 vector 이므로 각의만큼 axis 가 존재하고 모든 axis 에 대해 moving 해야 함이다.

즉, 단원으로 볼때



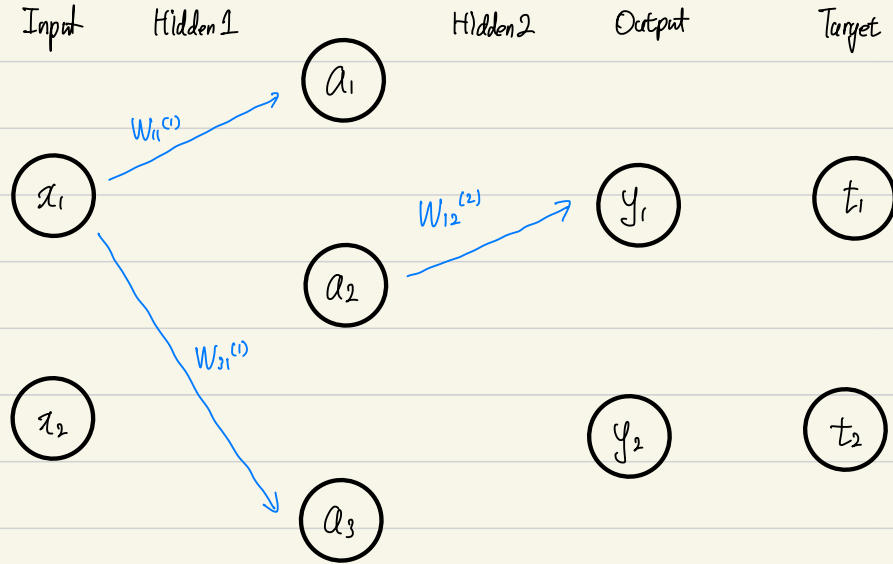
만일 진행한다면



은 밖이 못 가므로

모든 axis 에 대하여 parameter moving을 한다.

3. Simple fully connected classification model.



input $-0 = 2$, class : 3

no bias

$$\text{sigmoid} = \frac{1}{1 + e^{-x}}$$

Cross entropy

$$\text{softmax} = \frac{\exp(y_i)}{\sum \exp(y_i)}$$

Feed forward

let $S(x)$: sigmoid, $P(x)$: softmax

$$a_1 = S(w_{11}^{(1)}x_1 + w_{12}^{(1)}x_2)$$

$$y_1 = P(w_{11}^{(2)}a_1 + w_{12}^{(2)}a_2 + w_{13}^{(2)}a_3)$$

$$a_2 = S(w_{21}^{(1)}x_1 + w_{22}^{(1)}x_2)$$

$$y_2 = P(w_{21}^{(2)}a_1 + w_{22}^{(2)}a_2 + w_{23}^{(2)}a_3)$$

$$a_3 = S(w_{31}^{(1)}x_1 + w_{32}^{(1)}x_2)$$

$$\text{loss} = - \sum_{i \in C} t_i \log(y_i) = -t_1 \log y_1 - t_2 \log y_2$$

또는 binary 이니까 $-t_1 \log x_1 - (1-t_1) \log (1-x_1)$

이제 목표는 loss 를 계산 시켜주는 W 를 찾는 것이다.

즉, 이를 $\text{loss} = f(W)$ 라고 생각하면 됨. 또한 $w^{(1)}$ $w^{(2)}$ 다 순서를 생각하여 밑과 의미상

따라서 모든 $w_{ij}^{(k)}$ 에 대하여 ∇loss 를 계산하면 됨

gradient. important * : weights 만이 variable 다.다.

$W_2^{(2)}$ 의 gradient 를 구해봅시다.

by definition
$$\nabla = \lim_{h \rightarrow 0} \frac{\text{loss}(W_2^{(2)} + h, \dots) - \text{loss}(W_2^{(2)}, \dots)}{h}$$

실제로 code 에서는 h 를 충분히 작은 수로 설정하면 됩니다.

하지만 위 과정은 시간이 상당히 오래 걸립니다. 또한 h 에 대한 오차도 있습니다.

또는 도함수를 직접 계산해 보려면 deep 이면 n nodes가 1000 이면

4. Chain Rule

loss function & \approx layer, softmax, sigmoid 등 여러 함수를 연결하여 사용할 수 있다.

이를 이용하여 gradient를 쉽게 구할 수 있다.

$$\text{let } h(x) = f(g(x)) \quad , \quad y = g(x)$$

$$h'(x) = \frac{dh}{dx} = \frac{dh}{dy} \frac{dy}{dx} = \frac{dh}{dy} \times \frac{dy}{dx} \quad \text{구하다.}$$

$$\approx, \quad x=a \text{ at point of } \text{gradient} \quad f'(a) = \left. \frac{dh}{dy} \right|_{y=b} \times \left. \frac{dy}{dx} \right|_{x=a} \quad \text{where } b = g(a)$$

$$\text{e.g.) } f(x) = x^2 + 2x + 1 = g(h(x)) \quad \text{where } g(x) = x^2 \quad , \quad h(x) = x+1 \quad \nearrow y = h(x)$$

$$\text{gradient at } x=2 \quad \text{i) } f'(x)|_{x=2} = 2(x+1)|_{x=2} = 6$$

$$\begin{aligned} \text{ii) } f'(x)|_{x=2} &= \left. \frac{df}{dx} \right|_{x=2} = \left. \frac{df}{dy} \right|_{y=3} \times \left. \frac{dy}{dx} \right|_{x=2} \\ &= \left. \frac{dy^2}{dy} \right|_{y=3} \times \left. \frac{d(x+1)}{dx} \right|_{x=2} = 2y|_{y=3} \times 1|_{x=2} = 6 \end{aligned}$$

이를 위하여 $W_{ij}^{(k)}$ 에 대한 gradient를 계산해 보겠습니다.

예를 들어 $\frac{\partial \text{Loss}}{\partial W_{12}^{(1)}}$ 라면 softmax를 중간에 매개 변수로. 따라서 이러한 함수를 미리 외워둡시다.

Sigmoid and softmax

{	Sigmoid	$S(x) = \frac{1}{1 + \exp(-x)}$	→ 해당 value 만으로 대응
	softmax	$P(x_i) = \frac{\exp(x_i)}{\sum \exp(x_i)}$	→ 모든 output을 대응

보통 Sigmoid는 Activation softmax는 K개의 class 각각 probability를 주는 용도로 사용

또한 Sigmoid는 Binary 다. 이를 K개의 class로 확장하면 softmax가 된다. 이를 확장시키는 logit 부위 양자화 합니다.

Logit. = log + odds

odds = $\frac{\text{성공}}{\text{실패}}$ ratio 이라 logit $L = \ln \frac{p}{1-p}$ \Rightarrow 즉, odds는 확률 1 이고 logit은 확률 0 이 된다. \Rightarrow logistic regression 이다.

Sigmoid model의 output에 probability p 를 주고 싶다면 logit을 이용하면 된다.

즉, $L = \ln \frac{p}{1-p}$ 이어서 $L \Rightarrow$ output 이라고 하여 p 에 대해 정리하면 $p = \frac{1}{1 + \exp(-L)}$ 이 된다.

∴ Sigmoid : $f(x) = \frac{1}{1 + e^{-x}}$

Softmax 이는 binary 인 Sigmoid를 multi-class로 확장한 것이다.

Induction

binary case) $e^{t_1} = \frac{p_1}{1-p_1} = \frac{p_1}{p_2}$ where $p_i = P(Y=C_i | x)$ x : given data.

then for k classes let $i=1, 2, \dots, k$ $e^{t_i} := \frac{p_i}{p_k}$, thus $\sum_{i=1}^{k-1} e^{t_i} = \sum_{i=1}^{k-1} \frac{p_i}{p_k} = \frac{1}{p_k} \sum_{i=1}^{k-1} p_i$

and because $\sum_{i=1}^k p_i = 1$ $p_k = 1 - \sum_{i=1}^{k-1} p_i$ ∴ $\sum_{i=1}^{k-1} e^{t_i} = \frac{1-p_k}{p_k}$

∴ $p_k = \frac{1}{\sum_{i=1}^{k-1} e^{t_i} + 1}$, also as we defined $p_i = e^{t_i} p_k = \frac{e^{t_i}}{\sum_{i=1}^{k-1} e^{t_i} + 1}$, because $e^{t_k} = \frac{p_k}{p_k} = 1$

∴ $p_i = \frac{e^{t_i}}{\sum_{i=1}^k e^{t_i}}$

$$1) \text{ Sigmoid} \quad S(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{e^x+1} \quad \Rightarrow \quad \frac{\partial S}{\partial x} = \frac{e^x(1+e^x) - e^{2x}}{(1+e^x)^2} = \frac{e^x}{(1+e^x)^2} = \frac{1 \times e^x}{(1+e^x)(1+e^x)} = S(x) \{1-S(x)\}$$

$$2) \text{ Softmax} \quad P(x_j) = \frac{\exp(x_j)}{\sum \exp(x_i)} \quad \Rightarrow \quad \frac{\partial P}{\partial x_j} = \frac{1}{\{\sum \exp(x_i)\}^2} \left\{ \exp(x_j) \right\} \left\{ \sum \exp(x_i) - \exp(x_j) \right\}$$

~~~~~~ 2-항상식 case 사용.

$$= \frac{\exp(x_j)}{\sum \exp(x_i)} \left\{ \frac{\sum \exp(x_i) - \exp(x_j)}{\sum \exp(x_i)} \right\} = \frac{\exp(x_j)}{\sum \exp(x_i)} \left\{ 1 - \frac{\exp(x_j)}{\sum \exp(x_i)} \right\} = P(x_j) \{1 - P(x_j)\}$$

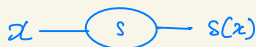
이제 끝 더 생각해  $P(x_j)$  을  $x_u$  에 대해서 미분하는 경우가 있다.

$$\frac{\partial P}{\partial x_u} = \frac{\partial}{\partial x_u} \left\{ \frac{\exp(x_j)}{\sum \exp(x_i)} \right\} = \exp(x_j) \frac{\partial}{\partial x_u} \left\{ \frac{1}{\sum \exp(x_i)} \right\} = \exp(x_j) \left\{ \frac{-\exp(x_u)}{\{\sum \exp(x_i)\}^2} \right\}$$

$$= - \frac{\exp(x_j)}{\sum \exp(x_i)} \frac{\exp(x_u)}{\sum \exp(x_i)} = -P(x_j)P(x_u)$$

$$\therefore \frac{\partial P(x_j)}{\partial x_u} = \begin{cases} P(x_j) \{1 - P(x_j)\} & \text{where } j=u \\ -P(x_j)P(x_u) & \text{where } j \neq u \end{cases} \quad \Leftrightarrow \quad P(x_j) \{I_{(u=j)} - P(x_u)\}$$

여기서 미분할 때 사용하는 것



↑  
여기 미분하는 데 사용해서 결과로 얻어.

999) cross-entropy.  $-\sum_{i=1}^c t_i \log(y_i)$

$$\frac{\partial}{\partial y_j} \left( -\sum_{i=1}^c t_i \log(y_i) \right) = \frac{\partial}{\partial y_j} \left( -t_j \log(y_j) \right) = -t_j \frac{1}{y_j}$$

1000) softmax with cross-entropy

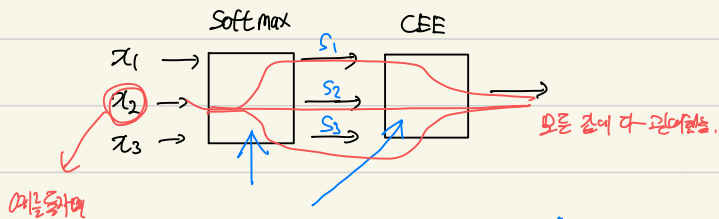
둘이 같이 다하면 상당히 이해하기 쉽고 99) 999) 의 결과를 이용하여 모든 작업 증명해서도 가능할 것이다.

1) 위의 결과 이용 즉, let  $f(x_j) = \text{CEE}(\text{softmax}(x_j))$ , let  $\text{CEE} = L(x)$ ,  $\text{softmax} = p(x)$ ,  $t_j = \text{target}$

$$\Rightarrow \frac{\partial}{\partial x_j} f(x_j) = \frac{\partial f(x_j)}{\partial p(x_j)} \times \frac{\partial p(x_j)}{\partial x_j}$$

Wrong!!

⇒ 이렇게 하면 틀릴 이유



두가지 matrix 이용하면 결과에 도달함.

모든 값이 다 관련됨.

다 됨. 즉, x1에 대한 gradient를 구하려면 s1, s2, s3 을 모두 고려해야 함.

2213 권이 이 방법을 이용해 하면 (by Chain Rule)

$$\frac{\partial}{\partial x_j} f(x_j) = \sum_{l=1}^k \frac{\partial f(x_j)}{\partial P(x_l)} \times \frac{\partial P(x_l)}{\partial x_j} = - \sum_{l=1}^k \frac{t_l}{P(x_l)} \times P(x_l) \{ I_{(l=j)} - P(x_j) \} = - \sum_{l=1}^k t_l \{ I_{(l=j)} - P(x_j) \}$$

$$= \sum_{l=1}^k t_l P(x_j) - \sum_{l=1}^k t_l I_{(l=j)} = P(x_j) \sum_{l=1}^k t_l - \sum_{l=1}^k t_l I_{(l=j)} = P(x_j) - t_j$$

$$\therefore \frac{\partial}{\partial x_j} f(x_j) = P(x_j) - t_j = y_j - t_j \quad \text{where } y_j: \text{output.} \quad \text{보통 이걸지 안하는 여를 일삼는다.}$$

한번이 라는 방법.

$$a_1 \cdots a_K \rightarrow \boxed{\text{softmax}} \rightarrow y_1 \cdots y_K \rightarrow \boxed{\text{CBE}} \rightarrow L(a_1 \cdots a_K) \quad \xrightarrow{\log \text{ loss}}$$

$$L(a_1 \cdots a_K, t_1 \cdots t_K) = - \sum_{i=1}^K t_i \log(y_i) = - \sum_{i=1}^K t_i \times \log \frac{\exp(a_i)}{\sum_{j=1}^K \exp(a_j)} = - \sum_{i=1}^K t_i a_i + \underbrace{\sum_{i=1}^K t_i \log \left\{ \sum_{j=1}^K \exp(a_j) \right\}}_1$$

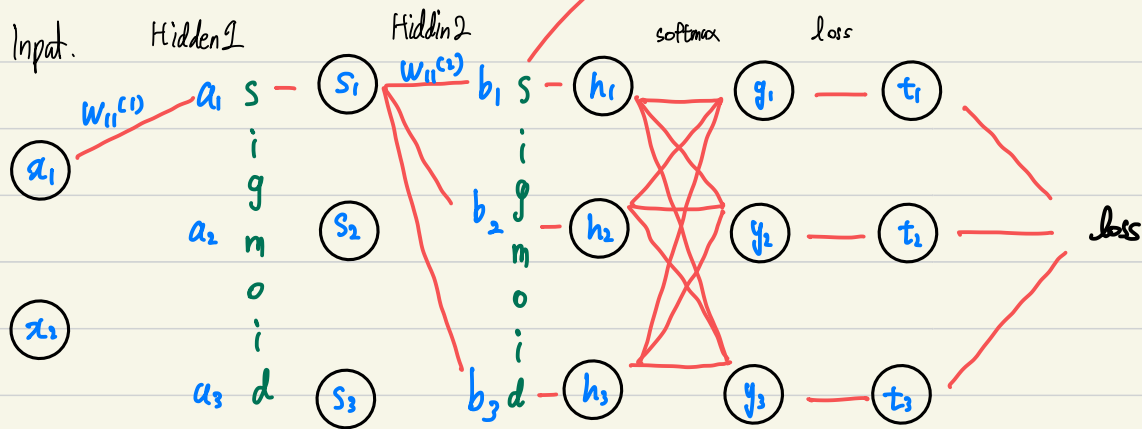
$$\therefore \frac{\partial L}{\partial a_i} = - t_i + \frac{\exp(a_i)}{\sum_{j=1}^K \exp(a_j)} = y_i - t_i$$

$$\therefore \frac{\partial L}{\partial a} = y - t$$

위에서 구한 것은 중간 과정일 뿐입니다. 이제 진짜 필요한  $w_{11}^{(1)}$  을 구해보겠습니다.

다시 처음 그림에서  $w_{11}^{(1)}$  의 gradient 를 구하기 위한 graph

→ 실제 modeling 단계는 없어야 함





상자만 들어올리면

$$\frac{\partial L}{\partial w_{11}^{(1)}} = \frac{\partial L}{\partial s_1} \times \frac{\partial s_1}{\partial a_1} \times \frac{\partial a_1}{\partial w_{11}^{(1)}} = \frac{\partial L}{\partial s_1} \times \{s_1(1-s_1)\} \times \{x_1\}$$

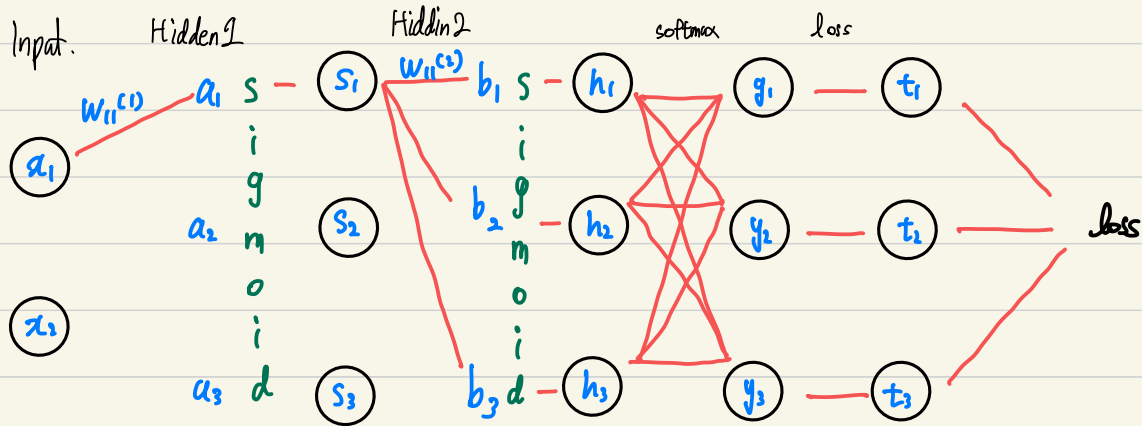
$$\frac{\partial L}{\partial s_1} = \sum_{i=1}^3 \frac{\partial L}{\partial h_i} \times \frac{\partial h_i}{\partial b_i} \times \frac{\partial b_i}{\partial s_1} = \sum_{i=1}^3 \frac{\partial L}{\partial h_i} \times \{h_i(1-h_i)\} \times \{w_{i1}^{(2)}\}$$

$$\frac{\partial L}{\partial h_i} = y_i - t_i$$

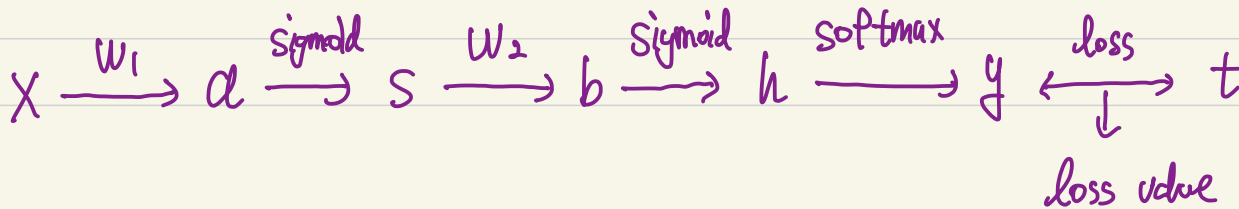
$$\therefore \frac{\partial L}{\partial w_{11}^{(1)}} = \sum_{i=1}^3 \{y_i - t_i\} \times \{h_i(1-h_i)\} \times \{w_{i1}^{(2)}\} \times \{s_1(1-s_1)\} \times \{x_1\}$$

## 5. Matrix form

이러한 구조를 Matrix로 나타내



$\Rightarrow$  by matrix form 이걸로 Linear layer bias == True



Feed forward

S: Sigmoid

P: softmax

$$a = XW_1 + b_1$$

matrix는 항상 form 을 생각해야 함.

$$S = S(a)$$

$$b = SW_2 + b_2$$

$$X: 1 \times 2$$

$$W_1: 2 \times 3$$

$$b_1: 1 \times 3$$

$$\Rightarrow$$

$$a: 1 \times 3$$

$$h = S(b)$$

$$S: 1 \times 3$$

$$W_2: 3 \times 3$$

$$b_2: 1 \times 3$$

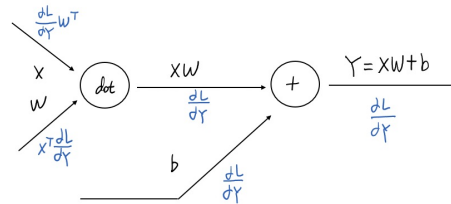
$$\Rightarrow$$

$$b: 1 \times 3$$

$$y = P(h)$$

## 역전파 (chain Rule)

Affine



i.e) gradient at  $w, b$  :  $\frac{dL}{dW} = X^T \frac{dL}{dY}$   
 $\frac{dL}{db} = \frac{dL}{dY}$  , 이진함으로 돌려주는 때:  $\frac{dL}{dx} = \frac{dL}{dY} w^T$

이 과정 즉, 중간에  $Y$ 로 끼임 과정을  
parameterize 하고 생각하면 된다.

## Back propagation

위에서 했던 방법과 모질리 같은 방법이다. (예) 영접법같은 Gauss-elimination 으로 풀고자 할 때)

위에서 했던 방법을 다시 생각해 보면

$X \xrightarrow{W_1} Y \xrightarrow{W_2} Z$  이걸  $\frac{\partial Z}{\partial W_1}$  을 구하기 위해  $Y$ 를 매개변수로 보고 chain rule 을 이용하였다.

즉  $\frac{\partial Z}{\partial W_1} = \frac{\partial Z}{\partial Y} \frac{\partial Y}{\partial W_1}$  을 하였다. 이제  $W_2$  입장에서 보면 미분은 역방향으로 단순히 지나감을 뿐이다.

이걸 matrix form 의 미분은 효율적 계산 graph 를 이용한다. 또한 크게 두개의 미분어 계산되는데 첫번째는

해당 layer의 parameter 에 대한 미분이고 두번째는 이전층으로 흘러들어가는 미분이다.

i) Affine  $Y = XW + b$

$$\frac{\partial L}{\partial W} = \frac{\partial Y}{\partial W} \times \frac{\partial L}{\partial Y} = X^T \frac{\partial L}{\partial Y}$$

$$\frac{\partial L}{\partial b} = \frac{\partial Y}{\partial b} \times \frac{\partial L}{\partial Y} = \frac{\partial L}{\partial Y}$$

}  $\Rightarrow$  해당 layer에서 각각 미분한다.

$$\frac{\partial L}{\partial X} = \frac{\partial L}{\partial Y} W^T$$

이걸 이전층  $X$ 로 흘러들어가는 미분이다.