

2. Neuron Modeling

2-1. Modeling a Neuron

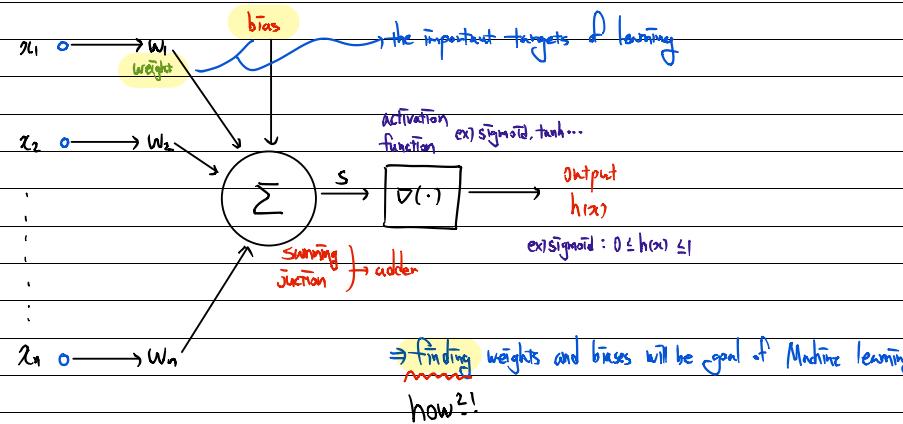
- three basic elements

1. synapses (weights)

2. adder (input vector \rightarrow scalar : $\underline{x} \rightarrow \underline{w}^T \underline{x}$)

→ 우리가 signal을 구하려면 weighted sum ($\underline{w}^T \underline{x}$)를 해야함!

3. activation function (nonlinear)



2.2 gradient descent

- a general technique for minimizing differentiable function

$\Rightarrow \min_{\theta} J(\theta)$: parameter θ 를 변화시킬 때 어떤 값 $J(\theta)$ 를 가장 작게 만드는 방법 ($\because J(\theta)$: cost function)

(Objective function : 학습과정에서 optimize 하려고 하는 function)

Cost function : 전체 dataset에 대해 계산된 loss들의 Avg or sum

모델 학습 시 일반적으로 minimizing 하는 function

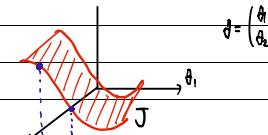
Loss function : 하나의 sample에 대해 모델의 $f(x) - \hat{f}(x)$: 허리를 측정하는 function)

- Idea : $J(\theta)$ 는 parameter space 안의 surface

1. start from somewhere on J

2. roll down the surface, decreasing J step by step

• two things to decide at each step!!



2.2 Gradient descent

- two things to decide at each step!!

① which direction?

⇒ 우리의 cost function이 decreasing 하는 방향!

방향? → $\nabla J(\theta)$: gradient는 경사가 가장 빠르게 증가하는 방향

→ 즉! $-\nabla$ 로 direction 결정하자!

$\therefore -\nabla J(\theta) := \Delta(\theta)$ (현재 위치 θ 에서 $\Delta(\theta)$ 로 몇걸자?)

② how much?

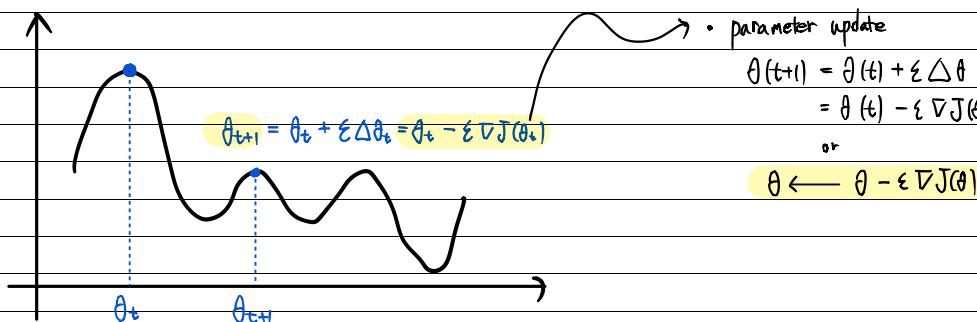
⇒ $N \cdot M \rightarrow f''(\frac{\text{곡률}}{\text{정보})}$

∴ $f'' \Rightarrow$ 곡률大, 경사한 곳 step short

곡률小, 경사한 곳 step long

⇒ but in our case... f'' 를 구하는 건 불가능..

∴ ϵ (learning rate) (hyperparameter)



2.3 Gradient Descent Algorithm

algorithm 1 gradient descent

- initialize θ (초기화)
- while stopping criterion not met do
- sample m examples: $\mathbb{X}_m = \{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$
- compute gradient estimate: $\hat{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_{i=1}^m L(y^{(i)}, \hat{y}^{(i)})$
- apply update: $\theta \leftarrow \theta - \epsilon \hat{g}$
- end while

ex) simple linear regression

$$S = \frac{1}{n} \sum_{i=1}^n (y_i - w^T x_i)^2$$

$$= \frac{1}{n} \sum_{i=1}^n \epsilon^2$$

▷ m forward props

▷ ϵ : learning rate

$$\hat{g} = \frac{dS}{dW}$$

three variants

① $m=1$ (random \hat{g} ?!) sample \mathbb{X}_1 !)

⇒ stochastic GD

② $1 < m < N$: minibatch SGD

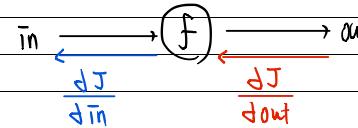
③ $m=N$: batch gradient descent.

2.4 Backprop Demystified

1. Single gate backprop

$$\text{out} = f(\bar{in})$$

*Output → Input
(gradient)*



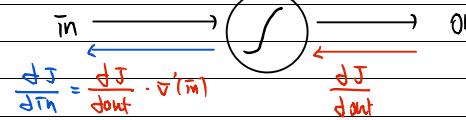
$$\frac{dJ}{d\bar{in}} = \frac{dJ}{dout} \cdot \frac{dout}{d\bar{in}}$$

$$= \frac{dJ}{dout} \cdot \frac{df(\bar{in})}{d\bar{in}}$$

Output local gradient
gradient gradient

2. Sigmoid

$$\text{out} = \sigma(\bar{in}) = \frac{1}{1+e^{-\bar{in}}}$$



$$\sigma(x) = \frac{1}{1+e^{-x}} = (1+e^{-x})^{-1}$$

$$\frac{d\sigma}{dx} = \frac{e^{-x}}{(1+e^{-x})^2}$$

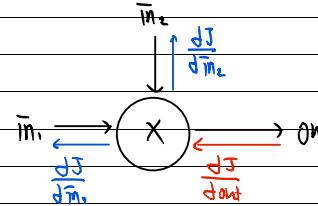
$$1 - \sigma(x) = \frac{e^{-x}}{1+e^{-x}}$$

$$\frac{dJ}{d\bar{in}} = \frac{dJ}{dout} \cdot \frac{dout}{d\bar{in}}$$

$$= \frac{dJ}{dout} \cdot \frac{d\sigma(\bar{in})}{d\bar{in}} = \frac{dJ}{dout} \cdot \sigma'(\bar{in}) = \frac{dJ}{dout} \cdot \sigma(1-\sigma)$$

3. Multiplication

$$\text{out} = \bar{in}_1 \cdot \bar{in}_2$$



$$\frac{dJ}{d\bar{in}_1} = \frac{dJ}{dout} \cdot \frac{dout}{d\bar{in}_1}$$

$$= \frac{dJ}{dout} \cdot \frac{\bar{in}_1 \cdot \bar{in}_2}{\bar{in}_1}$$

$$= \frac{dJ}{dout} \cdot \bar{in}_2$$

$$\frac{dJ}{d\bar{in}_2} = \frac{dJ}{dout} \cdot \frac{dout}{d\bar{in}_2}$$

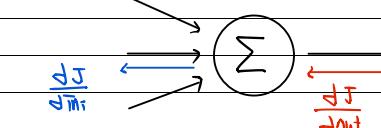
$$= \frac{dJ}{dout} \cdot \frac{\bar{in}_1 \cdot \bar{in}_2}{\bar{in}_2}$$

$$= \frac{dJ}{dout} \cdot \bar{in}_1$$

$$= \frac{dJ}{dout} \cdot \bar{in}_1$$

4. Summation

$$\text{out} = \sum_i \bar{in}_i = \bar{in}_1 + \bar{in}_2 + \dots + \bar{in}_r$$



$$\frac{dJ}{d\bar{in}_1} = \frac{dJ}{dout} \cdot \frac{dout}{d\bar{in}_1}$$

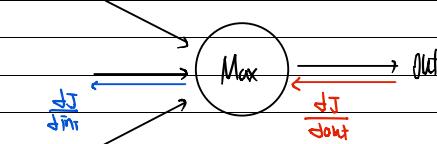
$$= \frac{dJ}{dout} \cdot \frac{\bar{in}_1 + \bar{in}_2 + \dots + \bar{in}_r}{\bar{in}_1}$$

$$= \frac{dJ}{dout} \cdot \frac{1}{r}$$

$$= \frac{dJ}{dout}$$

5. Max (CNN: maxpooling)

$$\text{out} = \max_i \{\bar{in}_i\}$$



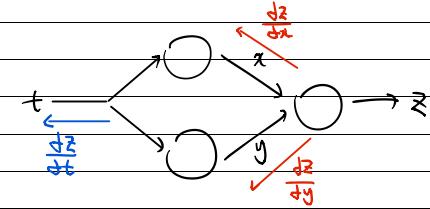
$$\frac{dJ}{d\bar{in}_i} = \frac{dJ}{dout} \cdot \frac{dout}{d\bar{in}_i}$$

$$= \frac{dJ}{dout} \cdot \frac{d\max(\bar{in}_i)}{d\bar{in}_i}$$

$$\frac{d\max(\bar{in}_i)}{d\bar{in}_i} = \begin{cases} 1 & \text{if } \bar{in}_i \text{ is max} \\ 0 & \text{o.w.} \end{cases}$$

⇒ Max (forward) ⇔ Max (backprop)

6. Backprop through fanout



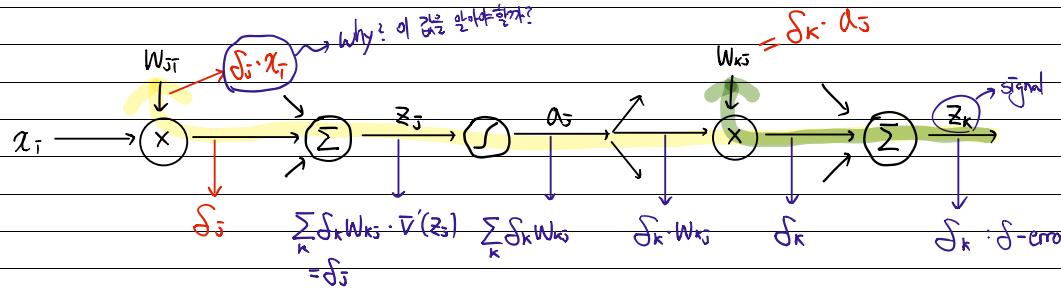
$$\frac{dJ}{dt} = \frac{dJ}{dz} \cdot \frac{dz}{dt} + \frac{dJ}{dz} \cdot \frac{dz}{dy}$$

⇒ fanout (forward) ⇔ Sum (backprop)

⇒ Sum (forward) ⇔ fanout (backprop)

example !!

• Compute $\frac{\partial J}{\partial w_{ji}}$



• parameter update

$$w_{ji} \leftarrow w_{ji} - \eta \nabla$$

을 알아내기 예문!

$$\Rightarrow w_{ji}^{(t+1)} \leftarrow w_{ji}^{(t)} - \eta \cdot \frac{\partial J}{\partial w_{ji}}$$

$$= w_{ji}^{(t)} - \eta \cdot (\delta_j \cdot x_i)$$

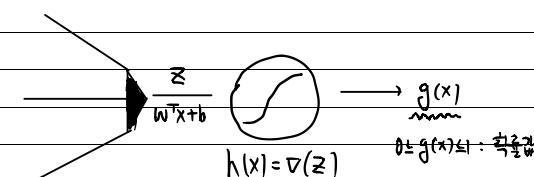
$$= w_{ji}^{(t)} - \eta \cdot \left(\sum_k \delta_k \cdot w_{kj} \right) \cdot v'(z_j) \cdot x_i$$

2.5 logistic regression

• logistic regression (in early neural nets: logistic regression unit = Neuron)

: outputs probability of a binary response. ex) 0 vs 1 / dead or alive
return soft labels (probability) \Rightarrow allow uncertainty

→ output: real (like regression) but bounded (like classification, probability)



• problem

$$\text{given : } x \in \mathbb{R}^n$$

$$\text{Want : } \hat{y} = p(y=1 | x)$$

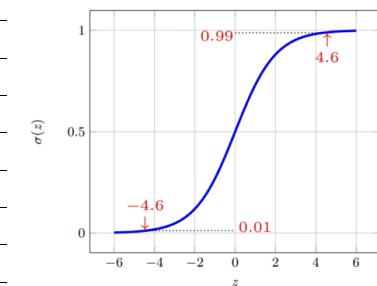
• model

$$\hat{y} = v(w^T x + b), \quad w \in \mathbb{R}^n, b \in \mathbb{R}$$

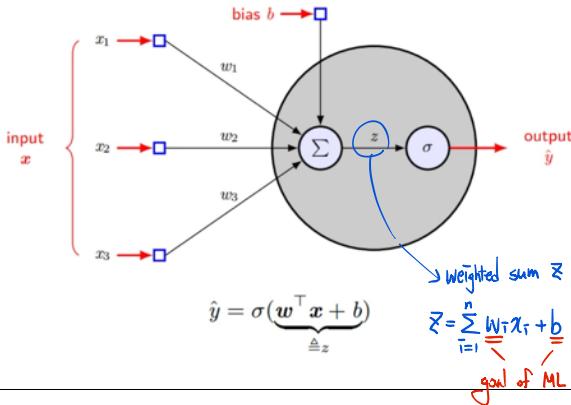
• sigmoid

$$v(z) = \frac{1}{1 + e^{-z}}$$

$$v'(z) = v(z)(1 - v(z))$$



- Computational graph



- Probabilistic interpretation

- given training set drawn from P(data)

$$X = \{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$$

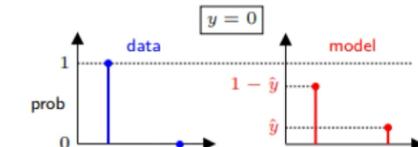
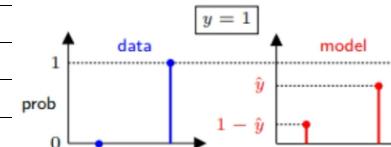
Ex) 49 4명 유무

→ consider label $y^{(i)} \in \{0, 1\}$ as hard probability

$$y \equiv P(y=1|x) \Rightarrow P(y=0|x) = 1 - P(y=1|x) = 1-y$$

then! observed pmf $p_{\text{data}} \in \{y \equiv P(y=1|x), 1-y \equiv P(y=0|x)\}$

fitted pmf $p_{\text{model}} \in \{\hat{y}, 1-\hat{y}\}$



- Loss (pointwise)

- We use log loss

$$L(y, \hat{y}) = -y \log \hat{y} - (1-y) \log (1-\hat{y})$$

If $y=1 \Rightarrow L(y, \hat{y}) = -\log \hat{y}, \hat{y} \rightarrow 1, L(y, \hat{y}) \rightarrow 0$
 If $y=0 \Rightarrow L(y, \hat{y}) = -\log(1-\hat{y}), \hat{y} \rightarrow 0, L(y, \hat{y}) \rightarrow 0$

★ $y = P(y=1|x)$ ★

- Cost function

$$J(w, b) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})]$$

for update parameter w & b , we need to compute $\frac{\partial J}{\partial w}, \frac{\partial J}{\partial b}$
 $(\because w_j \leftarrow w_j - \eta \nabla \frac{\partial J}{\partial w})$

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial z} \cdot \frac{\partial z}{\partial w} = (\hat{y}^{(i)} - y^{(i)}) \cdot (x_1^{(i)}) = \frac{1}{m} \sum (\hat{y}^{(i)} - y^{(i)}) x_1^{(i)} = \frac{1}{m} \sum (\nabla(w^T x^{(i)} + b) - y^{(i)}) x_1^{(i)}$$

$$\frac{\partial J}{\partial b} = \frac{1}{m} \sum (\nabla(w^T x^{(i)} + b) - y^{(i)}) \cdot 1$$

2.6 Training by Backprop

Training a neuron

- find w and b that minimize cost function $J(w, b)$
- use iterative optimization i.e. gradient descent

forward → output → output vs label → error → delta → gradient → backward

Minibatch!

$$\bar{z} = \bar{x}^T \bar{w} + b$$

$$\begin{bmatrix} w_1 & w_2 \\ b \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + b = w_1 x_1 + w_2 x_2 + b$$

$$\nabla (w_1 x_1 + w_2 x_2 + b) = a - \hat{y}$$

Training logistic regression by backprop

- minimize cost function by gradient descent.

$$J(w, b) = -\frac{1}{m} \sum (y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)}))$$

- repeat over training examples

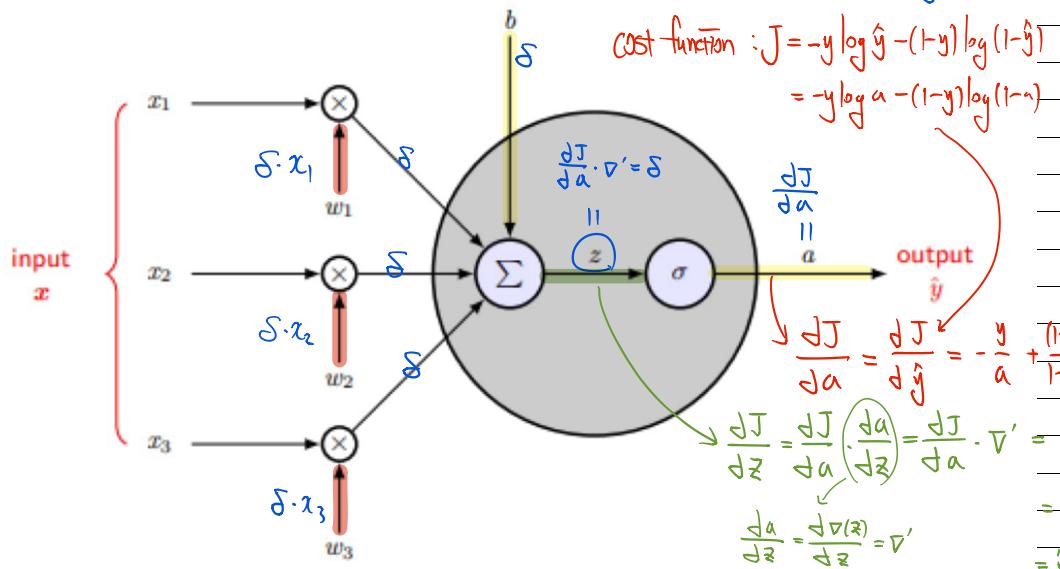
ε : learning rate (hyperparameter)

$$w \leftarrow w - \varepsilon \nabla J_w$$

$$b \leftarrow b - \varepsilon \nabla J_b$$

$$\begin{aligned} \frac{\partial J}{\partial w_1} &= \frac{1}{m} \sum x_1^{(i)} \delta^{(i)} \\ \frac{\partial J}{\partial w_2} &= \frac{1}{m} \sum x_2^{(i)} \delta^{(i)} \end{aligned}$$

forward prop:



$$\begin{aligned} z &= w^T x + b \\ a &= \sigma(w^T x + b) = \hat{y} \end{aligned}$$

$$\text{cost function: } J = -y \log \hat{y} - (1-y) \log (1-\hat{y})$$

$$= -y \log a - (1-y) \log (1-a)$$

- SGD equation (1 example)

$$\Rightarrow J(w, b) = L(y, a) = -y \log a - (1-y) \log (1-a)$$

$$\frac{\partial J}{\partial w_i} = x_i \cdot \delta = x_i \cdot \frac{\partial J}{\partial z} = x_i(a-y)$$

$$\frac{\partial J}{\partial b} = \delta = \frac{\partial J}{\partial z} = (a-y)$$

$$\Rightarrow \text{apply update} \quad w_i \leftarrow w_i - \varepsilon \frac{\partial J}{\partial w_i} = w_i - \varepsilon x_i \cdot \delta = w_i - \varepsilon x_i(a-y)$$

$$b \leftarrow b - \varepsilon \frac{\partial J}{\partial b} = b - \varepsilon (a-y)$$

$$= \frac{(a-y)}{a} + \frac{(1-y)}{1-a}$$

$$= a-y = \delta \quad (\text{delta error})$$

$$= \hat{y} - y$$

Arranging minibatch

- two options to arrange m examples

► in columns

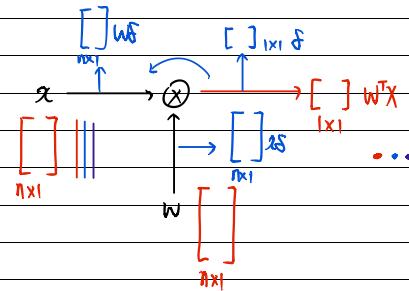
$$X_{(1)} = \begin{bmatrix} x_1^{(1)} & x_1^{(2)} \\ x_2^{(1)} & \dots & x_2^{(n)} \\ x_3^{(1)} & & x_3^{(n)} \\ x_n^{(1)} & & x_n^{(n)} \end{bmatrix} = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(n)} \end{bmatrix}$$

► in rows (design matrix)

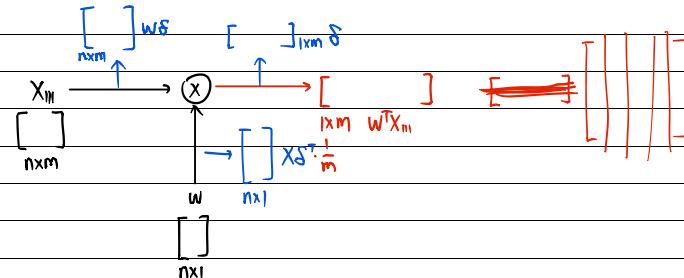
$$X = \begin{bmatrix} \cdots x^{(1)\top} \cdots \\ \vdots \\ \cdots x^{(m)\top} \cdots \end{bmatrix}$$

Weighted Sum

- 1 example (column)

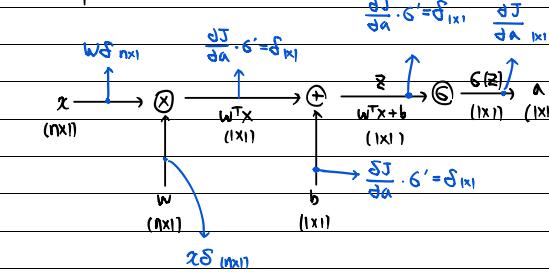


- Size-m minibatch (column)

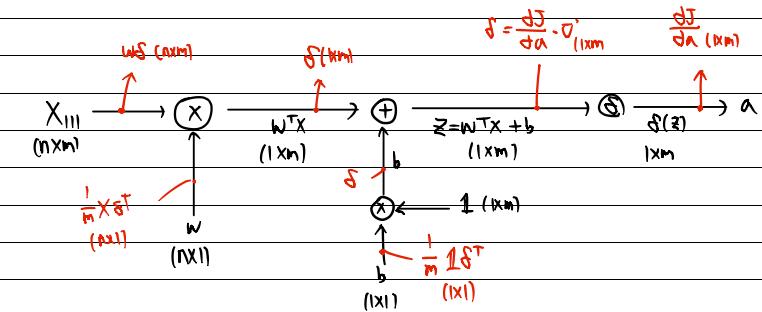


Forward / backward prop

- 1 example



- Size m minibatch



$$\delta = \frac{\partial J}{\partial a} \cdot \sigma'(z) \quad \frac{\partial J}{\partial a} \quad (1 \times m)$$

$$\frac{\partial J}{\partial a} \quad (1 \times m)$$