

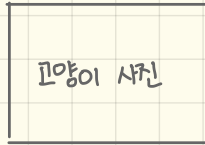
Neural Network and Deep Learning

[2-1. Logistic Regression as a Neural Network]

<Binary Classification>

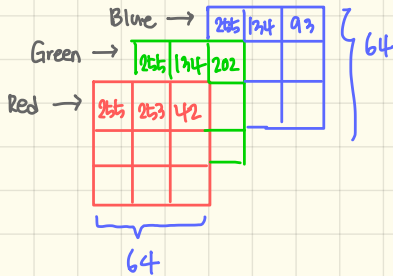
* Binary Classification: 그렇다 / 아니다 2개의 분류하는 것임, 결과가 '그렇다'이면 1임, '아니다'이면 0임 판별

ex)



→ 1 (cat) vs 0 (non cat)
y

x



$$X = \begin{bmatrix} 255 \\ 253 \\ 42 \\ \vdots \\ 255 \\ 124 \\ 202 \\ \vdots \\ 255 \\ 124 \\ 93 \end{bmatrix}$$

이므로, 전체치원: $64 \times 64 \times 3 = 12,288$

$$\therefore n = n_x = 12288$$

* notation

- $(x, y) : x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$

- m training example: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$M_{\text{train}} = m$ / $M_{\text{test}} = \# \text{test example}$

$$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix} \begin{matrix} \uparrow \\ n_x \\ \uparrow \end{matrix} \Rightarrow X \in \mathbb{R}^{n_x \times m}$$

X: shape = (n_x, m)

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix} \Rightarrow Y \in \mathbb{R}^{1 \times m}, Y \text{ shape} = (1, m)$$

<Logistic Regression>

* Logistic Regression: y 가 0 또는 1로 정해져 있는 이진 분류 문제에 사용되는 알고리즘

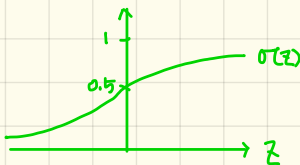
X : 입력특성, y : 주어진 x 에 해당하는 실제 값, \hat{y} : y 의 예측값 $= P(y=1|x)$

이때, $x \in \mathbb{R}^{n_x}$ 이고 parameter: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$ 이므로

$$\Rightarrow \text{Output } \hat{y} = \sigma(w^T x + b) = \sigma(z)$$

↑ $\hat{y} = w^T x + b$ 로 계산시, $0 \leq \hat{y} \leq 1$ 의 조건을 만족하지 않을 수 있으므로
시그모이드 함수를 통하여 0과 1의 범위로 맞춰준다.

$$* \text{시그모이드 함수: } \sigma(z) = \frac{1}{1 + e^{-z}}$$



$$\therefore \hat{y} = \sigma(w^T x + b), \text{ where } \sigma(z) = \frac{1}{1 + e^{-z}}$$

given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$, want $\hat{y}^{(i)} \approx y^{(i)}$

<Logistic Regression Cost Function>

* 손실함수(Loss Function): 하나의 입력특성(x)에 대한 실제값(y)과 예측값(\hat{y})의 오차를 계산하는 함수

- 일반적으로 손실함수는 보통 $L(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$ 으로 사용하나, 로지스틱 회귀는 지역최소값에 빠질 수 있으므로 다른 손실함수를 사용한다.

- 로지스틱 회귀의 손실함수: $L(\hat{y}, y) = -y \log \hat{y} + (1-y) \log (1-\hat{y})$

i) If $y=1$: $L(\hat{y}, y) = -\log \hat{y} \Rightarrow$ want $\log \hat{y}$ large \therefore want \hat{y} large ≈ 1

ii) If $y=0$: $L(\hat{y}, y) = -\log (1-\hat{y}) \Rightarrow$ want $\log (1-\hat{y})$ large \therefore want \hat{y} small ≈ 0

* 비용함수(Cost Function): 하나의 입력에 대한 오차를 계산하는 함수가 손실함수이며, 모든 입력에 대한 오차를 계산하는 함수를 비용 함수라고 한다. 즉, 비용함수는 모든 입력에 대한 손실함수의 평균

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

★ 중요한 점: 비용함수의 값이 작아지도록 w 와 b 를 찾는 것이 목표! ★

< Gradient Descent >

* 경사하강법: 비용함수를 최소화 만드는 w 와 b 를 찾는 것므로, 임의의 점을 골라서 기울기를 따라 최적의 값으로 업데이트 함



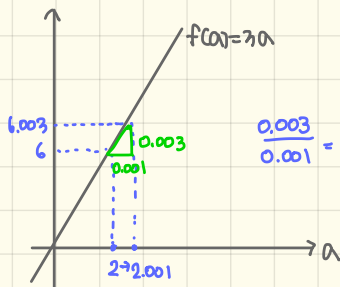
learning rate

$$- w : w - \alpha \frac{dJ(w, b)}{dw}$$

도함수 dw라고 함

$$- b : b - \alpha \frac{dJ(w, b)}{db}$$

< Derivatives (미분) >



i) $a=2 \rightarrow f(a)=6$

$a=2.001 \rightarrow f(a)=6.003$

$\frac{0.003}{0.001} = \frac{\text{height}}{\text{width}} = \text{slope (derivative) of } f(a) \text{ at } a=2 : 3$

ii) $a=5 \rightarrow f(a)=15$

$a=5.001 \rightarrow f(a)=15.003$

slope at $a=5 : 3$

* 도함수 (= slope): 변수 a 를 조금 변화시켰을 때, $f(a)$ 가 얼마나 변화하러 측정하는 것

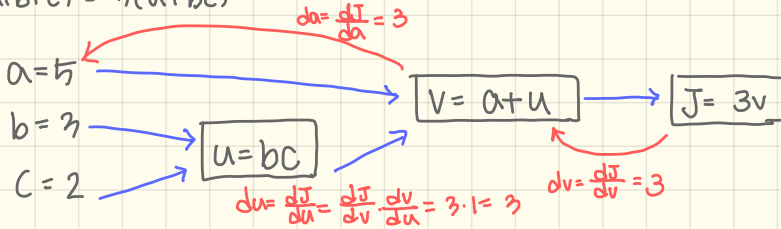
$$\frac{d}{da} f(a) = \frac{df(a)}{da} \quad \text{또} \quad f'(a)$$

< More Derivatives Examples >

함수 : $f(a)$	a^2	a^3	$\ln(a)$
도함수 : $\frac{d}{da} f(a)$	$2a$	$3a^2$	$\frac{1}{a}$

< Computation Graph >

ex) $J(a, b, c) = 3(a + bc)$



< Derivatives with a Computer Graph >

* 이론의 연쇄법칙: 합성함수의 도함수에 대한 공식

- 위의 예에서 $\frac{dJ}{da} = \frac{dJ}{dv} \cdot \frac{dv}{da} = 3 \cdot 1 = 3$ 으로 구할 수 있다.

- $\frac{dJ}{db} = \frac{dJ}{du} \cdot \frac{du}{db} = 3 \cdot c = 3 \cdot 2 = 6$ / $\frac{dJ}{dc} = \frac{dJ}{du} \cdot \frac{du}{dc} = 3 \cdot b = 9$

- 표기법: $d \text{ var} = \frac{d \text{ final output var}}{d \text{ var}}$

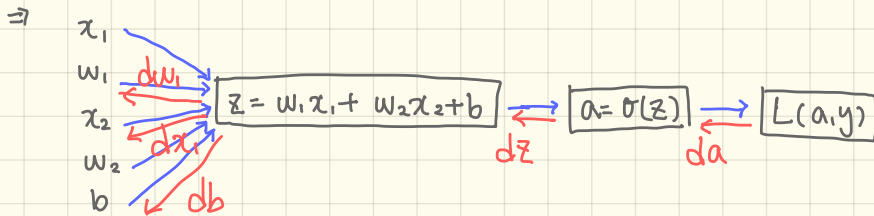
< Logistic Regression Gradient Descent >

$z = w^T x + b$

$\hat{y} = a = \sigma(z)$

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

z } Computer graph로 정리해보면,
(변위상 x_1, x_2 만 있다고 생각)



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

$$\pi) \quad dz = \frac{dL(a,y)}{dz} = \frac{dL(a,y)}{da} \cdot \frac{da}{dz} = \left\{ -\frac{y}{a} + \frac{1-y}{1-a} \right\} \{a(1-a)\}$$

$$= a-y$$

$$\bar{\pi}) \quad dw_1 = \frac{dL(a,y)}{dw_1} = x_1 dz$$

$$dw_2 = x_2 dz$$

$$db = dz$$

$$\left. \begin{array}{l} \Rightarrow \\ \end{array} \right\} \begin{array}{l} w_1 := w_1 - \alpha dw_1 \\ w_2 := w_2 - \alpha dw_2 \\ b := b - \alpha db \\ \text{또 업데이트 된다.} \end{array}$$

< Gradient Descent on m Examples >

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)}),$$

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

$$\text{If, } J=0 ; dw_1=0 ; dw_2=0 ; db=0$$

For $i=1$ to m ,

$$z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 += x_1^{(i)} dz^{(i)}$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

$\left. \begin{array}{l} \Rightarrow \\ \end{array} \right\} \text{특성 개수}(n) = 2 \text{ 라면 } J \text{ 점}$

[2-2. Python and Vectorization]

< Vectorization >

* 벡터화: 코딩에서 for문을 제거하는 기술

ex) $z = w^T x + b$

(Non-vectorization)

$$z = 0$$

for i in range($n-x$):

$$z += w[i] * x[i]$$

$$z += b$$

(Vectorization)

$$z = np.dot(w, x) + b$$

* SIMD (Single Instruction Multiple Data) : 병렬 프로세스의 한 종류로, 하나의 명령어로 여러개의 값을 동시에 처리하는 방식으로 벡터 연산을 가능하게 함

ex2) $v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$ $u = \begin{bmatrix} e^{v_1} \\ \vdots \\ e^{v_n} \end{bmatrix}$

(non-vectorization)

$$u = np.zeros((n,1))$$

for i in range(n):

$$u[i] = math.exp(v[i])$$

(vectorization)

$$u = np.exp(v)$$

< Vectorizing Logistic Regression >

$$z^{(1)} = w^T x^{(1)} + b$$

$$a^{(1)} = \sigma(z^{(1)})$$

$$z^{(2)} = w^T x^{(2)} + b$$

$$a^{(2)} = \sigma(z^{(2)})$$

$$z^{(n)} = w^T x^{(n)} + b$$

$$a^{(n)} = \sigma(z^{(n)})$$

...

$$\Rightarrow X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix}$$

$$[z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}] = w^T X + [b \ b \ \dots \ b]$$

$$= \underbrace{w^T x^{(1)} + b}_{= z^{(1)}} \cdot \underbrace{w^T x^{(2)} + b}_{= z^{(2)}} \cdot \dots \cdot \underbrace{w^T x^{(m)} + b}_{= z^{(m)}}$$

(vectorization coding)

$$z = np.dot(np.transpose(w), x) + b$$

< Vectorizing Logistic Regression's Gradient Descent

i) dz

$$dz^{(1)} = a^{(1)} - y^{(1)} \quad dz^{(2)} = a^{(2)} - y^{(2)} \quad \dots \quad dz^{(m)} = a^{(m)} - y^{(m)}$$

$$\Rightarrow dz = [dz^{(1)} \quad dz^{(2)} \quad \dots \quad dz^{(m)}]$$

$$A = [a^{(1)} \quad a^{(2)} \quad \dots \quad a^{(m)}]$$

$$Y = [y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(m)}]$$

$$\Rightarrow dz = A - Y$$
$$= [a^{(1)} - y^{(1)} \quad a^{(2)} - y^{(2)} \quad \dots \quad a^{(m)} - y^{(m)}]$$

ii) dw

$$dw = 0$$

$$dw \text{ += } x^{(1)} dz^{(1)}$$

$$dw \text{ += } x^{(2)} dz^{(2)}$$

$$\dots$$
$$dw \text{ += } x^{(m)} dz^{(m)}$$

$$dw \text{ /= } m$$

$$\Rightarrow dw = \frac{1}{m} X dz^T$$
$$= \frac{1}{m} \begin{bmatrix} x^{(1)} & \dots & x^{(m)} \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$
$$= \frac{1}{m} [x^{(1)} dz^{(1)} + \dots + x^{(m)} dz^{(m)}]$$

iii) db

$$db = 0$$

$$db \text{ += } dz^{(1)}$$

$$db \text{ += } dz^{(2)}$$

\dots

$$db \text{ += } dz^{(m)}$$

$$db \text{ /= } m$$

$$\Rightarrow db = \frac{1}{m} \sum_{i=1}^m dz^{(i)}$$
$$= \frac{1}{m} \text{np.sum}(dz)$$
$$= \text{np.sum}(dz)$$

non-vectorizing

$$J = 0, \quad dw_1 = 0, \quad dw_2 = 0, \quad db = 0$$

for $i = 1$ to m :

$$z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log(1 - a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 += x_1^{(i)} dz^{(i)} \quad \} \Rightarrow dw += x^{(i)} * dz^{(i)}$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

$$J = J/m, \quad dw_1 = dw_1/m, \quad dw_2 = dw_2/m$$

$$db = db/m$$

Vectorizing

$$Z = w^T X + b = \text{np.dot}(w^T, X) + b$$

$$A = \sigma(Z)$$

$$dz = A - Y$$

$$dw = \frac{1}{m} X dz^T$$

$$db = \frac{1}{m} \text{np.sum}(dz)$$

$$w = w - \alpha dw$$

$$b = b - \alpha db$$

< Broadcasting in Python >

ex1)

$$A = \begin{bmatrix} 56.0 & 0.0 & 4.4 & 68.0 \\ 1.2 & 104.0 & 52.0 & 8.0 \\ 1.8 & 175.0 & 99.0 & 0.9 \end{bmatrix} \rightarrow \text{cal} = A.\text{sum}(\text{axis}=0)$$

$$\text{percentage} = 100 * A / \text{cal}.\text{reshape}(1,4)$$

\uparrow (1,4) \uparrow (1,4)

⇒ (1,4) 행렬을 (1,4)로 나눔

ex2)

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \cancel{100} \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

확장

ex3)

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

(1,m) → (m,n)
확장

$$\begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}$$

ex4)

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

(m,1) → (m,n)
확장

$$\begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix}$$

∴ (m,n) $\begin{pmatrix} + \\ - \\ * \\ \div \end{pmatrix}$ (1,n) ~> (m,n)
 matrix $\begin{pmatrix} + \\ - \\ * \\ \div \end{pmatrix}$ (m,1) ~> (m,n) 또는 확장하여 계산

<Explanation of Logistic Regression Cost Function>

logistic regression cost function

$$\hat{y} = \sigma(w^T x + b) \quad \text{when } \sigma(z) = \frac{1}{1+e^{-z}}$$

$$\rightarrow \hat{y} = P(y=1|x)$$

$$\text{if } y=1 : P(y|x) = \hat{y}$$

$$\text{if } y=0 : P(y|x) = 1 - \hat{y}$$

$$\Rightarrow P(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)}$$

$$\text{if } y=1 : P(y|x) = \hat{y}^1 (1-\hat{y})^0 = \hat{y}$$

$$\text{if } y=0 : P(y|x) = \hat{y}^0 (1-\hat{y})^1 = 1-\hat{y}$$

코스트함수는 엔트로피가 함수이므로,

$$\log P(y|x) = \log (\hat{y}^y (1-\hat{y})^{(1-y)}) = \underbrace{y \log \hat{y} + (1-y) \log (1-\hat{y})}_{\text{최대화}} = -L(y, \hat{y}) : \underbrace{\text{손실함수}}_{\text{최대화}}$$

$$\therefore \text{손실함수 } L(\hat{y}, y) = -\log P(y|x) = -(\hat{y}^y (1-\hat{y})^{(1-y)})$$

m개의 훈련데이터를 학습시키면,

$$P(\text{labels of training set}) = \prod_{i=1}^m P(y^{(n)} | x^{(n)})$$

$$\xrightarrow{\text{로그}} \log P(\text{labels of training set}) = \log \prod_{i=1}^m P(y^{(n)} | x^{(n)}) = -\sum_{i=1}^m L(\hat{y}^{(n)}, y^{(n)})$$

$$\therefore \text{비용함수 } J(w, b) = -\log P(\text{labels of training set})$$

$$= \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(n)}, y^{(n)})$$