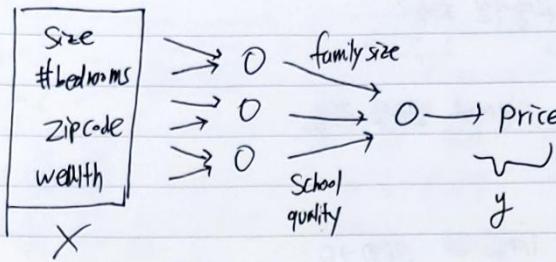


<Neural Network>

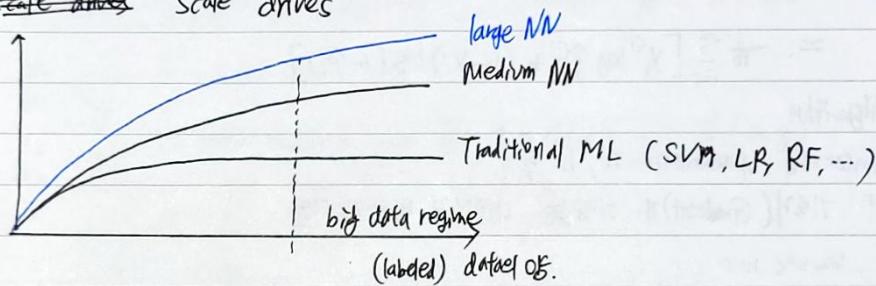


Input(x)	Output(y)	Application
Image	Object	Photo tagging ← CNN
Audio	Text	음성인식 ← t-series, RNN (for sequence data)
Image, radar info	다른 차의 위치	자율주행 ← Custom, Hybrid NN

▷ Structured data: 명확한 뜻이 정의됨 (size, # of bedrooms, price, -)

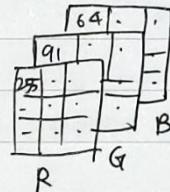
▷ Unstructured data: 명확한 뜻이 정의되지 않음. (→ pixel 값, 소리의 tone 등으로 변환해야 함).

▷ ~~Scale~~ Scale drives



<Binary off with LR> <2주차>

Cat image →



$$X = \begin{bmatrix} R \\ G \\ B \end{bmatrix}$$

전체 차원:
64x64x3 = 12288
↑↑↑
X는 Y의 RGB
모두 4열
 $N_x = 12288$

▷ 학습 데이터 표기

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

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

= M_{train}, M_{test} = # of test example.

$$X = \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(m)} \\ | & | & | & | \end{bmatrix}$$

← 신경망 구현시 원본 이미지 표기해야 하면.
= $n_x \times m$ matrix

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

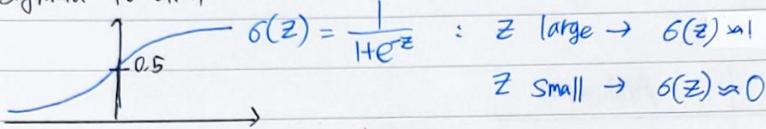
▷ Logistic Regression

Given X , want $\hat{y} = P(y=1|X)$: 고양이일 확률?

parameter: $W \in \mathbb{R}^{n \times d}$, $b \in \mathbb{R}$

Output $\hat{y} = \sigma(W^T X + b) \leftarrow (W^T X + b)$ 에 sigmoid 함수를 적용,
 $0 \leq \hat{y} \leq 1$ 로 막음.

• Sigmoid function



▷ Cost function of LR

at $\hat{y} = \sigma(W^T X + b)$ where $\sigma(z) = \frac{1}{1+e^{-z}}$,

• Loss function: $L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$ ← Cross-entropy loss

① If $y=1$: $L(\hat{y}, y) = -\log \hat{y}$ ← $\log \hat{y}$ 이 크게, \hat{y} 가 크게 (1에 최대한 가깝게)

② If $y=0$: $L(\hat{y}, y) = -\log(1-\hat{y})$ ← $\log(1-\hat{y})$ 이 크게, $1-\hat{y}$ 가 크게 (0에 최대한 가깝게)

• Cost function: $J(W, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$ Loss function의 평균.

최소화

$$= -\frac{1}{m} \sum_{i=1}^m [y_i \log \hat{y}_i + (1-y_i) \log(1-\hat{y}_i)]$$

▷ Gradient Descent Algorithm

- $J(w, b)$ 를 minimize하는 parameter w, b 찾기.
- Initial point w_0 초기 기울기 (Gradient)가 가장 높은 내리막길 방향으로 이동.

repeat { learning rate

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

↑ partial derivative
univariate case

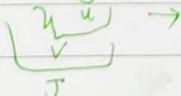
$$b := b - \alpha \frac{dJ(w, b)}{db}$$

• Chain Rules

$$a \xrightarrow{\text{affect}} v \rightarrow J \quad \text{일때, } \frac{dJ}{da} = \frac{dJ}{dv} \cdot \frac{dv}{da}$$

▷ 계산 그래프 (산출 그래프, Computational Graph)

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



$$\begin{aligned} u &= bc \\ v &= a + u \\ J &= 3v \end{aligned}$$

$$\begin{array}{c} a=5 \\ b=3 \\ c=2 \end{array} \xrightarrow{[u=bc]} \boxed{v=a+u} \xrightarrow{[J=3v]} \boxed{J=3v}$$

→ J를 계산하기 위해 이동

↳ J 계산, forward, left to right

• Derivative 계산 (\leftarrow backwards, right to left)

$$\begin{array}{l} a=5 \\ b=3 \\ c=2 \end{array} \xrightarrow{\quad} \boxed{U=b} \xrightarrow{\quad} \boxed{V=a+u} \xrightarrow{\quad} \boxed{J=3V}$$

$$\frac{dJ}{db} = \frac{dJ}{du} \cdot \frac{du}{db}$$

$$= 3 \cdot 1 = 3$$

$$J = 3V$$

$$V = 11 \rightarrow 11.001$$

$$\begin{aligned} "dc" &= \frac{dJ}{dv} \cdot \frac{dv}{du} \cdot \frac{du}{db} = 3 \cdot 1 \cdot c \\ &= 3c \quad \therefore \frac{dJ}{da} = 3 = "da" \quad \frac{dJ}{dv} = 3 = "dv" \end{aligned}$$

$$J = 33 \rightarrow 33.003$$

* backward propagation을 위해 코드를 쓰는 경우,

Final output (제일 신경) Variable이 있으면 (신경그래프의 마지막 노드, J)

$$\frac{d \text{Final output var}}{d \text{Var}} = "dvar", \text{ 최종 결과값 변수의 derivative}$$

▷ Logistic Regression을 위한 GDM (Single training model)

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

$$\hat{y} = a = \sigma(z) = \frac{1}{1+e^{-z}}$$

$$L(a, y) = -(y \log(a) + (1-y) \log(1-a)), \quad J(w, b) = \frac{1}{m} \sum L(a, y)$$

$$\begin{array}{l} x_1 \\ w_1 \\ x_2 \\ w_2 \\ b \end{array} \xrightarrow{\quad} z = w_1 x_1 + w_2 x_2 + b \rightarrow \hat{a} = a = \sigma(z) \rightarrow L(a, y)$$

② "dz" = $\frac{dL}{dz} = \frac{dL(a,y)}{dz}$

① $\frac{dL(a,y)}{da}$

• backward 계산

$"dA" = -\frac{y}{a} + \frac{1-y}{1-a}$

$"dz" = a - y$

$"dA" = \frac{da}{dz} \cdot da$

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

③ w, b 를 얼마나 바꿔야하는지 계산.

$$\frac{\partial L}{\partial w_1} = "dw_1" = x_1 \cdot dz, \quad \frac{\partial L}{\partial w_2} = dw_2 = x_2 \cdot dz, \quad \dots, \quad \frac{\partial L}{\partial b} = db = dz$$

④ Update

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

▷ Vectorization

$$Z = W^T x + b \text{에서}$$

$$W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \end{bmatrix}, \quad X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix} \quad W, X \in \mathbb{R}^{n_r}$$

- non-vectorized:

$$Z=0$$

for i in range(n_r)

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

$$Z += b \quad (\text{slow})$$

- Vectorized - numpy API

$$\cdot Z = \underbrace{\text{np.dot}(W, X) + b}_{W^T X}$$

$$\Leftrightarrow \text{np.outer}(U, V) = UV^T : ixj \text{ matrix}$$

$$\cdot U = e^V \quad \cdot U = \log V$$

$$\rightarrow U = \text{np.exp}(V) \quad \rightarrow U = \text{np.log}(V)$$

▷ Vectorizing LR

- M 개의 training example에서

$$(Z_1 = W^T X_1 + b, Z_2 = W^T X_2 + b, \dots, Z_m = W^T X_m + b) \quad X = \begin{bmatrix} \vdots & 1 & 1 & 1 \\ X_1 & | & X_2 & | & X_3 & \dots & | & X_m \\ \vdots & | & | & | & | & & | & | \end{bmatrix}, \quad W^T = [w_1, w_2, \dots, w_m], \quad b = [b, b, b, \dots, b]$$

$$\hookrightarrow [Z_1, Z_2, \dots, Z_m] = W^T X + [b, b, b, \dots, b]$$

$$\cdot Z = \text{np.dot}(W^T, X) + b \quad \cdot T: \text{전치 method.} \quad / A^*b: \text{'element-wise' 곱.}$$

$$\cdot A = [a_1, a_2, \dots, a_m] = \sigma(Z) \quad : 같은 크기의 행렬끼리만 적용 가능.$$

- M = np로 생성한 배열일 때

$$(N = \text{np.reshape}(M, 4, 25)) \quad 4 \times 25 \text{ 배열로 생성} \quad (2차원 배열 - (행, 열))$$

$$(N = M.reshape(4, 25)) \quad / ' - 1 ': 형태에 맞게 자동 설정$$

ex) Numpy에서 빠르게 생성

$$\text{np.array}([1, 2, 3]) \quad // 1차원 배열$$

$$\text{np.array}([[1, 2, 3]]) \quad // 2차원 배열 - Row vector$$

$$\text{np.array}([[1], [2], [3]]) \quad // 2차원 배열, column vector$$

▷ Broadcasting → 형태를 맞추기 위해 copy됨.

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

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 400 & 500 & 600 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 404 & 505 & 606 \end{bmatrix}$$

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

* , 뼈대생성시 rank 1 array는 만들지 말것.

- np.random.rand(5) → (5,1)
- ↓
- shape: (5,) ↓
- (5,1)
- (rank 1 array)

• np.random.rand(a,b,c)

: 입력된 숫자 차원의 random number array를 생성.

- assert(a.shape == (5,1)) ← dimension을 정확히 모을 경우.

- 3차원 배열: .reshape(행, 행, 열)

- np.dot(a,b), a = m×n matrix

$$m(n \times k) = m(k) \rightarrow \text{dot} = m \times k \text{ matrix}$$

- algorithm input으로 image의 뼈대 생성시, (length*height*3, 1)의 형태로 reshaping 필요.

- .shape method는 tuple을 반환.

- 3차원 array에서 shape[0] = 층수

shape[1] = 행수

shape[2] = 열수

- 2차원: shape[0] = 행수

shape[1] = 열수

- axis = 0 : 세로방향 (column) 계산 : each column 계산

axis = 1 : 가로방향 (row) 계산 : each row 계산

(3차원 데이터에서는 axis=0, 1, 2는 각각 층, 행, 열 방향 계산)

▷ Jupyter 실습

① EDA) train_x : image of (height, width, 3)

↓ ↓
num_px num_px

- train_set_x_orig: ndarray, shape of (M_train, height, width, 3) // dimension, shape 찾기

$$\begin{matrix} X \\ (a, b, c, d) \end{matrix} \longrightarrow X_{\text{flatten}} = X_{\text{reshape}}(X_{\text{shape}}[0], -1).T \quad // \text{reshaping}$$

$\begin{matrix} \uparrow \\ a \end{matrix}, \begin{matrix} \uparrow \\ b \cdot c \cdot d \end{matrix}, \begin{matrix} \uparrow \\ T \end{matrix} \text{ transpose}$

- standardize

$$- \frac{X - \bar{X}}{\sigma}$$

- Picture dataset에는 대신 pixel channel의 차이값인 255로 나누는 것이 편함.

② Building NN

1. model의 구조 정의

2. parameter 초기화

3. Loop : • current loss - forward propagation

• current gradient - back propagation

• Update parameters - GD.

- np.zeros(shape), shape = (0,0) 형태로 대입할 것.
- isinstance(x, float) - 자료형 검사, True or False 반환
 int
 str
 list
 dict
 :

• Forward / Backward Propagation

① Forward

- get X

$$\cdot \text{compute } A = \sigma(W^T X + b) \quad (= (a_1, a_2, \dots, a_m))$$

$$\cdot \text{calc cost function } J = -\frac{1}{m} \left\{ \sum_{i=1}^m y_i \log(a_i) + (1-y_i) \log(1-a_i) \right\}$$

② Backward

$$\cdot dW = \frac{\partial J}{\partial W} = \frac{1}{m} X(A-Y)^T \quad // (1/m) \cdot \text{np.dot}(X, ((A-Y).T))$$

$$\cdot db = \frac{\partial J}{\partial b} = \frac{1}{m} \sum_{i=1}^m (a_i - y_i) \quad // (1/m) \cdot \text{np.sum}(A-Y)$$

• Update

$$\cdot W = W - dW * \text{learning rate}$$

$$b = b - db * \text{learning rate}$$

▷ 사용된 함수

$S \leftarrow \text{Sigmoid}(z)$

$X_train_shape[0] = n_x$

$w, b \leftarrow \cdot \text{Initialize_with_zeros(dim)} : (dim, 1)$ 의 0행렬 W, $b=0$ 로 초기화

$\text{grads}, \text{cost} \leftarrow \cdot \text{propagate}(w, b, X, Y) : \text{Cost function 계산, 이에 따른 } dw \text{와 } db \text{ 계산. } \rightarrow \text{grads} = \{dw, db\}$

$\text{params}, \text{grads}, \text{costs} \leftarrow \cdot \text{Optimize}(w, b, X, Y, \text{num_iter}, \text{learning_rate}) : w \text{와 } b \text{ Update}$

• propagate에서 grad, cost 뿐만

• w, b update

• params (w, b)의 dict 와 grads(dw, db)의 dict, costs 반환

$Y_{\text{prediction}} \leftarrow \cdot \text{predict}(w, b, X) : \text{최종 } (w, b) \text{ 를 사용해서 } X \text{ 의 클래스 예측.}$

└ W = $n_x \times 1$ matrix

X: $n_x \times m$

A: $W^T X$ 과 차운 동일 $\rightarrow 1 \times m$ matrix

* Cost 함수 사용화 = 얼마나 잘 학습했는지 판단.

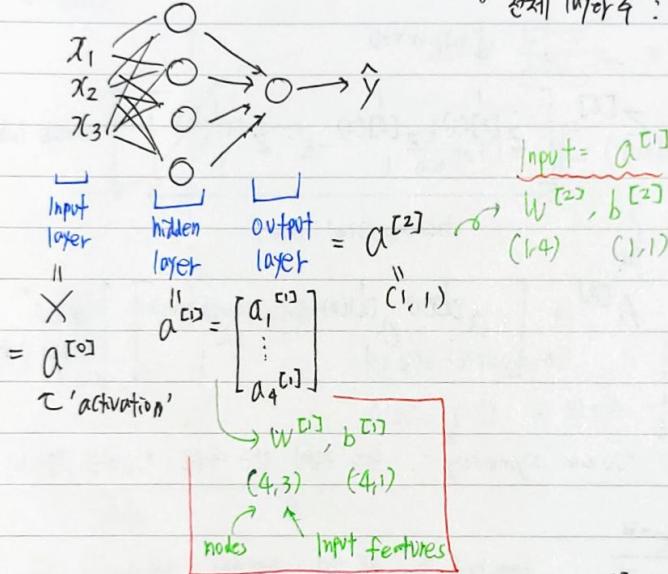
[3주차. Shallow neural network]

$$\begin{array}{c} x \\ \backslash \\ W^{[0]} \\ / \\ b^{[0]} \end{array} \quad (z^{[0]} = W^{[0]}x + b^{[0]}) \rightarrow (a^{[0]} = \sigma(z^{[0]})) \Rightarrow z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \rightarrow a^{[2]} \rightarrow L(a^{[2]}, y)$$

$dW^{[2]} \leftarrow \frac{\partial L}{\partial z^{[2]}} \quad db^{[2]} \leftarrow \frac{\partial L}{\partial b^{[2]}}$

> NN Representation

• 전체 layer 수 : hidden # + Output #



> Computing Output

• hidden layered node

$$Z = W^T x + b \quad a = \sigma(Z)$$

$$\begin{aligned} Z_1^{[0]} &= W_1^{[0]T} x + b_1^{[0]} \\ a_1^{[0]} &= \sigma(Z_1^{[0]}) \quad] a_i^{[0]} \leftarrow \# \text{layer} \\ Z_2^{[0]} &= W_2^{[0]T} x + b_2^{[0]} \\ a_2^{[0]} &= \sigma(Z_2^{[0]}) \quad] a_i^{[0]} \leftarrow \text{ith node} \end{aligned}$$

• Vectorizing

$$Z^{[0]} = \begin{bmatrix} -W_1^{[0]T} \\ -W_2^{[0]T} \\ -W_3^{[0]T} \\ -W_4^{[0]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[0]} \\ \vdots \\ b_4^{[0]} \end{bmatrix} = \begin{bmatrix} W_1^{[0]T} x + b_1^{[0]} \\ \vdots \\ W_4^{[0]T} x + b_4^{[0]} \end{bmatrix} = \begin{bmatrix} Z_1^{[0]} \\ \vdots \\ Z_4^{[0]} \end{bmatrix} \xrightarrow{\sigma} \begin{bmatrix} a_1^{[0]} \\ \vdots \\ a_4^{[0]} \end{bmatrix}$$

$W^{[0]}$ $b^{[0]}$

(4,3) (3,1) (4,1) (4,1) (4,1)

$$Z^{[2]} = W^{[2]}(a^{[1]}) + b^{[2]} \xrightarrow{\sigma} a^{[2]} = \hat{y}$$

(1,1) (1,4) (4,1) (1,1) (1,1)

- multiple training examples 경우에는

$$\begin{aligned} X &\rightarrow a^{[2]} = \hat{y} \\ \hookrightarrow X^{(1)} & \quad a^{[2](1)} \\ ; & \quad ; \\ X^{(m)} & \quad a^{[2](m)} \end{aligned}$$

for $i=1$ to m :

$$\begin{aligned} z^{[1](i)} &= W^{[1]}x^{(i)} + b^{[1]} \\ a^{[1](i)} &= \sigma(z^{[1](i)}) \\ z^{[2](i)} &= W^{[2]}a^{[1](i)} \cdot b^{[2]} \\ a^{[2](i)} &= \sigma(z^{[2](i)}) \\ \rightarrow z^{[2]} &= W^{[2]}X + b^{[2]} \\ A^{[2]} &= \sigma(z^{[2]}) \end{aligned}$$

$$W^{[2]} = \begin{bmatrix} \equiv & \dots & \equiv \end{bmatrix}, 1\text{마다 다른 값을 가짐.}$$

$$X = \begin{bmatrix} x^{(1)}_1 & x^{(1)}_2 & \dots & x^{(1)}_m \\ | & | & \dots & | \\ x^{(2)}_1 & x^{(2)}_2 & \dots & x^{(2)}_m \end{bmatrix}_{(n_x, m)}, \text{features}$$

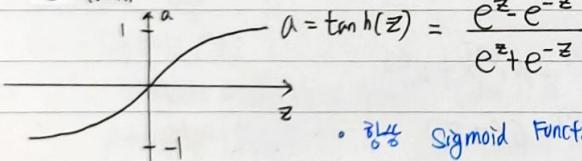
$$z = wx + b$$

$$\begin{aligned} z^{[L]} &= \begin{bmatrix} z^{[L](1)} & z^{[L](2)} & \dots & z^{[L](m)} \end{bmatrix} \\ &\quad \leftarrow \text{training setel Index} \quad \rightarrow \text{nodes index} \\ A^{[L]} &= \begin{bmatrix} a^{[L](1)}_1 & a^{[L](2)}_1 & \dots & a^{[L](m)}_1 \end{bmatrix} \quad \downarrow \# \text{ of hidden unit} \end{aligned}$$

Activation Function

$$a = g(z), g \text{는 non-linear 함수}$$

① tanh



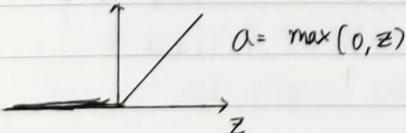
• 향상 Sigmoid Function보다 좋은 성능을 보임.

: 값이 $-1 \sim +1$ 사이에 위치하면서, 활성화수의 평균값이 hidden layer에서 계산될 때 0에 가깝기 때문.

- 데이터를 0으로 centering 할 때 유용함.
- 하지만 출력층의 경우, \hat{y} 가 0~1사이에 위치해야하기 때문에 binary classification의 output layer에서는 Sigmoid 사용.

② Relu

- Sigmoid / tanh 함수에서 그가 매우 크거나 작은 경우에 기울기가 0이 가기 위험을 보임.



* hidden layer에서는 ReLU, Output layer에서는 Sigmoid 사용 (binary clf)

$$\begin{aligned} \cdot dW &= dz \cdot \frac{\partial z}{\partial w} \\ \cdot db &= dz \cdot \frac{\partial z}{\partial b} \end{aligned}$$

• Formula ✘

forward propagation)

$$z^{[1]} = W^{[0]}x + b^{[0]}$$

$$A^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = W^{[1]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(z^{[2]})$$

$$(z^{[L]} = W^{[L]}A^{[L-1]} + b^{[L]})$$

$$(A^{[L]} = g^{[L]}(z^{[L]}))$$

Back propagation)

$$dz^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{m} \cdot dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} \cdot \text{np.sum}(dZ^{[2]}, \text{axis}=1, \text{keepdims=True})$$

$$dZ^{[0]} = \underbrace{W^{[0]T} dZ^{[2]}}_{(n^{[0]}, m)} * \underbrace{g^{[1]'}(z^{[1]})}_{(n^{[0]}, m)} // \text{element-wise product}$$

$$dW^{[0]} = \frac{1}{m} dZ^{[0]} X^T$$

$$db^{[0]} = \frac{1}{m} \text{np.sum}(dZ^{[0]}, \text{axis}=1, \text{keepdims=True})$$

• 단일 그래프 (single example)

$$\begin{array}{c} x \\ \xrightarrow{W^{[0]} \atop b^{[0]}} \boxed{z^{[1]} = W^{[0]}x + b^{[0]}} \rightarrow \boxed{a^{[1]} = g(z^{[1]})} \rightarrow \boxed{z^{[2]} = W^{[1]}a^{[1]} + b^{[2]}} \rightarrow \boxed{a^{[2]} = g(z^{[2]})} \rightarrow \boxed{J(a^{[2]}, y)} \end{array}$$

$$dZ^{[1]} = W^{[1]T} dZ^{[2]} * g'(z^{[1]})$$

$$\begin{matrix} W^{[2]} \\ b^{[2]} \end{matrix} \downarrow$$

$$dZ^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dZ^{[2]} a^{[1]T}$$

← / m 은 multi-example 일 때.

$$\text{dim } [W^{[2]} = (n^{[2]}, n^{[1]})]$$

$$b^{[2]} = (n^{[2]}, 1)$$

$$[z^{[2]}, dZ^{[2]} = (n^{[2]}, m)]$$

$$b^{[1]} = (n^{[1]}, 1)$$

$$[W^{[1]} = (n^{[1]}, n^{[0]})]$$

$$(n^{[1]}, m) = (n^{[1]}, n^{[2]}). (n^{[2]}, m) * (n^{[2]}, m)$$

$$dW^{[0]} = dZ^{[0]} \cdot X^T$$

$$db^{[0]} = dZ^{[0]}$$

▷ Random Initialization

• Weight - 랜덤 생성, 아주 작은 값으로 initialize / 0 일 경우 모든 hidden Unit이 symmetric 이 되므로 오류.

Weight $W^{[0]} = \text{np.random.rand}((n^{[1]}, n^{[0]})) * 0.01, n^{[0]} = n_x$

bias $b^{[0]} = \text{np.zeros}((n^{[1]}, 1))$ \hookrightarrow layer의 개수에 따라 달라짐.

$w^{[0]}, b^{[0]}, \dots$ 도 동일하게 Initialize.

→ 'break symmetry'로 뉴런이 같은 것을
제거지 않도록 설정.

* $a_4^{[2]}$: 2nd layer의 4th neuron의 activation Output

$a^{[2](4)}$: 2nd layer의 4th training example의 activation Output

▷ 실습

- def layer_sizes(X, Y)
 - return: n_x, n_h, n_y
- initialize_parameters(n_x, n_h, n_y) : $W_1 \sim b_2$ 정의
 - return: parameters = { }
- forward_propagation(X, parameters) : A2와 cache(Z, A의 dictionary)
 - return: A2, cache = { }
- Compute_cost(A2, Y, parameters)
 - return: cost
- backward_propagation(parameters, cache, X, Y) : dW, db의 dict, grad 계산.
 - return: grads = { }
- Update_parameters(parameters, grads, learning_rate) : parameter dict를 update.
 - return: parameters = { }
- NN_model(X, Y, n_h, num_iterations)
 - Initialize parameters
 - loop
 - Forward prop
 - Cost computing
 - back propagation
 - Update params with GD
 - return: parameters = { }
- predict(parameters, X) : 선출된 parameter를 forward propagation으로 A2, cache 선출 후 \hat{Y} 계산.

〈4주차, Deep Neural Network〉

▷ Notation

$L = 4$ (# layers)

$n^{[l]}$ = # Units in layer l

$$\alpha^{[l]} = g^{[l]}(z^{[l]}) \quad , \quad \alpha^{[0]} = X, \alpha^{[L]} = \hat{y}$$

- layer가 깊어질수록 더 복잡한 feature의 연산 수행.

▷ Forward Propagation

- 1부터 L -th layer까지 for loop 사용,

$$\begin{cases} z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]} & \rightarrow \text{'cache'로 저장, 나중의 backward step에 쓰임.} \\ A^{[l]} = g^{[l]}(z^{[l]}) \end{cases}$$

▷ Parameter dimension

$$X : (n^{[0]}, 1)$$

$$\begin{matrix} z^{[1]} = W^{[1]} X + b^{[1]} \\ (n^{[0]}, 1) \quad \downarrow \quad \downarrow \\ (n^{[0]}, n^{[1]}) \quad (n^{[0]}, 1) \quad (n^{[1]}, 1) \end{matrix}$$

$$Z^{[1]} = (n^{[1]}, m) \rightarrow \alpha^{[1]} = (n^{[1]}, m)$$

$$W^{[1]} = (n^{[0]}, n^{[1]})$$

$$W^{[l]} = (n^{[l]}, n^{[l-1]}) \quad b^{[l]} = (n^{[l]}, 1)$$

broadcasted

dZ, dA : same dimension as Z, A

* Shallow, many hidden nodes < Deep, less hidden nodes.

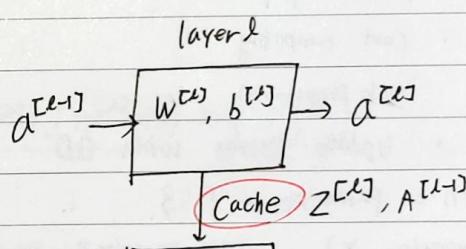
▷ Forward & Backward Functions

$$\text{Layer } l : W^{[l]}, b^{[l]}$$

Forward:

- Input: $a^{[l-1]}$

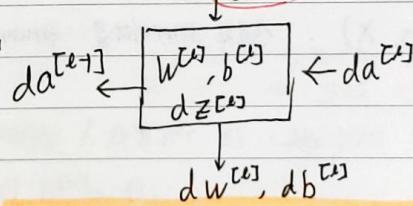
- Output: $\alpha^{[l]}$, 'cache' $Z^{[l]}$



Backward:

- Input: $da^{[l]}$, Cache($Z^{[l]}$)

- Output: $da^{[l-1]}$, $dW^{[l]}$, $db^{[l]}$



$$\begin{cases} dz^{[l]} = da^{[l]} * g^{[l]}'(Z^{[l]}) \\ dW^{[l]} = dZ^{[l]} \cdot A^{[l-1]T} \\ db^{[l]} = dZ^{[l]} \\ da^{[l-1]} = W^{[l]T} \cdot dZ^{[l]} \end{cases}$$

(vectorized)

$$dz^{[l]} = dA^{[l]} * g^{[l]}'(Z^{[l]}) = W^{[l+1]T} \cdot dZ^{[l+1]} + \dots$$

$$dW^{[l]} = \frac{1}{m} dZ^{[l]} \cdot A^{[l-1]T}$$

$$db^{[l]} = \frac{1}{m} \cdot \text{np.sum}(dZ^{[l]}, axis=1, keepdims=True)$$

$$da^{[l-1]} = W^{[l]T} \cdot dZ^{[l]}$$

▷ Parameters: W, b

Hyper parameters:

- | |
|--|
| α
iter
hidden layers, # hidden units
activation function |
| |

▷ 실습

- #1. . initial_parameters_deep(layer_dims) : W, b 에 대한 dict, parameter return.
- . model : [Linear part (L-1 layer)+Relu] + Sigmoid part (Output layer)
- . linear_forward(A, w, b) : return Z, cache = (A, w, b) L-1 까지, 0 ~ L-1
- . linear_activation_forward(A_prev, w, b, activation)
 - helper func: sigmoid, relu : input Z, return A, activation-cache
 - return A, cache (= linear-cache, activation-cache)
- . L_model_forward(X, parameters) : linear-activation forward with relu n-1 번 진행
 → with sigmoid 1번 진행.
 - $l=1$ to $l=L-1$ 까지 $A_{\text{prev}}=A$, linear_activ_forward(A_{prev} , parameters['w l '], p l ['b l '], act=relu)
 - $l=L$ 까지 linear_activ_forward(A , \sim , activ=sigmoid) 진행 → AL, caches 반환
- . linear_backward(dZ, cache) : return dA_prev, dW, db
 ↳ (A_{prev} , W, b)
- . linear_activation_backward(dA, cache, activation) : return dA_prev, dW, db
 - * cache는 linear-cache, activation-cache를 나누어서 대입
 - * linear, activation 부분 한꺼번에 backward 하는 경우.
- . Lmodel_backward(AL, Y, caches)
 - . dAL 계산
 - . $L \rightarrow L-1$ layer: sigmoid
 - . $L-2 \rightarrow 0$: Relu
 - . $dA^{[l]}$ input → Output $dA^{[l-1]}, dW^{[l]}, db^{[l]}$
 - . Update_parameters(parameters, grads, learning_rate) : return parameters