

Deep Learning

1

1. Neural networks and Deep Learning

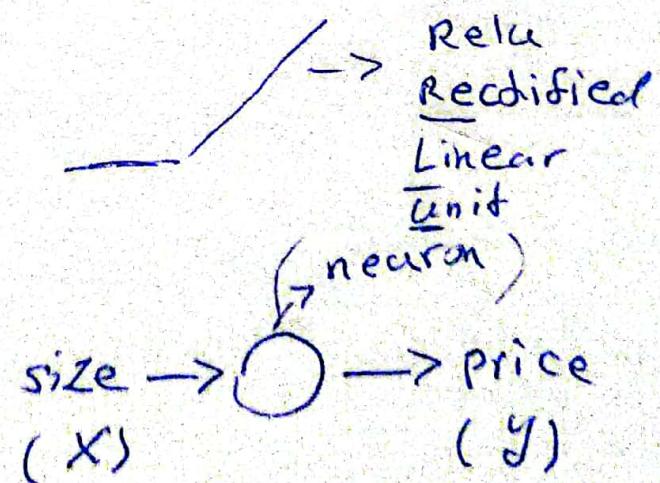
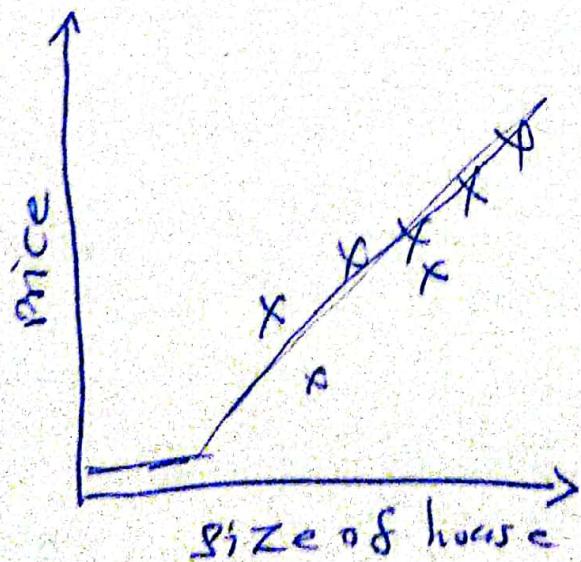
2. Implementing Deep Neural Networks: Hyperparameter tuning,
Regularization and Optimization.

3. Structuring Your Machine Learning project

4. Convolutional Neural Networks

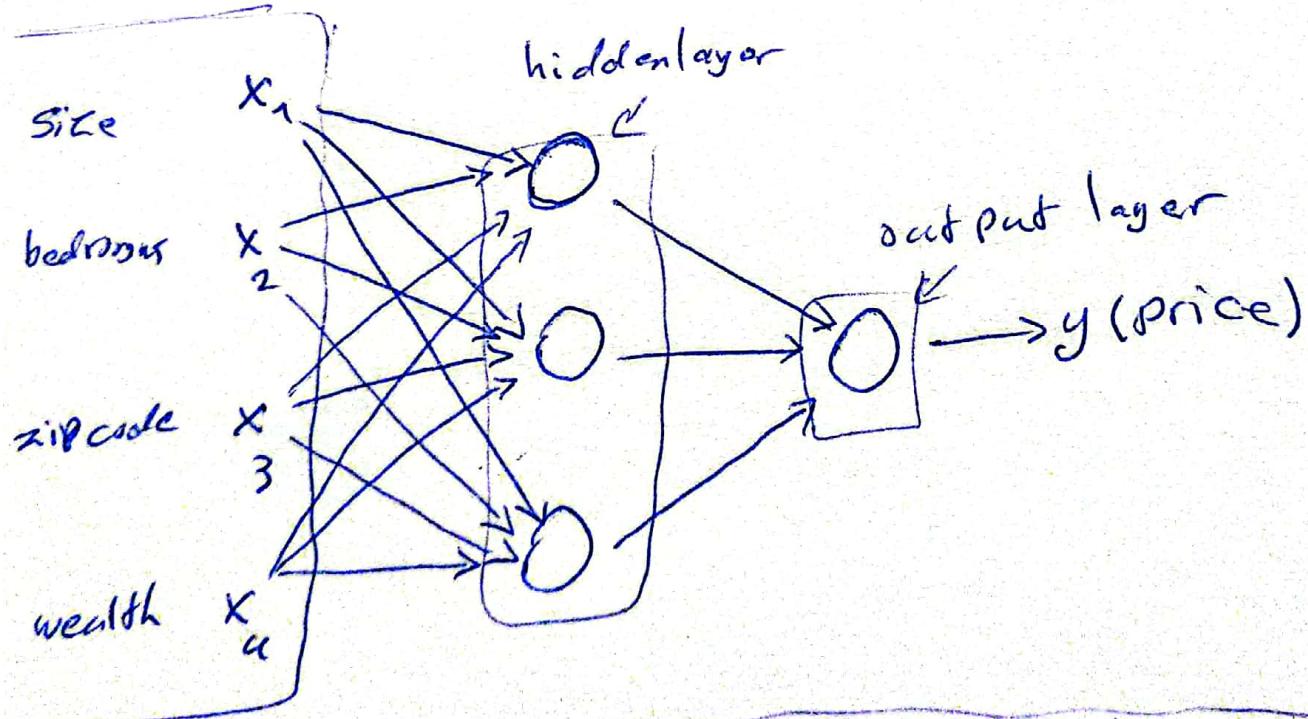
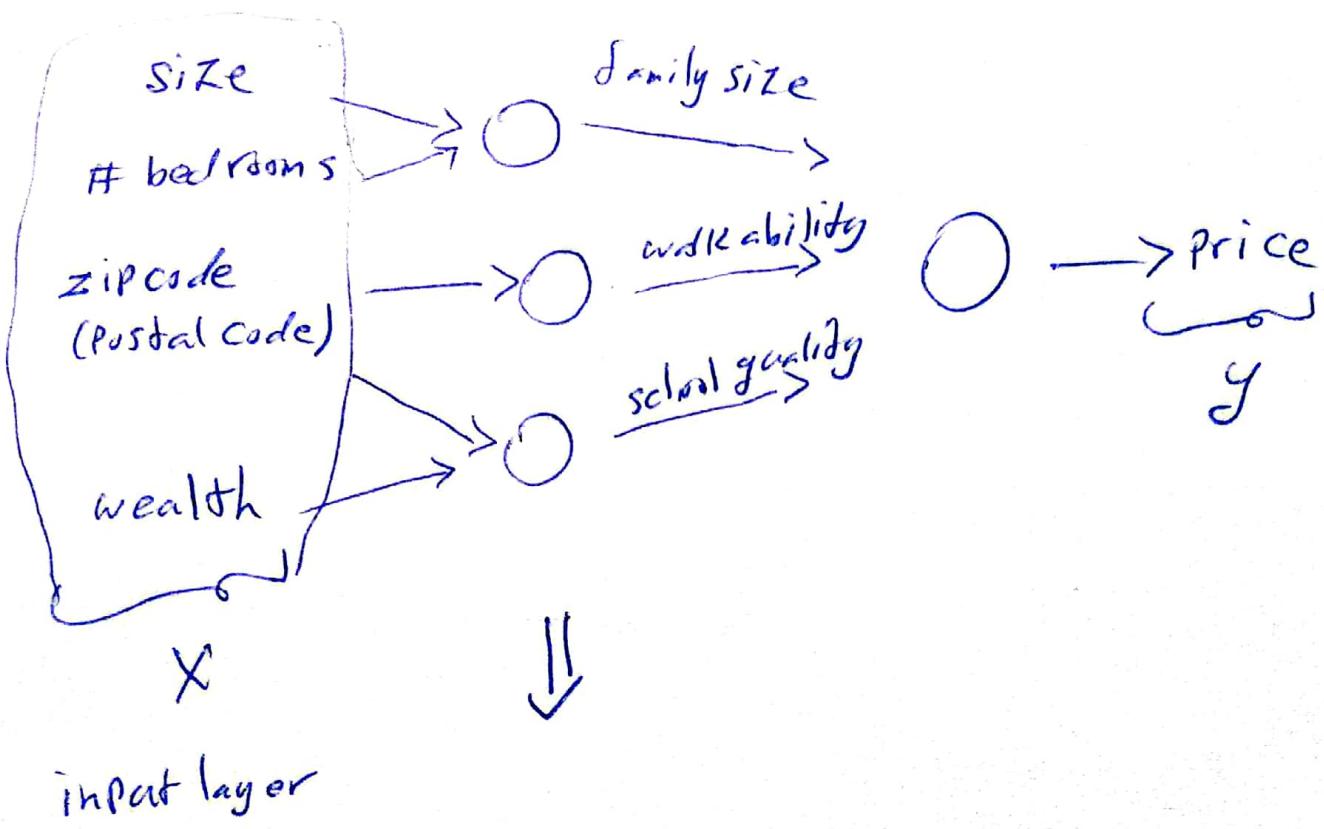
5. Natural Language Processing : Building sequence
models (RNNs, LSTM)

what is a neural network?



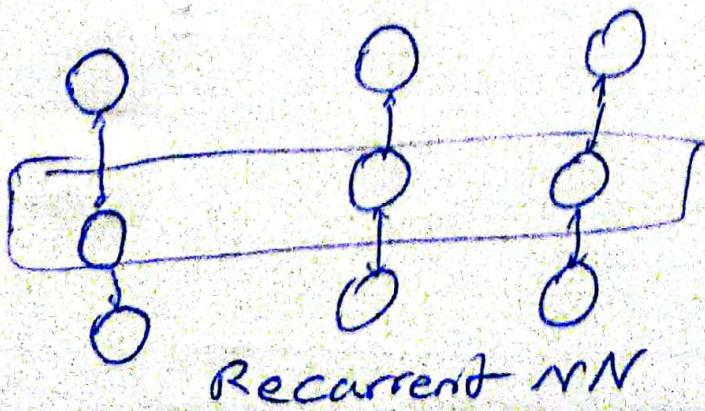
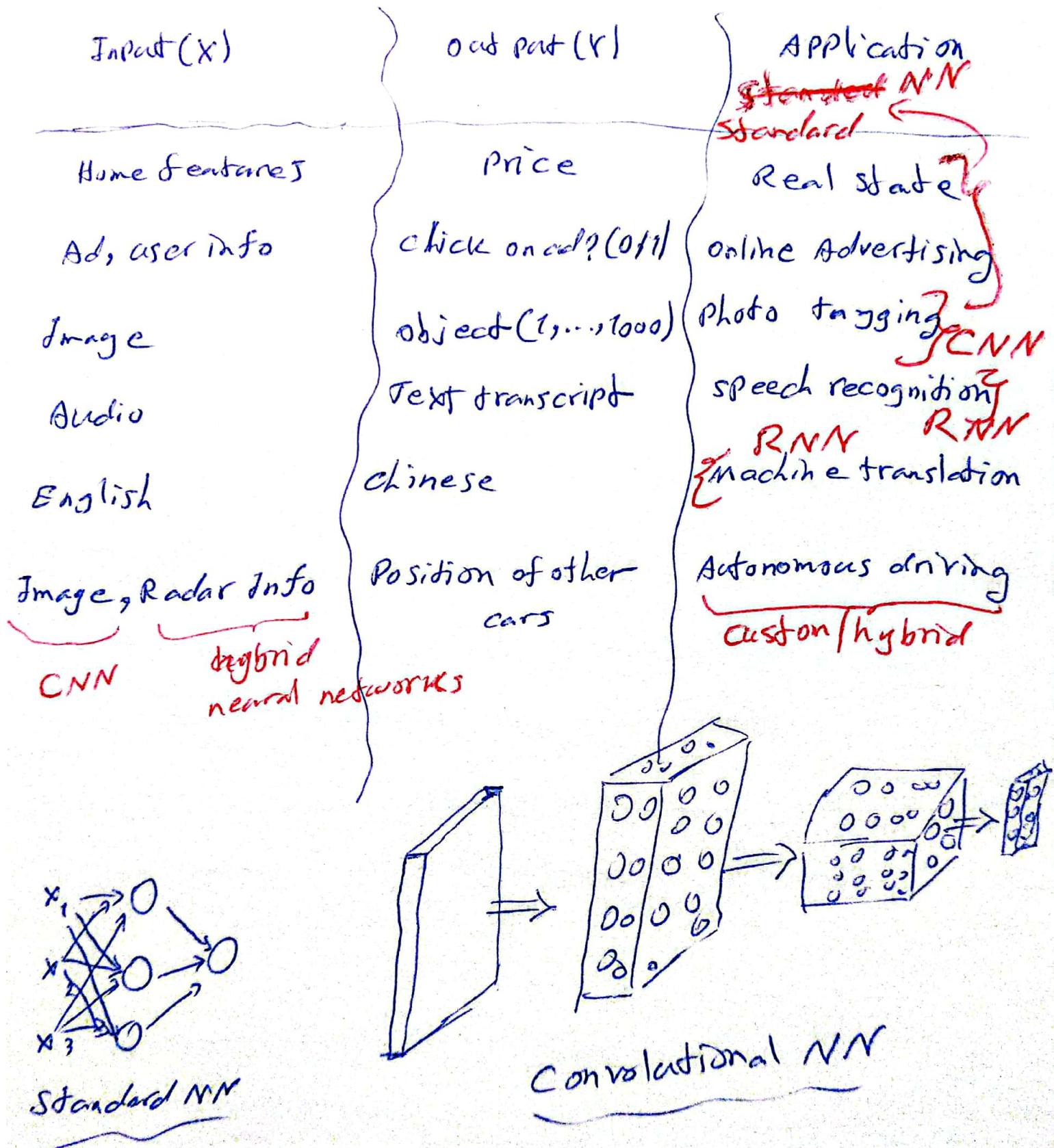
Housing price prediction

(2)



Supervised Learning with Neural Networks

Supervised learning examples (deep learning) ③



Supervised Learning

9

Structured data

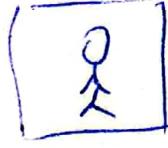
size	# feature	...	price
2024	?	:	400
1000	3	:	300
:	:	:	:
3000	a	:	500

↓
databases

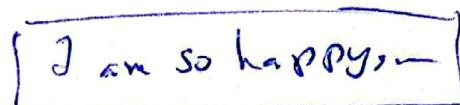
unstructured data



Audio



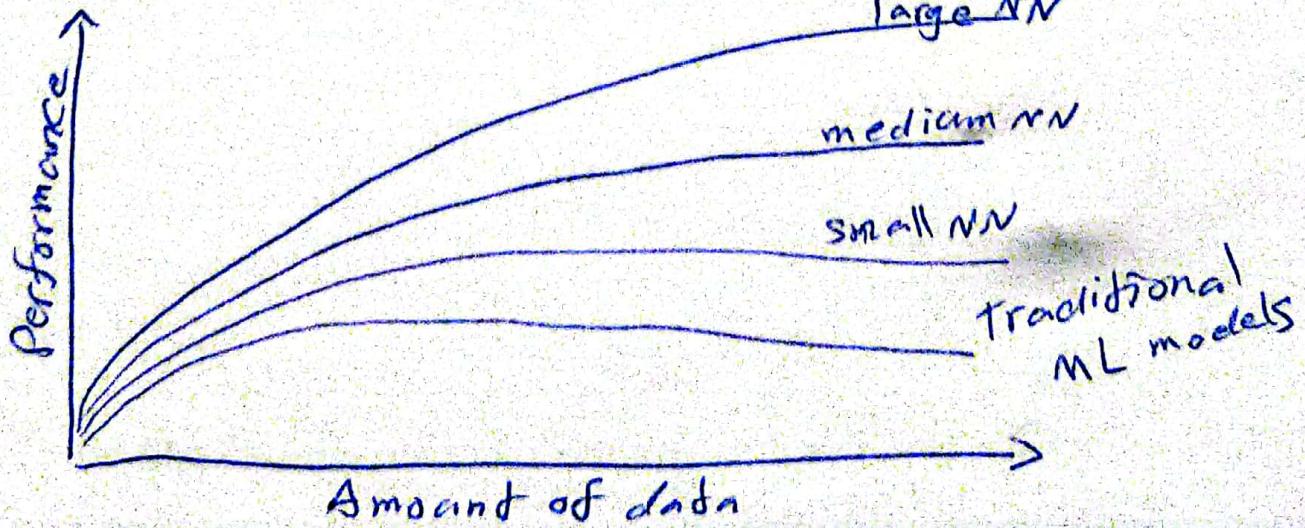
Image



Text

historically computers (classical ML) was better in understanding Structured data , but human are really good in understanding unstructured data , tanks deep learning computer nowadays are really better in understanding unstructured data in compare to the past.

is
why deep learning recently ~~takes~~ taking off?

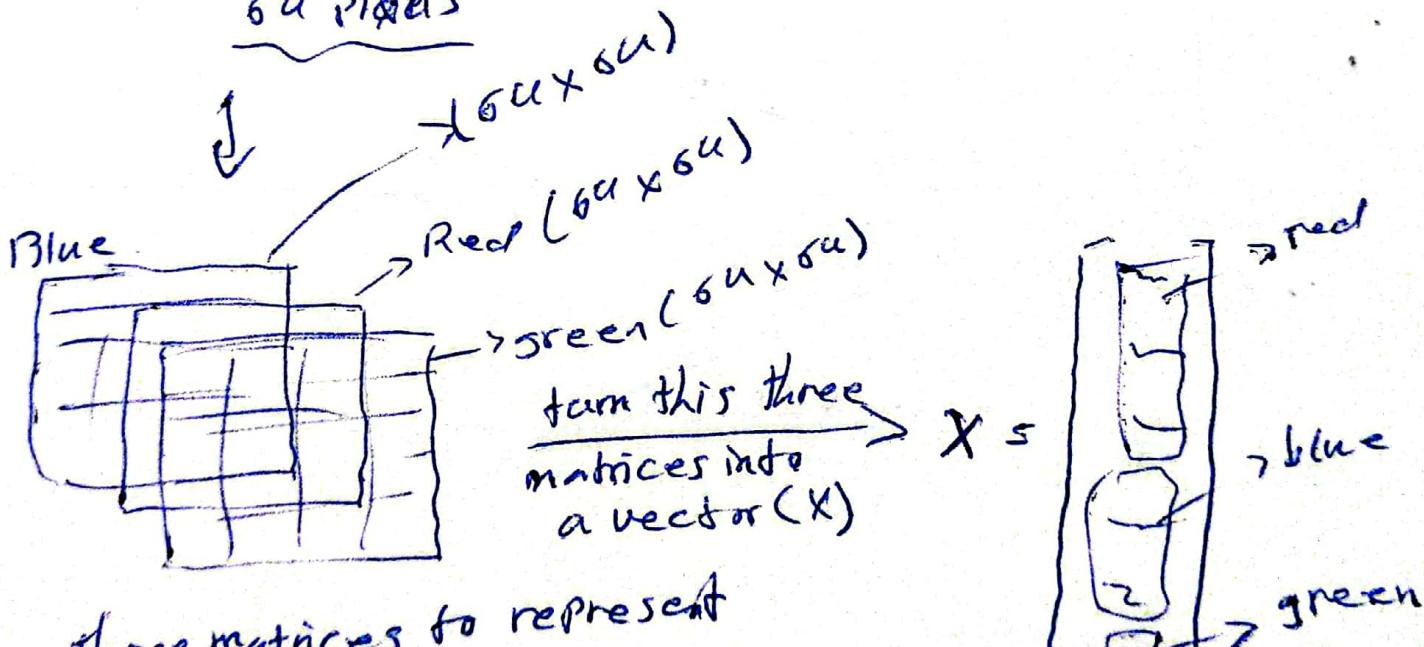
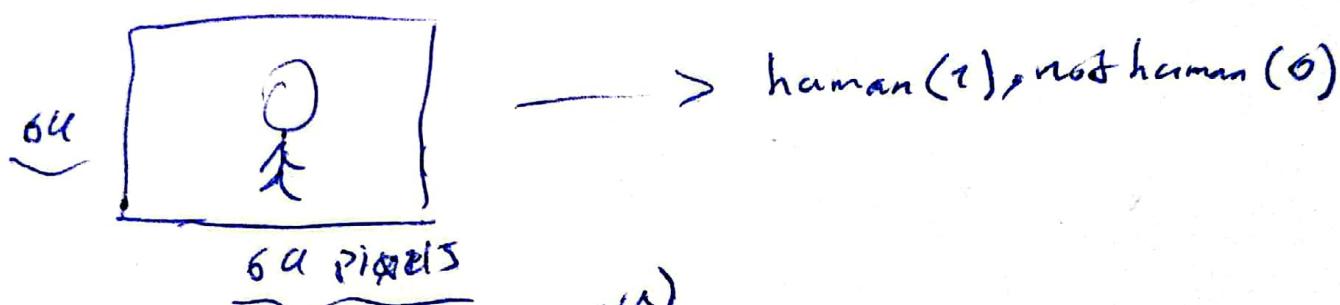


Binary classification

logistic regression) an algorithm to ~~classifies~~ a ~~binary problem~~ for binary classification.

Example:

an image



three matrices to represent
the three main colors in the
picture.

$$64 \times 64 \times 3 = 12288$$

~~so~~ -
so the dimension of
 X is: $n_x = (12288, 1)$

notation

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

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

$$m = m_{train} + m_{test}$$

$$m \rightarrow \text{number of examples}$$

$$X = \begin{bmatrix} \cdot & | & | \\ x^{(1)} & x^{(2)} & \dots x^{(m)} \\ | & | & | \end{bmatrix} \quad n_x \rightarrow \text{number of features}$$

stacking all m (examples) into ~~the~~ X matrix

so: $X = \mathbb{R}^{(n_x \times m)}$

$X.shape = (n_x, m)$

number of
features

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

$$y \in \mathbb{R}^{1 \times m}$$

$$y.shape = (1, m)$$

logistic regression

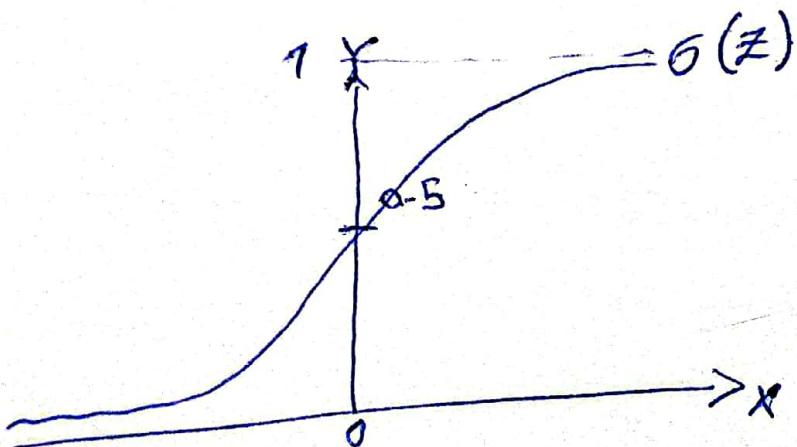
Given X , want $\hat{y} = P(y=1|X)$

$$X \in \mathbb{R}^{n_x}$$

parameters; $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

$$\underline{Z = w^T \times X + b}$$

output $\hat{y} = \sigma(w^T \times X + b)$



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

if Z is large then

$$\sigma(Z) \approx \frac{1}{1+0} \approx 1$$

if Z is small then

$$\sigma(Z) \approx \frac{1}{1+\infty} \approx 0$$

cost function for logistic regression

~~loss~~ ~~error~~ function

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

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

if $y=0$ then; $L(\hat{y}, y) = -\log(1-\hat{y})$

cost function (logistic regression)

(8)

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

$$= -\frac{1}{2} \sum_{i=1}^m \left[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)}) \right]$$

- what is the difference between the cost function and the loss function?

the loss function computes the error for a single training example, the cost function is the average of the loss functions of the entire training set.

Gradient descent

logistic regression
algorithm

Recap: $\hat{y} = \sigma(w^T \cdot x + b)$, $\sigma(z) = \frac{1}{1+e^{-z}}$

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

want to find w, b that minimize $J(w, b)$

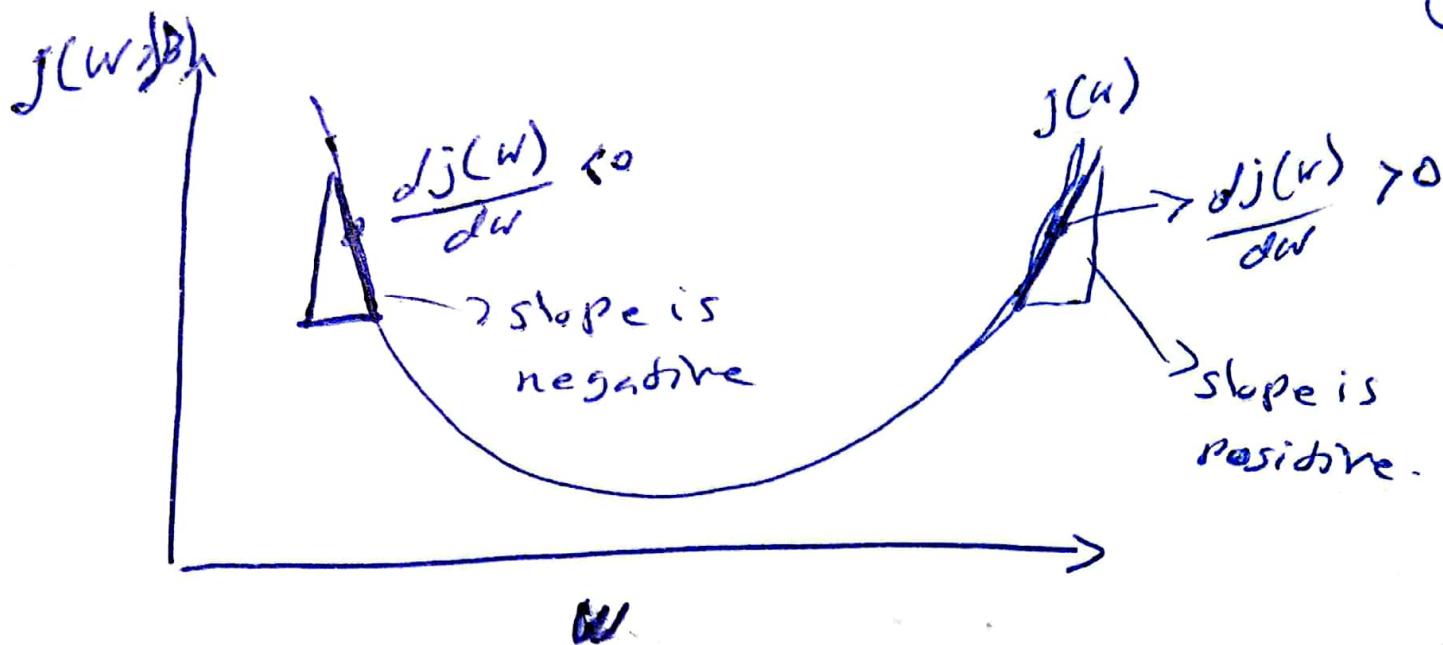
Repeat {

$$w := w - \alpha \frac{\partial J(w, b)}{\partial w}$$

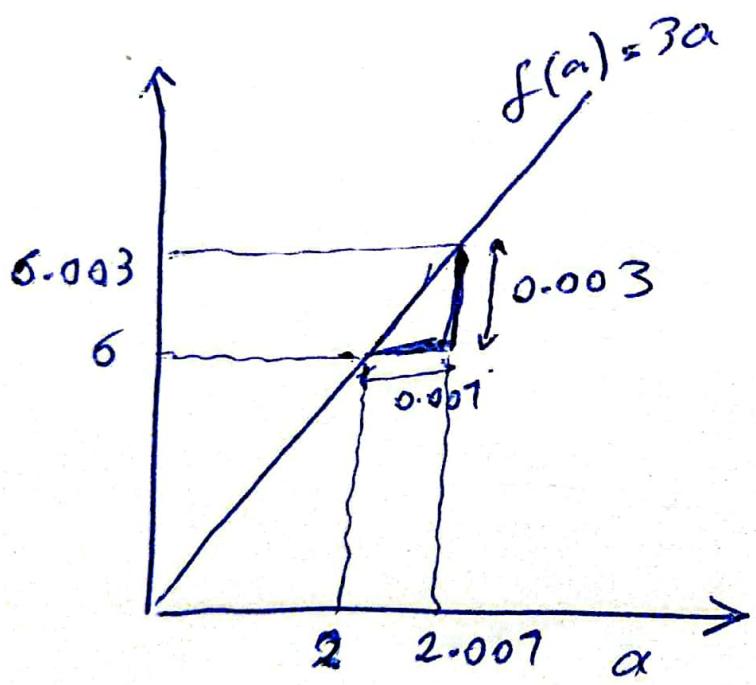
$$b := b - \alpha \frac{\partial J(w, b)}{\partial b}$$

}

9



derivative



$$a=2 \quad f(a)=6$$

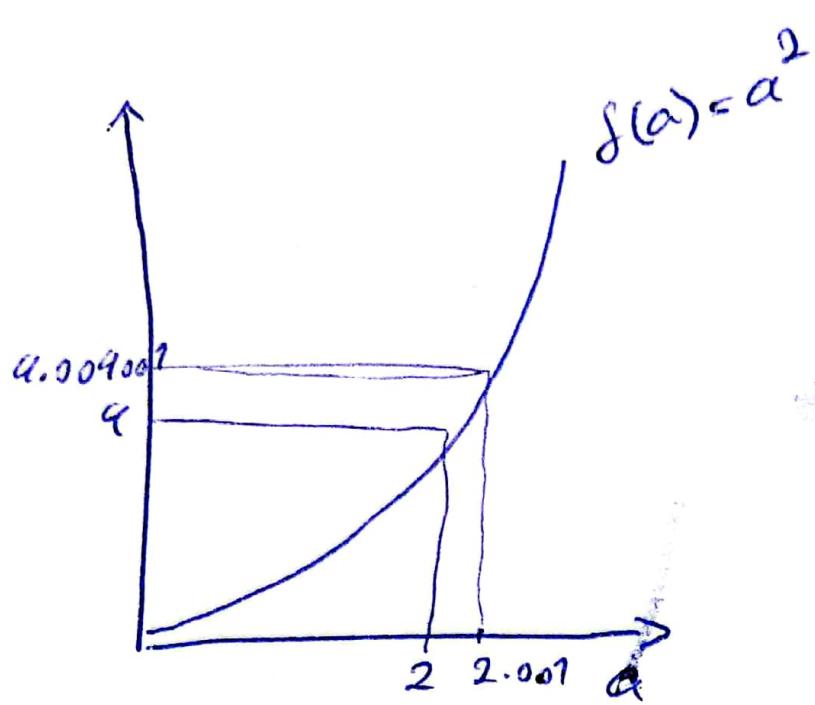
$$a=2.007 \quad f(a)=6.003$$

slope (derivative) of $f(a)$ at $a=2 = 3$

$$\text{slope} = \frac{\text{height}}{\text{width}} = \frac{0.003}{0.007}$$

 $= 3$

$$\boxed{\frac{df(a)}{da} = 3}$$



$$a=2 \quad f(a)=4$$

$$a=2.001 \quad f(a) \approx 4.004007$$

slope (derivative) of $f(a)$
at $a=2 = 4$

$$\frac{d f(a)}{da} = 4 \quad \text{when } a=2$$

$$a=5 \quad f(a)=25$$

$$a=5.001 \quad f(a) \approx 25.010$$

$$\frac{d f(a)}{da} = \frac{d}{da} a^2 = 2a$$

$$f'(a) = a^3 \quad \frac{d}{da} f(a) = 3a^2$$

$$f(a) = \log(a)$$

$\ln(a)$

$$\frac{d f(a)}{da} = \frac{1}{a}$$

computation graph

exp:

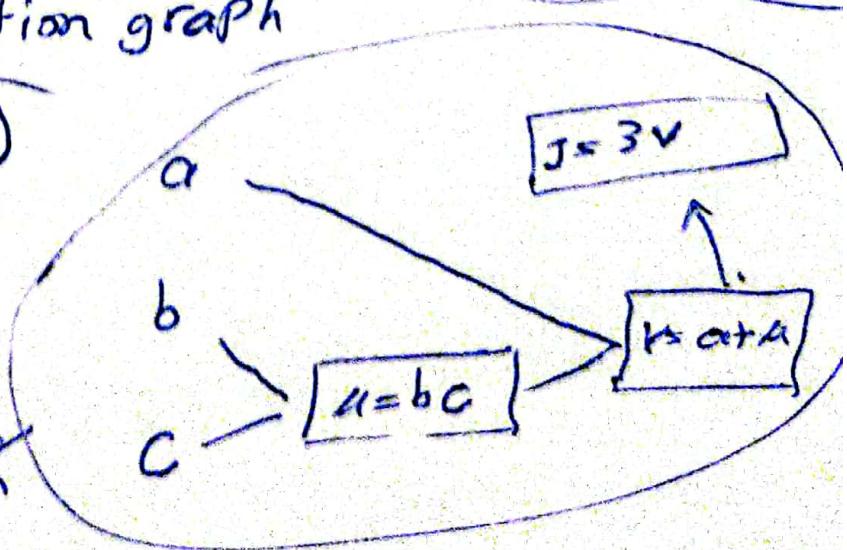
$$j(a, b, c) = 3(\underbrace{a + b \cdot c}_{u})$$

$$u = b \cdot c$$

$$v = a \cdot u$$

$$j = 3v$$

the computation graph



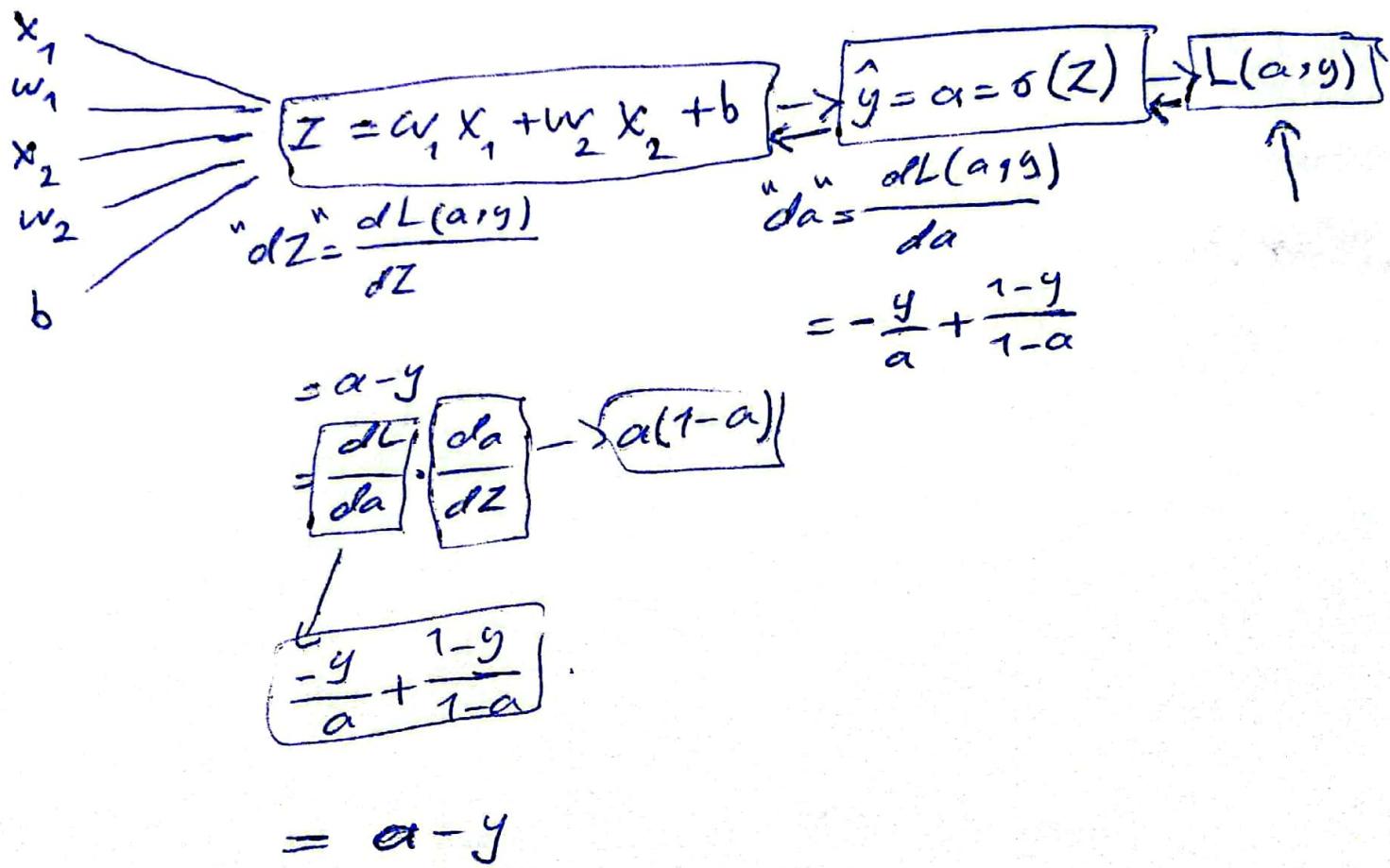
derivative with a computation graph

Logistic regression Gradient descent

$$Z = w^T \cdot X + b$$

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

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



$$\frac{\partial L}{\partial w_1} = dL \cdot dZ = x_1 \cdot dZ, \quad dL \cdot dZ = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial Z} = \frac{\partial L}{\partial a} \cdot \sigma'(Z) = \frac{\partial L}{\partial a} \cdot (a(1-a))$$

$$w_1 := w_1 - \alpha \frac{\partial L}{\partial w_1}$$

$$w_2 := w_2 - \alpha \frac{\partial L}{\partial w_2}$$

$$b := b - \alpha \frac{\partial L}{\partial b}$$

Gradient Descent on m example

72

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

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

$$\underline{\frac{d w_1}{d w_1}}^{(i)}, \underline{\frac{d w_2}{d w_2}}^{(i)}, \underline{\frac{d b}{d b}}^{(i)} \quad (x^{(i)}, y^{(i)})$$

$$\underline{\frac{\partial J(w, b)}{\partial w_1}} = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial w_1} L(a^{(i)}, g^{(i)})}_{\partial a^{(i)} / \partial w_1 = (x^{(i)}, y^{(i)})}$$

$$j=0, \frac{d w_1}{d w_1}=0, \frac{d w_2}{d w_2}=0 \Rightarrow \frac{d b}{d b}=0$$

for $i=1 \rightarrow m$,

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

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

$$J \leftarrow - \left[y^{(1)} \log a^{(1)} + (1 - y^{(1)}) \log (1 - a^{(1)}) \right]$$

$$\begin{aligned} \frac{d w_1}{d w_1} &= x_1^{(i)} \cdot d z^{(i)} \\ \frac{d w_2}{d w_2} &= x_2^{(i)} \cdot d z^{(i)} \end{aligned} \quad \left. \right\} n=2$$

$$d b \leftarrow d I^{(i)}$$

$$|j|=m ; |d w_1|=m ; |d w_2|=m ; |d b|=m$$

$$\frac{\partial w_1}{\partial w_1} = \frac{\partial j}{\partial w_1} ; \quad \frac{\partial w_2}{\partial w_2} = \frac{\partial j}{\partial w_2} ; \quad \frac{\partial b}{\partial b} = \frac{\partial j}{\partial b}$$

$$w_1 := w_1 - \alpha \frac{\partial w_1}{\partial w_1}$$

$$w_2 := w_2 - \alpha \frac{\partial w_2}{\partial w_2}$$

$$b := b - \alpha \frac{\partial b}{\partial b}$$

what is vectorization?

$$z = w^T x + b; \quad w = \begin{bmatrix} : \\ : \\ : \end{bmatrix} \quad x = \begin{bmatrix} : \\ : \\ : \end{bmatrix} \quad w \in \mathbb{R}^{n_x} \\ x \in \mathbb{R}^{n_x}$$

$$z = \underbrace{np.\text{dot}(w, x)}_{w^T \cdot x} + b$$

Parallelization: both CPU and GPU have

Parallelization in structure called SIMD -

Stand for Single Instruction multiple Data.

- In neural network implementation, whenever possible
avoid explicit for-loops.

logistic regression derivatives with for loops:

$$J = 0 ; \quad d\omega_1 = 0 ; \quad d\omega_2 = 0 ; \quad db = 0$$

for $i = 1$ to m :

$$Z^{(i)} = w^T \cdot x^{(i)} + b$$

$$\alpha^{(i)} = \sigma(Z^{(i)}) \quad (1 - \alpha^{(i)})$$

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

$$dZ^{(i)} = \alpha^{(i)} - y^{(i)}$$

$$d\omega_1 += x_1^{(i)} \cdot dZ^{(i)}$$

$$d\omega_2 += x_2^{(i)} \cdot dZ^{(i)}$$

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

$$J = J/m ; \quad d\omega_1 = d\omega_1/m ; \quad d\omega_2 = d\omega_2/m ; \quad db = db/m$$

Vectorization of logistic regression

$$X = \begin{bmatrix} | & | & | \\ X^{(1)} & X^{(2)} & \dots & X^{(m)} \\ | & | & | \end{bmatrix}_{(n \times m)}$$

$$[Z^{(1)}, Z^{(2)}, \dots, Z^{(m)}] = w^T \cdot X + [b, b, \dots, b]_{(1 \times m)}$$

$$= [w^{(1)}, w^{(2)}, \dots, w^{(m)}] \cdot \begin{bmatrix} | & | & | \\ X^{(1)} & X^{(2)} & \dots & X^{(m)} \\ | & | & | \end{bmatrix}_{(n \times m)} + [b, b, \dots, b]$$

$$= [w^T \cdot X + b, w^T \cdot X + b, w^T \cdot X + b, \dots, w^T \cdot X + b]_{(1 \times m)}$$

$$\Rightarrow Z = w^T \cdot X + b$$

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

$$A = \sigma(Z)$$

$$A = \frac{1}{1 + np \cdot \exp(Z)} ; \boxed{dZ = A - Y}$$

$$db = \frac{1}{m} \sum_{i=1}^m dZ^{(i)} \Rightarrow \frac{1}{m} \times np \cdot \text{sum}(dZ)$$

$$dW = \frac{1}{m} \times X \times dZ^T$$

$$= \frac{1}{m} \left[X^{(1)}, \dots, X^{(m)} \right] \begin{bmatrix} dZ^{(1)} \\ \vdots \\ dZ^{(m)} \end{bmatrix}$$

$$\leq \frac{1}{m} \left[X^{(1)} \frac{dZ^{(1)}}{dZ}, \dots, X^{(m)} \frac{dZ^{(m)}}{dZ} \right] \quad (m \times 1)$$

Implementing Logistic Regression (completely)

$$Z = w^T \cdot x + b; \quad Z = np.dot(w.T, x) + b$$

$$A = \sigma(Z); \quad A = \frac{1}{1 + np.exp(Z)}$$

$$dZ = A - Y$$

$$dW = \frac{1}{m} \cdot X \cdot dZ^T$$

$$db = \frac{1}{m} \cdot \text{sum}(dZ)$$

$$w := w - \alpha \cdot dW$$

$$b := b - \alpha \cdot db$$

a single iteration of logistic regression algorithm.

"Broadcasting in python"

- broadcasting is a technique, which can be used to speed up the code.

"example of broadcasting in Python"

(17)

$A = np.array([$

$$[1 \ 2 \ 5 \ 2]$$

$$[2 \ 3 \ 5 \ 3]$$

$$[3 \ 4 \ 5 \ 6]$$

])

$A.shape = (3 \times 4)$

$cal = A.sum(axis=0)$ # means: sum up all the entire column values in each column.



$$[6 \ 9 \ 15 \ 6]$$

percentage = A / cal # shape ~~(1, 4)~~



percentage.shape = (3×4) # Percentage of each value in the corresponding column.

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

broadcasting?

General principle of broadcasting

18

$$\text{matrix} \begin{pmatrix} m, n \end{pmatrix} \xrightarrow{\substack{+ \\ x}} \begin{pmatrix} 1, n \end{pmatrix} \rightsquigarrow (m, n)$$

$$\\ \xrightarrow{\substack{- \\ x}} \begin{pmatrix} m, 1 \end{pmatrix} \rightsquigarrow (m, n)$$

$$\text{matrix} \begin{pmatrix} m, 1 \end{pmatrix} + \underbrace{R}_{(1, 0)} \rightsquigarrow (m, 1)$$

$$\text{matrix} \begin{pmatrix} 1, m \end{pmatrix} + R \rightsquigarrow (1, m)$$

broadcasting in numpy

`np.sum(axis=0)` # sum up vertically

`np.sum(axis=1)` # sum up horizontally

Logistic regression cost function

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

interpret $\hat{y} = P(y=1 | X)$

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

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

$$P(y|X) = \hat{y}^y (1-\hat{y})^{(1-y)} \leftarrow$$

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

$\underbrace{\hat{y}}_{\text{if } y=1} \quad \underbrace{(1-\hat{y})^0}_{(1-\hat{y})^0 = 1} = 1$

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

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

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

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

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

soft

$$P(\text{labels in target set}) = \prod_{i=1}^m P(y^{(i)} | x^{(i)})$$

$$\log P(\dots) = \sum_{i=1}^m \underbrace{\log P(y^{(i)} | x^{(i)})}_{-L(\hat{y}^{(i)}, y^{(i)})}$$

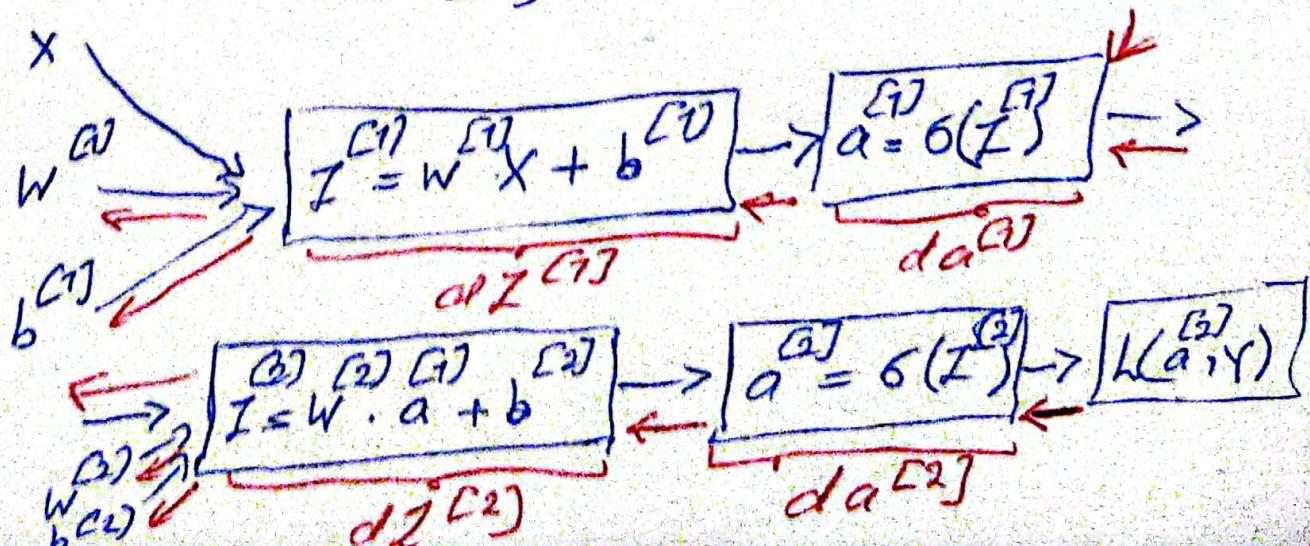
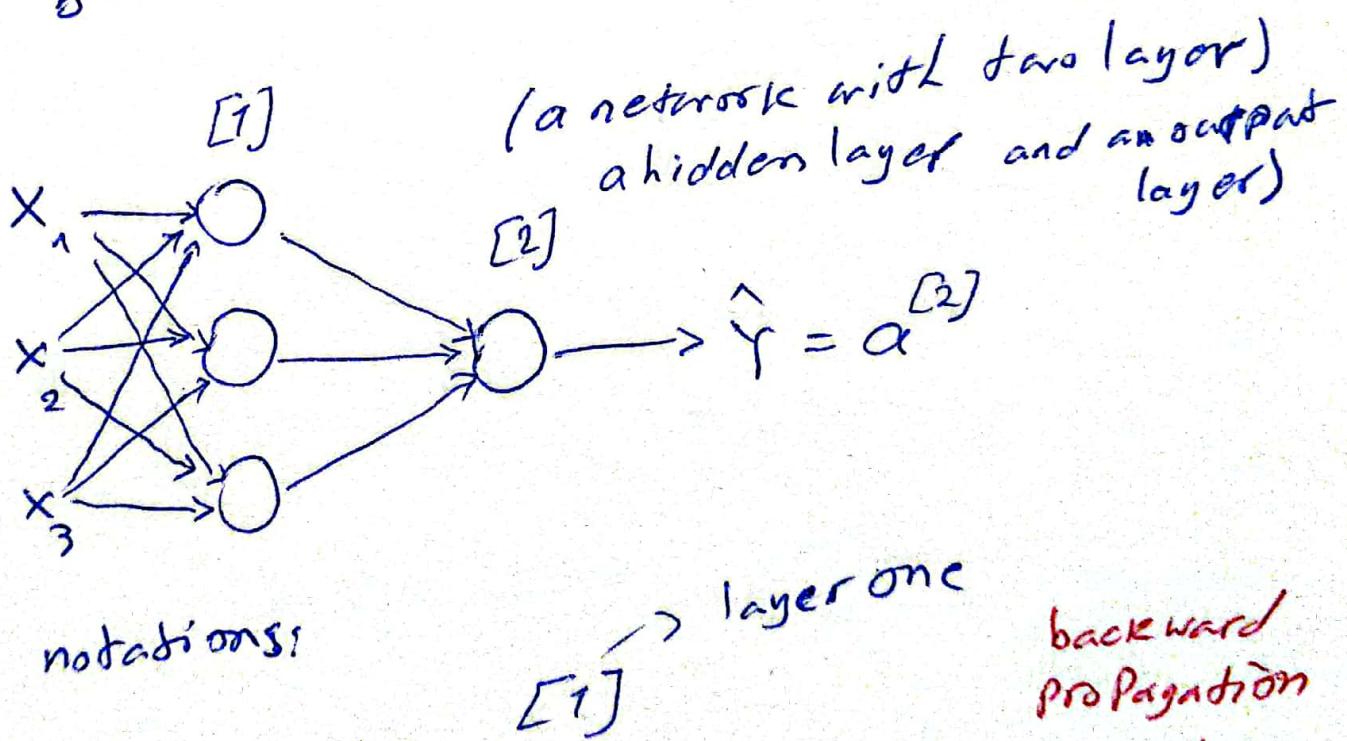
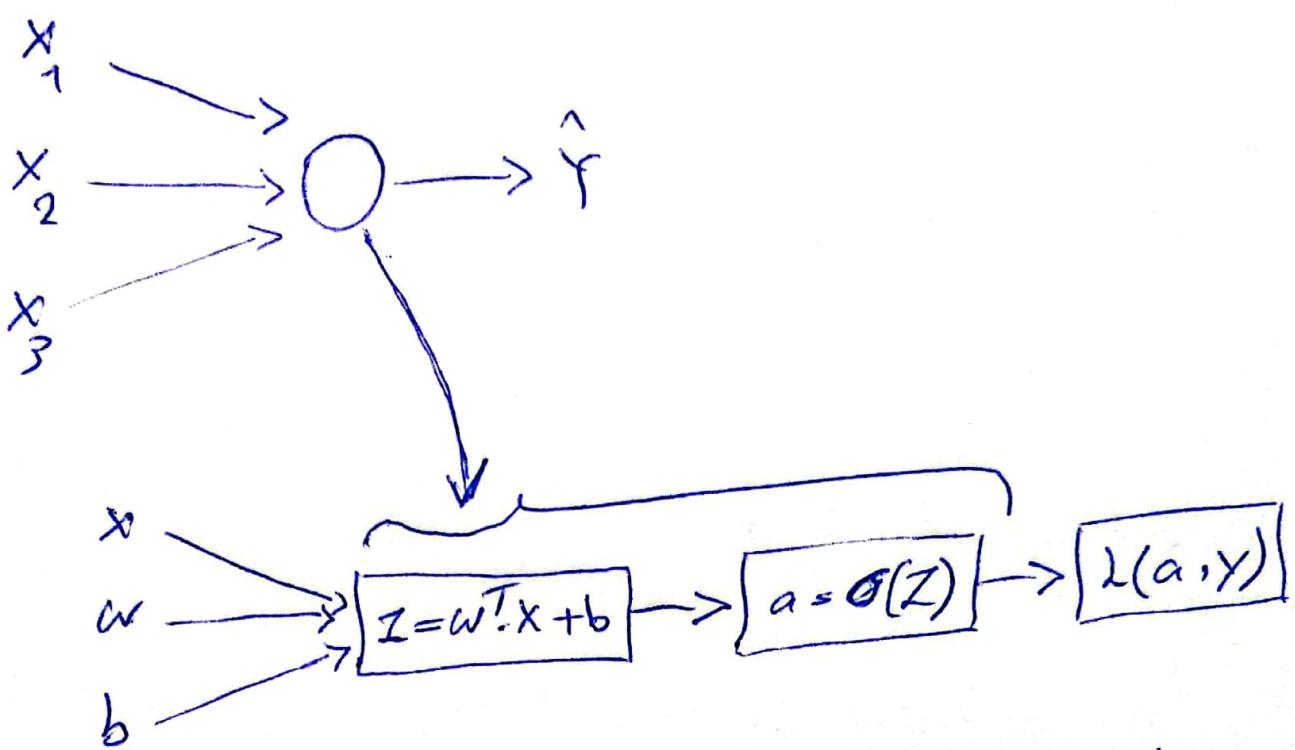
$$= - \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

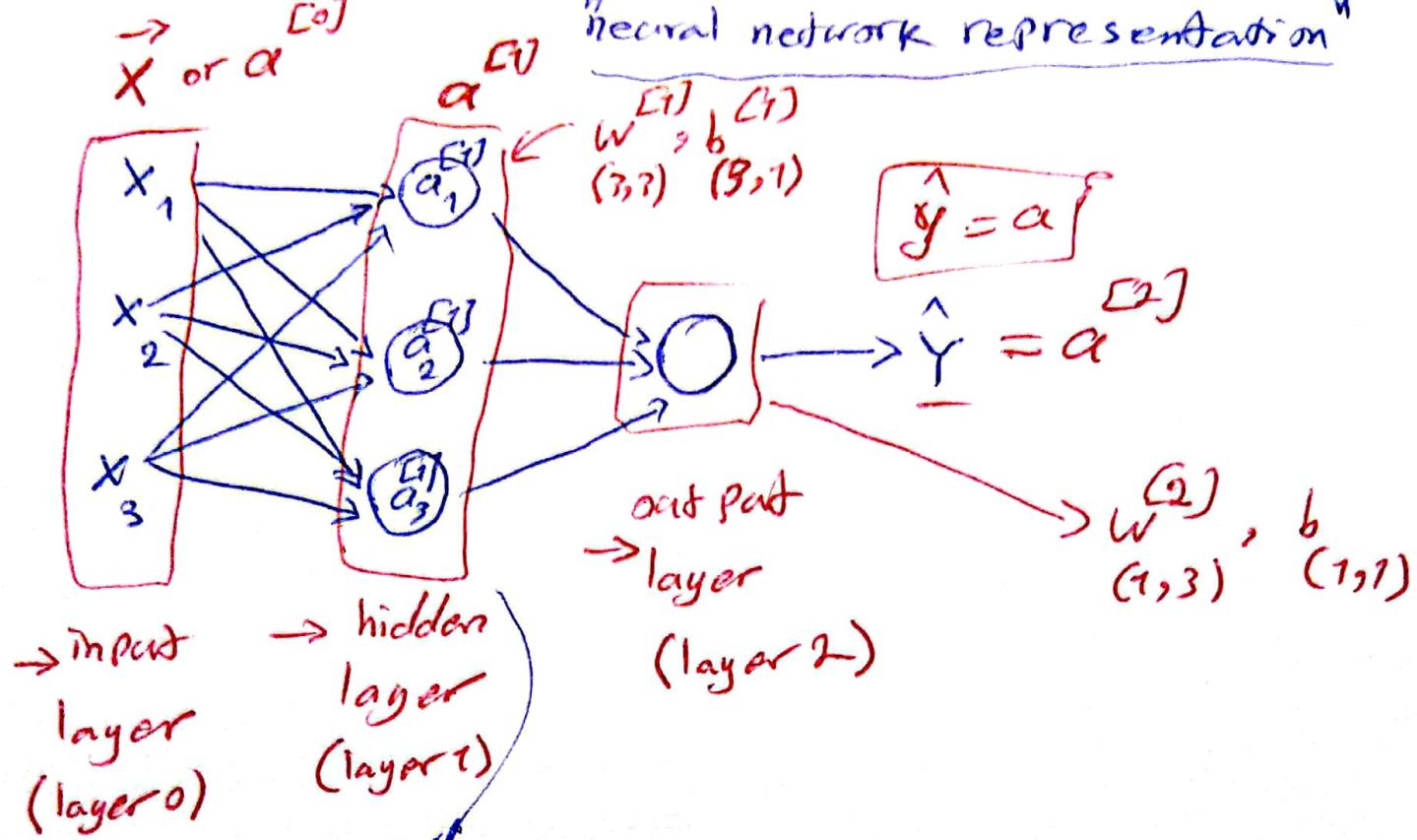
$$\text{cost } J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

(minimize)

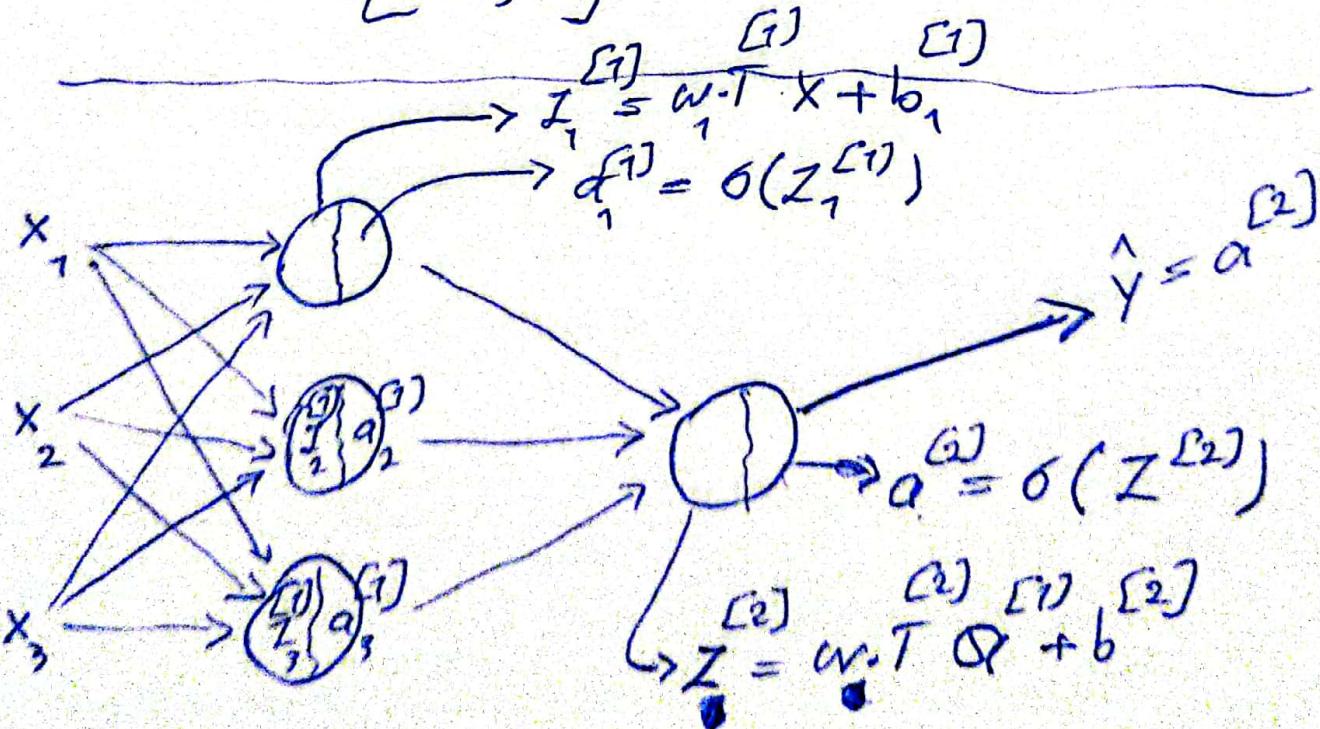
"Neural network overview"

(20)





$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \end{bmatrix}$$



hidden layer (1)

$$\begin{array}{l}
 I_1^{(G)} = w_1^{(G)T} X + b_1^{(G)}, \quad a_1^{(G)} = \sigma(Z_1^{(G)}) \\
 I_2^{(G)} = w_2^{(G)T} X + b_2^{(G)}, \quad a_2^{(G)} = \sigma(Z_2^{(G)}) \\
 I_3^{(G)} = w_3^{(G)T} X + b_3^{(G)}, \quad a_3^{(G)} = \sigma(Z_3^{(G)}) \\
 I_a^{(G)} = w_a^{(G)T} X + b_a^{(G)}, \quad a_a^{(G)} = \sigma(Z_a^{(G)})
 \end{array}$$

$$\left[\begin{array}{c} w_1^{(G)T} \\ w_2^{(G)T} \\ w_3^{(G)T} \\ w_a^{(G)T} \end{array} \right] \left[\begin{array}{c} X_1 \\ X_2 \\ X_3 \end{array} \right] + \left[\begin{array}{c} b_1^{(G)} \\ b_2^{(G)} \\ b_3^{(G)} \\ b_a^{(G)} \end{array} \right]_{(a,1)} = \left[\begin{array}{c} w_1^{(G)T} X + b_1^{(G)} \\ w_2^{(G)T} X + b_2^{(G)} \\ w_3^{(G)T} X + b_3^{(G)} \\ w_a^{(G)T} X + b_a^{(G)} \end{array} \right]$$

$$\Rightarrow \left[\begin{array}{c} I_1^{(G)} \\ I_2^{(G)} \\ I_3^{(G)} \\ I_a^{(G)} \end{array} \right] \rightarrow a^{(G)} = \left[\begin{array}{c} a_1^{(G)} \\ \vdots \\ a_a^{(G)} \end{array} \right] = \sigma(Z^{(G)})$$

Given input X : $\rightarrow n \text{ num units}$

$$Z^{(1)}_{(u,1)} = w^{(1)}_{(u,1)} X_{(3,1)} + b^{(1)}_{(4,1)}$$

$$a^{(1)}_{(u,1)} = g(Z^{(1)}_{(u,1)})$$

$$Z^{(2)}_{(1,1)} = w^{(2)}_{(1,1)} a^{(1)}_{(u,1)} + b^{(2)}_{(1,1)}$$

$$a^{(2)}_{(1,1)} = g(Z^{(2)}_{(1,1)})$$

"vectorizing across multiple examples"

~~layer~~
~~example~~
~~ith~~

$$\vec{X} = \begin{bmatrix} | & | & | \\ X^{(1)} & X^{(2)} & \dots & X^{(m)} \\ | & | & & | \\ (n \times m) \end{bmatrix} \rightarrow \text{num examples}$$

num features

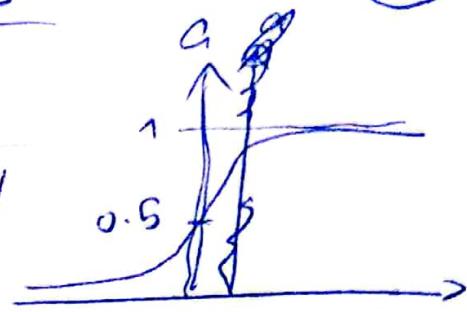
$$Z^{(1)} = \begin{bmatrix} |_{G(1)[1]} & |_{G(1)[2]} & \dots & |_{G(1)[m]} \\ Z^{(1)[1]} & Z^{(1)[2]} & \dots & Z^{(1)[m]} \\ | & | & & | \\ (H, m) \end{bmatrix} \rightarrow \text{num units}$$

$$A^{(1)} = \begin{bmatrix} |_{G(1)[1]} & |_{G(1)[2]} & \dots & |_{G(1)[m]} \\ a^{(1)}, a^{(1)[2]}, \dots, a^{(1)[m]} \\ | & | & & | \\ (H, m) \end{bmatrix} \rightarrow \text{num units}$$

"Activation Functions"

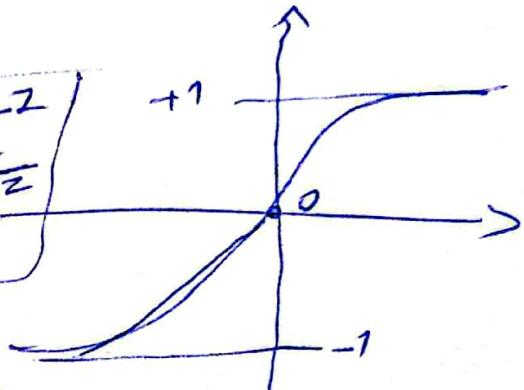
1. Sigmoid function

$$a = \frac{1}{1+e^{-z}}$$



2. tanh function

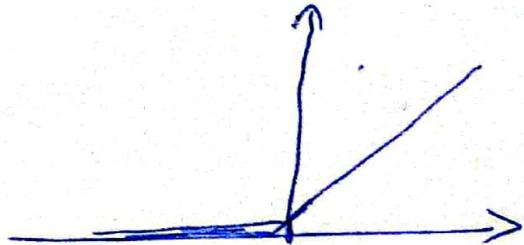
$$a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



almost always is a better choice

~~as~~ than sigmoid function

3. rectified linear unit (ReLU)

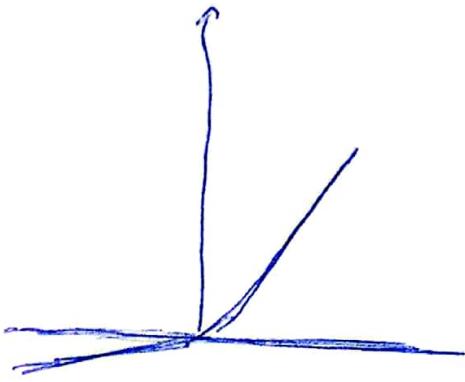


$$a = \max(0, z)$$

- if the output of NN ~~is~~ is binary, the output layer should ~~be~~ use sigmoid activation function and other layers (hidden layers), ~~not~~ should use ReLU or sometimes tanh

4. Leaky Relu

$$\alpha = \max(0.01z, z)$$



why should we use non-linear activation functions?

if does not matter how many hidden layer one uses in a NN, if ~~the~~ they do not have an activation function, the result is just the same as a regression model.

linear or logistic

Derivatives of activation functions

1. Sigmoid activation functions

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

$\frac{d}{dz} g(z) = \text{slope of } g(x) \text{ at } z$

$$g(z) = \alpha$$

$$\frac{d}{dz} g(z) = \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}} \right)$$

$$= g(z) (1 - g(z)) = \alpha(1 - \alpha)$$

2. tanh activation function

$$g(z) = \tanh(z)$$

$$= \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\frac{d}{dz} g(z) = 1 - (\tanh(z))^2$$

if $g(z) = a$

$$\text{so: } \frac{d}{dz} g(z) = 1 - a^2$$

3. Relu

$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \\ \text{undefined} & \text{if } z = 0 \end{cases} \quad ?$$

4. Leaky Relu

$$g(z) = \max(0.001z, z)$$

$$g'(z) = \begin{cases} 0.001 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Gradient Descent for NNs

Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$

$$n_x = n^{[0]}, n^{[1]}, n^{[2]} = 1$$

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

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

$$w^{[2]} \\ (n^{[2]}, n^{[1]})$$

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

cost function: $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]})$

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

Gradient Descent:

Repeat {

compute predictions $(\hat{y}^{(1)}, \dots, \hat{y}^{(m)})$

$$\frac{\partial J}{\partial w^{[1]}} = \frac{\partial J}{\partial w^{[1]}} , \frac{\partial J}{\partial b^{[1]}} = \frac{\partial J}{\partial b^{[1]}} , \dots \rightarrow$$

$$w^{(G)} = w^{(G)} - \alpha d w^{(G)}$$

$$b^{(G)} = b^{(G)} - \alpha d b^{(G)}$$

$$w^{(2)} = w^{(2)} - \alpha d w^{(2)}$$

$$b^{(2)} = b^{(2)} - \alpha d b^{(2)}$$

∫

"formulas for computing derivatives"

forward propagation:

$$Z^{(G)} = w^{(G)} X + b^{(G)}$$

$$A^{(G)} = g^{(G)}(Z^{(G)})$$

$$Z^{(2)} = w^{(2)} \cdot A^{(G)} + b^{(2)}$$

$$A^{(2)} = g^{(2)}(Z^{(2)})$$

Back propagation

$$dZ^{(2)} = A^{(2)} - Y$$

$$dw^{(G)} = \frac{1}{m} dZ^{(2)} \cdot A^{(G)T}$$

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

$$dZ^{(G)} = W^{(2)T} \cdot dZ^{(2)} \otimes g^{(G)'}(Z^{(G)})$$

elementwise product

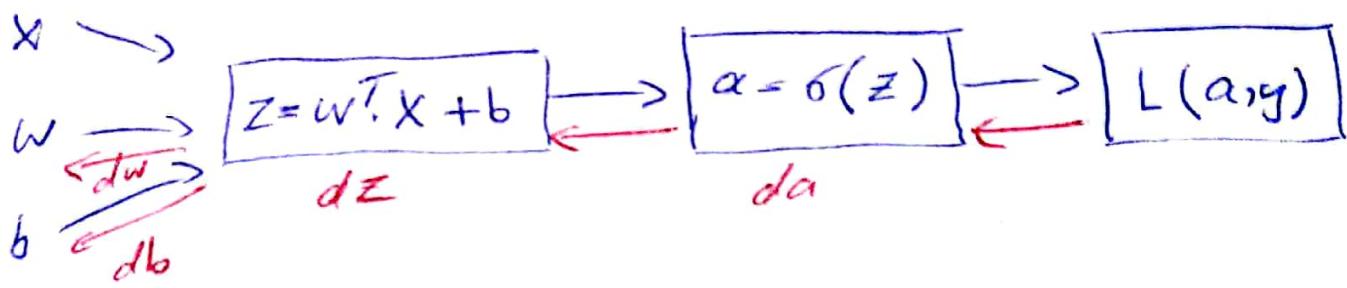
$$dw^{(G)} = \frac{1}{m} dZ^{(G)} \cdot X^T$$

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

prerand (n,)

Same dimension

logistic regression



$$L(\alpha, y) = -y \log \alpha - (1-y) \log(1-\alpha)$$

$$d\alpha = \frac{d}{d\alpha} L(\alpha, y)$$

$$= -\frac{y}{\alpha} + \frac{1-y}{1-\alpha}$$

$$dZ = \alpha - y$$

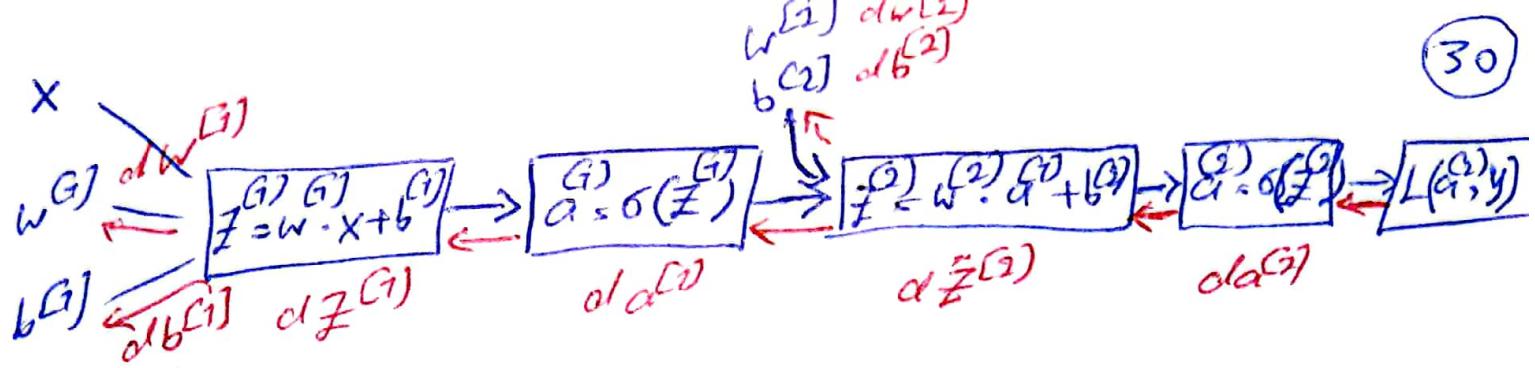
$$= d\alpha \cdot g'(Z) ; \quad g(Z) = \sigma(Z)$$

$$\frac{\partial L}{\partial Z} = \frac{\partial L}{\partial \alpha} \cdot \boxed{\frac{\partial \alpha}{\partial Z}} \rightarrow \frac{d}{dZ} g(Z) = g'(Z)$$

$$dZ = d\alpha \cdot g'(Z)$$

$$dw = dZ \cdot X$$

$$db = dZ$$



$$\frac{\partial Z^{(2)}}{\partial a^{(1)}} = \alpha - y$$

$$\frac{\partial w^{(2)}}{\partial Z^{(2)}} = \frac{\partial Z^{(2)}}{\partial Z^{(1)}} \cdot \alpha \cdot T$$

$$\frac{\partial b^{(2)}}{\partial Z^{(2)}} = \frac{\partial Z^{(2)}}{\partial Z^{(1)}}$$

~~" $\frac{\partial w}{\partial Z^{(2)}} = \frac{\partial Z^{(2)}}{\partial Z^{(1)}} \cdot R$ "~~

~~$\frac{\partial Z^{(2)}}{\partial a^{(1)}}$~~

$$\frac{\partial Z^{(2)}}{\partial a^{(1)}} = w^{(2)T} \cdot \frac{\partial Z^{(2)}}{\partial Z^{(1)}} \times g^{(2)'}(Z^{(2)})$$

\nearrow element wise product

~~$w^{(2)}$~~

$$\frac{\partial w^{(1)}}{\partial Z^{(1)}} = \frac{\partial Z^{(1)}}{\partial Z^{(0)}} \cdot X^T$$

$$\frac{\partial b^{(1)}}{\partial Z^{(1)}} = \frac{\partial Z^{(1)}}{\partial Z^{(0)}}$$

Summary of gradient descent

Back propagation for NN (two layers)

~~gradient~~

for one example

$$dZ^{(2)} = a^{(2)} - y$$

$$dW^{(2)} = dZ^{(2)} \cdot a^{(1)T}$$

$$db^{(2)} = dZ^{(2)}$$

$$dZ^{(1)} = W^{(2)T} dZ^{(2)} \times g^{(1)}(Z^{(1)})$$

elementwise product

$$dW^{(1)} = dZ^{(1)} X^T$$

$$db^{(1)} = dZ^{(1)}$$

vectorized (for many examples)

$$dZ^{(2)} = A^{(2)} - Y$$

$$dW^{(2)} = \frac{1}{m} dZ^{(2)} \cdot A^{(1)T}$$

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

$$dZ^{(1)} = W^{(2)T} dZ^{(2)} \times g^{(1)}(Z^{(1)})$$

$$dW^{(1)} = \frac{1}{m} dZ^{(1)} \cdot X^T$$

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

"Random initialization"

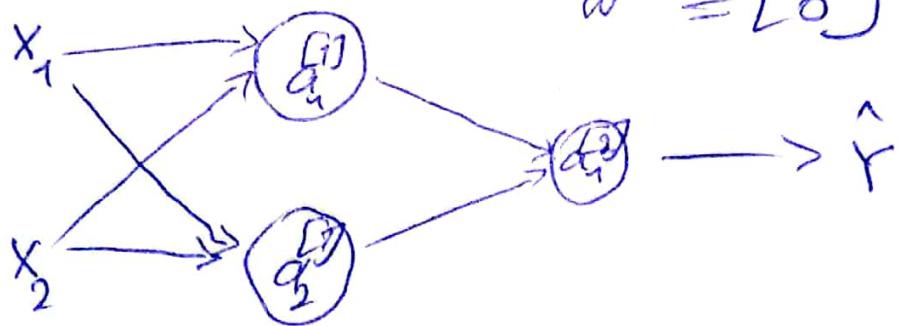
if the weights ($W^{(n)}$) in a neural network initialized to zero, the network won't work.

exp:

the symmetry breaking problem

(32)

$$w^{[2]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$



if: $w^{(1)} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$ $b^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

then: $a_1^{(1)} = a_2^{(1)}$ $dF_1^{(1)} = dF_2^{(2)}$

so the both units will be exactly the same, means they both calculate the same value and ~~it will after~~ they will compute the same value ~~all the time~~ in each iteration.

the solutions initializing the parameters randomly.

$$\underline{w^{(1)} = np.random.randn(2, 2) * 0.01}$$

$$\underline{b^{(1)} = np.zeros(2)}$$

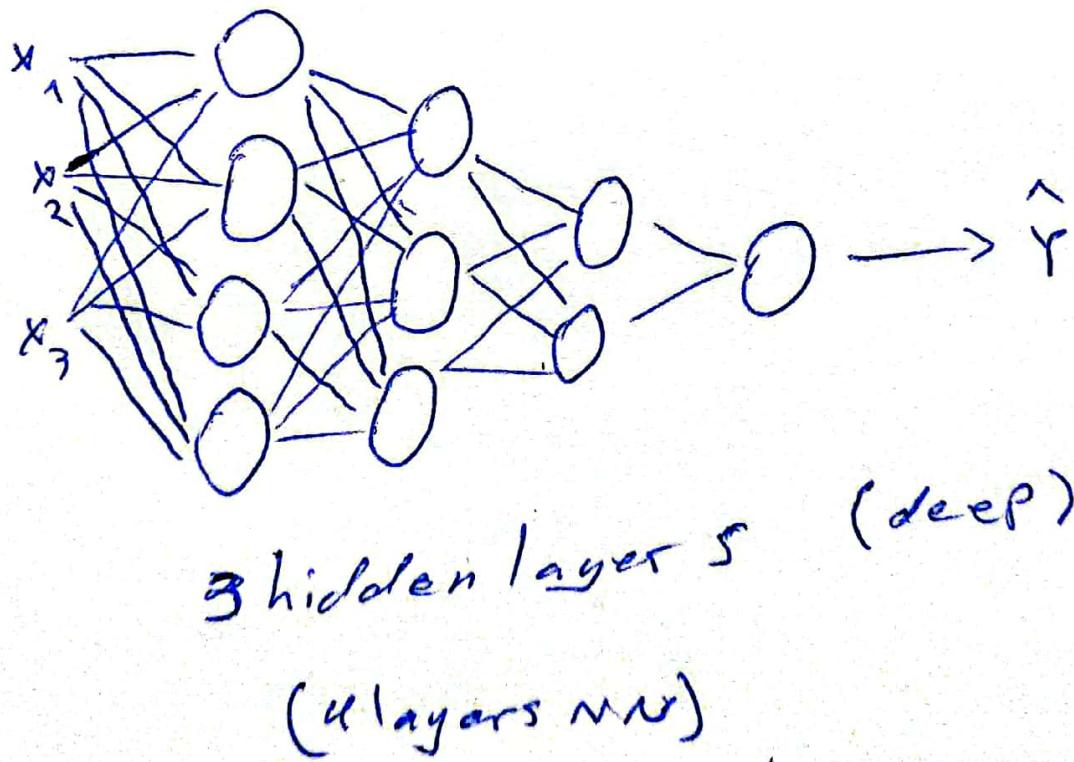
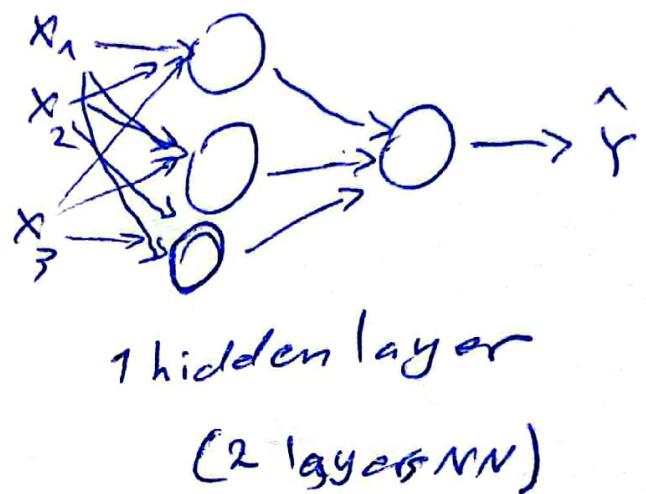
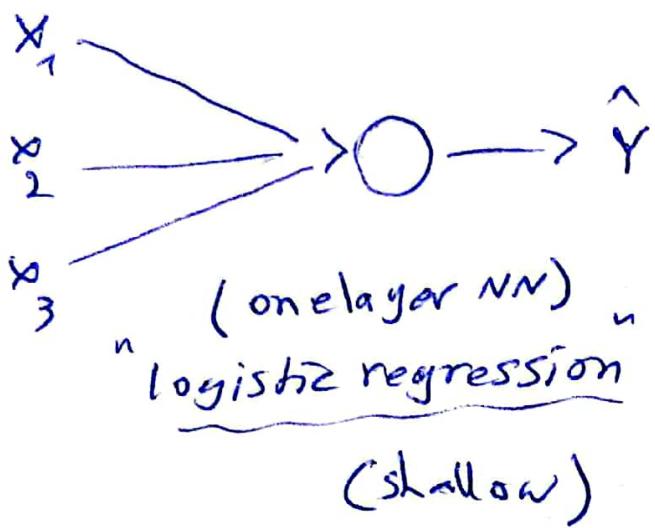
$$\underline{w^{(2)} = np.random.randn(2, 1) * 0.01}$$

$$\underline{b^{(2)} = np.zeros(1)}$$

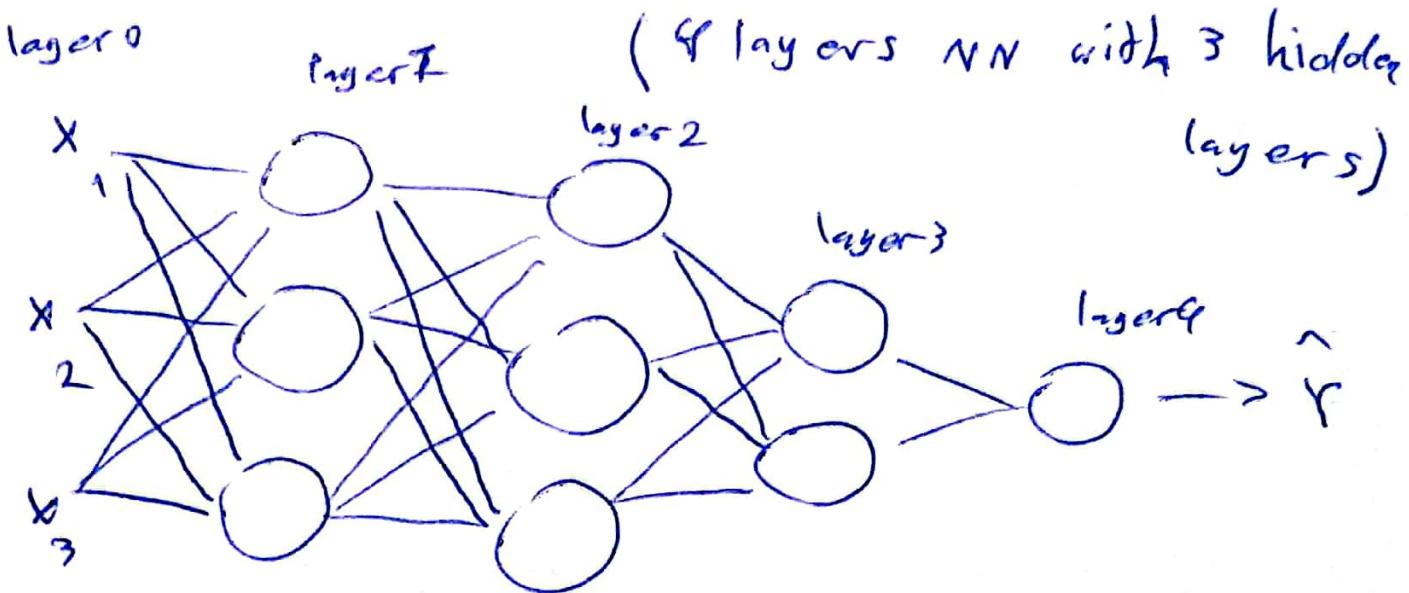
Deep L-layer Neural Network

33

what is a deep neural network?



Deep neural network notation



$$L = 4 \quad (\# \text{ layers})$$

$n^{[L]}$ = # units in layer L

$$n^{[1]} = 3; n^{[2]} = 3; n^{[3]} = 2; n^{[4]} = 1$$

$$n^{[0]} = n_x = 3$$

$a^{[L]}$ = activation function of layer L

$$\text{exp: } a^{[1]} = g^{[1]}(z^{[1]})$$

$$a^{[2]} = \max(z^{[2]}, 0)$$

$$x = a^{[0]}; \hat{y} = a^{[L]}$$

"forward propagation in a deep network"

x (a single example):
 → or $a^{[0]}$

$$x: z^{[1]} = w \cdot x + b^{[1]}$$

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

$$z^{[2]} = w \cdot a^{[1]} + b^{[2]}$$

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

...

$$z^{[q]} = w \cdot a^{[q-1]} + b^{[q]} = a$$

$$a^{[q]} = g^{[q]}(z^{[q]})$$

$$a^{[q]} = \hat{y}$$

General rule:

$$z^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = g^{[l]}(z^{[l]})$$

vectorized way

$$z^{[1]} = w^{[1]} \otimes x + b^{[1]}$$

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

$$z^{[2]} = w^{[2]} \otimes a^{[1]} + b^{[2]}$$

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

$$z^{[4]} = w^{[4]} \otimes a^{[3]} + b^{[4]}$$

$$a^{[4]} = g^{[4]}(z^{[4]}) = \hat{y}$$

General rule:

$$z^{[l]} = w^{[l]} \otimes a^{[l-1]} + b^{[l]}$$

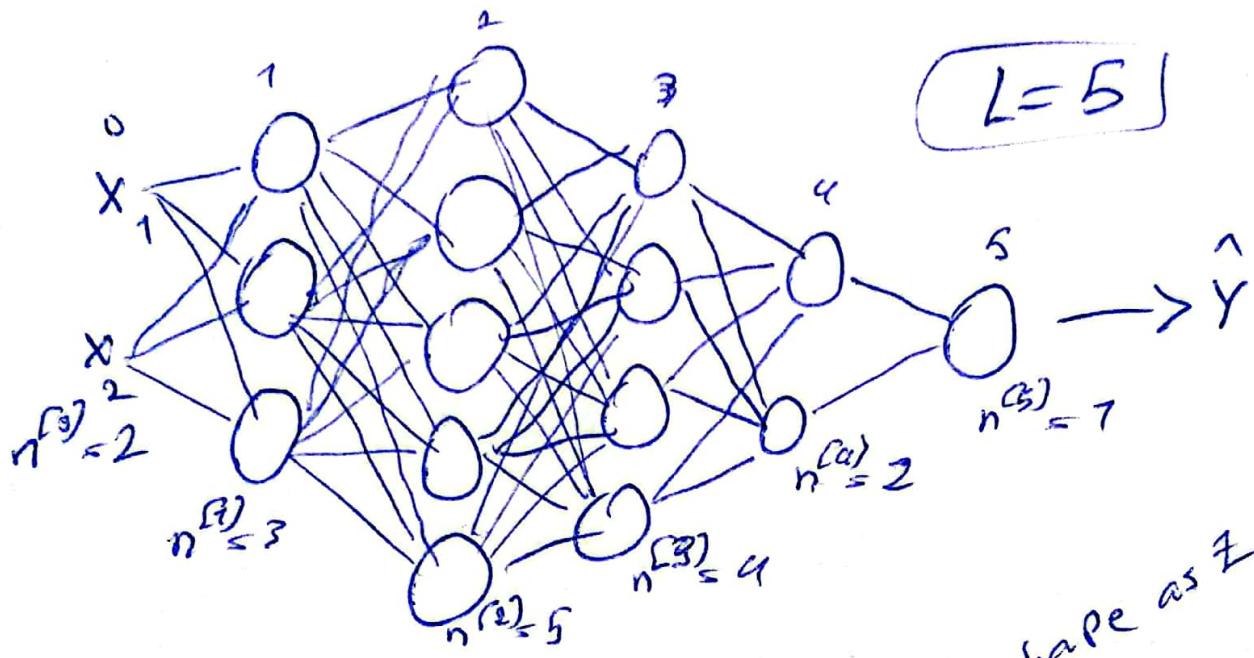
$$a^{[l]} = g^{[l]}(z^{[l]})$$

for $l = 1, 2, 3, 4$

Getting the matrix dimensions right

36

Parameters : $w^{[L]}$ and $b^{[L]}$



layer 1 ; $n^{[1]} = 3$

same shape as z

$$z^{[1]} = w^{[1]} \cdot x + b^{[1]}$$

\downarrow
 $(3,1)$ $(3,2)$ $(2,1)$ $(3+1)$
 $(n^{[0]}, 1)$ $(n^{[0]}, 1)$ $(n^{[0]}, 1)$ $(n^{[1]}, 1)$

In General:

$$w^{[L]} = (w^{[1]}, w^{[2]}, \dots, w^{[L-1]})$$

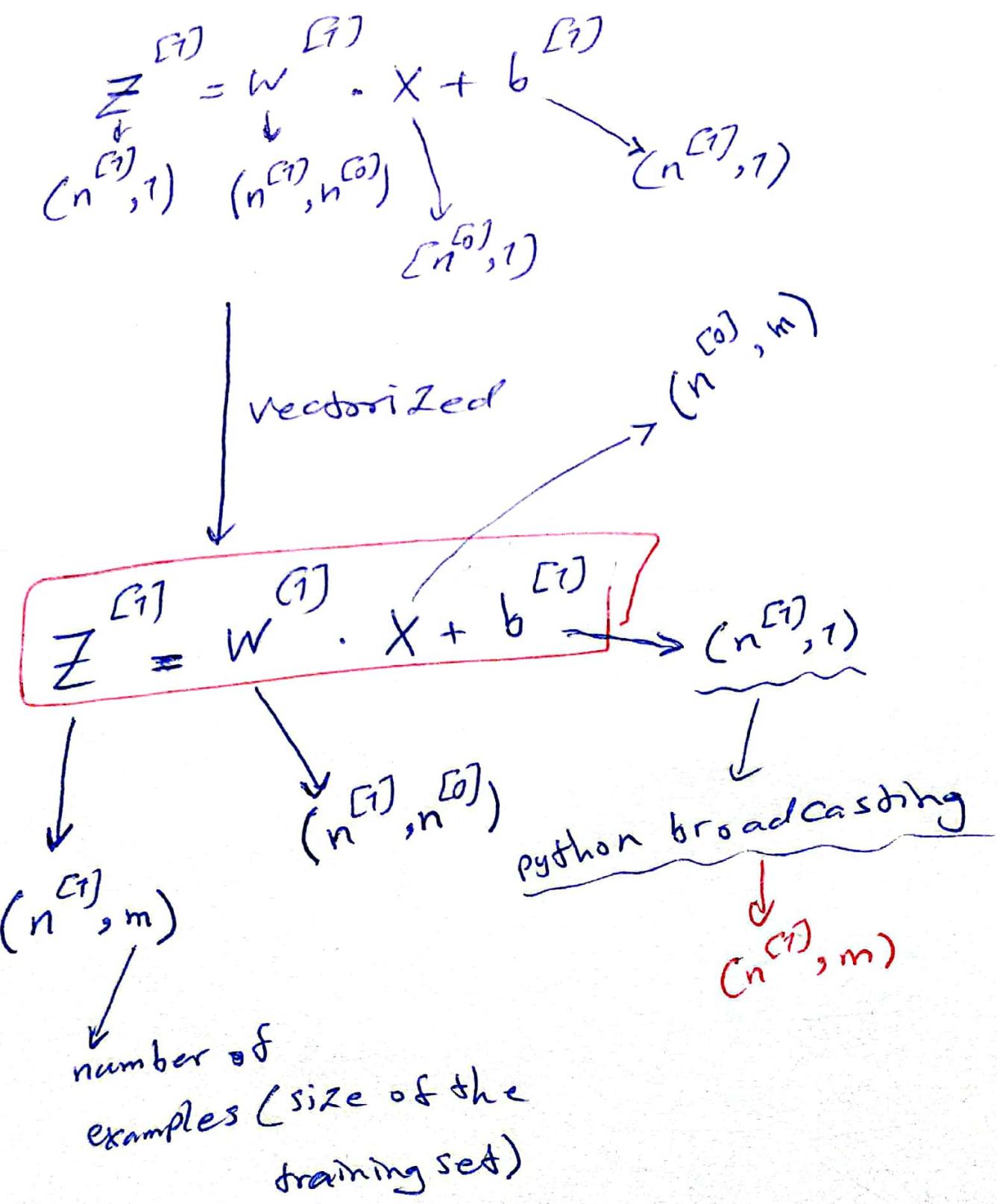
$$b^{[L]} = (n^{[L]}, 1)$$

$$dw^{[L]} = (w^{[1]}, w^{[2]}, \dots, w^{[L-1]})$$

$$[::] = [::; ::] \cdot [::]$$

$$z^{[L]} = a^{[L]} = (n^{[L]}, 1)$$

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



$$Z^{[l]} = A^{[l]} \cdot s \quad (n^{[l]}, m)$$

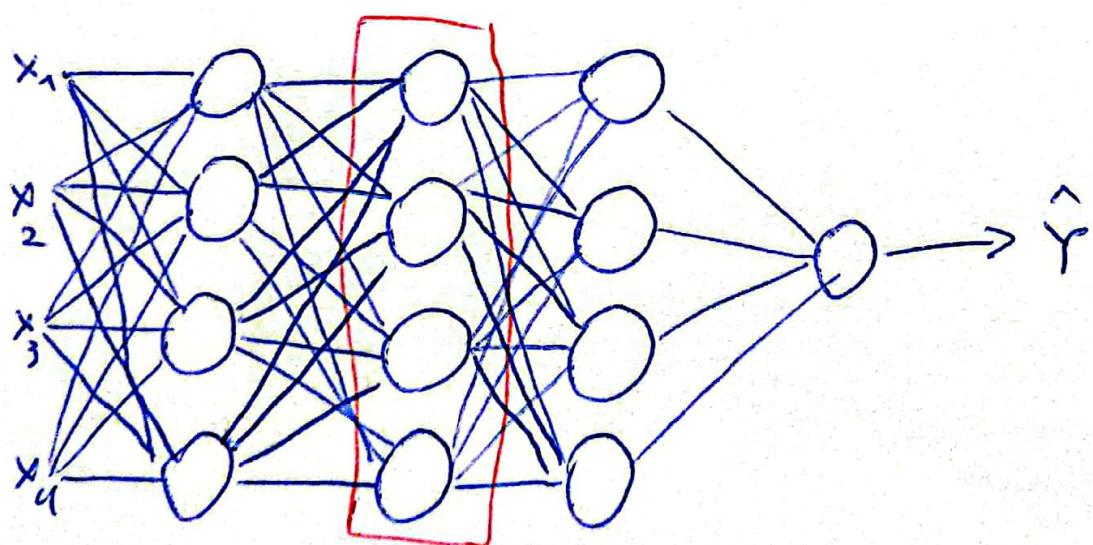
$$A^{[0]} \cdot s = X = (n^{[0]}, m)$$

$$dZ^{[l]} = dA^{[l]} \cdot s = (n^{[l]}, m)$$

why deep representations?

"Circuit theory" and deep learning
 number of units is relatively small
 informally, there are functions you can compute with
 a small L-layer deep neural network that shallower
 networks require exponentially more hidden units
 to compute.

"forward and backward functions"



layer L : $w^{(L)}; b^{(L)}$

forwards: input $a^{(L-1)}$; output $a^{(L)}$

$$Z^{(L)} = w^{(L)} \cdot a^{(L-1)} + b^{(L)}; \text{ cache } Z^{(L)}$$

$$a^{(L)} = g^{(L)}(Z^{(L)})$$

Backward: input: $\underline{da}^{(L)}$; output: $\underline{da}^{(L-1)}$

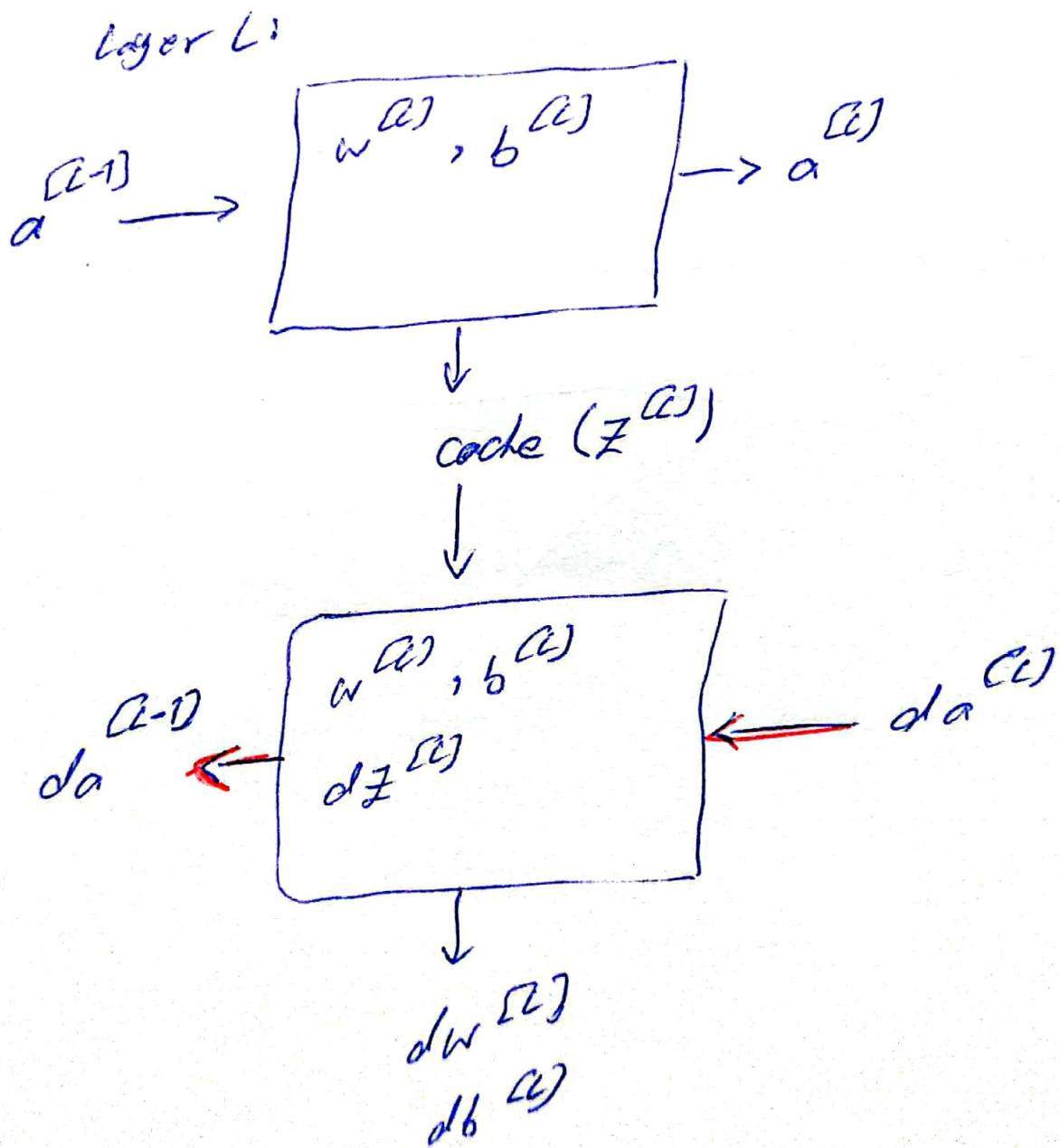
(39)

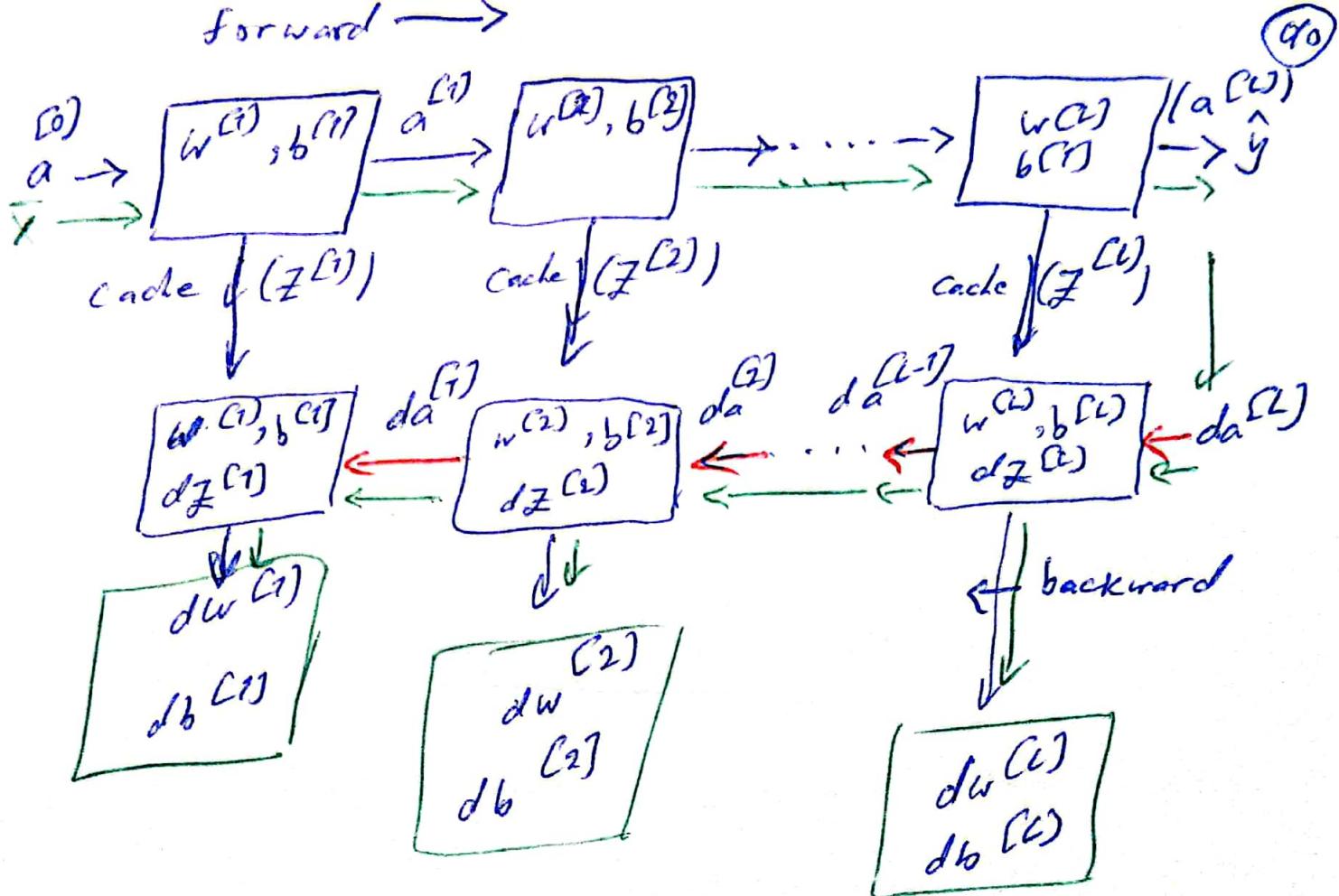
cache ($Z^{(L)}$)

~~dw~~

$dw^{(L)}$

$db^{(L)}$





$$w^{(l)} := w^{(l)} - \alpha dw^{(l)}$$

$$b^{(l)} := b^{(l)} - \alpha db^{(l)}$$

forward propagation for layer L

input: $a^{(L-1)}$

output: $a^{(L)}$; cache($Z^{(L)}$)

$$Z^{(L)} = w^{(L)} \cdot a^{(L-1)} + b^{(L)}$$

$$a^{(L)} = g^{(L)}(Z^{(L)})$$

vectorized

$$Z^{(L)} = w^{(L)} \cdot A^{(L-1)} + b^{(L)}$$

$$A^{(L)} = g^{(L)}(Z^{(L)})$$

$x = A^{(0)} \rightarrow \square \rightarrow \square \rightarrow \square \rightarrow \dots$

backward propagation for layer L:

Input $da^{[L]}$

Output $da^{[L-1]}, dw^{[L]}, db^{[L]}$

$$\delta Z^{[L]} = da^{[L]} \times g^{(L)'}(Z^{[L]})$$

$$dw^{[L]} = \delta Z^{[L]} \cdot a^{[L-1]T}$$

$$db^{[L]} = \delta Z^{[L]}$$

$$da^{[L-1]} = w^{[L]T} \cdot \delta Z^{[L]}$$

$$\delta Z^{[L]} = w^{[L-1]T} \cdot \delta Z^{[L+1]} \times g^{(L)'}(Z^{[L]})$$

vectorized version

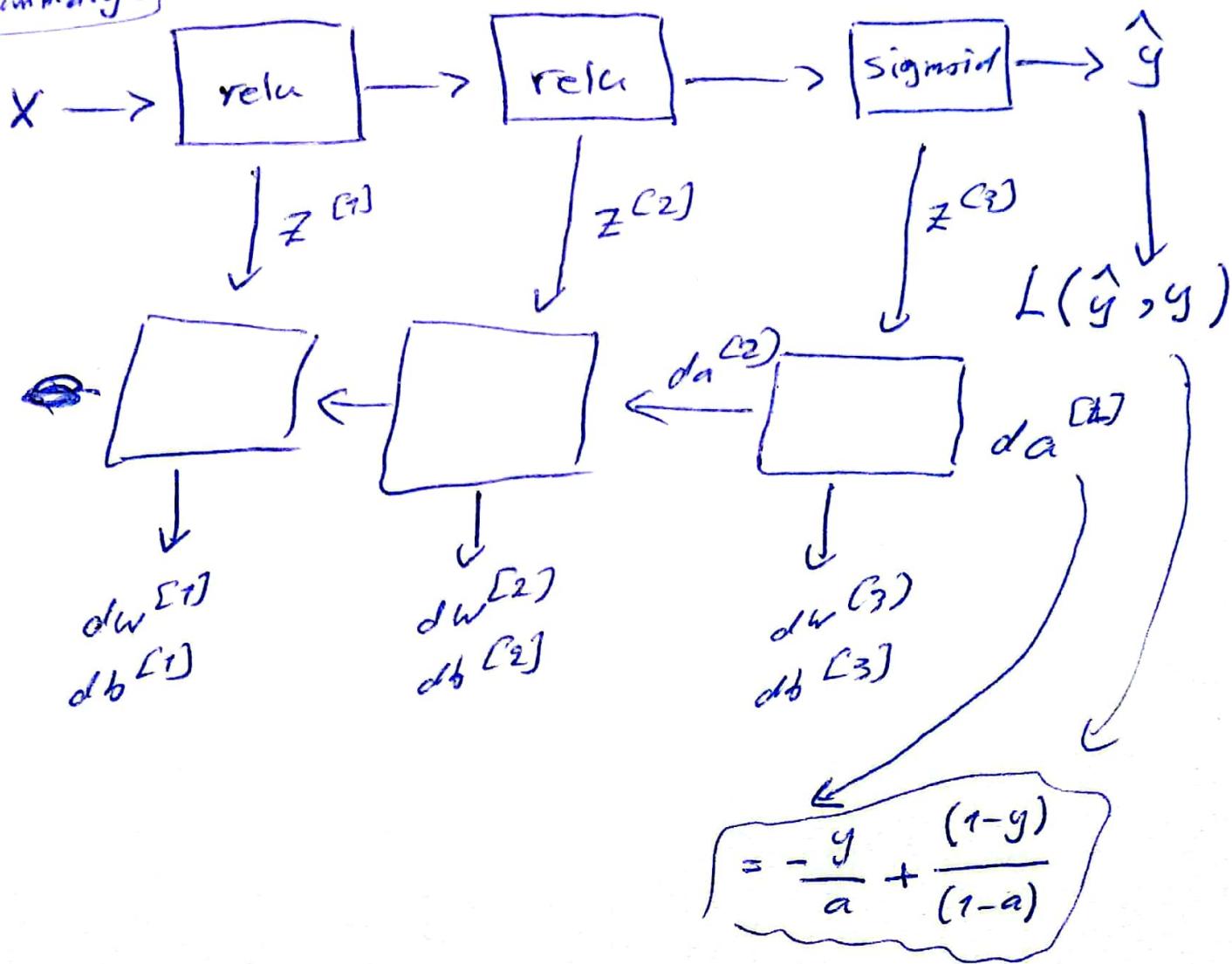
$$\delta Z^{[L]} = dA^{[L]} \times g^{(L)'}(Z^{[L]})$$

$$dw^{[L]} = \frac{1}{m} \delta Z^{[L]} \cdot A^{[L-1]T}$$

$$db^{[L]} = \frac{1}{m} \cdot np \cdot \text{sum}(\delta Z^{[L]}, \text{axis}=1, \text{keepdims=True})$$

$$dA^{[L-1]} = w^{[L]T} \cdot \delta Z^{[L]}$$

Summary:



Parameters vs. hyperParameters

Parameters: $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(L)}, b^{(L)}$

Hyper Parameters: learning rate (α)

iterations

hidden layers (L)

hidden units ($n^{(1)}, n^{(2)}, \dots$)

choice of activation function

momentum, minibatch size, regularization
parameters