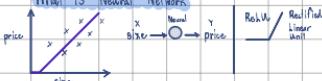


What is Neural Networks



Supervised learning with Neural network

train with label → has input and output final correct

Standard Neural Network → data with simple input output

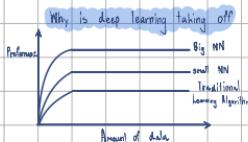
Convolutional Neural Network → image

Recurrent Neural Network → audio → text, translator

Custom/Hybrid Neural Network → mix

Structured Data → with columns or rows (spiral data)

Unstructured Data → Audio, Image, Text



Data → more training data

Computation → fast training (better hardware)

Algorithms → performance training

Binary Classification E.g. Input $\rightarrow \text{Output}$ by using 0 or 1

Logistic Regression

Input \rightarrow output \rightarrow $P(y=1|x) = \sigma(w^T x + b)$ (sigmoid function)

$$P(y=1|x) = \frac{1}{1+e^{-w^T x - b}}$$

Cost function $J(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$

$$J(\hat{y}, y) = -[y \log \hat{y} + (1-y) \log(1-\hat{y})]$$

If $y=1$: $J(\hat{y}, y) = -\log \hat{y} \rightarrow$ want $\log \hat{y}$ large

If $y=0$: $J(\hat{y}, y) = -\log(1-\hat{y}) \rightarrow$ want $\log(1-\hat{y})$ small

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

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

Gradient Descent → find best value of w, b to minimize $J(w, b)$ Global minimum

$$w = w - \alpha \frac{\partial J(\hat{y}, y)}{\partial w} \quad b = b - \alpha \frac{\partial J(\hat{y}, y)}{\partial b}$$

Computation Graph → explain why it organizes compute gradients or compute derivatives.

Derivatives with Computation Graphs

$$a = b$$

$$b = 3$$

$$c = 2$$

$$\frac{\partial (3a+3b)}{\partial a} = 3$$

$$\frac{\partial (3a+3b)}{\partial a} = 3 \quad \frac{\partial (3a+3b)}{\partial b} = 3$$

$$\frac{\partial (3a+3b)}{\partial a} = 3 \quad \frac{\partial (3a+3b)}{\partial b} = 3$$

$$\frac{\partial (3a+3b)}{\partial a} = 3 \quad \frac{\partial (3a+3b)}{\partial b} = 3$$

$$\frac{\partial (3a+3b)}{\partial a} = 3 \quad \frac{\partial (3a+3b)}{\partial b} = 3$$

$$\frac{\partial (3a+3b)}{\partial a} = 3 \quad \frac{\partial (3a+3b)}{\partial b} = 3$$

logistic Regression Gradient Descent (using for update parameter)

$$x \cdot w_0 + w_1 x + b \rightarrow \hat{y} = a = b(x) \rightarrow \hat{y}(a, x)$$

$$\begin{aligned} \frac{\partial J}{\partial w_0} &= \frac{\partial J(a, y)}{\partial w_0} = \frac{\partial a}{\partial w_0} = \frac{\partial a}{\partial x} \cdot \frac{\partial x}{\partial w_0} = a(1-a) \cdot 1 \\ \frac{\partial J}{\partial w_1} &= \frac{\partial J(a, y)}{\partial w_1} = \frac{\partial a}{\partial w_1} = a(1-a) \cdot x \end{aligned}$$

$$\begin{aligned} \frac{\partial J}{\partial w_0} &= \frac{\partial J}{\partial w_0} = x_1 \cdot \frac{\partial J}{\partial x_1} = \frac{\partial J}{\partial x_1} = x_1 \cdot d_n \\ \frac{\partial J}{\partial w_1} &= \frac{\partial J}{\partial w_1} = x_2 \cdot \frac{\partial J}{\partial x_2} = \frac{\partial J}{\partial x_2} = x_2 \cdot d_n \end{aligned}$$

Gradient Descent on m examples

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

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

$$\frac{\partial J}{\partial w_j} = \frac{1}{m} \sum_{i=1}^m \frac{\partial J}{\partial w_j}(a^{(i)}, y^{(i)}) = \frac{\partial J}{\partial w_j}(a^{(i)}, y^{(i)})$$

coding: $J = 0, dw_1 = 0, dw_2 = 0, db = 0$

for $i=1$ to m

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

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

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

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

$$\text{plus value gradient } dw_1 += y^{(i)} \cdot da^{(i)}$$

$$dw_2 += x^{(i)} \cdot da^{(i)}$$

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

Average Gradients

$$J = \frac{1}{m}, dw_1 = \frac{1}{m} \cdot dw_1, dw_2 = \frac{1}{m} \cdot dw_2$$

$$w_1 = w_1 - \alpha \cdot dw_1, w_2 = w_2 - \alpha \cdot dw_2$$

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

→ multiply more times (np.dot, np.sum) → SIMD - single instruction multiple data

Vectorization ($\mathbf{z} = \mathbf{w}^T \mathbf{x} + b$)

examples

$$\mathbf{y} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \rightarrow \mathbf{w}^T \mathbf{x} + b$$

another function

$\mathbf{w} = \text{np.exp}(\mathbf{V})$ | $\text{np.log}()$, $\text{np.abs}()$, np.minimum
with for loop
 $\mathbf{w} = \text{np.zeros}(n, 1)$

for i in range(n):

$\mathbf{w}[i] = \text{math.exp}(\mathbf{V}[i])$

Vectorizing Logistic Regression

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

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

$$\mathbf{z} = \mathbf{I}_m \mathbf{x}^{(1)} \quad \mathbf{x}^{(1)} \dots \mathbf{z}^{(m)} \mathbf{x} = \mathbf{I}_m \mathbf{x} + \mathbf{b} \quad b_1, b_2, \dots, b_d$$

call ii) Broadcasting = $\mathbf{w}^T \mathbf{x} + b = \mathbf{w}^T \mathbf{x} + b \dots \mathbf{w}^T \mathbf{x} + b$

$$\begin{aligned} \mathbf{z} &= \text{np.dot}(\mathbf{w}, \mathbf{x}) + b \quad \text{↑ broadcast} \\ \mathbf{a} &= \text{I}_{\mathbf{a}^{(1)} \dots \mathbf{a}^{(m)}} = \text{G}(\mathbf{z}) \end{aligned}$$

Vectorizing logistic Regression's Gradient Computation

$$\mathbf{d}\mathbf{z}^{(i)} = \mathbf{a}^{(i)} - \mathbf{y}^{(i)}, \text{np.sum}(\mathbf{d}\mathbf{z}^{(i)}) \cdot \mathbf{a}^{(i)} - \mathbf{y}^{(i)}$$

$$\mathbf{d}\mathbf{z} = \mathbf{I}_m \mathbf{d}\mathbf{z}^{(1)} \dots \mathbf{d}\mathbf{z}^{(m)}$$

$$\mathbf{A} = \mathbf{I}_{\mathbf{d}\mathbf{z}^{(1)} \dots \mathbf{d}\mathbf{z}^{(m)}} \quad \mathbf{Y} = \mathbf{I}_{\mathbf{y}^{(1)} \dots \mathbf{y}^{(m)}}$$

$$\begin{aligned} \mathbf{d}\mathbf{z} &= \mathbf{A} \cdot \mathbf{Y} = \mathbf{I}_{\mathbf{d}\mathbf{z}^{(1)} \dots \mathbf{d}\mathbf{z}^{(m)}} \cdot \mathbf{Y} = \mathbf{I}_{\mathbf{d}\mathbf{z}^{(1)} \dots \mathbf{d}\mathbf{z}^{(m)} \cdot \mathbf{y}^{(1)} \dots \mathbf{y}^{(m)}} \\ \text{sum } \mathbf{d}\mathbf{z} &= 0 \quad \mathbf{d}\mathbf{b} = 0 \quad \mathbf{d}\mathbf{b} = 0 \quad \mathbf{d}\mathbf{b} = 0 \quad \mathbf{d}\mathbf{b} = 0 \\ \mathbf{d}\mathbf{w} &= \mathbf{X}^T \cdot \mathbf{d}\mathbf{z} = \mathbf{X}^T \cdot \mathbf{A} \cdot \mathbf{Y} = \mathbf{X}^T \cdot \mathbf{d}\mathbf{z}^{(1)} \dots \mathbf{d}\mathbf{z}^{(m)} \\ \mathbf{d}\mathbf{w} &= \mathbf{I}_m \mathbf{d}\mathbf{z}^{(1)} \dots \mathbf{d}\mathbf{z}^{(m)} \quad \mathbf{d}\mathbf{b} = \mathbf{d}\mathbf{z}^{(1)} \dots \mathbf{d}\mathbf{z}^{(m)} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{d}\mathbf{w} &= \mathbf{I}_m \quad \mathbf{d}\mathbf{b} = \mathbf{I}_m \end{aligned}$$

Broadcasting in Python \Rightarrow make code run faster

$\mathbf{A} = \text{np.array}([3, 2, 1], \text{ndmin}=2)$ → sum of num in any columns

$\mathbf{a} = \mathbf{A} \cdot \text{sum}(\mathbf{A}, \text{axis}=0)$ → sum of num in any rows by axis

percentage = $100 \cdot \mathbf{A} / (\mathbf{A}. \text{reshape}((3, 1)))$

Example

$$\begin{aligned} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \cdot \begin{bmatrix} 10 & 11 & 12 \\ 13 & 14 & 15 \\ 16 & 17 & 18 \end{bmatrix} &= \begin{bmatrix} 101 & 102 & 103 \\ 134 & 135 & 136 \\ 167 & 168 & 169 \end{bmatrix} \\ \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \cdot \begin{bmatrix} 10 & 11 & 12 \\ 13 & 14 & 15 \\ 16 & 17 & 18 \end{bmatrix} &= \begin{bmatrix} 101 & 102 & 103 \\ 134 & 135 & 136 \\ 167 & 168 & 169 \end{bmatrix} \\ \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \cdot \begin{bmatrix} 10 & 11 & 12 \\ 13 & 14 & 15 \\ 16 & 17 & 18 \end{bmatrix} &= \begin{bmatrix} 101 & 102 & 103 \\ 134 & 135 & 136 \\ 167 & 168 & 169 \end{bmatrix} \end{aligned}$$

A node python/numpy vectors

→ return a Transpose

$\mathbf{a} = \text{np.random.randn}(5, 1); \text{print}(\mathbf{a}, \mathbf{T})$

print(np.dot(a, a.T)) → transpose value

$\mathbf{a} = \text{np.random.randn}(5, 1);$

print(np.dot(a, a.T)) → product value

↑ check condition

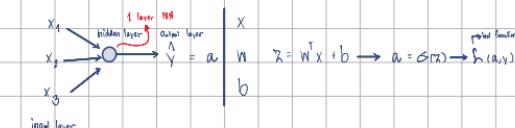
assert(a.shape == (5, 1))

Explain of logistic regression cost function

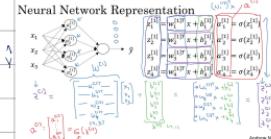
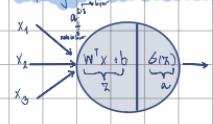
$$\begin{aligned} \mathbf{y} &= \sigma(\mathbf{w}^T \mathbf{x} + b) \quad ; \quad \mathcal{L}(\mathbf{y}) = \frac{1}{2} \mathbf{y}^2 \\ \mathbf{y} &= \sigma(\mathbf{y} = 1) \end{aligned}$$

$$\begin{aligned} \text{if } \mathbf{y} = 1 : \quad p(\mathbf{y} | \mathbf{x}) &= \mathbf{y} \\ \mathbf{y} = 0 : \quad p(\mathbf{y} | \mathbf{x}) &= 1 - \mathbf{y} \end{aligned} \quad \left. \begin{array}{l} \mathbf{p}(\mathbf{y} | \mathbf{x}) \\ \mathbf{p}(\mathbf{y} | \mathbf{x}) \end{array} \right\} \mathcal{L}(\mathbf{y})$$

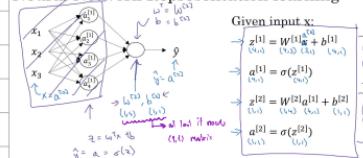
Neural Network Overview



Computing Neural Network Output

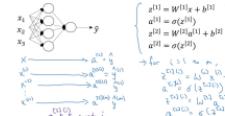


Neural Network Representation learning



Vectorize Across Multiple Examples

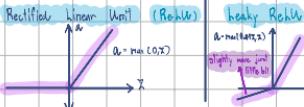
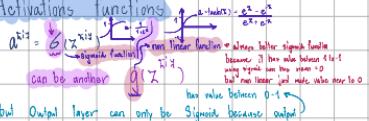
Vectorizing across multiple examples



Vectorizing across multiple examples



Activations Functions



Why Non-linear Activation function
e.g. decision boundary
→ using for learn image, sounds, languages

Derivatives of activation functions

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

Slope: $\frac{d}{dx} g(x) = \frac{1}{1 + e^{-x}} \cdot \left(-\frac{1}{1 + e^{-x}} \right) = g(x)(1 - g(x))$

$x=10, g(x) \approx 1, \frac{d}{dx} g(x) \approx (1)(1-0) \approx 0$

$x=-10, g(x) \approx 0, \frac{d}{dx} g(x) \approx (0)(1-0) \approx 0$

$x=0, g(x) \approx \frac{1}{2}, \frac{d}{dx} g(x) \approx \frac{1}{2}(1-\frac{1}{2}) = \frac{1}{4}$

Tanh activation function: $g(x) = \tanh(x)$

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

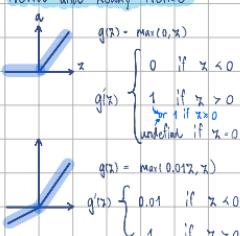
$$\frac{d}{dx} g(x) = 1 - (\tanh(x))^2$$

$x=10, \tanh(x) \approx 1, g'(x) \approx 0$

$x=-10, \tanh(x) \approx -1, g'(x) \approx 0$

$x=0, \tanh(x)=0, g'(x)=1$

ReLUs and Leaky ReLUs



Gradient Descent for Neural networks

Parameters $\rightarrow w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}$
 $(w^{(1)}, b^{(1)}) \in \mathbb{R}^{n^{(1)} \times n^{(2)}} \quad (w^{(2)}, b^{(2)}) \in \mathbb{R}^{n^{(2)} \times 1}$

Cost function $\rightarrow J(w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}) = \frac{1}{m} \sum_{i=1}^m L(y_i, \hat{y}_i)$

Gradient Descent \rightarrow loop to compute gradients ($\frac{\partial J}{\partial w^{(1)}}, \dots, \frac{\partial J}{\partial b^{(2)}}$)

$$\frac{\partial J}{\partial w^{(1)}} = \frac{\partial J}{\partial w^{(1)}} \cdot \frac{\partial w^{(1)}}{\partial w^{(1)}} = \frac{\partial J}{\partial w^{(1)}}$$

$$w^{(1)} = \frac{\partial J}{\partial w^{(1)}} \cdot -\alpha \cdot \frac{\partial J}{\partial w^{(1)}}$$

$$b^{(1)} = b^{(1)} - \alpha \cdot \frac{\partial J}{\partial b^{(1)}}$$

Formulas for computing derivatives

Forward Propagation:

$$z^{(1)} = W^{(1)} X + b^{(1)}$$

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

$$z^{(2)} = W^{(2)} a^{(1)} + b^{(2)}$$

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

for which sigmoid function

Backpropagation (Calculate gradient descent)

$$\frac{\partial J}{\partial z^{(2)}} = a^{(2)} - y^{(2)}$$

(how different value between prediction and target)

$$\frac{\partial J}{\partial w^{(2)}} = \frac{\partial J}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial w^{(2)}}$$

gradient descent rule: $\frac{\partial J}{\partial w^{(2)}} = \frac{\partial J}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial w^{(2)}}$

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

np.sum($\frac{\partial J}{\partial z^{(2)}}, \text{axis}=1, \text{keepdims=True}$)

$$\frac{\partial J}{\partial z^{(1)}} = (W^{(2)})^\top \frac{\partial J}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial z^{(1)}}$$

derivative of activation function

$$\frac{\partial J}{\partial w^{(1)}} = \frac{1}{2} \frac{\partial J}{\partial z^{(1)}}$$

$\frac{\partial J}{\partial w^{(1)}} = \frac{1}{2} \frac{\partial J}{\partial z^{(1)}} \cdot \frac{\partial z^{(1)}}{\partial w^{(1)}}$

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

np.sum($\frac{\partial J}{\partial z^{(1)}}, \text{axis}=1, \text{keepdims=True}$)

Random Initialization

both hidden unit are still computing exactly the same function

What happen if you initialize weight to zero → model still same same function

How to do:

- $w^{(1)} = np.random.randn(2, 3)$ → make weight matrix
- $b^{(1)} = np.zeros((2, 1))$ → make bias vector
- $w^{(1)} = \dots, b^{(1)} = 0$ → make bias vector

Deep k-layer Network

Shallow \rightarrow 1 NN, Deep \Rightarrow multiple NN

Deep neural network notation

$$x_1 \rightarrow \text{Hidden} \rightarrow \text{Hidden} \rightarrow \text{Hidden} \rightarrow y$$

$$z_1, z_2, z_3, \dots, z_k$$

$$a_1, a_2, a_3, \dots, a_k$$

$$w^{(1)}, w^{(2)}, w^{(3)}, \dots, w^{(k)}$$

$$n^{(1)}, n^{(2)}, n^{(3)}, \dots, n^{(k)}$$

n_{out} = k

$$h = k \text{ (H layers)}$$

$$n^{(k)} = H \text{ units in layer } k$$

$$n^{(k)} = n_j$$

z_k = activations in layer k

$$a_k = f(z_k), w \text{ & } b = \text{weights for } z_k$$

Forward Propagation in a Deep learning

$$\begin{array}{c} \text{Input } x_1 \\ \text{Input } x_2 \\ \text{Input } x_3 \end{array} \xrightarrow{\text{Neuron 1}} \text{Neuron 2} \xrightarrow{\text{Neuron 3}} \text{Neuron 4} \xrightarrow{\text{Neuron 5}}$$

$$x^{(1)} = W^{(1)} x + b^{(1)}$$

$$x^{(2)} = W^{(2)} x^{(1)} + b^{(2)}$$

$$a^{(1)} = g^{(1)}(x^{(1)})$$

$$a^{(2)} = g^{(2)}(x^{(2)})$$

Output layer

$$a^{(3)} = g^{(3)}(x^{(2)}) = y$$

vectorize

$$x^{(1)} = W^{(1)} x + b^{(1)}$$

$$A^{(1)} = g^{(1)}(x^{(1)})$$

$$\begin{pmatrix} 1 & x_1 & x_2 & x_3 \end{pmatrix} \begin{pmatrix} 1 & a^{(1)} \\ A^{(1)} & 1 \end{pmatrix} = \begin{pmatrix} 1 & x_1 & x_2 & x_3 \end{pmatrix}$$

$$A^{(2)} = W^{(2)} A^{(1)} + b^{(2)}$$

$$A^{(3)} = g^{(3)}(A^{(2)}) = Y$$

$$Y = g(A^{(3)}) = Y$$

Getting your matrix dimensions right

$$x \rightarrow 3N \rightarrow 3N \rightarrow MN \rightarrow 2N \rightarrow 1N \rightarrow Y$$

Differentiator $W^{(1)}, b^{(1)}$

$$x^{(1)} = W^{(1)} x + b^{(1)}$$

$$(3,1) \quad (3,2) \quad (2,1) \quad (2,1)$$

$$W^{(1)} = (n^{(1)}, n^{(2)})$$

$$x^{(2)} = W^{(2)} x^{(1)} + b^{(2)}$$

$$(5,1) \quad (5,3) \quad (3,1) \quad (5,1)$$

$$x^{(3)} = W^{(3)} x^{(2)} + b^{(3)}$$

$$(5,1) \quad (5,1) \quad (1,1) \quad (1,1)$$

$$Y = g^{(3)}(x^{(3)})$$

$$\text{In my Neuron 3 having different } \hat{x} \rightarrow 4 \text{ Neurons} \rightarrow 3 \text{ Neurons} \rightarrow 1 \text{ Neuron} \rightarrow Y$$

Circuit theory and deep learning

There are functions you can compute with a

"small" h-layer deep neural network that shallow networks require exponentially more hidden layers

units to compute

\rightarrow complexity $\approx O(\log n)$

Efficient Deep Network $y = x_1 \text{XOR } x_2 \text{XOR } \dots \text{XOR } x_n \text{XOR }$

$$x_1 \rightarrow \text{XOR}$$

$$x_2 \rightarrow \text{XOR}$$

$$x_3 \rightarrow \text{XOR}$$

$$\vdots$$

$$\rightarrow \text{complexity } \approx O(2^n)$$

Inefficient Shallow Network

$$\begin{array}{c} x_1 \rightarrow 0 \\ x_2 \rightarrow 0 \\ x_3 \rightarrow 0 \\ x_4 \rightarrow 0 \end{array} \rightarrow \text{OR} \rightarrow Y$$

Building Blocks of a Deep neural network and propagation

Forward and backward functions

layer l: $W^{(l)}, b^{(l)}$

$$\begin{array}{c} \text{1) Forward: Input } x^{(l-1)} \text{, Output } a^{(l)} \\ \text{Calculation: } a^{(l)} = W^{(l)} x^{(l-1)} + b^{(l)} \\ x^{(l)} = W^{(l)} a^{(l)} + b^{(l)}; a^{(l)} = g^{(l)}(x^{(l)}) \end{array}$$

$$\begin{array}{c} \text{2) Backward: Input } x^{(l-1)}, \text{ Output } da^{(l-1)}, da^{(l)}, db^{(l)} \\ \text{Calculation: } da^{(l-1)} = da^{(l)} \cdot g'(x^{(l)}) \\ da^{(l)} = da^{(l)} \cdot a^{(l-1)} \\ db^{(l)} = da^{(l)} \\ da^{(l-1)} = W^{(l)} a^{(l-1)} \cdot da^{(l)} + g'(x^{(l)}) da^{(l)} \\ da^{(l)} = W^{(l+1)T} da^{(l+1)} + g'(x^{(l)}) da^{(l)} \end{array}$$

Vectorized version

$$\begin{array}{c} da^{(l-1)} = da^{(l)} \cdot g'(x^{(l)}) \\ da^{(l)} = da^{(l)} \cdot a^{(l-1)} \\ db^{(l)} = da^{(l)} \\ da^{(l-1)} = W^{(l)} a^{(l-1)} \cdot da^{(l)} + g'(x^{(l)}) da^{(l)} \\ da^{(l)} = W^{(l+1)T} da^{(l+1)} + g'(x^{(l)}) da^{(l)} \end{array}$$

$$\begin{array}{c} da^{(l-1)} = da^{(l)} \cdot g'(x^{(l)}) \\ da^{(l)} = da^{(l)} \cdot a^{(l-1)} \\ db^{(l)} = da^{(l)} \\ da^{(l-1)} = W^{(l)} a^{(l-1)} \cdot da^{(l)} + g'(x^{(l)}) da^{(l)} \\ da^{(l)} = W^{(l+1)T} da^{(l+1)} + g'(x^{(l)}) da^{(l)} \end{array}$$

da^{(l-1)} = da^{(l)} \cdot g'(x^{(l)}) da^{(l)} = 1/M \sum_i (da^{(l)}_i, axis=1, keepdims=True)

$$\begin{array}{c} \text{Summary: } X \rightarrow \text{ReLU} \rightarrow \text{ReLU} \rightarrow \text{sigmoid} \rightarrow Y \rightarrow \text{compute loss} \\ \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\ x^{(1)} \quad x^{(2)} \quad x^{(3)} \quad x^{(4)} \end{array}$$

$$\begin{array}{c} \text{backwards iteration: } \left\{ \begin{array}{c} da^{(4)} \\ da^{(3)} \\ da^{(2)} \\ da^{(1)} \end{array} \right\} \left\{ \begin{array}{c} da^{(4)} \\ da^{(3)} \\ da^{(2)} \\ da^{(1)} \end{array} \right\} \left\{ \begin{array}{c} da^{(4)} \\ da^{(3)} \\ da^{(2)} \\ da^{(1)} \end{array} \right\} \left\{ \begin{array}{c} da^{(4)} \\ da^{(3)} \\ da^{(2)} \\ da^{(1)} \end{array} \right\} \\ \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\ \frac{da^{(4)}}{db^{(4)}} \quad \frac{da^{(3)}}{db^{(3)}} \quad \frac{da^{(2)}}{db^{(2)}} \quad \frac{da^{(1)}}{db^{(1)}} \end{array}$$

$$\begin{array}{c} \text{Parameters vs Hyperparameters} \\ \text{Hyperparameter} \Rightarrow \text{e.g. learning rate, # iterations, # hidden layer, ...} \\ \text{Parameters} \Rightarrow \text{weights, biases} \end{array}$$

$$\frac{da^{(4)}}{db^{(4)}} = -\frac{2}{M} \sum_{i=1}^M (1-y_i) y_i$$