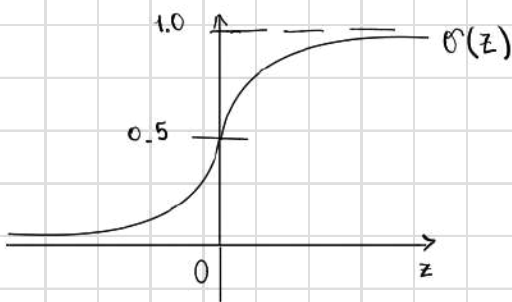


Logistic Regression

- Input: $x \in \mathbb{R}^{n \times 1}$ ("n" dimensional vector)
- Parameters: $w \in \mathbb{R}^{n \times 1}$, $b \in \mathbb{R}$.
- Output: $\hat{y} = P(y=1|x)$, $0 \leq \hat{y} \leq 1$.

$\hat{y} = w^T x + b$ doesn't make sense for logistic regression because \hat{y} can't be a large number. So we use the sigmoid function in this case: $\hat{y} = \sigma(w^T x + b) \rightarrow w^T x + b = z \rightarrow \hat{y} = \sigma(z)$



$$\sigma(z) = \frac{1}{1 + e^{-z}}; \quad \lim_{z \rightarrow \infty} \sigma(z) = 1; \quad \lim_{z \rightarrow -\infty} \sigma(z) = 0.$$

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

Loss Function For Logistic Regression (single error)

We don't use $L(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$ in L.R. because the optimization will not converge

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

• If $y=0$: $L(\hat{y}, y) = -\log \hat{y} \leftarrow$ want $\log(\hat{y})$ large, so want \hat{y} large

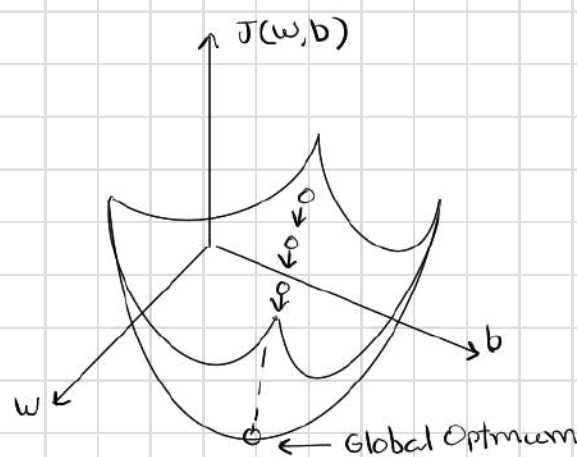
• If $y=1$: $L(\hat{y}, y) = -\log(1-\hat{y}) \leftarrow$ want $\log(1-\hat{y})$ large, so want \hat{y} small

Cost Function (average error)

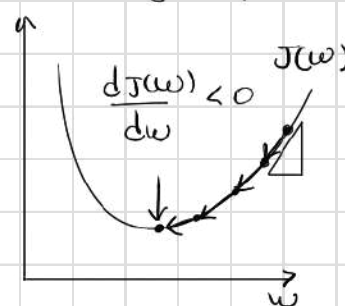
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \cdot \log \hat{y}^{(i)} + (1-y^{(i)}) \cdot \log(1-\hat{y}^{(i)})]$$

Gradient Descent

Want to find w, b to minimize $J(w, b)$.



Considering only w



Repeat until converged:

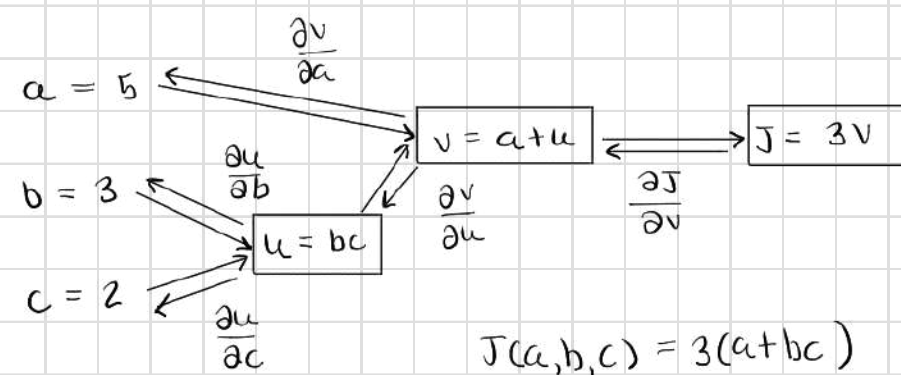
$$\{w := w - \alpha \frac{dJ(w)}{dw}\}$$

Considering multiple parameters, we use "partial derivative".

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

Computation Graph

To find partial derivatives we're going to consider the propagation throught the computational graph.



$$\frac{\partial J}{\partial u} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial u} = 3 \times 1 = 3 \quad \frac{\partial J}{\partial v} = 3$$

$$\frac{\partial J}{\partial a} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial a} = 3 \times 1 = 3$$

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial u} \cdot \frac{\partial u}{\partial b} = 3 \times c = 6$$

$$\frac{\partial J}{\partial c} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial u} \cdot \frac{\partial u}{\partial c} = 3 \times b = 9$$

Gradient Descent on "m" Example

$$z = w^T x + b, \hat{y} = a = \sigma(z); L(a, y) = -(y \log(a) + (1-y) \log(1-a))$$

$$\frac{\partial a}{\partial z} = a(1-a) \quad \frac{\partial L(a, y)}{\partial a} = -\frac{y}{a} + \frac{1-y}{1-a} \quad \frac{\partial L(a, y)}{\partial w_1} = x_1 \cdot \frac{\partial L(a, y)}{\partial z}$$

$$\frac{\partial L(a, y)}{\partial w_2} = x_2 \cdot \frac{\partial L(a, y)}{\partial z}$$

$$\frac{\partial L(a, y)}{\partial b} = \frac{\partial L(a, y)}{\partial z}$$

$$\frac{\partial L(a, y)}{\partial z} = \frac{\partial L(a, y)}{\partial a} \cdot \frac{\partial a}{\partial z} = \left(-\frac{y}{a} + \frac{1-y}{1-a}\right) a(1-a) = a - y$$

Algorithm:

$$J = 0; dw_1 = 0; dw_2 = 0; db = 0$$

For $i = 1$ to M :

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

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

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

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

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

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

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

$$J = J/M; dw_1 = dw_1/M; dw_2 = dw_2/M; db = db/M;$$

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

vectorized Implementation

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

$$A = \sigma(Z)$$

$$dZ = A - Y$$

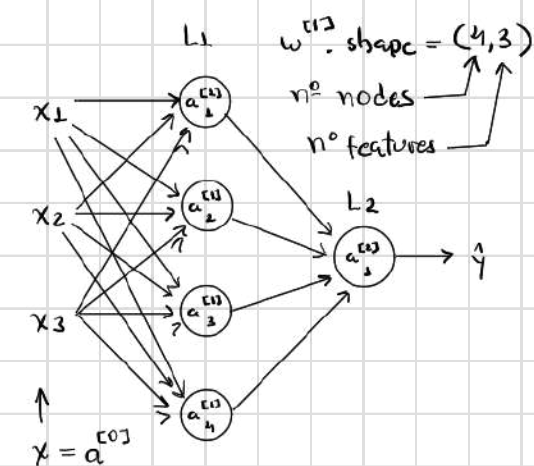
$$dw = \frac{1}{M} X dZ^T$$

$$db = \frac{1}{M} dZ$$

$$w := w - \alpha dw$$

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

Neural Network



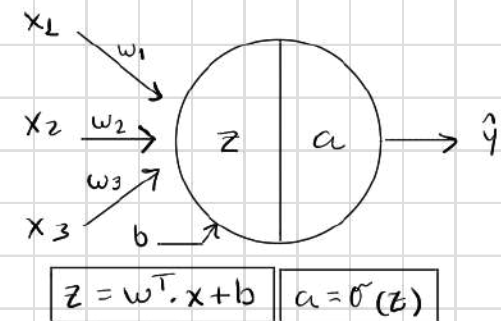
$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}, a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}, a_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}, a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}, a_4^{[1]} = \sigma(z_4^{[1]})$$

$$z_1^{[2]} = w_1^{[2]T} x + b_1^{[2]}, a_1^{[2]} = \sigma(z_1^{[2]})$$



vectorizing across multiple examples

for $i = 1$ to m :

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

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

$$z^{[2]}(i) = w^{[2]T} a^{[1]}(i) + b^{[2]}$$

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

$$Z^{[1]} = \begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} w_1^{[1]T} \cdot x + b_1^{[1]} \\ w_2^{[1]T} \cdot x + b_2^{[1]} \\ w_3^{[1]T} \cdot x + b_3^{[1]} \\ w_4^{[1]T} \cdot x + b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

$$Z^{[2]} = \begin{bmatrix} w_1^{[2]T} \end{bmatrix} \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix} + \begin{bmatrix} b_1^{[2]} \end{bmatrix} = \begin{bmatrix} w_1^{[2]T} \cdot a + b_1^{[2]} \end{bmatrix} = \begin{bmatrix} z_1^{[2]} \end{bmatrix}$$

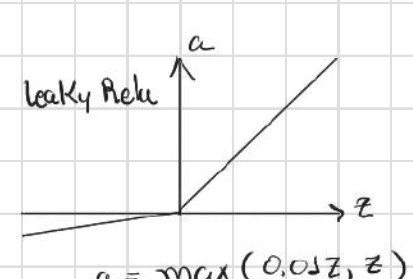
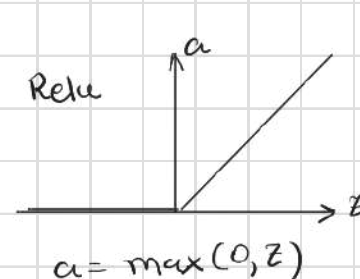
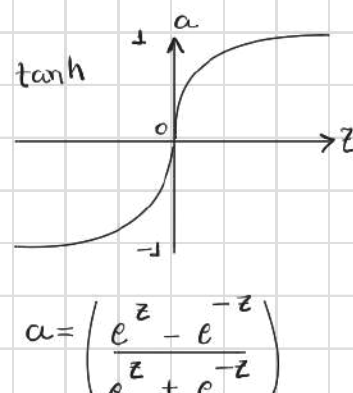
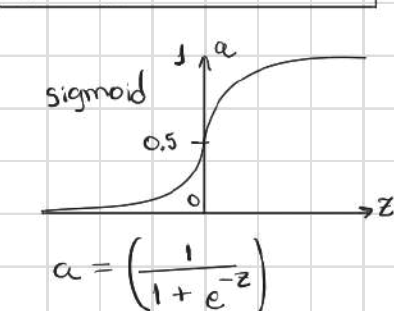
$$Z^{[1]} = W^{[1]} X + b^{[1]}$$

$$A^{[1]} = \sigma(Z^{[1]})$$

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

$$A^{[2]} = \sigma(Z^{[2]})$$

Activation Functions



Gradient Descent For Neural Networks

In this case, $g(z)$ represents a generic activation function.

Forward propagation

Back propagation

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

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

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

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

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

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

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

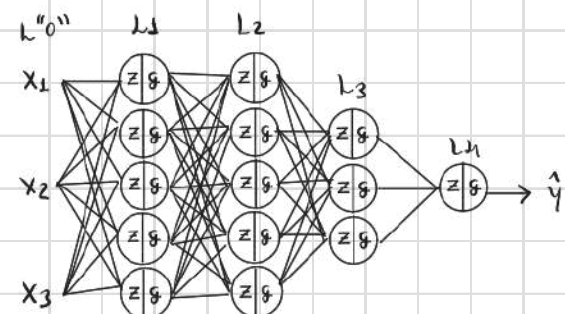
$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dz^{[1]} x^T \quad \text{elementwise product}$$

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

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

Deep Learning Notation



$$x = a^{[0]}$$

L = number of layers

 $n^{[l]}$ = number of units in layer "l" $a^{[l]}$ = activations in layer "l"

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

$$W^{[l]}, b^{[l]} = \text{weights for } z^{[l]}$$

Forward propagation

$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$$

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

Vectorized vectorization

$$Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$$

$$A^{[l]} = g^{[l]}(Z^{[l]})$$

$$\hat{y} = g^{[L]}(Z^{[L]}) = A^{[L]}$$

Parameters Dimensions

$$Z^{[l]} \rightarrow (n^{[l]}, 1)$$

$$W^{[l]} \rightarrow (n^{[l]}, n^{[l-1]})$$

$$A^{[l]} \rightarrow (n^{[l]}, 1)$$

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

$$dW^{[l]} \rightarrow (n^{[l]}, n^{[l-1]})$$

$$db^{[l]} \rightarrow (n^{[l]}, 1)$$

Notes that:

$$z^{[l]}, a^{[l]} \rightarrow (n^{[l]}, 1)$$

$$Z^{[l]}, A^{[l]} \rightarrow (n^{[l]}, m)$$

$$\text{for } l=0, A^{[0]} = X \rightarrow (n^{[0]}, m)$$

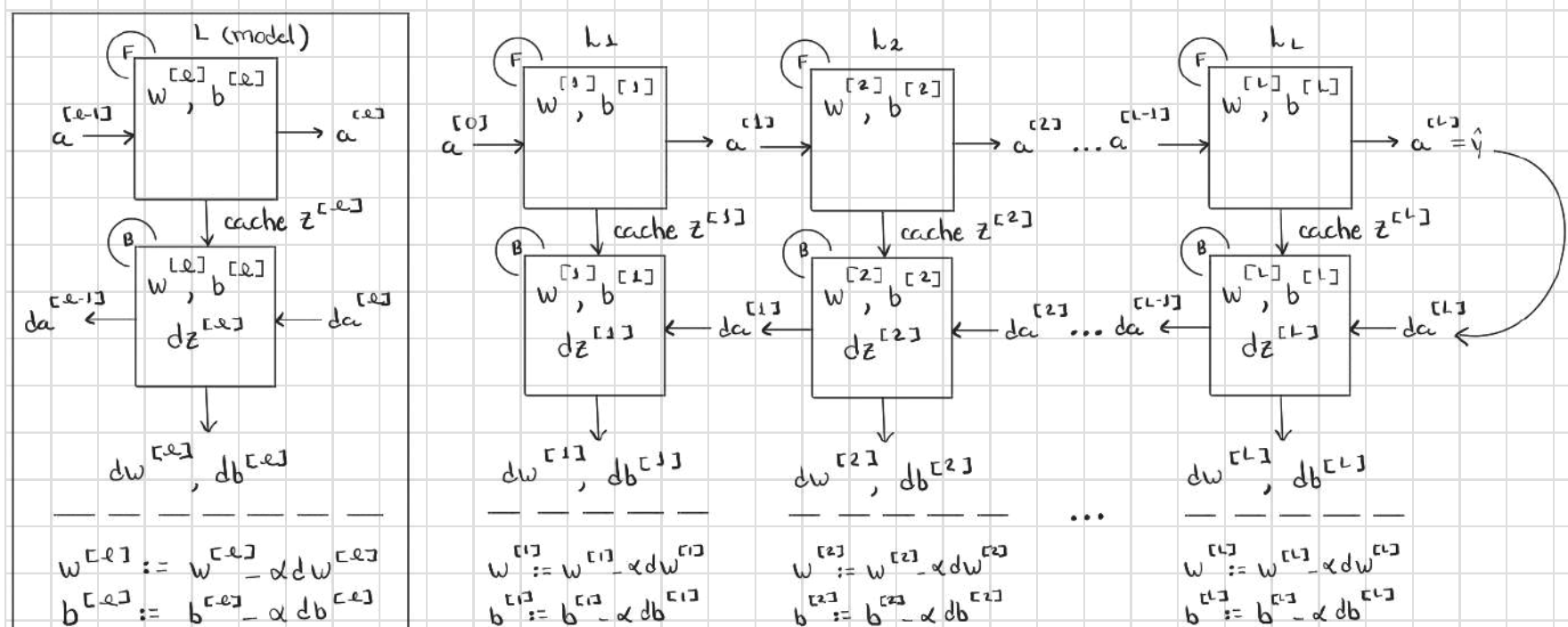
$$dz^{[l]}, dA^{[l]} \rightarrow (n^{[l]}, m)$$

$$Z^{[1]} = W^{[1]} \cdot X + b^{[1]}$$

$(n^{[0]}, m) \quad (n^{[1]}, n^{[0]}) \quad (n^{[0]}, m) \quad (n^{[1]}, m)$

python broadcast

Forward and Backward function blocks

Forward: Input $a^{[l-1]}$, output $a^{[l]}$, cache $z^{[l]}$. Backward: Input $da^{[l]}$, output $da^{[l-1]}$, $dW^{[l]}$, $db^{[l]}$.

ANNs and the human neuron

