# NEURAL NETWORKS & DEEP LEARNING

By: Mohamed Aziz Tousli

# BINARY CLASSIFICATION

- $(x, y), x \in R^{n_x}, y \in \{0,1\}$
- $m$: #Training examples
- $m_{test}$: #Test examples
- $X = \begin{bmatrix} \vdots & \cdots & \vdots \\ x^{(1)} & \cdots & x^{(m)} \\ \vdots & \cdots & \vdots \end{bmatrix} \in (n_x, m)$
- $Y = \begin{bmatrix} y^{(1)} & \cdots & y^{(m)} \end{bmatrix} \in (1, m)$

**Logistic regression**:

- $\hat{y} = P(y = 1|x) = \sigma(z + b); z = w^T \rightarrow \hat{Y} = \sigma(\Theta^T X)$

- Parameters: $w \in R^{n_x}, b \in R \rightarrow \Theta = \begin{bmatrix} \theta_0 = b \\ \theta_1 \\ \vdots \\ \theta_{n_x} \end{bmatrix}$ } w

- PS: $X_0 = 1, x \in R^{n_x} + 1$
- Sigmoid function: $\sigma(z) = \frac{1}{1+e^{-z}}$

**Loss error function**: $L(\hat{y}, y) = -y\log\hat{y} + (1 - y)\log(1 - \hat{y}) = \frac{1}{2}(\hat{y} - y)^2$

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

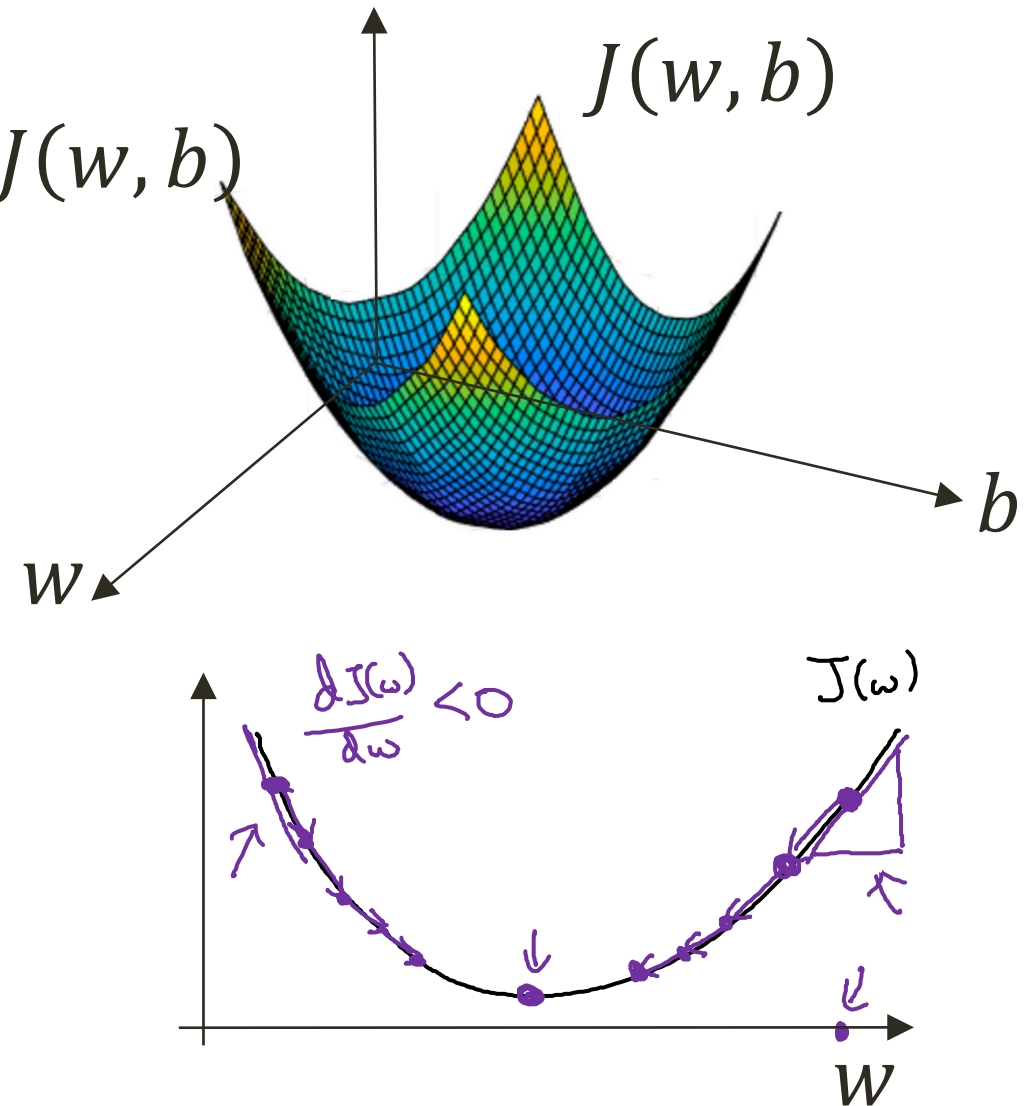# GRADIENT DESCENT

Want to find $w, b$ that minimize $J(w, b)$
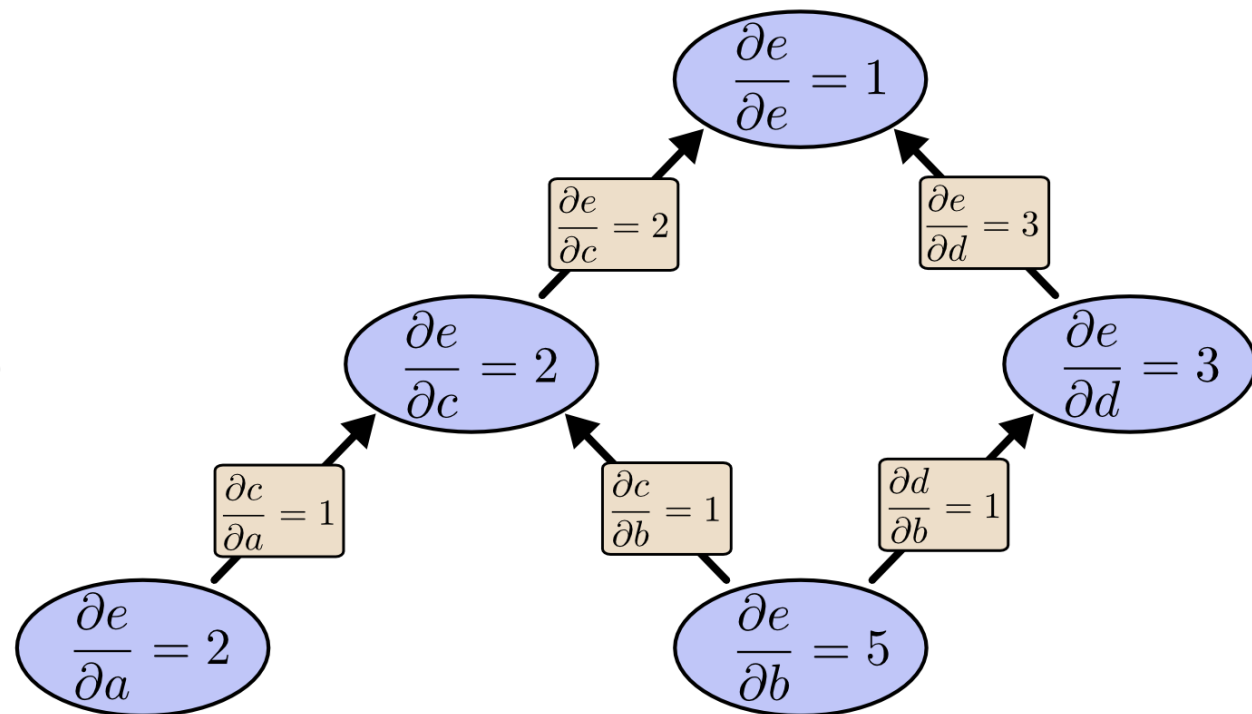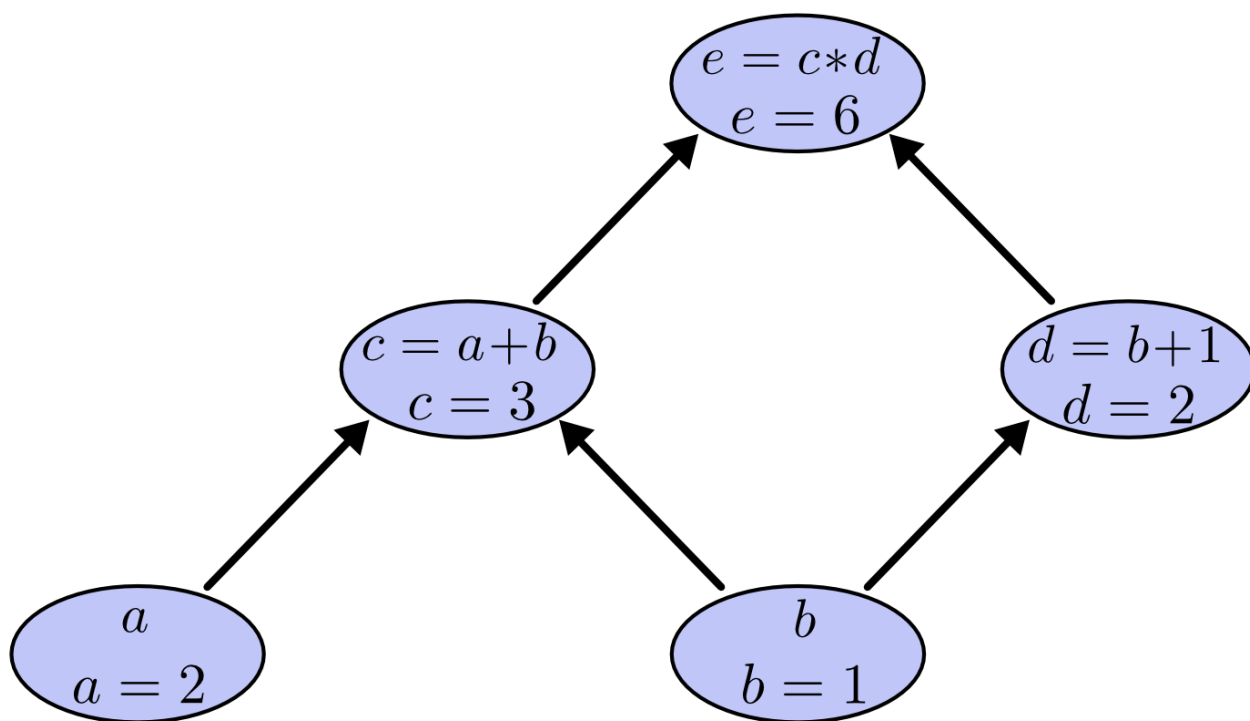
Repeat {

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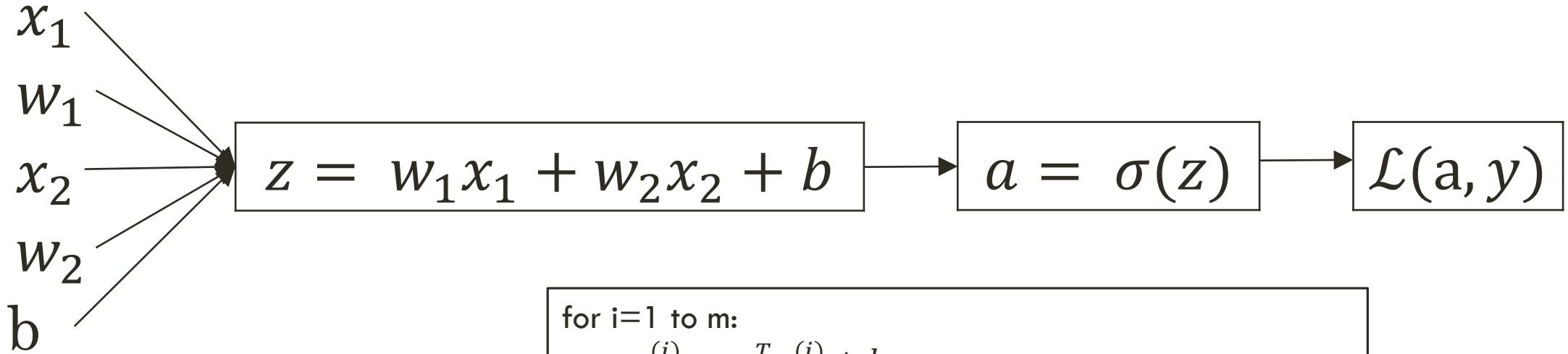
$$b := b - \alpha \frac{dJ(w, b)}{db}$$

}

$\alpha = learning\ rate$

$dw = \dfrac{dJ(w, b)}{dw}$

$db = \dfrac{dJ(w, b)}{db}$



$J(w, b)$

$\dfrac{dJ(\omega)}{d\omega} < 0$

$J(\omega)$

$w$

# COMPUTATION GRAPH

# LOGISTIC REGRESSION DERIVATIES

$x_1$

$w_1$

$x_2$

$w_2$

$b$

$$z = w_1 x_1 + w_2 x_2 + b$$

$$a = \sigma(z)$$

$$\mathcal{L}(a, y)$$

$$dz = \frac{dL}{dz} = a - y$$

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

$$dw = x . dz$$

$$db = dz$$

for i=1 to m:

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

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

$\quad J += -[y^{(i)} log a^{(i)} + (l1 - y^{(i)}) log(1 - a^{(i)})$

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

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

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

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

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

# VECTORIZATION

**Non vectorized:** For loop → Slow

**Vectorized:** Tensor calculations → Fast
- Whenever possible, avoid explicit for-loops
- $u = Av \rightarrow u = np.dot(A, v)$
- $w_1, w_2 \dots \rightarrow one\ W$

Python
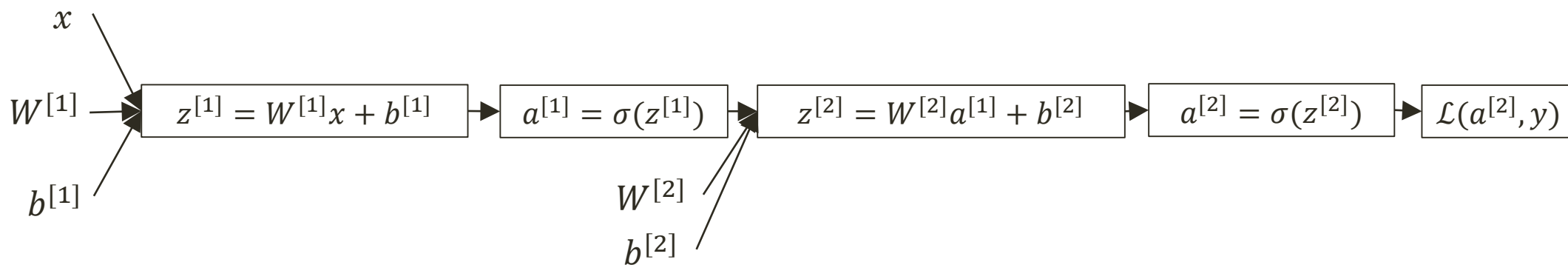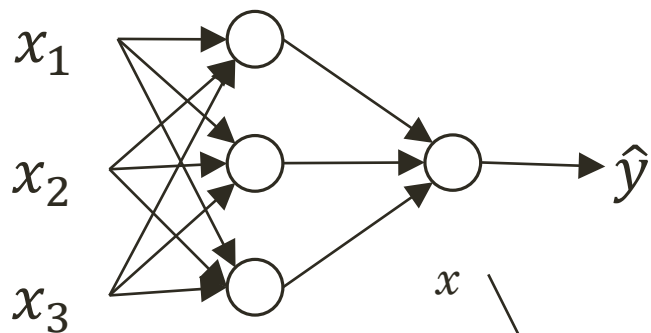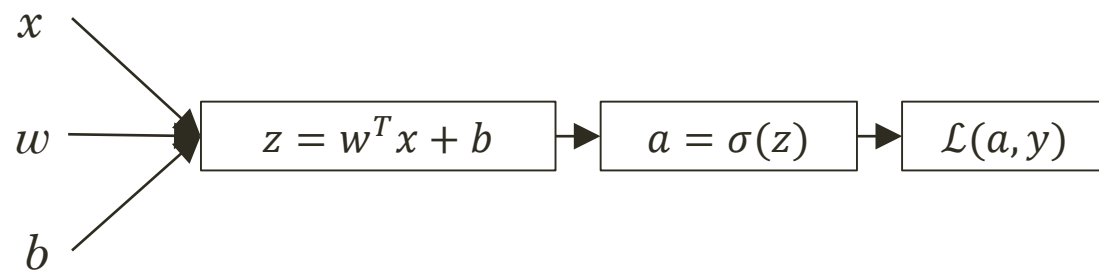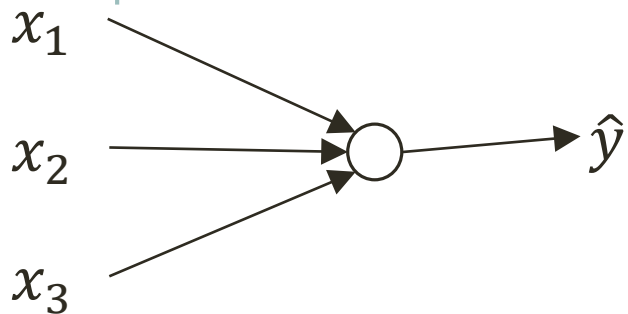Broadcasting

**Vectorizing logistic regression:**
- $Z = \left[z^{(1)} \dots z^{(n+1)}\right] = w^T X + [b \dots b] = \left[w^T x^{(1)} + b \dots w^T x^{(2)} + b\right] = \boldsymbol{np.dot(w.T, X) + b}$
- $A = \left[a^{(1)} \dots a^{(m)}\right] = \boldsymbol{\sigma(Z)}$
- $dZ^{[L]} = \boldsymbol{A}^{[L]} - \boldsymbol{Y} = \left[a^{(1)} - y^{(1)} \dots a^{(m)} - y^{(m)}\right]; dZ^{[l]} = w^{[l+1]T} dZ^{[l+1]} * g^{[l]'}\left(Z^{[l]}\right)$
- $db^{[l]} = \frac{1}{n}\boldsymbol{np.sum\left(dZ^{[l]}, axis = 1, keepdims = True\right)}, dw^{[l]} = \frac{1}{n}\boldsymbol{dZ^{[l]}A^{[l-1]T}}$
- $w := \boldsymbol{w - \alpha dw};\ b := \boldsymbol{b - \alpha db}$

# NEURAL NETWORK

$x_1$

$x_2$

$x_3$

$\hat{y}$

$x$

$w$

$b$

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

$$a = \sigma(z)$$

$$\mathcal{L}(a, y)$$

$x_1$

$x_2$

$x_3$

$\hat{y}$

$x$

$W^{[1]}$

$b^{[1]}$

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

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

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

$W^{[2]}$

$b^{[2]}$

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

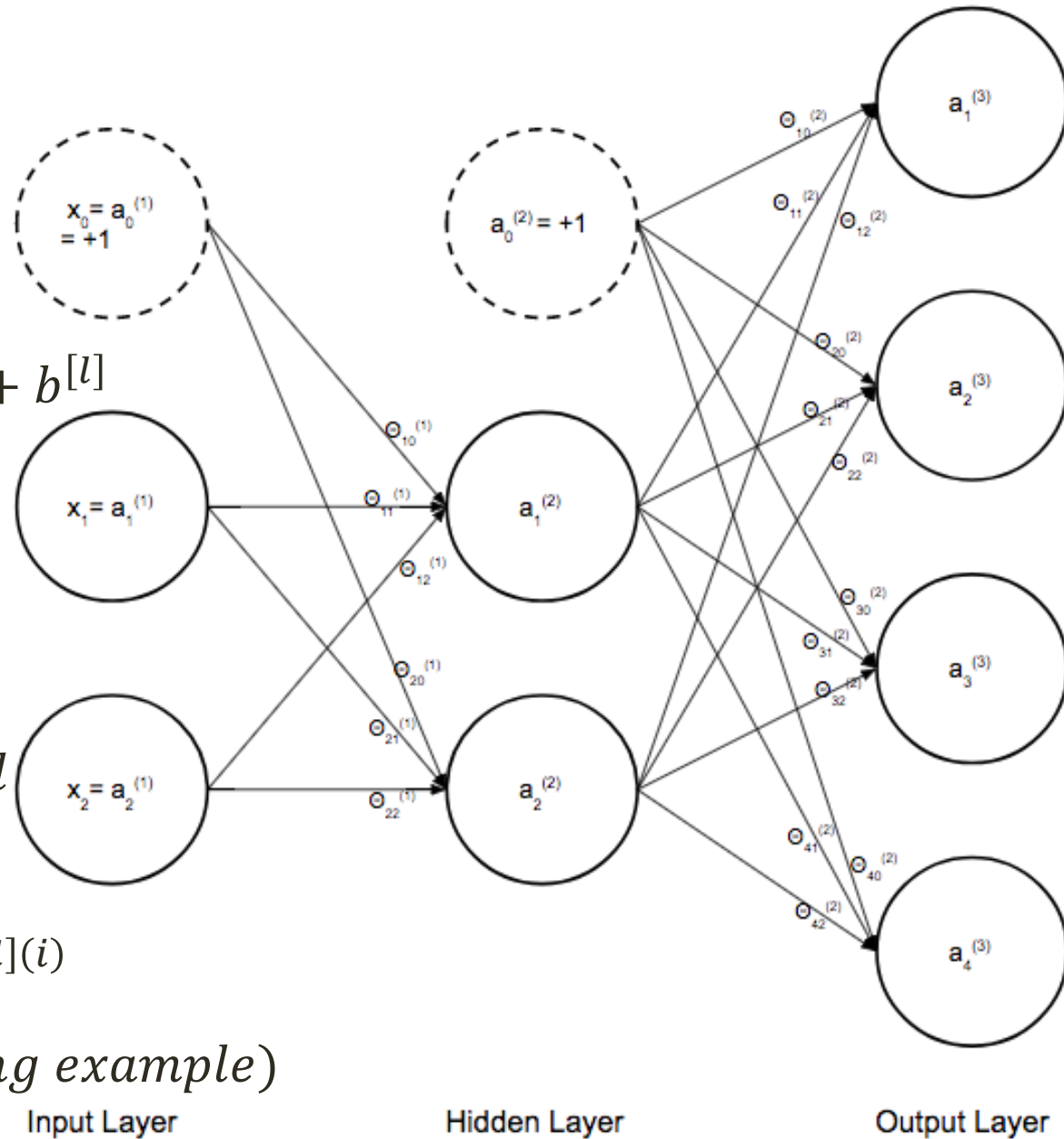$$\mathcal{L}(a^{[2]}, y)$$

# NEURAL NETWORK REPRESENTATION

$$z^{[l]} = \begin{bmatrix} \vdots \\ z_i^{[l]} \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots \\ w_i^{[l]T}X + b_i^{[l]} \\ \vdots \end{bmatrix} = W^{[l]}x + b^{[l]}$$

$$a^{[l]} = \begin{bmatrix} \vdots \\ a_i^{[l]} \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots \\ \sigma\left(z_i^{[l]}\right) \\ \vdots \end{bmatrix} = \sigma(z^{[l]})$$

$a_j^{[l](i)} = activation\ in\ node\ j\ in\ layer\ l$

$in\ example\ i$

➜ $for\ i = 1\ to\ m: Calculate\ z^{[l](i)}, a^{[l](i)}$

**PS:** $Z^{[l]}, A^{[l]} \in (\#hidden\ units, \#training\ example)$



Input Layer          Hidden Layer          Output Layer

# MORE ON VECTORIZED IMPLEMENTATION

$$z^{[1](1)} = w^{[1]} x^{(1)} + \cancel{b^{[1]}}^{\uparrow 0} \quad , \quad z^{[1](2)} = w^{[1]} x^{(2)} + \cancel{b^{[1]}}^{\uparrow 0} \quad , \quad z^{[1](3)} = w^{[1]} x^{(3)} + \cancel{b^{[1]}}^{\uparrow 0}$$

$$w^{[1]} = \begin{bmatrix} \rule[0.5ex]{2em}{0.4pt} \\ \rule[0.5ex]{2em}{0.4pt} \\ \rule[0.5ex]{2em}{0.4pt} \\ \rule[0.5ex]{2em}{0.4pt} \end{bmatrix} \qquad w^{[1]} x^{(1)} = \begin{bmatrix} \bullet \\ \vdots \\ \bullet \\ \vdots \\ \bullet \end{bmatrix} \qquad w^{[1]} x^{(2)} = \begin{bmatrix} \bullet \\ \vdots \\ \bullet \\ \vdots \\ \bullet \end{bmatrix} \qquad w^{[1]} x^{(3)} = \begin{bmatrix} \bullet \\ \vdots \\ \bullet \\ \vdots \\ \bullet \end{bmatrix}$$

$$z^{[1]} = w^{[1]} X + b^{[1]}$$

$$w^{[1]} \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} \cdots \\ | & | & | \end{bmatrix} = \begin{bmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{bmatrix} = \begin{bmatrix} | & | & | \\ z^{[1](1)} & z^{[1](2)} & z^{[1](3)} \cdots \\ | & | & | \end{bmatrix} = z^{[1]}$$

$$+ b^{[1]} \qquad + b^{[1]} \qquad + b^{[1]}$$

$$w^{[1]} x^{(1)} = z^{[1](1)}$$

# ACTIVATION FUNCTIONS

| Name | Plot | Equation | Derivative |
|------|------|----------|------------|
| Sigmoid |  | $f(x) = \sigma(x) = \dfrac{1}{1 + e^{-x}}$ | $f'(x) = f(x)(1 - f(x))$ |
| Tanh |  | $f(x) = \tanh(x) = \dfrac{(e^x - e^{-x})}{(e^x + e^{-x})}$ | $f'(x) = 1 - f(x)^2$ |
| Rectified Linear Unit (relu) |  | $f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$ | $f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$ |
| Leaky Rectified Linear Unit (Leaky relu) |  | $f(x) = \begin{cases} 0.01x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$ | $f'(x) = \begin{cases} 0.01 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$ |

# SUMMARY:
## FORWARD AND BACKWARD PROPAGATION

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$
$$A^{[1]} = g^{[1]}(Z^{[1]})$$
$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$
$$A^{[2]} = g^{[2]}(Z^{[2]})$$
$$\vdots$$
$$A^{[L]} = g^{[L]}(Z^{[L]}) = \hat{Y}$$

$$dZ^{[L]} = A^{[L]} - Y$$
$$dW^{[L]} = \frac{1}{m} dZ^{[L]} A^{[L]^T}$$
$$db^{[L]} = \frac{1}{m} np.\text{sum}(dZ^{[L]}, axis = 1, keepdims = True)$$
$$dZ^{[L-1]} = dW^{[L]^T} dZ^{[L]} g'^{[L]}(Z^{[L-1]})$$
$$\vdots$$
$$dZ^{[1]} = dW^{[L]^T} dZ^{[2]} g'^{[1]}(Z^{[1]})$$
$$dW^{[1]} = \frac{1}{m} dZ^{[1]} A^{[1]^T}$$
$$db^{[1]} = \frac{1}{m} np.\text{sum}(dZ^{[1]}, axis = 1, keepdims = True)$$

# WEIGHT INITIALIZATION AND DIMENSIONALITY
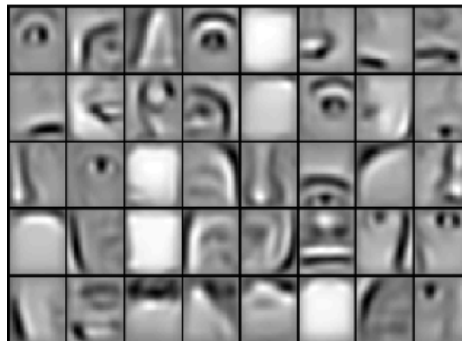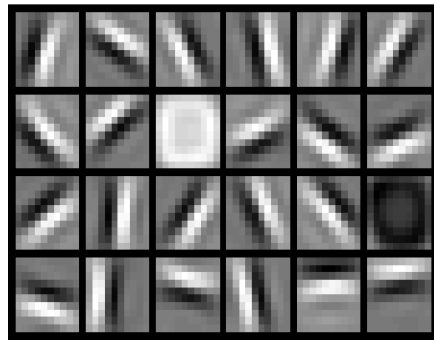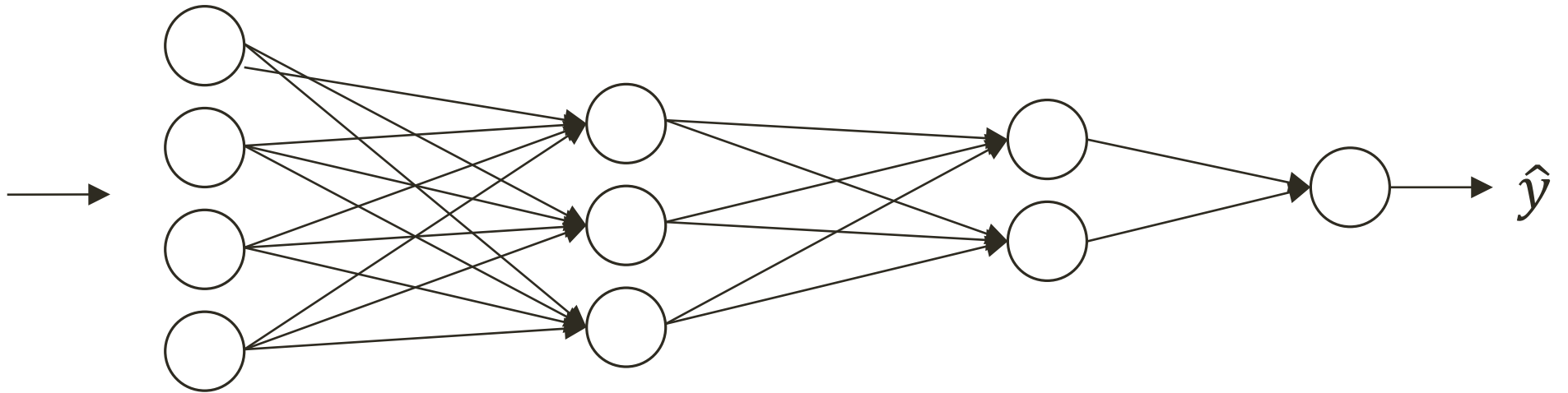
Zero initialization → Activations will be the same

Random initialization → **Good**

PS: b can be initialized to zero

<u>Notations</u>:
- $n^{[l]} = \#units\ in\ layer\ l$
- $W^{[l]}, dW^{[l]} \in \left(n^{[l]}, n^{[l-1]}\right)$
- $A^{[l]}, Z^{[l]}, b^{[l]}, db^{[l]} \in \left(n^{[l]}, 1\right)$
- $A^{[l]}, dA^{[l]}, Z^{[l]}, dZ^{[l]} \in \left(n^{[l]}, m\right)$

# DEEP REPRESENTATION

# HYPERPARAMETERS

**Parameters**: $W, b$

**Hyperparameters**: $Learning\ rate\ \alpha, \#iterations, \#hidden\ layer\ L,$

$\#hidden\ units, choice\ of\ activation\ functions$

**PS**: Applying deep learning is a very empirical process

Idea

Code

Experiment