1. Suppose that the network has $L$ layers.

   Suppose that layer $\ell$, for $\ell = 1, \ldots, L$, contains $n_\ell$ neurons, the network maps from $\mathbb{R}^{n_1} \to \mathbb{R}^{n_L}$

   Let $W^{[\ell]} \in \mathbb{R}^{n_\ell \times n_{\ell-1}}$ to denote the matrix of weights at layer $\ell$.

   Let $b^{[\ell]} \in \mathbb{R}^{n_\ell}$ is the vector of biases for layer $\ell$.

   Def (3.1) & (3.2) given that :

   (3.1): $a^{[1]} = x \in \mathbb{R}^{n_1}$

   (3.2): $a^{[\ell]} = \sigma(W^{[\ell]} a^{[\ell-1]} + b^{[\ell]}) \in \mathbb{R}^{n_\ell}$  for $\ell = 2, \ldots, L$

   (5.2) defined $z^{[\ell]} = W^{[\ell]} a^{[\ell-1]} + b^{[\ell]}$, so (3.2) can be written as $a^{[\ell]} = \sigma(z^{[\ell]})$, for $\ell = 2, \ldots, L$

   Also, we let $D^{[\ell]} \in \mathbb{R}^{n_\ell \times n_\ell}$ denote the diagonal matrix with $(i,i)$ entry given by $\sigma'(z_i^{[\ell]})$.

   i.e. $D^{[\ell]} = \mathrm{diag}(\sigma'(z^{[\ell]}))$  to avoid the Hadamard product.

   By Chain rule, we have $\nabla a^{[L]}(x) = \dfrac{\partial a^{[L]}}{\partial x} = \dfrac{\partial a^{[L]}}{\partial a^{[1]}} = \dfrac{\partial a^{[L]}}{\partial a^{[L-1]}} \cdot \dfrac{\partial a^{[L-1]}}{\partial a^{[L-2]}} \cdots \dfrac{\partial a^{[2]}}{\partial a^{[1]}}$ —— ①

   Then, $\dfrac{\partial a^{[\ell]}}{\partial a^{[\ell-1]}} = \dfrac{\partial \sigma(z^{[\ell]})}{\partial z^{[\ell]}} \cdot \dfrac{\partial z^{[\ell]}}{\partial a^{[\ell-1]}} = \dfrac{\partial \sigma(z^{[\ell]})}{\partial z^{[\ell]}} \cdot \dfrac{\partial (W^{[\ell]} a^{[\ell-1]} + b^{[\ell]})}{\partial a^{[\ell-1]}} = D^{[\ell]} W^{[\ell]}$  for $\ell = 2, \ldots, L$ —— ②

   By ① and ②, we have

   $$\nabla a^{[L]}(x) = \left(D^{[L]} W^{[L]}\right)_{n_L \times n_{L-1}} \cdot \left(D^{[L-1]} W^{[L-1]}\right)_{n_{L-1} \times n_{L-2}} \cdots \left(D^{[2]} W^{[2]}\right)_{n_2 \times n_1}$$

   Since $n_L = 1$, $\nabla a^{[L]}(x) \in \mathbb{R}^{1 \times n_1}$

   The pseudocode for this algorithm :

```
Input : network weights {W[ℓ], b[ℓ]}ℓ=2^L , σ(x) and σ'(x) , input x = a[1]
Output : grad_x = ∂a[L]/∂x
# Step 1 : Forward pass ( calculate and store z[ℓ] and a[ℓ] )
a[1] = x
for ℓ = 2 to L
    z[ℓ] = W[ℓ] @ a[ℓ-1] + b[ℓ]      # z[ℓ] ∈ ℝ^nℓ    # @ : 矩陣乘法
    a[ℓ] = σ(z[ℓ])                    # a[ℓ] ∈ ℝ^nℓ
    Store z[ℓ] and a[ℓ]
# Step 2 : Backward pass
D[L] = σ'(z[L])          # scalar , since nL = 1
G = D[L] @ W[L]          # G ∈ ℝ^{1×n_{L-1}}
for ℓ = L-1 down to 2
    D[ℓ] = diag(σ'(z[ℓ]))
    G = G @ (D[ℓ] @ W[ℓ])
# Result
grad_x = G
return grad_x      # grad_x is ∇a[L](x) ∈ ℝ^{1×n_1}
```

2. 我們課堂中所舉例的神經網路模型都是比較小的數字，像是4層、5層而已，那如果太多層會發生什麼事？有辦法找到"最佳層數"嗎？

2. 我們課堂中所舉例的神經網路模型都是比較小的數字，像是4層、5層而已，那如果太多層會發生什麼事？有辦法找到"最佳層數"嗎？