

1.

$a^{(1)} = x \in \mathbb{R}^n$ is a vector.

Let $z^{(1)} = W^{(1)} a^{(1-1)} + b^{(1)}$, which is linear part for a .

Then $a^{(1)} = \sigma(z^{(1)})$ is operating a activation function σ such as sigmoid, ReLU on $z^{(1)}$, which is non-linear part for a .

These operation are to perform linear transformation and non-linear activation on the input x layer by layer until the last layer L , and obtain the output $a^{(L)}$.

So, if we want to get $\nabla_x a^{(L)}$, we can use chain rule to get $\nabla_x a^{(L)} = \frac{\partial a^{(L)}}{\partial a^{(L-1)}} \cdot \frac{\partial a^{(L-1)}}{\partial a^{(L-2)}} \cdot \dots \cdot \frac{\partial a^{(1)}}{\partial x}$.

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

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

Then, $\frac{\partial a^{(1)}}{\partial a^{(1-1)}} = \frac{d\sigma(z^{(1)})}{dz^{(1)}} \cdot W^{(1)}$

Define $\delta^{(1)} := \frac{\partial a^{(L)}}{\partial a^{(1)}} \in \mathbb{R}^n$

From the output layer:

$$\delta^{(L)} = \sigma'(z^{(L)}) = (W^{(L)})^T \delta^{(L)}$$

$$\delta^{(l)} = (W^{(l+1)})^T \delta^{(l+1)} \circ \sigma'(z^{(l)}) \text{ for } l = L-1, L-2, \dots, 2, \text{ where } \circ \text{ is Hadamard product.}$$

Then we will get $\nabla_x a^{(L)} = (W^{(2)})^T \delta^{(2)}$ (We don't need to multiply $\sigma'(z^{(1)})$ since the first layer has no activation)

Algorithm:

$n_L = 1 \rightarrow \delta^{(L)}$ is a scalar

delta = sigma_prime($z^{(L)}$)

for l in $(L-1)$ to 2 :

delta = ($W^{(l+1)}$)^T * delta * sigma_prime($z^{(l)}$)

grad_x = $W^{(2)}$ ^T * delta

return grad_x