

Auto differentiation

$$f(a) = 4a^2 - 3$$

$$f'(a) = 8a$$

$$f'(1.8) = 14.4$$

slope of $f(a)$
at $a = 1.8$

$$f(1.8) = 4 \cdot (1.8)^2 - 3 = 9.96$$

What should we do to
'a' to reduce $f(a)$?

$$f(1.7) = 4 \cdot (1.7)^2 - 3 = 8.96$$

we have a specific value

$$f(a, k) = k \cdot a + 17$$

and we want to pick
new value for 'a' such
that $f(a, k)$ is reduced

$$\frac{\partial f(a, k)}{\partial a} = k$$

2-Layer Network

$$\text{Layer 1} \begin{cases} z^{(1)} = A^{(0)} W^{(1)T} + b^{(1)} \\ A^{(1)} = \sigma(z^{(1)}) \end{cases}$$

↑ sigmoid activation

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

$$A^{(2)} = \sigma(z^{(2)})$$

↑ We don't have to use the same activation function for each layer

$$A^{(2)} = \sigma \left(\underbrace{A^{(1)} W^{(2)T}} + b^{(2)} \right)$$

$$= \sigma \left(\sigma(z^{(1)}) W^{(2)T} + b^{(2)} \right)$$

$$\sigma \left(\left(\sigma \left(\sigma \left(A^{(0)} \underbrace{W^{(1)T}} + b^{(1)} \right) W^{(2)T} + b^{(2)} \right) \right) W^{(3)T} + b^{(3)} \right)$$

$$\mathcal{L}(\hat{y}, y) = \frac{1}{2} \| \hat{y} - y \|^2$$

(1) How does $\frac{\partial \mathcal{L}}{\partial W^{(1)}}$ change if we add hidden layers?

(2) How does $\frac{\partial \mathcal{L}}{\partial W^{(1)}}$ change if we change an activation function?

(3) How does $\frac{\partial \mathcal{L}}{\partial W^{(2)}}$ change if we change the loss function?

Auto diff / Auto grad

