

Lecture 12 Backprop & Improving Neural Networks

1. Backpropagation

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

$$\text{with } L^{(i)} = -[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})]$$

$$\text{update: } w^{[l]} = w^{[l]} - \alpha \frac{\partial J}{\partial w^{[l]}}$$

- $\frac{\partial J}{\partial w^{[l]}} \times$, instead: $\frac{\partial L}{\partial w^{[l]}} = -[y^{(i)} \frac{\partial L}{\partial w^{[l]}} \cdot \log \sigma(w^{[l]} a^{[l]} + b^{[l]}) + (1-y^{(i)}) \frac{\partial L}{\partial w^{[l]}} \log (1-\sigma(\cdot))]$

$$= -y^{(i)} \frac{1}{\sigma(a^{[l]})} a^{[l]} \cdot (1-a^{[l]}) \cdot a^{[l]T} + (1-y^{(i)}) \frac{1}{1-\sigma(a^{[l]})} a^{[l]} (1-a^{[l]}) \cdot a^{[l]T}$$

$$\therefore \frac{\partial}{\partial w^{[l]}} (w^{[l]} a^{[l]} + b^{[l]})^2$$

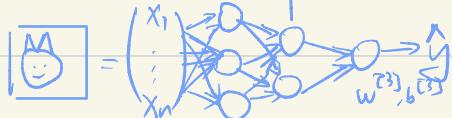
$$= -y^{(i)} (1-a^{[l]}) \cdot a^{[l]T} + (1-y^{(i)}) a^{[l]T}$$

$$= -y^{(i)} \cdot a^{[l]T} + y^{(i)} a^{[l]T} + a^{[l]T} - y^{(i)} a^{[l]T}$$

$$= -(y^{(i)} - a^{[l]T}) a^{[l]T}$$

for vectorization

The cat example:



$$\therefore \frac{\partial}{\partial w^{[l]}} (w^{[l]} a^{[l]} + b^{[l]})^2$$

$$\begin{matrix} (1,2) \\ (1,1) \end{matrix} \quad \begin{matrix} (2,1) \\ (2,2) \end{matrix}$$

$$\Rightarrow \frac{\partial J}{\partial w^{(3)}} = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} - a^{(3)}) a^{(2)T}$$

$$\cdot \frac{\partial L}{\partial w^{(3)}} = \frac{\partial L}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a^{(2)}} \cdot \frac{\partial a^{(2)}}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial w^{(2)}}$$

$$\begin{aligned} \frac{\partial L}{\partial w^{(3)}} &= (a^{(3)} - y) \cdot a^{(2)T} \\ a^{(2)T} &= \frac{\partial z^{(2)}}{\partial w^{(2)}} \end{aligned}$$

$$\therefore \frac{\partial L}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(3)}} = \frac{\partial L}{\partial z^{(3)}} = \frac{\partial L}{\partial w^{(3)}} \div \frac{\partial z^{(3)}}{\partial w^{(3)}}$$

$$\Rightarrow \frac{\partial L}{\partial w^{(3)}} = (a^{(3)} - y) \cdot w^{(3)} \cdot a^{(2)} \underbrace{(1 - a^{(2)})}_{(2,1)} \cdot a^{(1)T}$$

element wise product:

$$w^{(3)T} \boxed{*} a^{(2)} \underbrace{(1 - a^{(2)})}_{(2,1)} \cdot \underbrace{(a^{(3)} - y) a^{(1)T}}_{(1,3)}$$

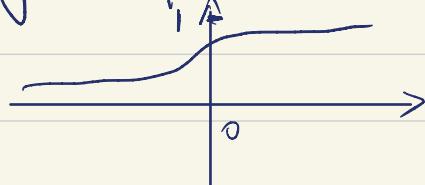
when same matrix size,

$$\begin{bmatrix} 1 \\ 2 \end{bmatrix} * \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 \\ 4 \end{bmatrix}$$

2. Improve NNs

1) Activation functions

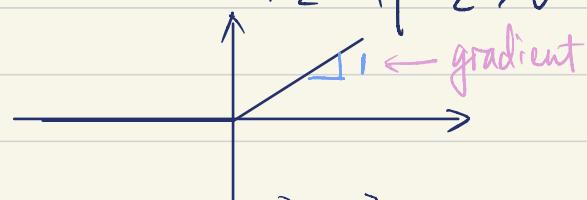
D Sigmoid function



$$\sigma(z) = \frac{1}{1+e^{-z}}$$

$$\sigma'(z) = \sigma(z) \cdot (1 - \sigma(z))$$

$$② \text{ReLU}(z) = \begin{cases} 0 & \text{if } z \leq 0 \\ z & \text{if } z > 0 \end{cases}$$



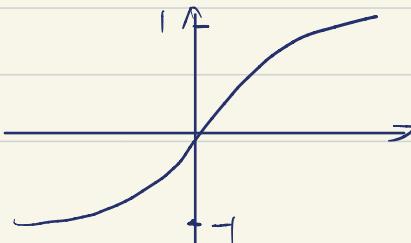
Most used.

The other two has low gradients when z is $\ll 0$ or $\gg 0$

$$\text{ReLU}'(z) = \begin{cases} z & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \end{cases}$$

$$③ \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\tanh'(z) = 1 - \tanh^2(z)$$



- Why do we need activation function?

example: $\begin{matrix} 0 & 0 \\ 0 & 0 \\ 0 \end{matrix}$ activations: Id. $z \mapsto z$

$$\begin{aligned}\hat{y} &= a^{[3]} = z^{[3]} = w^{[3]} a^{[2]} + b^{[3]} \\ &= w^{[3]} (w^{[2]} a^{[1]} + b^{[2]}) + b^{[3]} = w^{[3]} w^{[2]} (w^{[1]} x + b^{[1]}) + w^{[3]} b^{[2]} + b^{[3]} \\ &= Wx + B\end{aligned}$$

$$\text{with: } W = w^{[3]} w^{[2]} w^{[1]}$$

$$B = w^{[3]}. w^{[2]}. b^{[1]} + w^{[3]} b^{[2]} + b^{[3]}$$

\Rightarrow without activation function \rightarrow linear combination

2) Initialization Method

- Normalizing your input

$$\begin{aligned}M &= \frac{1}{m} \sum_{i=1}^m x^{(i)} \quad \Rightarrow \quad x = \frac{x - M}{S} \\ S^2 &= \frac{1}{m} \sum_{i=1}^m (x^{(i)})^2\end{aligned}$$

Before:



After:

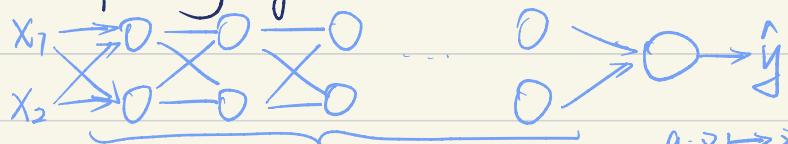


directly towards the middle

* The M and σ are computed on the training data. And these M and σ have to be used on the test set as well.

- Vanishing / Exploding gradients

example:



$$g: z \mapsto z, b=0$$

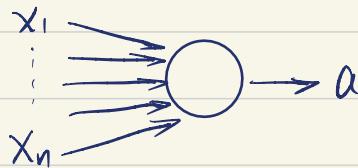
10 layers

$$\hat{y} = w^{[L]} \cdot a^{[L-1]} = w^{[L]} \cdot w^{[L-1]} \cdot a^{[L-2]} = w^{[L]} \cdot \underbrace{w^{[L-1]} \cdots w^{[1]}}_{10 \text{ layers}} \cdot x$$

$$w^{[1]} = \begin{pmatrix} 1.5 & 0 \\ 0 & 1.5 \end{pmatrix} \quad \xrightarrow{\text{very large}} \quad \begin{pmatrix} 2 & 2 \\ 1.5 & 0 \\ 0 & 1.5 \end{pmatrix}$$

very large

example with one neuron:



$$a = g(z)$$

$$z = w_1 x_1 + \dots + w_n x_n$$

large $n \rightarrow$ small $w_i \Rightarrow w_i \sim \frac{1}{n}$

for sigmoid: $w^{(l)} = \text{np.random.randn(shape)} * \text{np.sqrt}(\frac{1}{n^{(l-1)}})$

for ReLU: $w^{(l)} = \text{np.random.randn(shape)} * \text{np.sqrt}(\frac{2}{n^{(l-1)}})$

for tanh: $w^{(l)} \sim \sqrt{\frac{1}{n^{(l-1)}}}$ (Xavier Initialization)

He Initialization: $w^{(l)} \sim \sqrt{\frac{2}{n^{(l)} + n^{(l-1)}}}$ (often used)

3) Optimization

$$X = (x^{(1)}, x^{(2)}, \dots, x^{(m)})$$

$$Y = (y^{(1)}, y^{(2)}, \dots, y^{(m)})$$

$$X = (x^{\{1\}}, x^{\{2\}}, \dots, x^{\{T\}})$$

$$x^{(1)} \dots x^{(1000)}$$

$$Y = (y^{\{1\}}, \dots, y^{\{m\}})$$

Algorithm for mini-batch gradient descent:

For iteration $t = 1, \dots$:

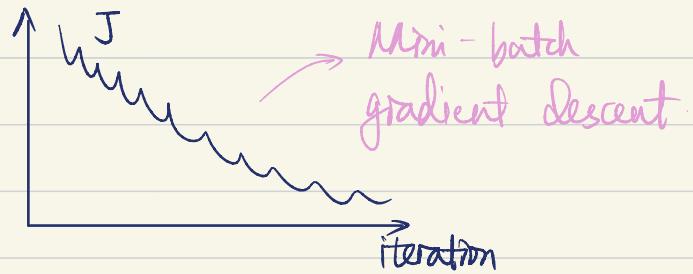
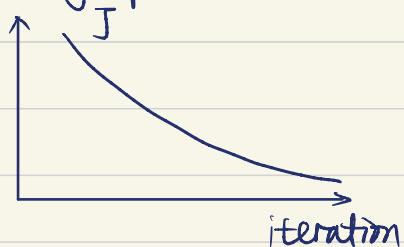
Select batch (x^{t+3}, y^{t+3})

Do Forward propagation $\Rightarrow J = \frac{1}{1000} \sum_{i=1}^{1000} L^{(i)}$

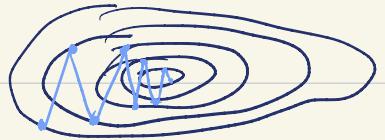
Backprop Batch

Update $w^{(l)}, b^{(l)}$

Show in graph:



- Gradient Descent + Momentum algorithm.



↑ move
smaller

← move larger

weight

look at
past updates

velocity : tracks the direction that we should take regarding the current update and also the past update

$$\begin{cases} V = BV + (1-B) \frac{\partial L}{\partial W} \\ W = W - \alpha V \end{cases}$$

the avg. of vertical is smaller (different direction)
the avg. of horizontal is larger (same direction)