

Vanishing or exploding gradients

=> very deep networks with L layers

$$\bar{y} = w^L w^{L-1} \dots w^2 w^1 x$$

$$w > 1$$

$$w < 1$$

$$w = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$$

$$w = \begin{bmatrix} 0.8 & 0 \\ 0 & 0.8 \end{bmatrix}$$

\bar{y} will be very large number

\bar{y} will be very small number

Hard to train

Weight initialization

n : # of input features

w : weight matrices

$$\text{variance}(w) = \frac{1}{n} \leftarrow \text{can be hyperparameter}$$

Relu activation

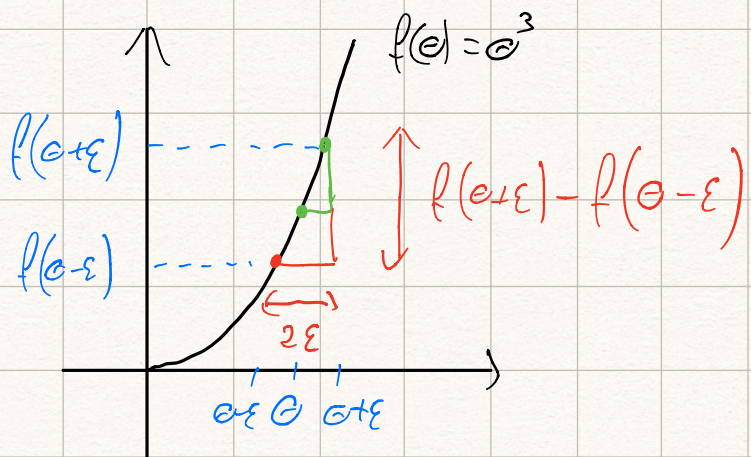
$$w^l = n.p. \text{ random} \cdot \text{randn}(\text{shape}(w)) * n.p. \cdot \text{sqr}t\left(\frac{2}{n^{l-1}}\right)$$

use 2 with relu
 n input features from layer $l-1$

Tanh activation

$$\text{sqr}t\left(\frac{1}{n^{l-1}}\right) \text{ "xavier initialization"}$$

Numerical approximation of gradients



$$f'(a) = \lim_{\epsilon \rightarrow 0} \frac{f(a+\epsilon) - f(a-\epsilon)}{2\epsilon}$$

$$\epsilon := 0,01 \Rightarrow \frac{1,01^3 - 0,99^3}{2 \cdot 0,01} \approx 3,0001$$

$$a := 1 \Rightarrow$$

$$f'(a) = 3 \cdot a^2 = 3$$

approximation error

very close

with one side approximation:

$$\frac{f(a+\epsilon) - f(a)}{\epsilon} \Rightarrow \frac{1,01^3 - 1^3}{0,01} \approx \frac{0,0303}{0,01} = 3,03$$

approximation error

\Rightarrow using two side calculation

\Rightarrow more accurate approximation

gradient checking

Reshape $w^1, b^1, w^2, b^2 \dots w^L, b^L$ into a big vector θ

• $J(w^1, b^1 \dots w^L, b^L) = J(\theta)$ concatenate

Reshape $dw^1, db^1 \dots dw^L, db^L$ into a big vector $d\theta$

Is $d\theta$ the gradient of J ?



for each i :

$$d\theta_{\text{approx}}[i] = \frac{J(\theta_1 \dots \theta_{i+\epsilon} \dots \theta_n) - J(\theta_1 \dots \theta_{i-\epsilon} \dots \theta_n)}{2\epsilon}$$

$$\approx d\theta[i] = \frac{\partial J}{\partial \theta_i}$$

Is $d\theta_{\text{approx}} \approx d\theta$?

if $\epsilon = 10^{-7}$

check : $\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2}$

10^{-7} great
 \vdots
 10^{-5} maybe
 \vdots
 10^{-3} worry

Don't use gradient checking in training only to debug
doesn't work with dropout