## numerical approximation of gradients

$f(\theta) = \theta^3$



$\varepsilon = 0.01$

approximation of gradient $\approx \dfrac{f(\theta+\varepsilon) - f(\theta-\varepsilon)}{2\varepsilon}$

$f'(\theta) = \lim_{\theta \to 0} \dfrac{f(\theta+\varepsilon) - f(\theta-\varepsilon)}{2\varepsilon}$

the error of this approximation is $O(\varepsilon^2)$

$\dfrac{f(\theta+\varepsilon) - f(\theta)}{\varepsilon} \longrightarrow$ error is $O(\varepsilon)$

---

## gradient checking of a NN

① take $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots$
and reshape into a big vector $\theta$

$w^{[1]}, w^{[2]}, \dots, w^{[L]}$ ; reshape them into vectors
concatenate them with $b^{[1]}, b^{[2]}, \dots, b^{[L]}$

② take $dw^{[1]}, db^{[1]}, \dots$
reshape them into a big vector $d\theta$

③ $J(\theta) = J(\theta_1, \theta_2, \theta_3, \dots \theta_i \dots)$ ; $eps = 10^{-7}$

for each $i$

$d\theta_{approx}[i] = \dfrac{J(\theta_1, \theta_2, \dots, \theta_i+\varepsilon, \dots) - J(\theta_1, \theta_i-\varepsilon}{2\varepsilon}$

check $\dfrac{\|d\theta_{approx} - d\theta\|_2}{\|d\theta_{approx}\|_2 + \|d\theta\|_2}$ $\approx 10^{-7}$ great

---

practical tips

1/ don't use gradient checking in training
use it only to debug

2/ if algorithm fails grad check, look at components
to identify the bug

3/ remember to add the derivative of the regularization
term

4/ doesn't work with dropout

## normalizing training sets

mean $= \mu = \frac{1}{m} \sum_{i=1}^{m} x^{[i]}$

variance $= \sigma^2 = \frac{1}{m} \sum_{i=1}^{m} x^{[i]} ** 2$

$X := \frac{X - \mu}{\sigma}$

NB: use same $\mu$ and $\sigma$ to normalize the test set

## why normalization ?

we will need less # iterations to converge to the minimum

$\Rightarrow$ it will easier and faster to train

⚠️ if the features are in the same scale, no need to do this step

## Vanishing gradients / exploding gradients

suppose that: $g^{[\ell]}(z^{[\ell]}) = z^{[\ell]}$

$b^{[\ell]} = 0$      for all $\ell$

$\Rightarrow \hat{y} = w^{[L]} \cdots \cdots w^{[1]} X$

if we suppose $w^{[\ell]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$ for $\ell$ in $1 \ldots 2 -$

$\Rightarrow \hat{y} = w^{[L]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} X$

$\simeq 1.5^{L-1} X$

if we have a very deep NN :

$\hat{y}$ will explode

if we suppose $w^{[\ell]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$

$\hat{y} = w^{[L]} \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}^{L-1} X$

$\simeq 0.5^{L-1} X$

if we have a very deep NN

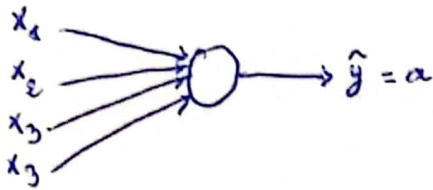activations will decrease exponentially

Conclusion

if weights $> 1$ : we will increase exponentially

if weights $< 1$ : we will decrease exponentially

solution for vanishing/exploding

better choice for parameter initialization

### 1/ single neuron



suppose $b = 0$

$$z = \sum_{i=1}^{n} w_i x_i$$

we want $z$ not too large not too small

if $n$ is large : we want small $w_i$

$\Rightarrow$ set $\text{Var}(w_i) = \dfrac{1}{n}$ or $\dfrac{2}{n}$ if we are using a relu activation

$w^{[\ell]} = np.random.randn(\text{Shape}) \times np.sqrt\left(\dfrac{1}{n^{[\ell-1]}}\right)$

nR : activation : $\text{var}(w_i) = \sqrt{\dfrac{2}{n^{[\ell-1]}}}$ ; He initialization

or

$\text{var}(w_i) = \sqrt{\dfrac{2}{n^{[\ell-1]} + n^{[\ell]}}}$

This variance can be another hyperparameter to fine tune

① Random initialization is used to break symmetry and make sure different hidden units can learn different things

② He initialization works well with ReLU activation