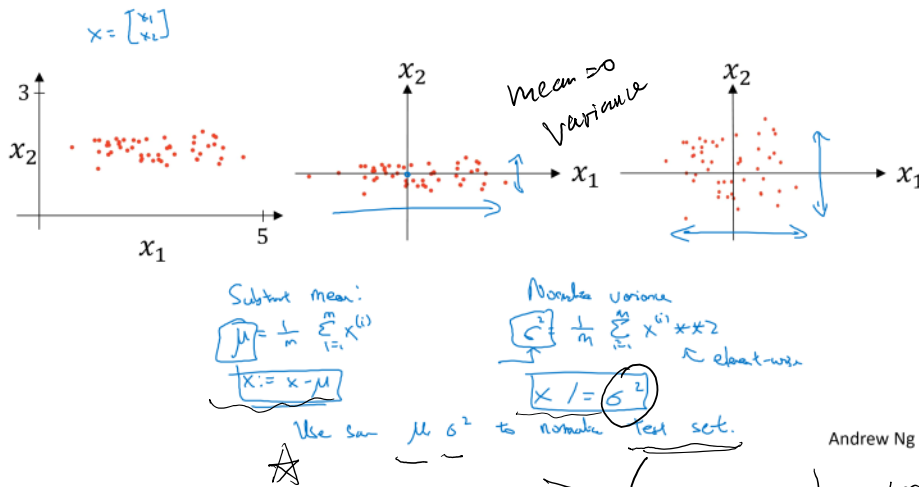


deeplearning.ai

## Setting up your optimization problem

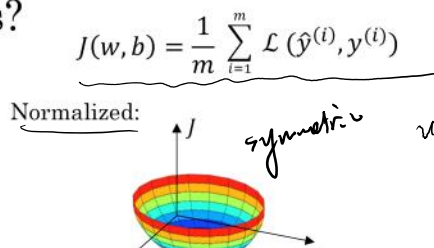
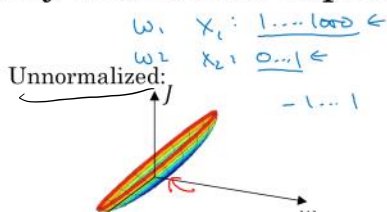
### Normalizing inputs

#### Normalizing training sets

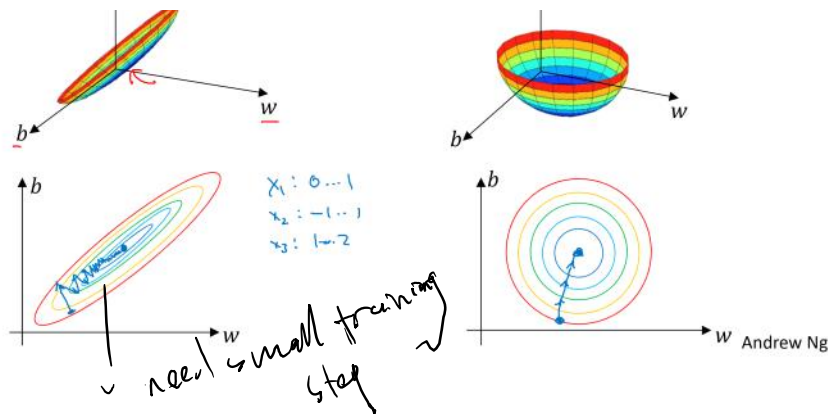


Don't normalize training & test set differently  
want same transformation

#### Why normalize inputs?



unnormalize + feature of different scales  
 $w_1, w_2$  different value  
elongated bowl



$w_1, w_2$  different value  
 elongated bowl <sup>2</sup>/<sub>sig</sub>.

In practice,  $w$  higher dim  
 this two-dim plot shows the  
 intuition

if feature similar scale,  $w$  not  
 be important but this won't  
 do any harm

vanish exploding gradient



deeplearning.ai

Setting up your  
 optimization problem

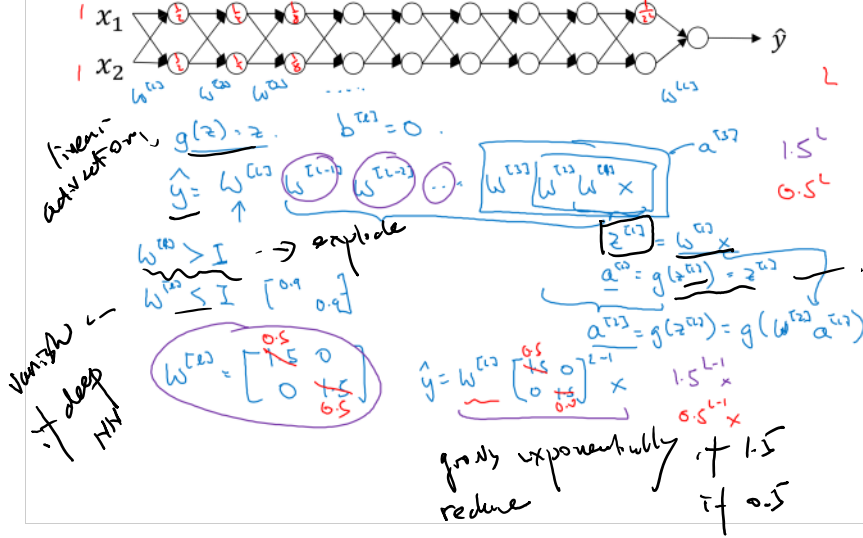
Vanishing/exploding  
 gradients

one big problem

# Vanishing/exploding gradients

example  
L=150

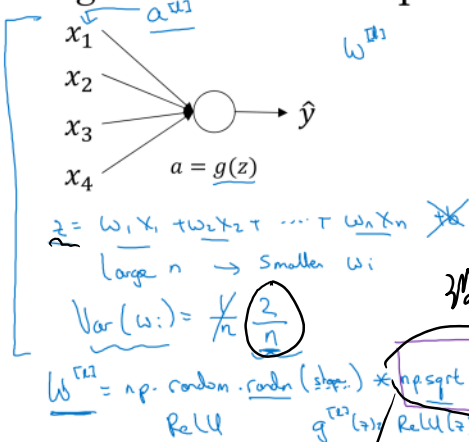
→ vanishing  
exploding gradient



Andrew Ng

Solution

## Single neuron example



variant

Other variants:

tanh

Xavier initialization

Xavier Initialization

$\frac{2}{n^{(l-1)} + n^{(l)}}$

2nd RN is 0

feature, variance  $\rightarrow 1$

can add hyper-param for this

Reduce the vanishing/exploding gradient

not too much  $> 1$ .

not too much  $< 1$

Andrew Ng

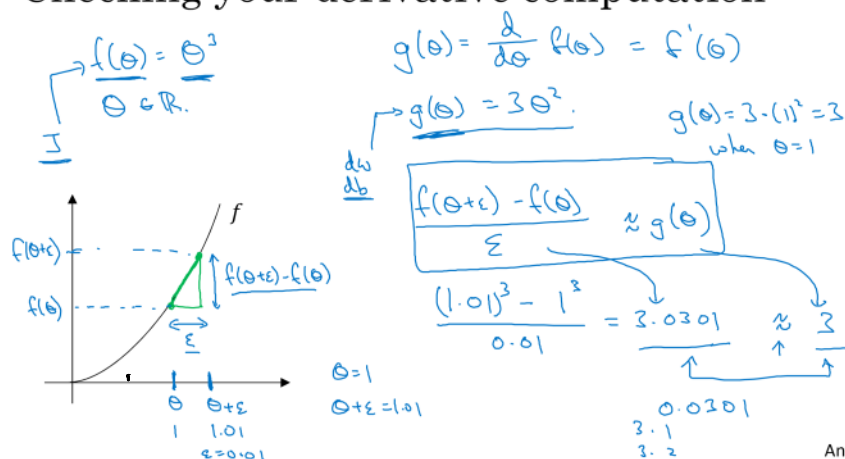


deeplearning.ai

## Setting up your optimization problem

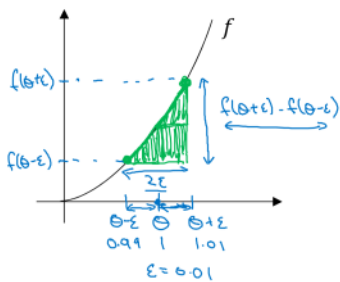
### Numerical approximation of gradients

#### Checking your derivative computation



## Checking your derivative computation

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



$$\left[ \frac{f(\theta+\epsilon) - f(\theta-\epsilon)}{2\epsilon} \approx g(\theta) \right]$$

$$\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.0001 \approx 3$$

$$g(\theta) = 3\theta^2 = 3$$

approx error: 0.0001  
(prev slide: 3.0301, error: 0.03)

$$f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta+\epsilon) - f(\theta-\epsilon)}{2\epsilon} \quad \begin{matrix} \mathcal{O}(\epsilon^2) \\ 0.01 \\ 0.0001 \end{matrix} \quad \left| \quad \frac{f(\theta+\epsilon) - f(\theta)}{\epsilon} \quad \begin{matrix} \text{error: } \mathcal{O}(\epsilon) \\ 0.01 \end{matrix} \right.$$

Andrew Ng

We would rather use the two-sided difference

numerical approximation

Just calculus.

Two sided.

for gradient checking

gradient check



deeplearning.ai

## Setting up your optimization problem

## Gradient Checking

$J$  is the loss of  $J(\theta)$ :

Take  $qW_{[1]}, qP_{[1]}, \dots, qW_{[L]}, qP_{[L]}$  and reshape into a big vector  $q\theta$ .

$$J(m_{[1]}, p_{[1]}, \dots, m_{[L]}, p_{[L]}) = J(\theta)$$

Take  $W_{[1]}, P_{[1]}, \dots, W_{[L]}, P_{[L]}$  and reshape into a big vector  $\theta$ .

## Gradient check for a neural network

### Gradient checking (Grad check)

$$J(\theta) = J(\theta_1, \theta_2, \dots)$$

for each  $i$ :

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

Where to get  $\epsilon$ ? still the prop thing?

$$\approx \frac{\partial J}{\partial \theta_i}$$

$$d\theta_{approx} \approx d\theta$$

Check

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

$\epsilon = 10^{-7}$

$$\approx \frac{10^{-7}}{10^{-5}} \leftarrow \text{great!}$$

$$\rightarrow 10^{-3} \leftarrow \text{worry!}$$

normalize it into a ratio.

Andrew Ng

why  $\epsilon$ ?

How to define whether two vectors are close? Euclidean distance

$\rightarrow$  if big value need debug

## Gradient checking implementation notes

- Don't use in training – only to debug → slow computation
- If algorithm fails grad check, look at components to try to identify bug.

- Remember regularization.

- Doesn't work with dropout.

- Run at random initialization; perhaps again after some training.

$$\frac{\partial \mathcal{O}_{\text{train}}[i]}{\partial \theta} \leftrightarrow \frac{\partial \mathcal{O}[i]}{\partial \theta}$$

$$J(\theta) = \frac{1}{n} \sum_i \ell(y^{(i)}, \theta) + \frac{\lambda}{2n} \sum_i \|\theta^{(i)}\|_2^2$$

$\partial \theta = \text{grad of } J \text{ wrt. } \theta$

$J$  keep prob = 1.0  
 → check before turning dropout  
 run grad check after some iterations

梯度检查只适用于调试  
 不适用于训练。  
 正则化项  
 正则化函数

Andrew Ng