

# Week 1c - Setting up your optimization problem

Friday, August 14, 2020 11:25 PM

## ① Normalizing inputs

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

i. Subtract mean to set mean to zero  
center data

$$\mu = \frac{1}{m} \sum_{i=1}^n x^{(i)}$$

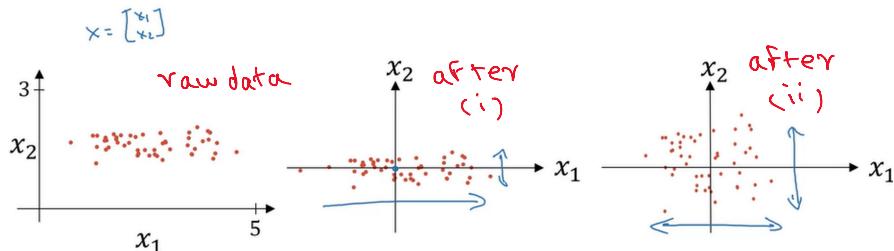
$$x := x - \mu$$

ii. Normalize variance

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^n x^{(i)} * 2$$

$$x := \sigma$$

iii. Use same  $\mu$  &  $\sigma$  to normalize test set



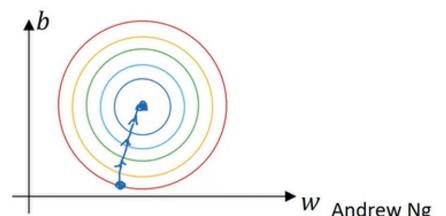
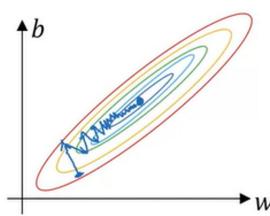
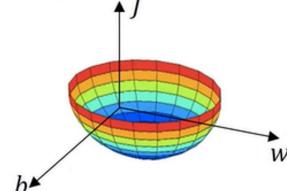
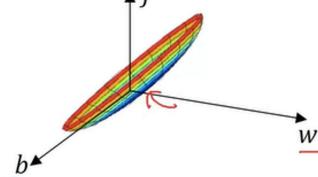
→ why normalize inputs

Why normalize inputs?

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Unnormalized:  $w_1 x_1: 1 \dots 1000$   $w_2 x_2: 0 \dots$   $\rightarrow$  unmatched scale hurts learning rate

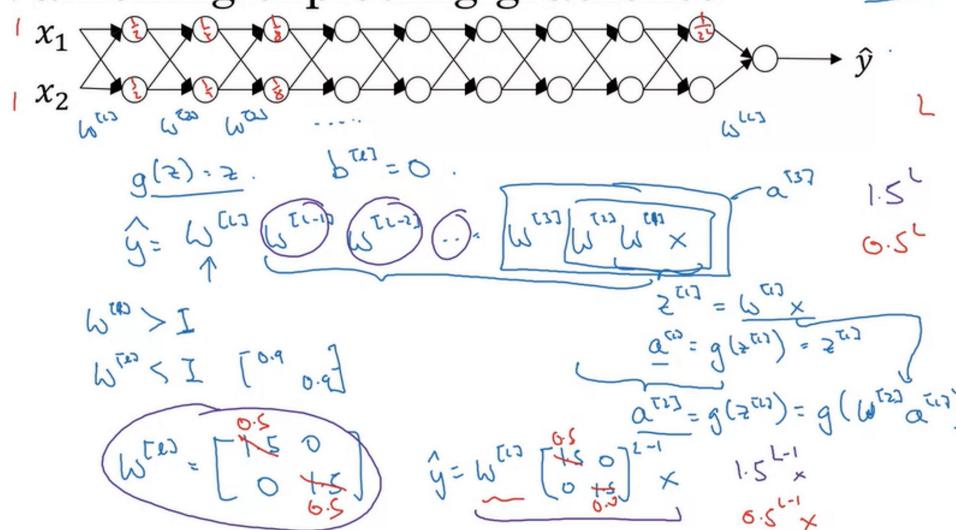
Normalized:



## ② Vanishing / Exploding gradients

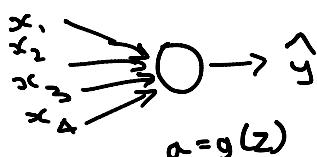
- When training deep NNs, gradients might become too small or very large
- If  $W^{[L]} > I$  (Identity matrix),  
the final value  $\hat{y}$  is in the proportion of  $W^{[L-1]} \dots W^{[1]}$  [exponentially high]
- If  $W^{[L]} < I$   
the final value  $\hat{y}$  is very low

### Vanishing/exploding gradients



Andrew Ng

### ③ Weight Initialization for Deep Neural Networks



$$z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

large  $n \rightarrow$  small  $w_i$

→ we could set  $\text{var}(w_i) = 1/n$

$$w^{[L]} = np.random.rand(\text{shape}) * np.sqrt\left(\frac{2}{n^{L-1}}\right)$$

$\frac{2}{n}$  works better for ReLU

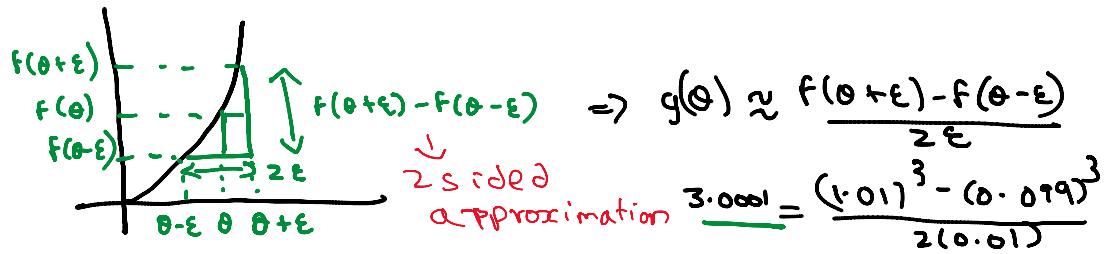
→ other variants

- $\tanh \rightarrow \sqrt{\frac{1}{n^{L-1}}} w^{[L]}$  works better (Xavier initialization)
- $\sqrt{\frac{2}{n^{L-1} + n^{L-2}}}$  is also used

### ④ Numerical approximation of gradients

→ Checking your derivative computation

$$\text{If } F(\theta) = \theta^3$$



1-sided approx

$$g(\theta) = 3\theta^2 = 3 \\ \therefore \text{error} = 0.0001$$

$$\Rightarrow g(\theta) \approx \frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \rightarrow \boxed{\text{order of error is } \epsilon}$$

$$\Rightarrow \underline{3.0301} \approx \frac{(1.01)^3 - (1)^3}{0.01} \Rightarrow \text{error} = 0.0301$$

• 2-sided approximation runs slower but is worth the improvement in accuracy

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

↓ order of error is in order of  $\epsilon^2$

## ⑤ Gradient Checking

i. Take  $w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}$  and reshape it into  $\theta$   
 $\rightarrow J(w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}) = J(\theta)$

ii. Take  $dW^{(1)}, db^{(1)}, \dots, dW^{(L)}, db^{(L)}$  and reshape it into  $d\theta$   
 Is  $d\theta$  gradient of  $\theta$ ?

$\rightarrow$  Gradient Checking (grad check)

- For each  $i$ :

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

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

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

- check if  $\delta\theta_{\text{approx}} \approx \delta\theta$

$$\Rightarrow \frac{\|\delta\theta_{\text{approx}} - \delta\theta\|_2}{\|\delta\theta_{\text{approx}}\|_2 + \|\delta\theta\|_2}$$

$\epsilon = 10^{-7}$  is generally used

↳ if O/P is  $10^{-7} \rightarrow$  great  
 $10^{-5} \rightarrow$  double check  
 $10^{-3} \rightarrow$  worry

## ⑥ Gradient checking Implementation Notes

- Use grad check only in testing / debug as it is very slow to be used in training
- If algorithm fails grad check, look at components to try
- Remember regularization
- Doesn't work with dropout  
Use  $\text{keep-prob} = 1.0$  then do grad-check & then set keep-prob to desired value
- Run at random initialization ; perhaps again after some training