**Week 1-3 Problem Set-up (Normalization, Initialization, G checking)**

## Normalizing training sets

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

Subtract mean:
$$\mu = \frac{1}{m} \sum_{i=1}^{m} x^{(i)}$$

$$x := x - \mu$$

Normalize variance:
$$\sigma^2 = \frac{1}{m} \sum_{i=1}^{m} x^{(i)} **2 \quad \leftarrow \text{element-wise}$$

$$x /= \sigma^2$$

Use same $\mu$ $\sigma^2$ to normalize test set.
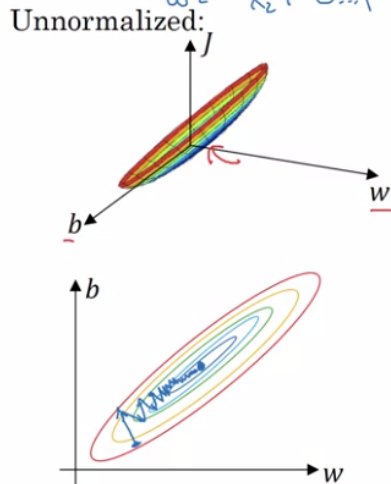
Andrew Ng

also do with test set, (divide by sigma rather than sigma^2)

Why normalize?

# Why normalize inputs?

$w_1, \quad x_i : 1 \cdots 1000$

$w_2 \quad x_2 : 0 \cdots 1$

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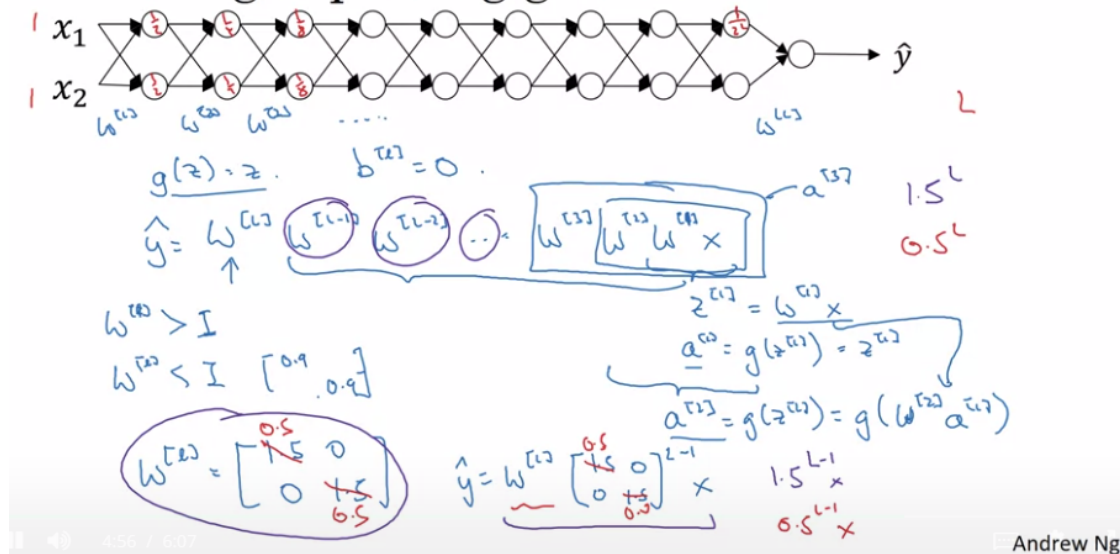
Unnormalized:

Normalized:

if you're running gradient descent on the cost function like the one on the left, then you might have to use a very **small learning rate** because if you're here that gradient descent might need a lot of steps to oscillate back and forth before it finally finds its way to the minimum. Whereas if you have a more spherical contours, then wherever you start gradient descent can pretty much go straight to the minimum. You can take much **larger steps** with gradient descent rather than needing to oscillate

around like like the picture on the left.
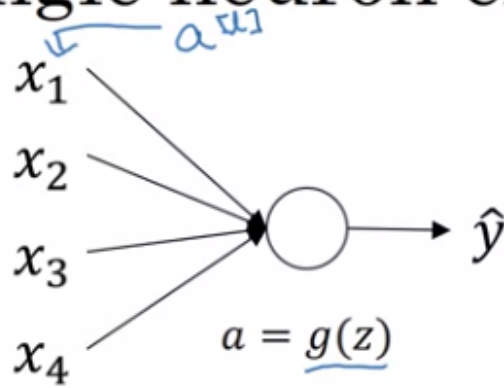
---

# Vanishing / Exploding Gradients



(dW will also have such extreme problems)

# Weight Initialization for NN

# Single neuron example

$x_1$

$x_2$     $a^{[l]}$

$x_3$     $\hat{y}$

$x_4$     $a = g(z)$

$$z = w_1 x_1 + w_2 x_2 + \cdots + w_n x_n \cancel{+ b}$$

$$\text{large } n \rightarrow \text{Smaller } w_i$$

$$\text{Var}(w_i) = \frac{1}{n} \quad \frac{2}{n}$$

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

ReLU       $g^{[l]}(z) = \text{ReLU}(z)$

ReLU: He Initialization, sqrt(2/n^([l-1))

Xavier Initialization for tanh

① the   Vorarts:

tanh   $\sqrt{\dfrac{1}{n^{[l-1]}}}$
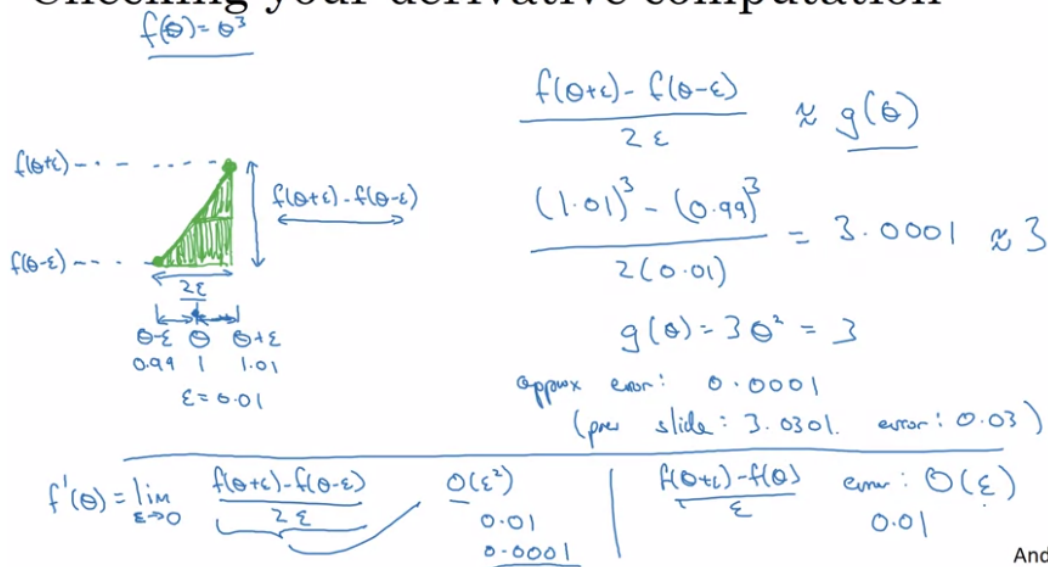
Xavier initialization

**Let's derive Xavier Initialization now, step by step.**

Our full derivation gives us the following initialization rule, which we apply to all weights:

$$W_{i,j}^{[l]} = \mathcal{N}\left(0, \frac{1}{n^{[l-1]}}\right)$$

# Gradient Checking

## Checking your derivative computation

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



$$\frac{f(\theta+\varepsilon) - f(\theta-\varepsilon)}{2\varepsilon} \approx g(\theta)$$

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

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

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

$$f'(\theta) = \lim_{\varepsilon \to 0} \frac{f(\theta+\varepsilon) - f(\theta-\varepsilon)}{2\varepsilon} \qquad \frac{O(\varepsilon^2)}{0.01} \qquad \frac{f(\theta+\varepsilon) - f(\theta)}{\varepsilon} \qquad \text{error: } O(\varepsilon)$$

$0.0001 \qquad\qquad 0.01$

Andrew Ng

## Gradient check for a neural network

Take $W^{[1]}, b^{[1]}, \ldots, W^{[L]}, b^{[L]}$ and reshape into a big vector $\theta$.

Concatenate

$$J(W^{[1]}, b^{[1]}, \ldots, W^{[L]}, b^{[L]}) = J(\theta)$$

Take $dW^{[1]}, db^{[1]}, \ldots, dW^{[L]}, db^{[L]}$ and reshape into a big vector $d\theta$.

concatenate

Is $d\theta$ the gradient of $J(\theta)$

## Gradient checking (Grad check)

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

for each $i$:

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

$$\approx \underline{d\theta[i]} = \frac{\partial J}{\partial \theta_i} \qquad \Bigg| \qquad d\theta_{appr} \overset{?}{\approx} d\theta$$

$i$

Check $\quad \dfrac{\| d\theta_{appr} - d\theta \|_2}{\to \| d\theta_{appr} \|_2 + \| d\theta \|_2}$

$$\approx 10^{-7} - \text{great!}$$
$$10^{-5}$$
$$\varepsilon = 10^{-7} \qquad \to 10^{-3} - \text{worry.}$$

$10^{-7} \to$ great / $10^{-5}$ / $10^{-3}$
-> wrong

Some other notes:
turn-off GC in training (only do this in debugging)

## Gradient checking implementation notes

- Don't use in training – only to debug $\qquad d\theta_{appr}[i] \longleftrightarrow \dfrac{d\theta[i]}{\uparrow}$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \uparrow \ \uparrow$

- If algorithm fails grad check, look at components to try to identify bug.
$\qquad\qquad db^{[l]}_{\vdash} \quad dw^{[l]}_{\vdash}$

- Remember regularization. $\qquad J(\theta) = \frac{1}{m} \sum_i \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_l \| w^{[l]} \|_F^2$
$\qquad\qquad\qquad\qquad\qquad\qquad d\theta = \text{grad of } J \text{ w.r.t. } \theta$

- Doesn't work with dropout. $\quad J \qquad \text{keep-prob} = 1.0$

- Run at random initialization; perhaps again after some training.
$\qquad W, b \approx 0$

it's not impossible that your implementation of gradient descent

is **correct when w and b are close to 0**, so at random initialization. But that as you run gradient descent and w and b become bigger, maybe your implementation of backprop is correct only when w and b is close to 0, but it gets more inaccurate when w and b become large. So one thing you could do, I don't do this very often, but one thing you could do is run grad check **at random initialization and then train the network for a while so that w and b have some time to wander away from 0,** from your small random initial values. And then run grad check again after you've trained for some number of iterations.