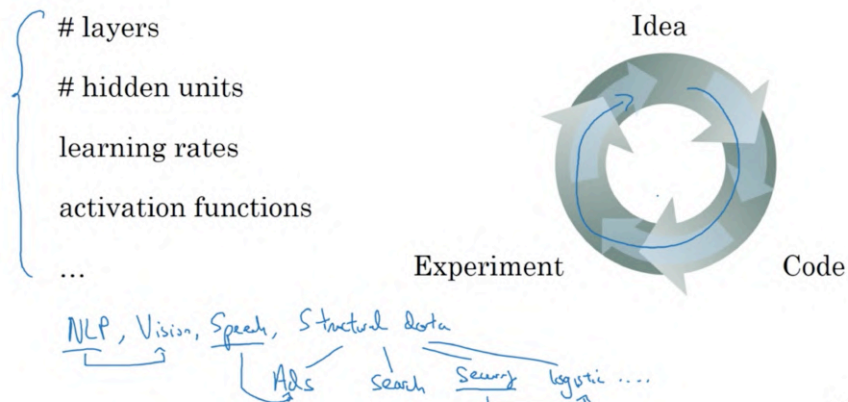
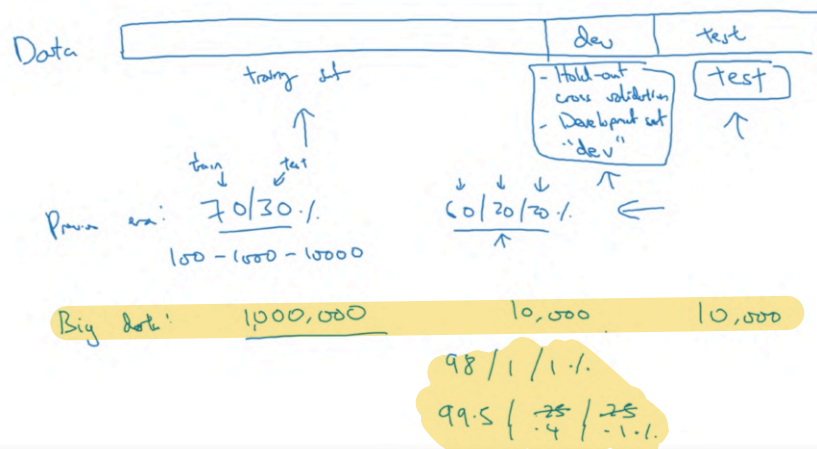


# Applied ML is a highly iterative process



Andrew Ng

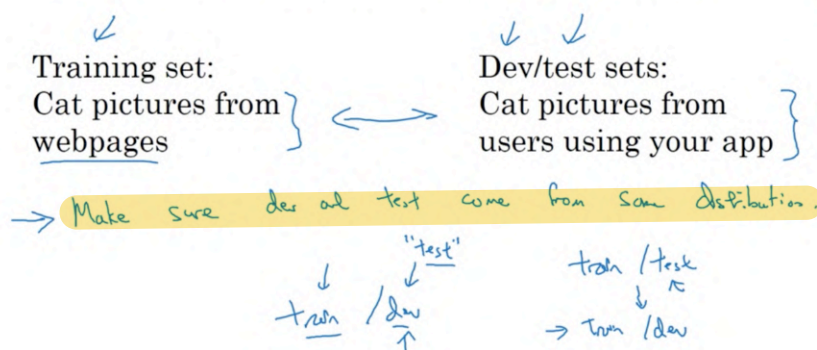
## Train/dev/test sets



Andrew Ng

## Mismatched train/test distribution

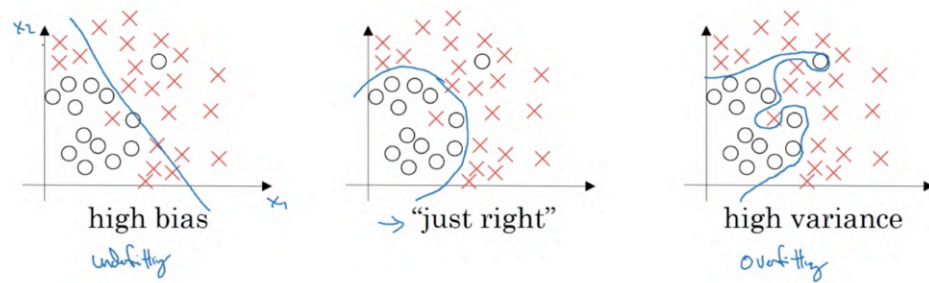
Corts



Not having a test set might be okay. (Only dev set.)

Andrew Ng

# Bias and Variance



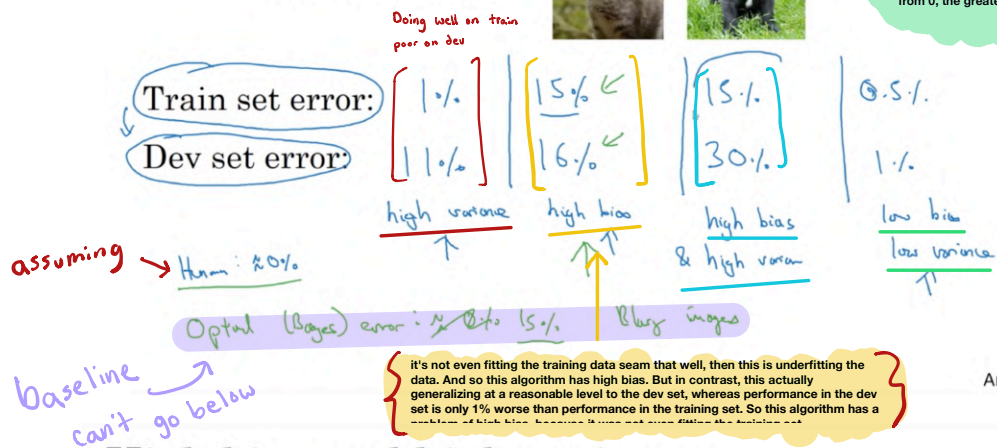
Andrew Ng

## Bias and Variance

### Cat classification

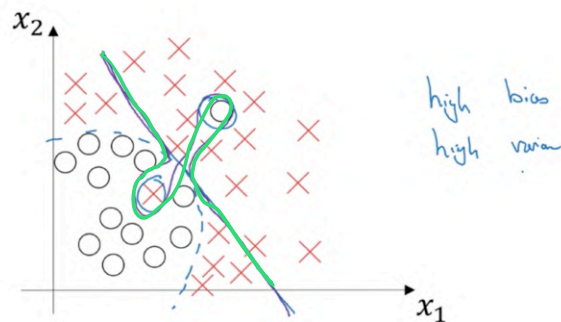


- Variance**
  - The greater the difference between train and dev set error, the greater the variance
- Bias**
  - The greater the difference between train and dev set errors from 0, the greater the bias



Andrew Ng

## High bias and high variance



Andrew Ng

## Basic recipe for machine learning



In the earlier era of machine learning, there used to be a lot of discussion on what is called the bias variance tradeoff. And the reason for that was that, for a lot of the things you could try, you could increase bias and reduce variance, or reduce bias and increase variance. But back in the pre-deep learning era, we didn't have many tools, we didn't have as many tools that just reduce bias or that just reduce variance without hurting the other one. But in the modern deep learning, big data era, so long as you can keep training a bigger network, and so long as you can keep getting more data, which isn't always the case for either of these, but if that's the case, then getting a bigger network almost always just reduces your bias without necessarily hurting your variance, so long as you regularize appropriately. And getting more data pretty much always reduces your variance and doesn't hurt your bias much. So what's really happened is that, with these two steps, the ability to train, pick a network, or get more data, we now have tools to drive down bias and just drive down bias, or drive down variance and just drive down variance, without really hurting the other thing that much. And I think this has been one of the big reasons that deep learning has been so useful for supervised learning, that there's much less of this tradeoff where you have to carefully balance bias and variance, but sometimes you just have more options for reducing bias or reducing variance without necessarily increasing the other one.

Andrew Ng

## Logistic regression

$\min_{w,b} J(w,b)$       $w \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$       $\lambda = \text{regularization parameter}$

$J(w,b) = \frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2$       $\lambda = \text{lambda}$       $\text{omit}$

Euclidean →  $L_2$  regularization      $\|w\|_2^2 = \sum_{j=1}^n w_j^2 = w^T w$       $w$  will be sparse

Manhattan →  $L_1$  regularization      $\frac{\lambda}{2m} \sum_{j=1}^n |w_j| = \frac{\lambda}{2m} \|w\|_1$

Please note that in the next video at 5:45, the Frobenius norm formula should be the following:

$$\|w^{[l]}\|^2 = \sum_{i=1}^{n^l} \sum_{j=1}^{n^{l-1}} (w_{ij}^{[l]})^2$$

The limit of summation of  $i$  should be from 1 to  $n^l$ .

The limit of summation of  $j$  should be from 1 to  $n^{l-1}$ .

(It's flipped in the video). The rows "i" of the matrix should be the number of neurons in the current layer  $n^l$ ;

whereas the columns "j" of the weight matrix should equal the number of neurons in the previous layer  $n^{l-1}$ .

Andrew Ng

## Neural network

$J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$

$\|w^{[l]}\|_F^2 = \sum_{i=1}^{n^{l+1}} \sum_{j=1}^{n^l} (w_{ij}^{[l]})^2$       $w: \begin{pmatrix} n^{l+1} & n^l \end{pmatrix}$

"Frobenius norm"      $\|\cdot\|_2^2$       $\|\cdot\|_F^2$

$dW = (\text{from backprop}) + \frac{\lambda}{m} W^{[l]}$       $\frac{\partial J}{\partial w^{[l]}} = dW^{[l]}$

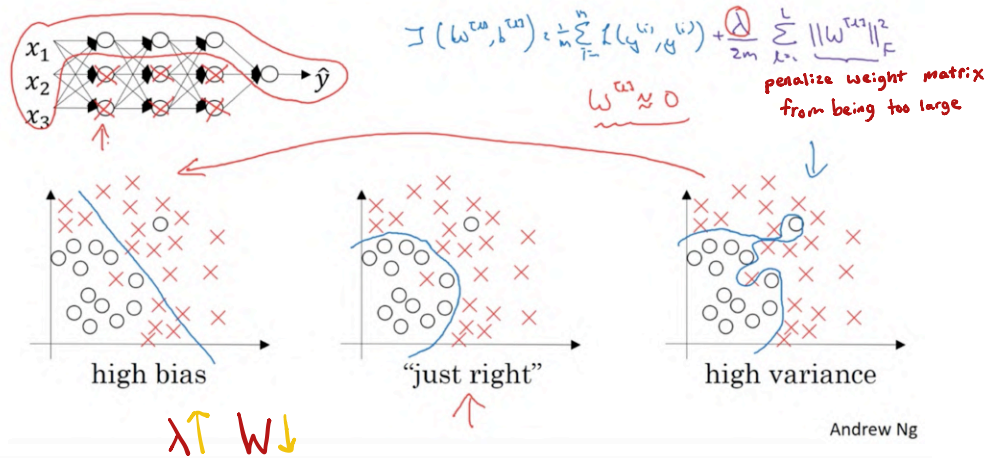
$\rightarrow W^{[l]} := W^{[l]} - \alpha dW^{[l]}$

"L2" → "Weight decay"      $W^{[l]} := W^{[l]} - \alpha \left[ (\text{from backprop}) + \frac{\lambda}{m} W^{[l]} \right]$

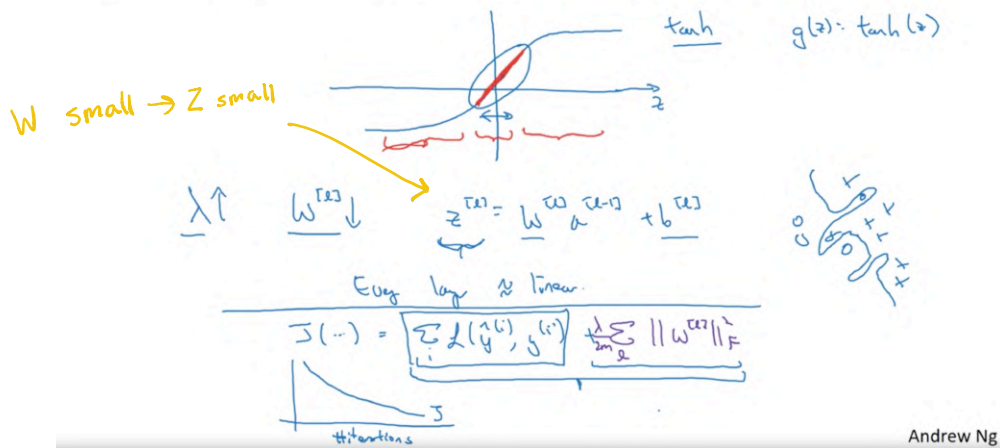
$\left[ \left( 1 - \frac{\alpha \lambda}{m} \right) W^{[l]} \right] = \left[ W^{[l]} - \left( \frac{\alpha \lambda}{m} \right) W^{[l]} \right] - \alpha (\text{from backprop})$

Andrew Ng

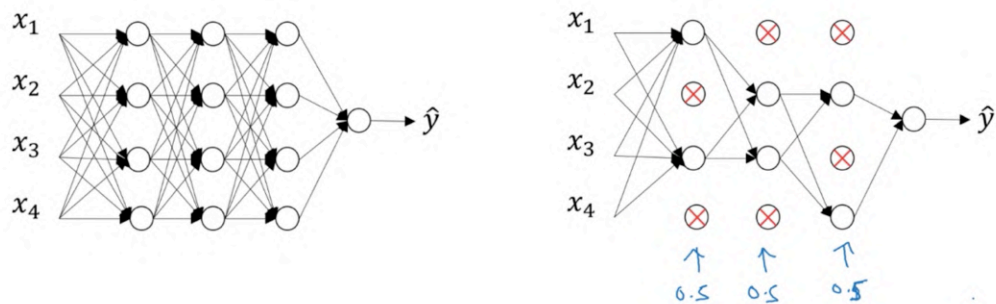
How does regularization prevent overfitting? (Variance)



How does regularization prevent overfitting?



Dropout regularization





## Implementing dropout ("Inverted dropout")

Illustrate with layer  $l=3$ .  $\text{keep-prob} = 0.8$   $0.2 \leftarrow \text{drop}$

$\rightarrow d3 = \text{np.random.rand}(a3.\text{shape}[0], a3.\text{shape}[1]) < \text{keep-prob}$

$a3 = \text{np.multiply}(a3, d3)$  #  $a3 \neq d3$ .

$\rightarrow a3 /= \text{keep-prob}$   $\leftarrow$

50 units.  $\rightarrow$  10 units shut off

$z^{(4)} = w^{(4)} \cdot a^{(3)} + b^{(4)}$

$\uparrow$   $\nwarrow$  reduced by 20%. Test

$/= 0.8$

In order to not reduce the expected value of  $z^{(4)}$ , what you do is you need to take this, and divide it by 0.8 because this will correct or just a bump that back up by roughly 20% that you need. So it's not changed the expected value of  $a3$ . And, so this line here is what's called the inverted dropout technique. And its effect is that, no matter what you set to keep-prob to, whether it's 0.8 or 0.9 or even one, if it's set to one then there's no dropout, because it's keeping everything or 0.5 or whatever, this inverted dropout technique, there is there is line to are due to the green box at dropping out. This makes test time easier because you have less of a scaling problem. By far the most common implementation of dropouts today as far as I know is inverted dropouts. I recommend you just implement this. But there were some early iterations of dropout that missed this divide by keep-prob line, and so at test time the average becomes more and more complicated

Andrew Ng

## Making predictions at test time

$$a^{(0)} = X$$

No drop out.

$$\begin{aligned} z^{(1)} &= w^{(1)} a^{(0)} + b^{(1)} \\ a^{(1)} &= g^{(1)}(z^{(1)}) \\ z^{(2)} &= w^{(2)} a^{(1)} + b^{(2)} \\ a^{(2)} &= \dots \end{aligned}$$

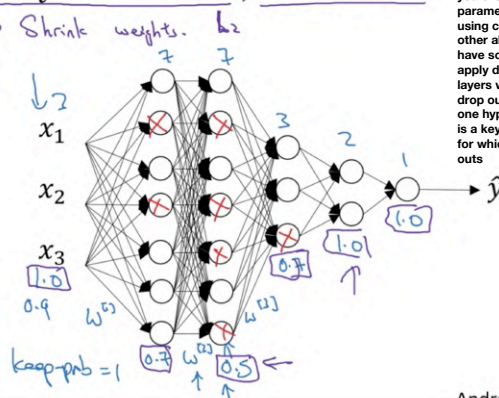
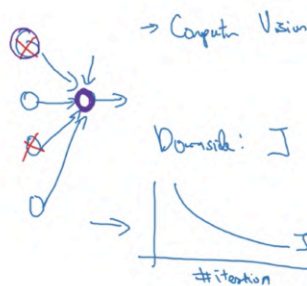
$/= \text{keep-prob}$

$\downarrow$   
 $\hat{y}$

Andrew Ng

## Why does drop-out work?

Intuition: Can't rely on any one feature, so have to spread out weights.  $\rightarrow$  Shrink weights.

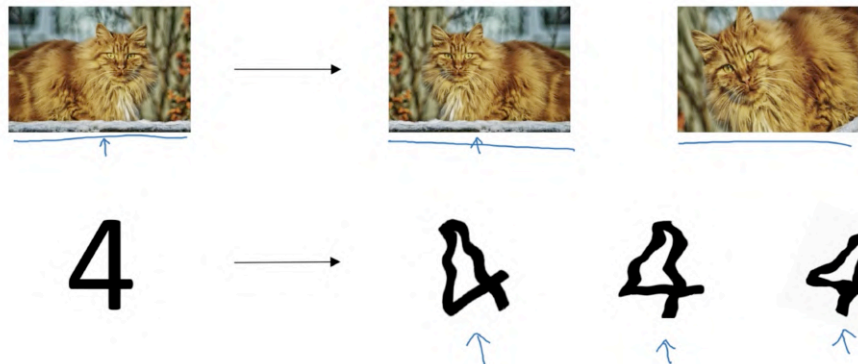


So just to summarize, if you're more worried about some layers overfitting than others, you can set a lower key prop for some layers than others.

The downside is, this gives you even more hyper parameters to search for using cross-validation. One other alternative might be to have some layers where you apply drop out and some layers where you don't apply drop out and then just have one hyper parameter, which is a key prop for the layers for which you do apply drop outs

Andrew Ng

## Data augmentation



Andrew Ng

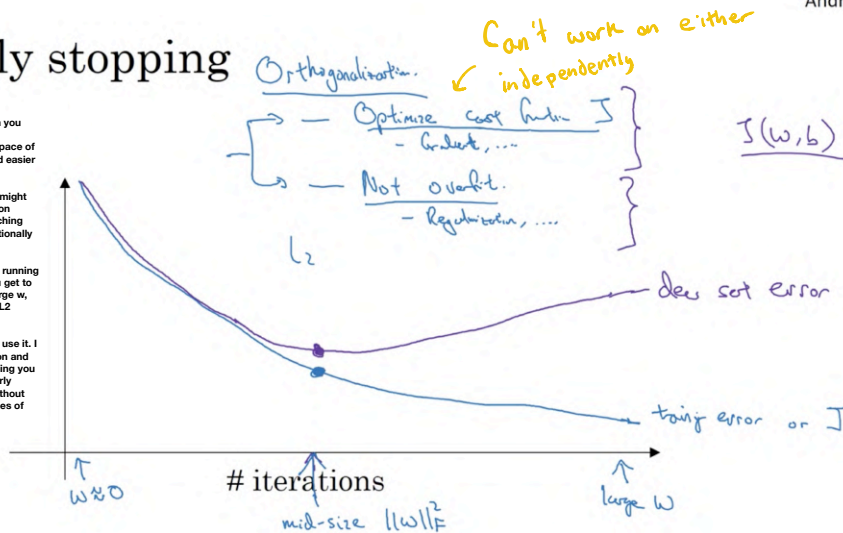
## Early stopping

Alternative is just use L2 regularization then you can just train the neural network as long as possible. I find that this makes the search space of hyper parameters easier to decompose, and easier to search over.

But the downside of this though is that you might have to try a lot of values of the regularization parameter lambda. And so this makes searching over many values of lambda more computationally expensive.

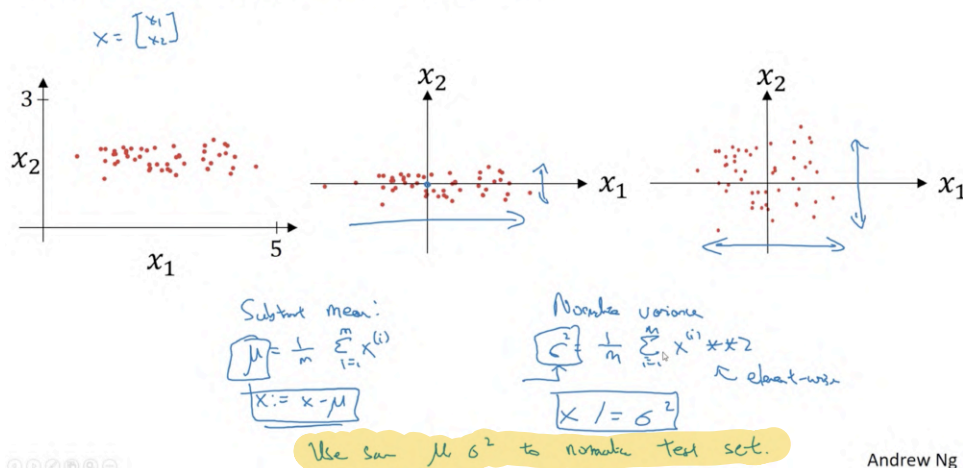
And the advantage of early stopping is that running the gradient descent process just once, you get to try out values of small  $w$ , mid-size  $w$ , and large  $w$ , without needing to try a lot of values of the L2 regularization hyperparameter lambda.

Despite it's disadvantages, many people do use it. I personally prefer to just use L2 regularization and try different values of lambda. That's assuming you can afford the computation to do so. But early stopping does let you get a similar effect without needing to explicitly try lots of different values of lambda.



Andrew Ng

## Normalizing training sets speeds up training

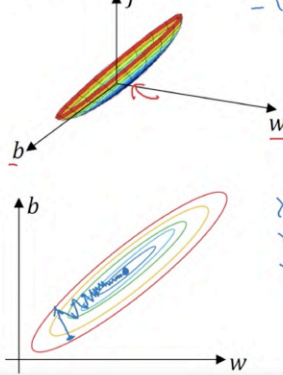


Andrew Ng

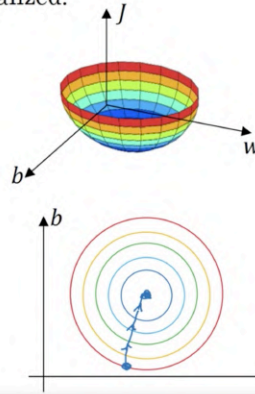
## Why normalize inputs?

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

Unnormalized:  
 $w_1: x_1: 1 \dots 1000 \leftarrow$   
 $w_2: x_2: 0 \dots 1 \leftarrow$   
 $-1 \dots 1$



Normalized:

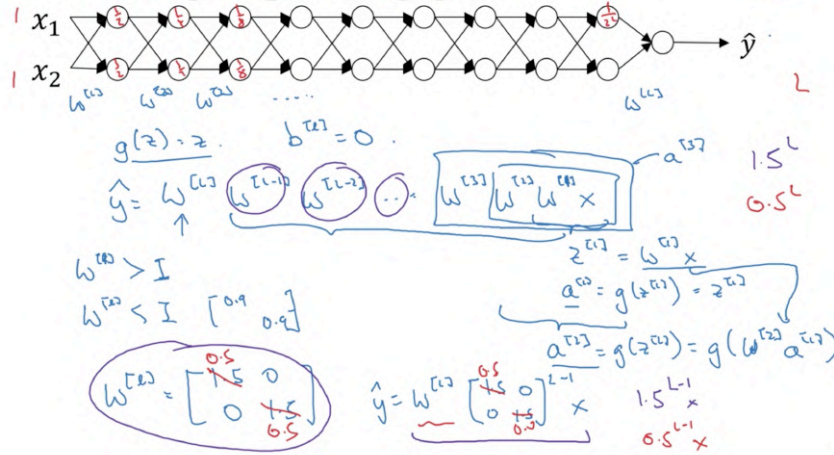


$x_1: 0 \dots 1$   
 $x_2: -1 \dots 1$   
 $x_3: -1 \dots 2$

Andrew Ng

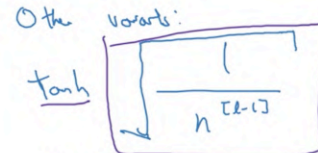
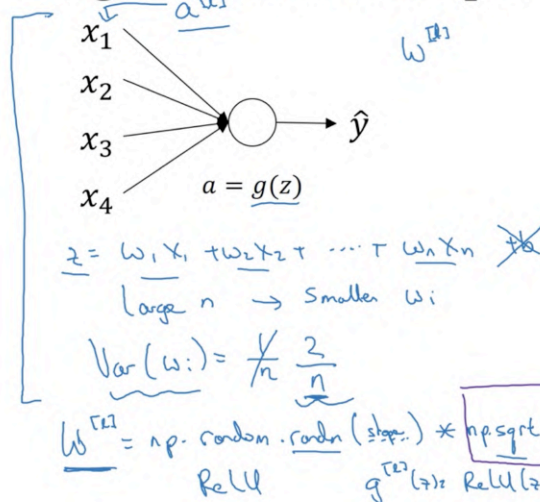
## Vanishing/exploding gradients

$L=150$



Andrew Ng

## Single neuron example



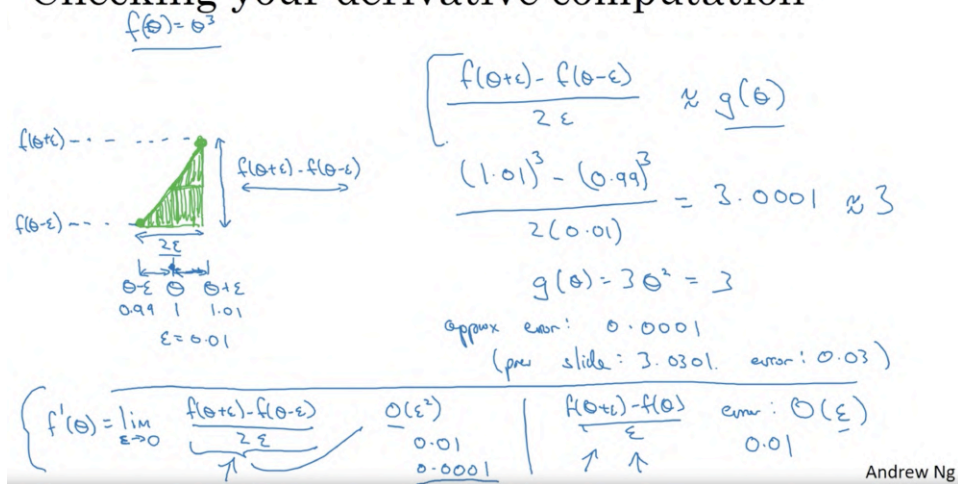
Other variants:  
 tanh  
 Xavier initialization  $\uparrow$

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

Andrew Ng



## Checking your derivative computation



## Gradient check for a neural network

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

concentrate

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

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

concentrate

Is  $d\theta$  the gradient of  $J$ .

Andrew Ng

## Gradient checking (Grad check)

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

for each  $i$ :

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

$$\approx \underline{d\theta[i]} = \frac{\partial J}{\partial \theta_i} \quad \left| \quad d\theta_{\text{approx}} \approx d\theta \right.$$

Check

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

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

$\approx \frac{10^{-7}}{10^{-5}} = 10^{-2} \leftarrow$

$\rightarrow 10^{-3} \leftarrow$

$10^{-7} - \text{great!}$

$10^{-3} - \text{worry}$

Andrew Ng



# Gradient checking implementation notes

- Don't use in training – only to debug

$$\frac{d\Theta_{\text{approx}}[i]}{\uparrow \uparrow} \longleftrightarrow \frac{d\Theta[i]}{\uparrow}$$

- If algorithm fails grad check, look at components to try to identify bug.

$$\frac{db^{[L]}}{\uparrow} \quad \frac{dw^{[L]}}{\uparrow}$$

- Remember regularization.

$$J(\Theta) = \frac{1}{n} \sum_i L(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2n} \sum_i \|w^{(i)}\|_F^2$$

$d\Theta = \text{gradient of } J \text{ wrt. } \Theta$

- Doesn't work with dropout.

$$J \quad \text{keep-prob} = 1.0$$

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

$$\underline{w, b} \approx 0$$