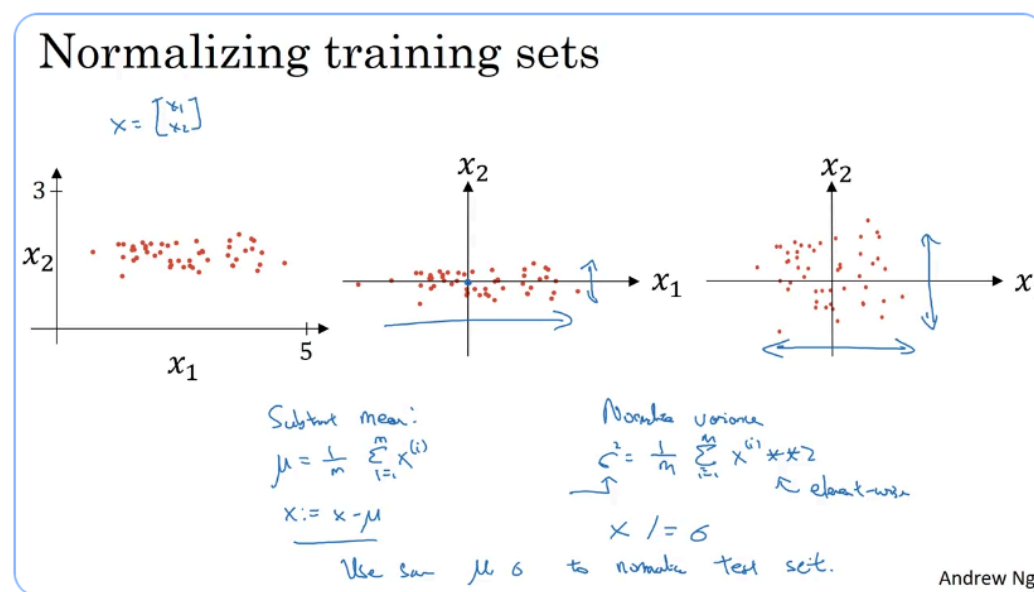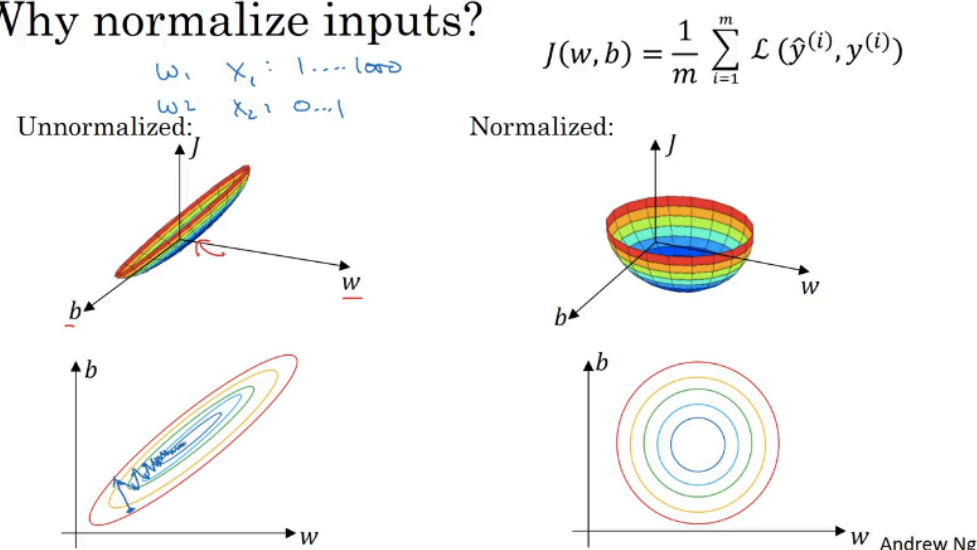# Setting Up your Optimization Problem



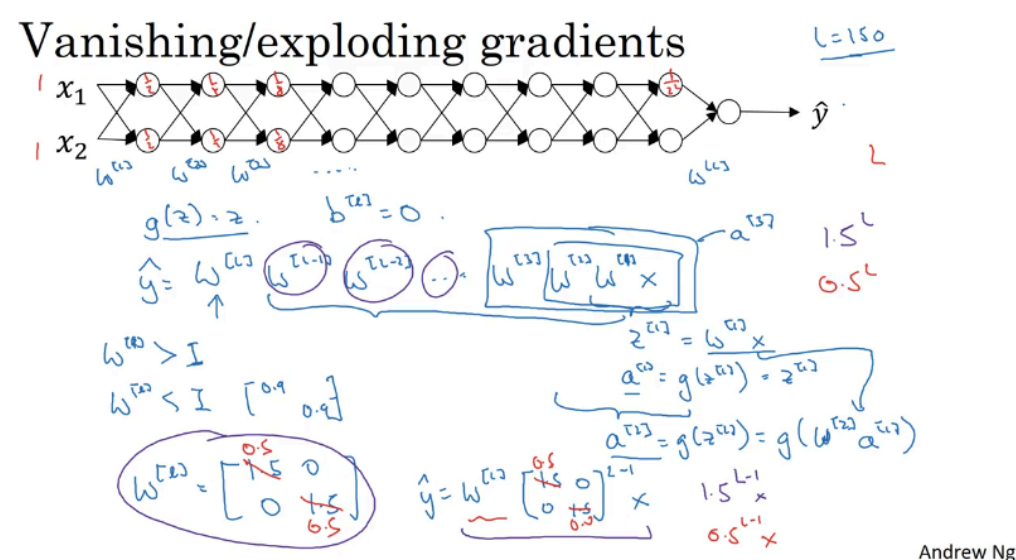If you use this to scale your training data, then use the same mu and sigma to normalize your test set.

In particular, you don't want to normalize the training set and a test set differently. Whatever this value is and whatever this value is, use them in these two formulas so that you scale your test set in exactly the same way rather than estimating mu and sigma squared separately on your training set and test set, because you want your data both training and test examples to go through the same transformation defined by the same Mu and Sigma squared calculated on your training data. → normalize both train and test set



Cost function after taking normalize → its easy for gradient descent to work

# Vanishing / Exploding Gradients

So the intuition I hope you can take away from this is that at the weights W, if they're all just a little bit bigger than one or just a little bit bigger than the identity matrix, then with a very deep network the activations can explode. And if W is just a little bit less than identity. So this maybe here's 0.9, 0.9, then you have a very deep network, the activations will decrease exponentially.
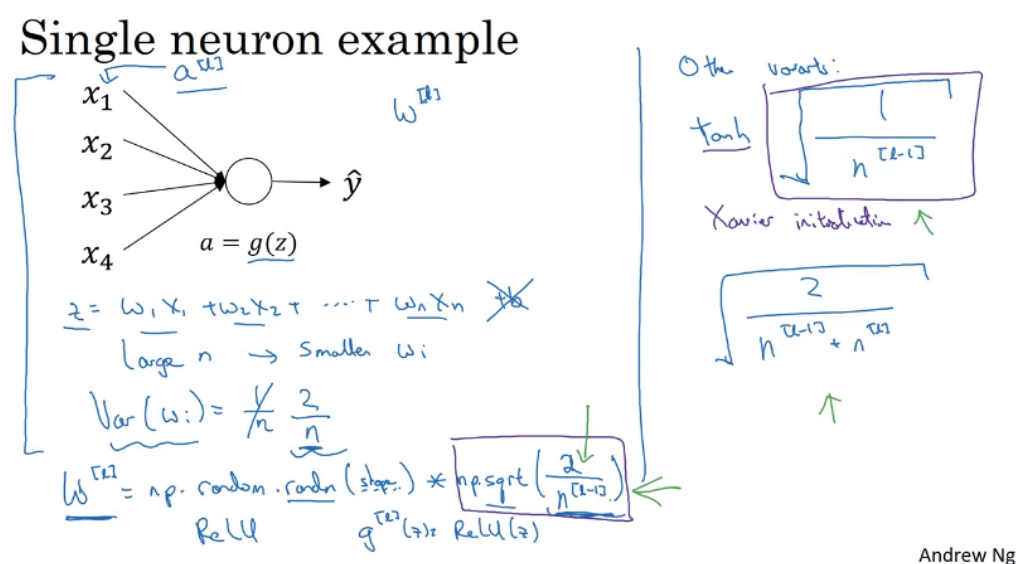
- if we choose W > I (identity matrix) → the weight could be very large (increase exponentially )

- if we choose W < I → the weight could be very small (decrease exponentially)

But with such a deep neural network, if your activations or gradients increase or decrease exponentially as a function of L, then these values could get really big or really small.

And this makes training difficult, especially if your gradients are exponentially smaller than L, then gradient descent will take tiny little steps. It will take a long time for gradient descent to learn anything.

To summarize, you've seen how deep networks suffer from the problems of vanishing or exploding gradients. In fact, for a long time this problem was a huge barrier to training deep neural networks. It turns out there's a partial solution that doesn't completely solve this problem but it helps a lot which is careful choice of how you initialize the weights.

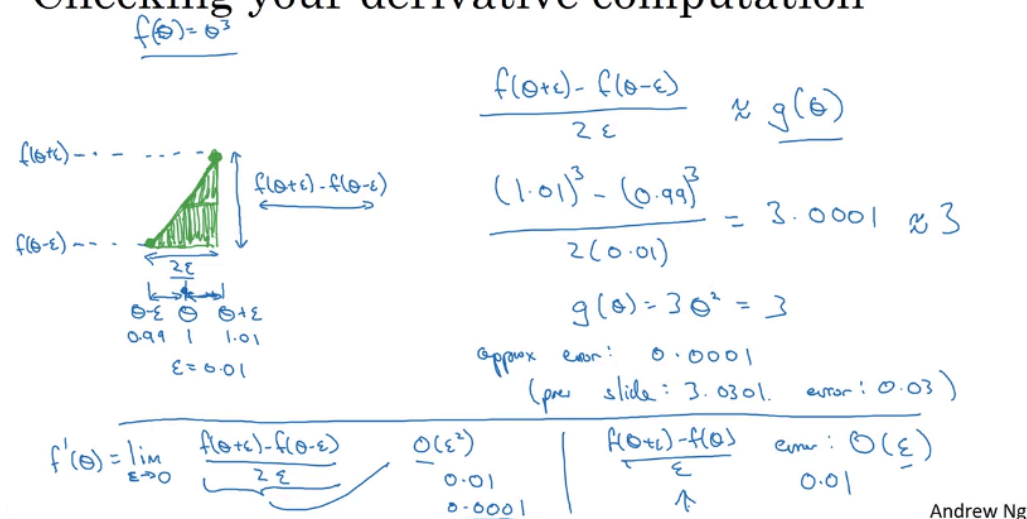# Weight Initialization for Deep Networks



Andrew Ng

 So in order to make z not blow up and not become too small, you notice that the larger n is, the smaller you want $w_i$ to be, right? Because z is the sum of the $w_i x_i$. And so if you're adding up a lot of these terms, you want each of these terms to be smaller. One reasonable thing to do would be to set the variance of $w$ to be equal to $\frac{1}{n}$, where n is the number of input features that's going into a neuron.

But for working a little bit better when we use:

- ReLU function, Variance of $w$ should be $\frac{2}{n}$, then we should initialize W with np.random.rand(shape) * np.sqrt($\frac{2}{n^{[l-1]}}$)

- tanh function: instead of using $\frac{2}{n^{[l-1]}}$ we could use $\frac{1}{n^{[l-1]}}$, this is called Xavier initalization, and there is another version $\frac{2}{n^{[l-1]}+n^{[l]}}$

# Numerical Approximation of Gradients



Andrew Ng

# Gradient Checking

Gradient checking (Grad check)    $J(\theta) = J(\theta_1, \theta_2, \dots)$

for each $i$:

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

$\approx d\theta[i] = \dfrac{\partial J}{\partial \theta_i}$  |  $d\theta_{approx} \overset{?}{\approx} d\theta$

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

$\varepsilon = 10^{-7}$    $10^{-5}$

$10^{-3}$ — worry.

Andrew Ng

check if our gradient is right:

using formula in the image

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

# Gradient Checking Implementation Notes

Gradient checking implementation notes

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

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

$db^{[l]}_{\ell} \quad dw^{[l]}_{\ell}$

$J(\theta) = \frac{1}{n}\sum \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m}\sum_l \|w^{[l]}\|_F^2$

- Remember regularization.

$d\theta = $ grad of $J$ wrt. $\theta$

- Doesn't work with dropout.    $J$    keep-prob = 1.0

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

$w, b \approx 0$

Andrew Ng

 There isn't an easy to compute cost function J that dropout is doing gradient descent on.

It turns out that dropout can be viewed as optimizing some cost function J, but it's cost function J defined by summing over all exponentially large subsets of nodes they could eliminate in any iteration. So the cost function J is very difficult to compute, and you're just sampling the cost function every time you eliminate different random subsets in those we use dropout. So it's difficult to use grad check to double check your computation with dropouts. So what I usually do is implement grad check without dropout. So if you want, you can set keep-prob and dropout to be equal to 1.0.

And then turn on dropout and hope that my implementation of dropout was correct.