

Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>

PA-1 : Parameters Initialisation 1-In general, initializing all the weights to zero results in the network failing to break symmetry. This means that every neuron



deeplearning.ai

Setting up your ML application

Train/dev/test sets

Applied ML is a highly iterative process

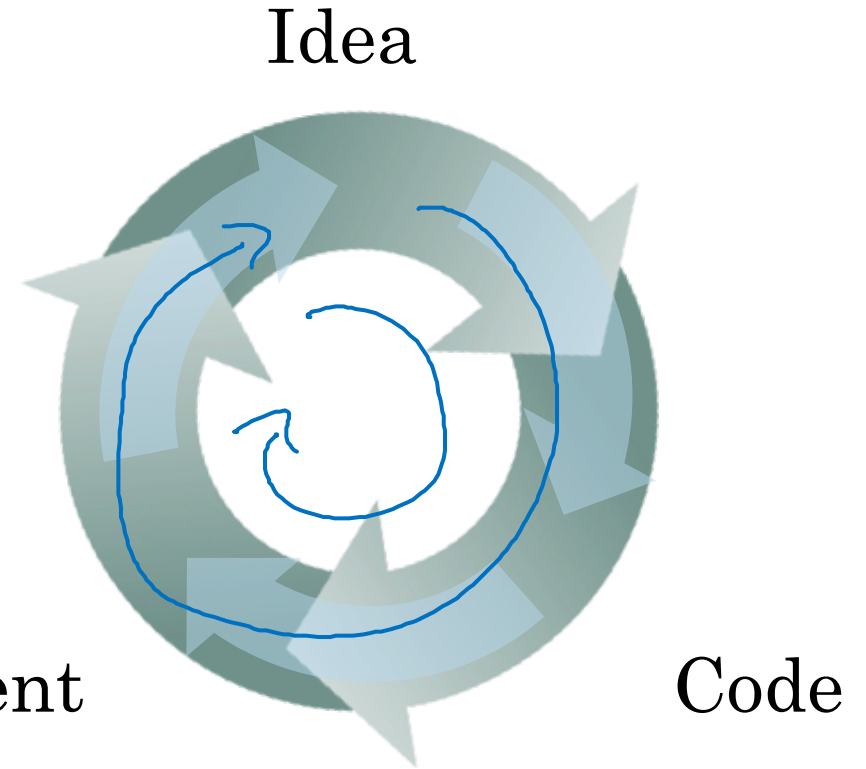
layers

hidden units

learning rates

activation functions

• • •



NLP, Vision, Speech, Structural Data

The diagram shows four terms at the top: "NLP", "Vision", "Speech", and "Structural Data".
- A bracket connects "NLP" and "Vision", with an arrow pointing from "Vision" towards the right.
- A bracket connects "Speech" and "Structural Data", with an arrow pointing from "Structural Data" towards the right.
- Below "Speech", there are three terms: "Ads", "Search", and "Security".
- An arrow points from "Speech" down to "Ads".
- Another arrow points from "Speech" down to "Search".
- A third arrow points from "Speech" down to "Security".
- A fourth arrow points from "Speech" down to "logistic ...".
- A bracket connects "Security" and "logistic ...", with an arrow pointing from "logistic ..." towards the right.

Train/dev/test sets



- Hold-out cross validation
- Development set "dev"

test



Previous era: $\frac{70}{30} \%$
100 - 1000 - 10000

$\frac{60}{20}{20} \%$

If you have small dataset then it is OK to split in the ratio of 60/20/20 and something like

Big data: $\frac{1,000,000}{}$

10,000

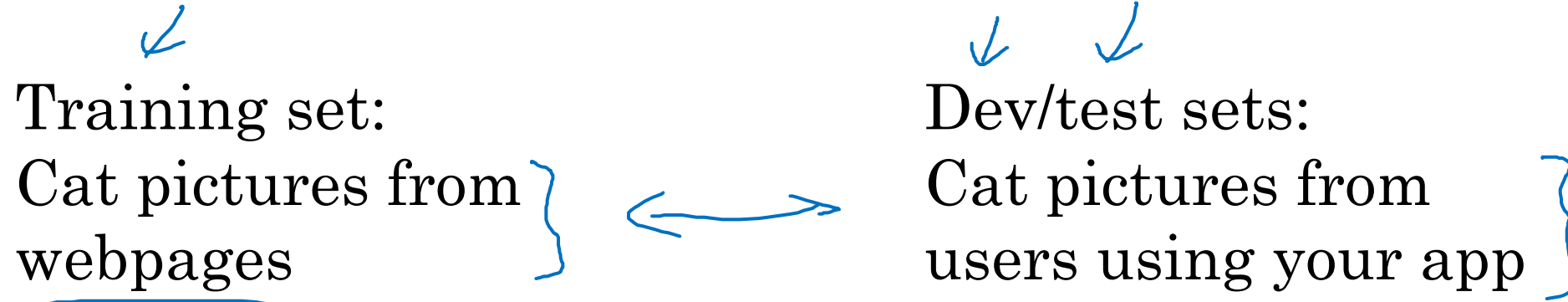
10,000

98 / 1 / 1 %

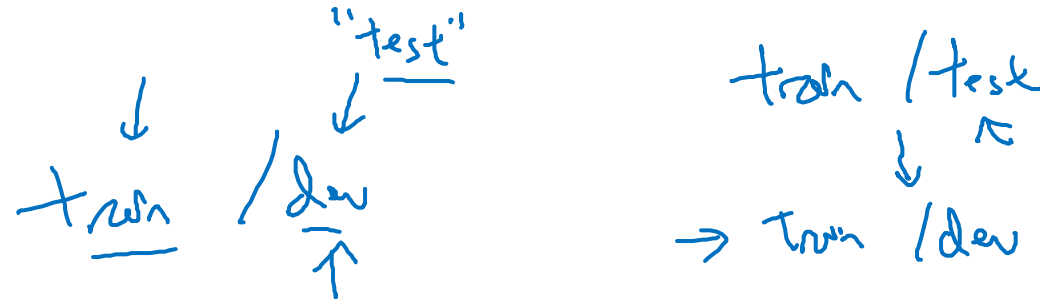
99.5 / .25 / .25
 .4 / .1 %

Mismatched train/test distribution

Certs



→ Make sure dev and test come from same distribution.



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

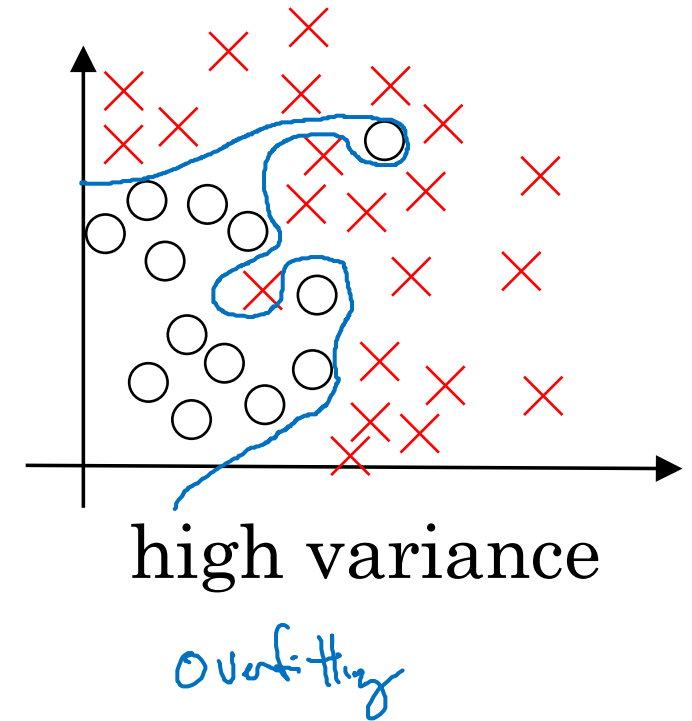
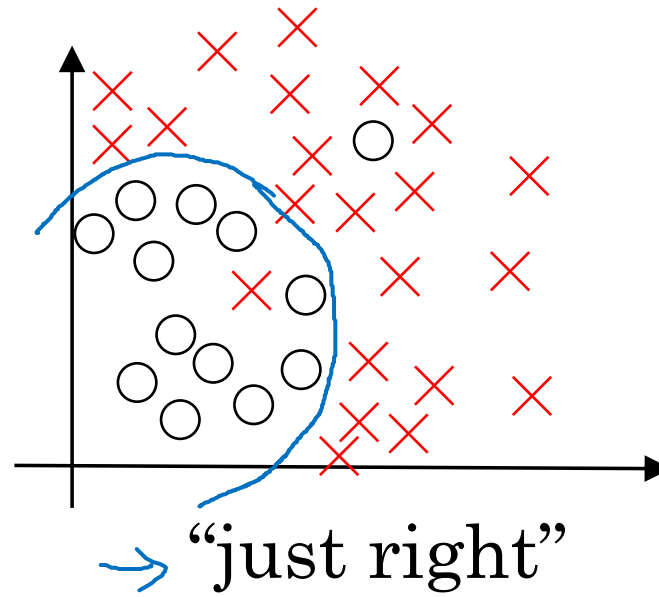
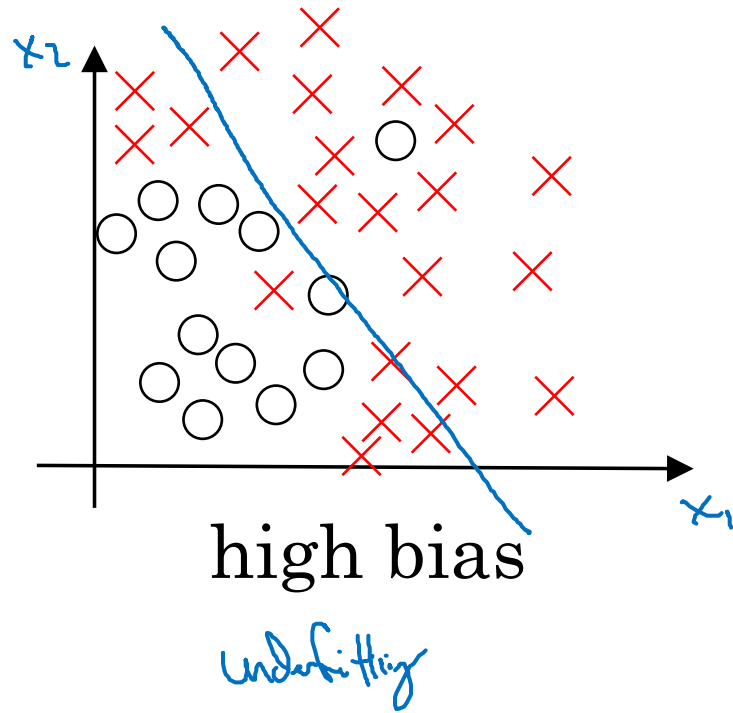


deeplearning.ai

Setting up your ML application

Bias/Variance

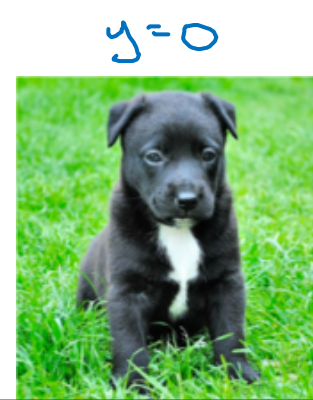
Bias and Variance



O and V paas paas hain

Bias and Variance

Cat classification



Train set error:

1%

15% \swarrow

15%

0.5%

Dev set error:

11%

16% \swarrow

30%

1%

high variance
 \uparrow

high bias
 \uparrow

high bias
& high variance

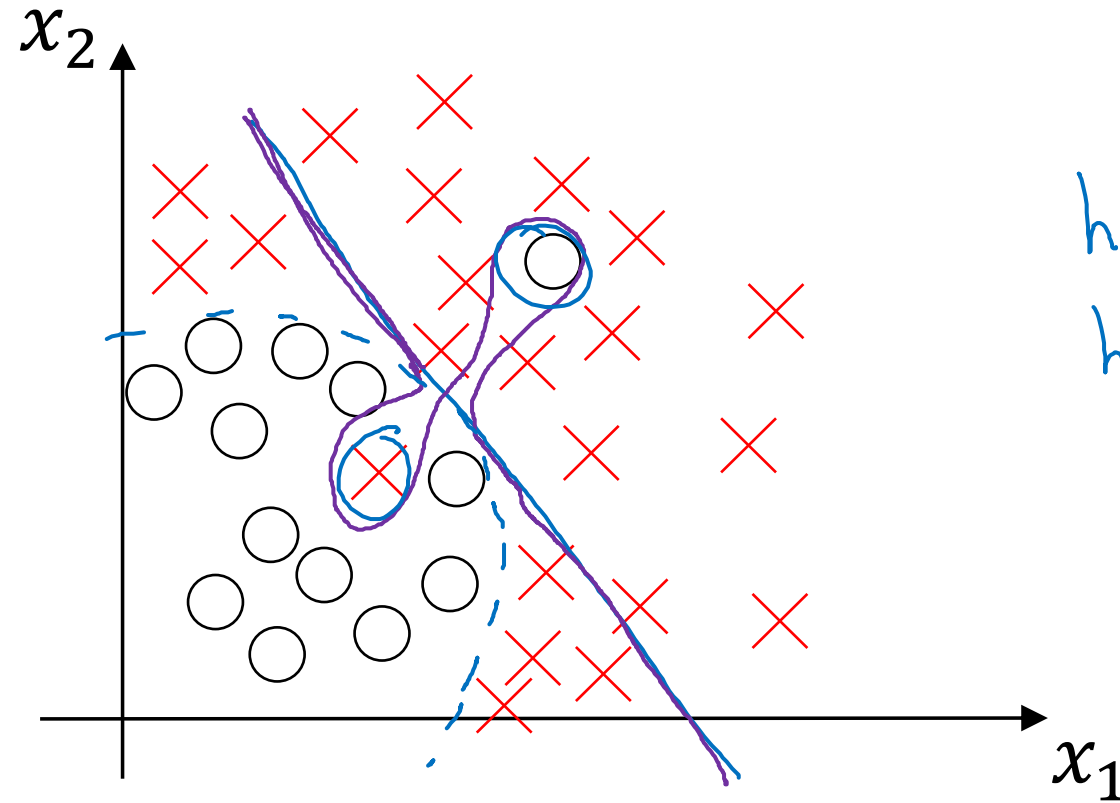
low bias
low variance
 \uparrow

Human: ~0%

Optimal (Bayes) error: ~~~0%~~ 15%

Blurry images

High bias and high variance



high bias
high variance

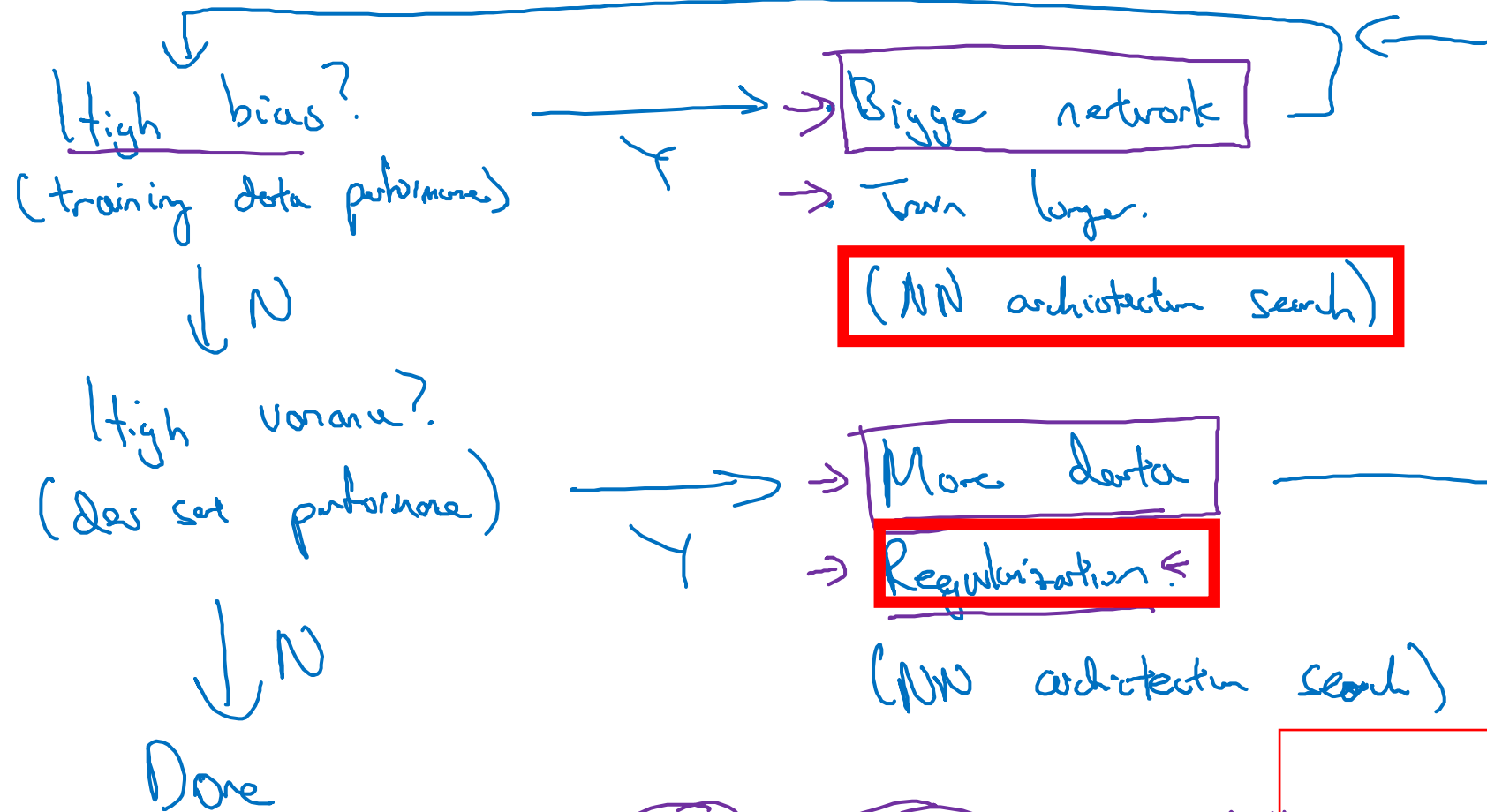


deeplearning.ai

Setting up your ML application

Basic “recipe” for machine learning

Basic recipe for machine learning



In earlier times there used to be trade-off but now w



deeplearning.ai

Regularizing your neural network

Regularization

Logistic regression

λ = regularization parameter

lambda

lambda

$$\min_{w,b} J(w,b)$$

$$\underline{w} \in \mathbb{R}^{n_x}, \underline{b} \in \mathbb{R}$$

$$J(w,b) = \underbrace{\frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, \hat{y}^{(i)})}_{\text{cost function}} + \frac{\lambda}{2m} \underbrace{\|w\|_2^2}_{\text{L2 regularization}}$$

~~$+\frac{\lambda}{2m} b^2$~~
omit

L_2 regularization $\underline{\|w\|_2^2} = \sum_{j=1}^{n_x} w_j^2 = w^T w \leftarrow$

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

w will be sparse

Neural network

$$\rightarrow J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(y^{(i)}, \hat{y}^{(i)})}_{\text{loss}} + \underbrace{\frac{\lambda}{2n} \sum_{l=1}^L \|w^{[l]}\|_F^2}_{\text{weight decay}}$$

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

$$w^{[l]}: \begin{matrix} n^{[l]} & n^{[l-1]} \\ \uparrow & \uparrow \end{matrix}$$

"Frobenius norm"

$$\|\cdot\|_2^2$$

$$\|\cdot\|_F^2$$

$$dw^{[l]} = \left[(\text{from backprop}) + \frac{\lambda}{n} w^{[l]} \right]$$

$$\rightarrow w^{[l]} := w^{[l]} - \alpha dw^{[l]}$$

$$\frac{\partial J}{\partial w^{[l]}} = dw^{[l]}$$

"Weight decay"

$$w^{[l]} := w^{[l]} - \alpha \left[(\text{from backprop}) + \frac{\lambda}{n} w^{[l]} \right]$$

$$= w^{[l]} - \frac{\alpha \lambda}{n} w^{[l]} - \alpha (\text{from backprop})$$

$$= \underbrace{\left(1 - \frac{\alpha \lambda}{n}\right)}_{\leq 1} \underbrace{w^{[l]}}_{\text{weight}} - \alpha (\text{from backprop})$$

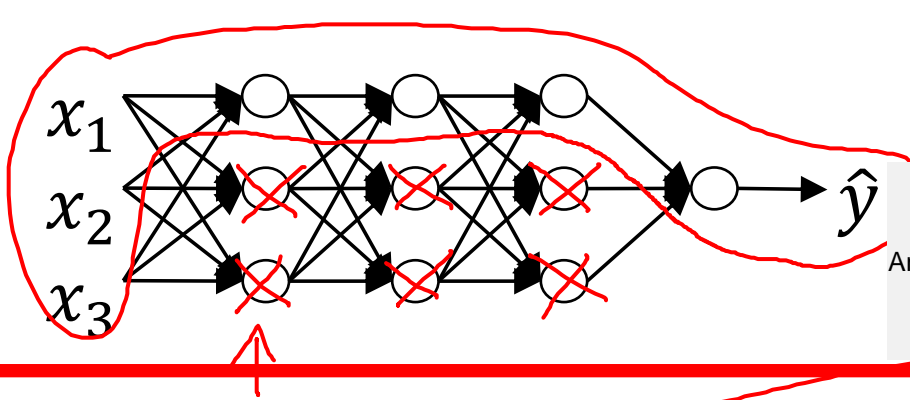


deeplearning.ai

Regularizing your neural network

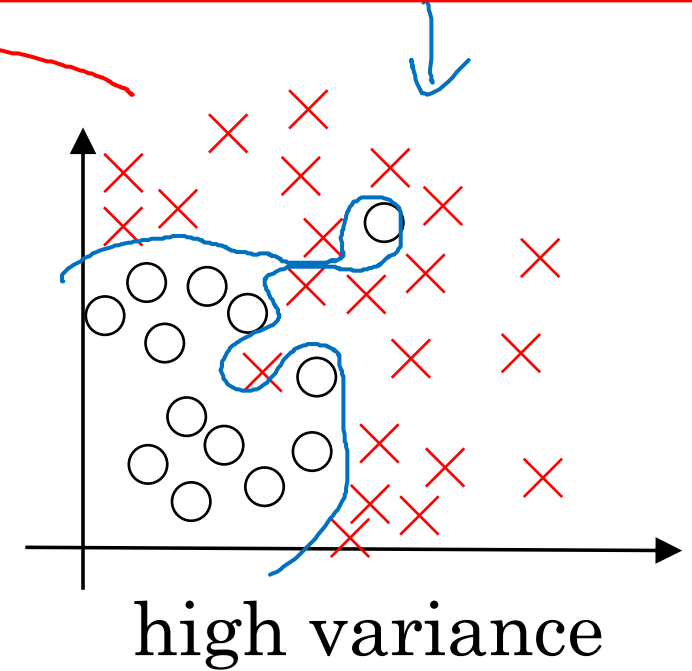
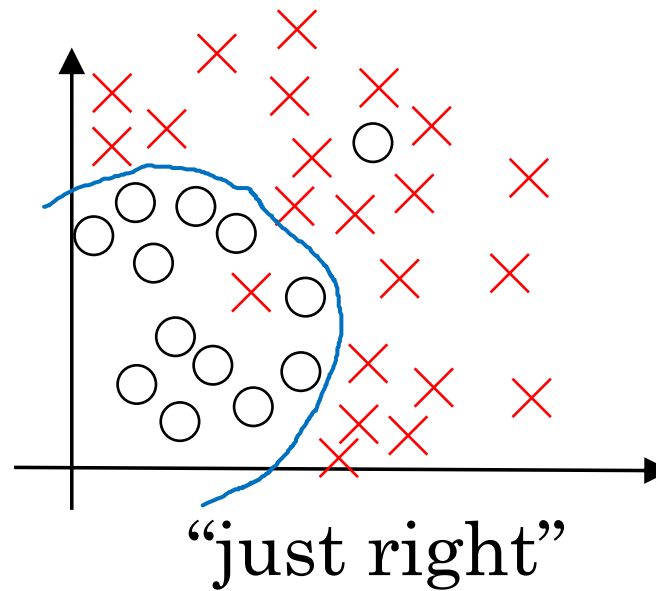
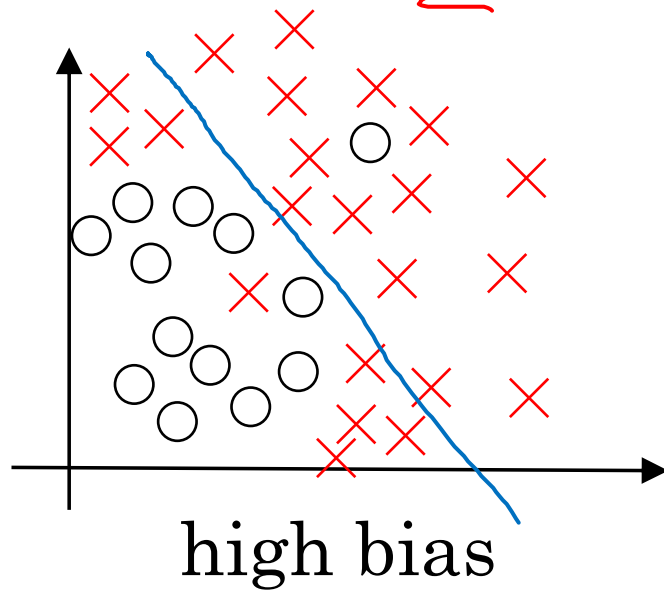
Why regularization reduces overfitting

How does regularization prevent overfitting?

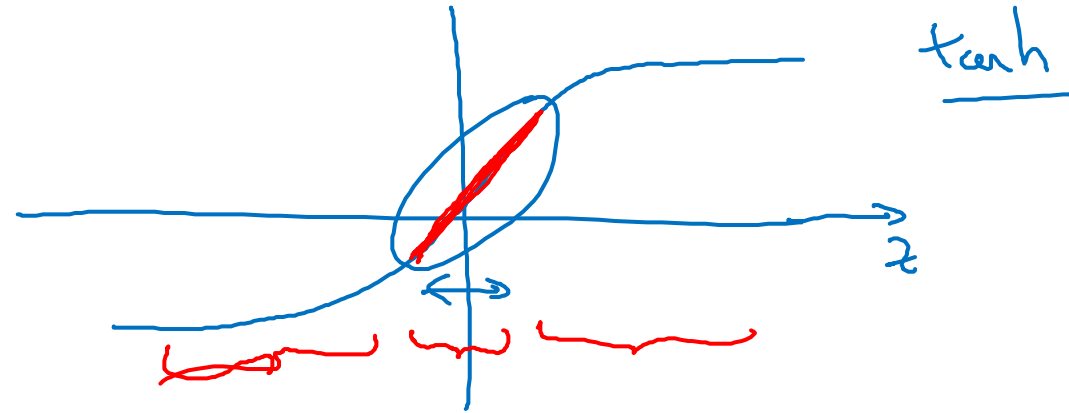


$$J(w^{(L)}, b^{(L)}) = \frac{1}{n} \sum_{i=1}^n l(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(L)}\|_F^2$$

And so what we did for regularization was add this extra term that penalizes the weight matrices from being too large. And we said that



How does regularization prevent overfitting?



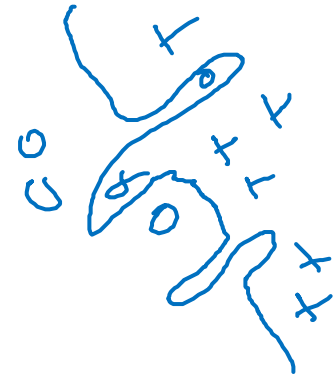
$$g(z) = \tanh(z)$$

$$\lambda \uparrow$$

$$W^{[L]} \downarrow$$

$$z^{[L]} = W^{[L]} a^{[L-1]} + b^{[L]}$$

Every layer \approx linear.



$$J(\dots) = \underbrace{\sum_i \mathcal{L}(\hat{y}^{(i)}, y^{(i)})}_{\text{training loss}} + \underbrace{\frac{\lambda}{2m} \sum_L \|W^{[L]}\|_F^2}_{\text{regularization term}}$$





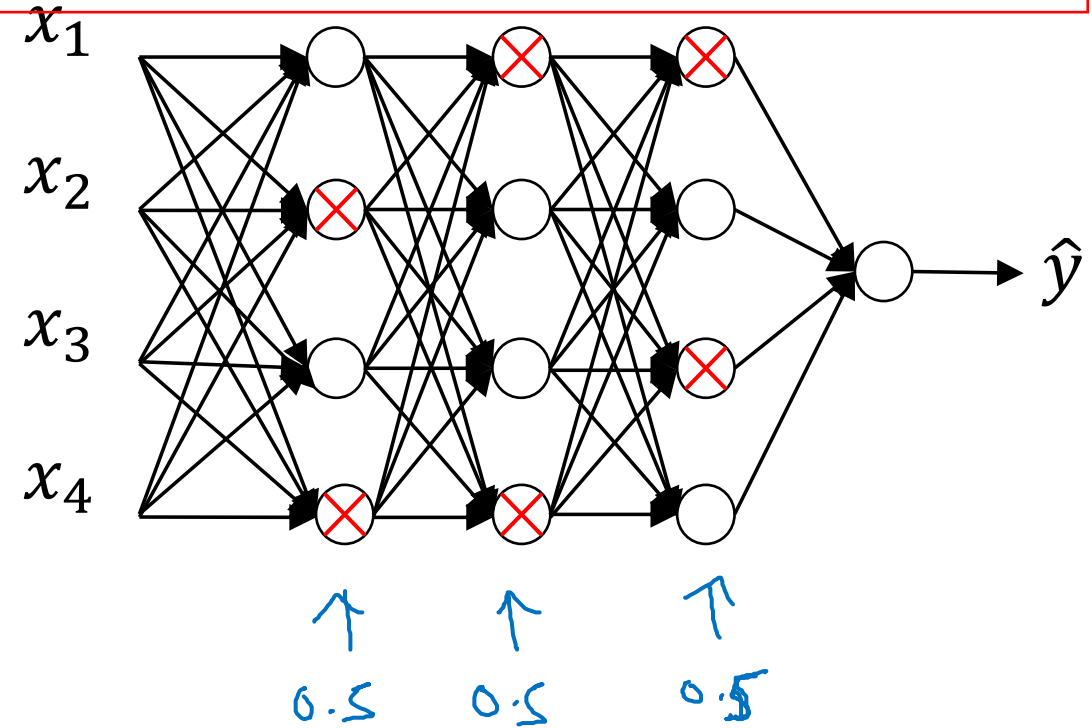
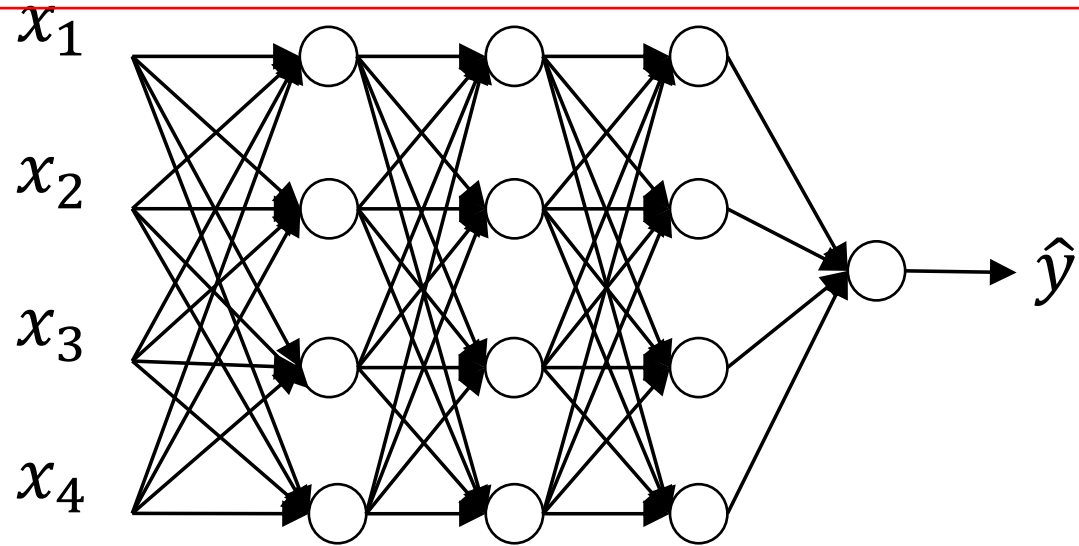
deeplearning.ai

Regularizing your neural network

Dropout regularization

Dropout regularization

what we do we associate a probability with each node denoting its probability of being dropped out. then we actually drop those nodes and do our training



Implementing dropout ("Inverted dropout")

Illustrate with layer $l=3$.

keep-prob = 0.8

0.2

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

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

$a3 \neq d3$.

→ $a3 /= \text{keep-prob}$ ←

50 units. \leadsto 10 units shut off

$$z^{[4]} = w^{[4]} \cdot \underbrace{a^{[3]}}_{\text{reduced by 20\%}} + b^{[4]}$$

↑

reduced by 20%

$/= 0.8$

Test

Making predictions at test time

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

No drop out.

No drop out at the testing time.

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

↓
↑
y

/= keep-prob



deeplearning.ai

Regularizing your neural network

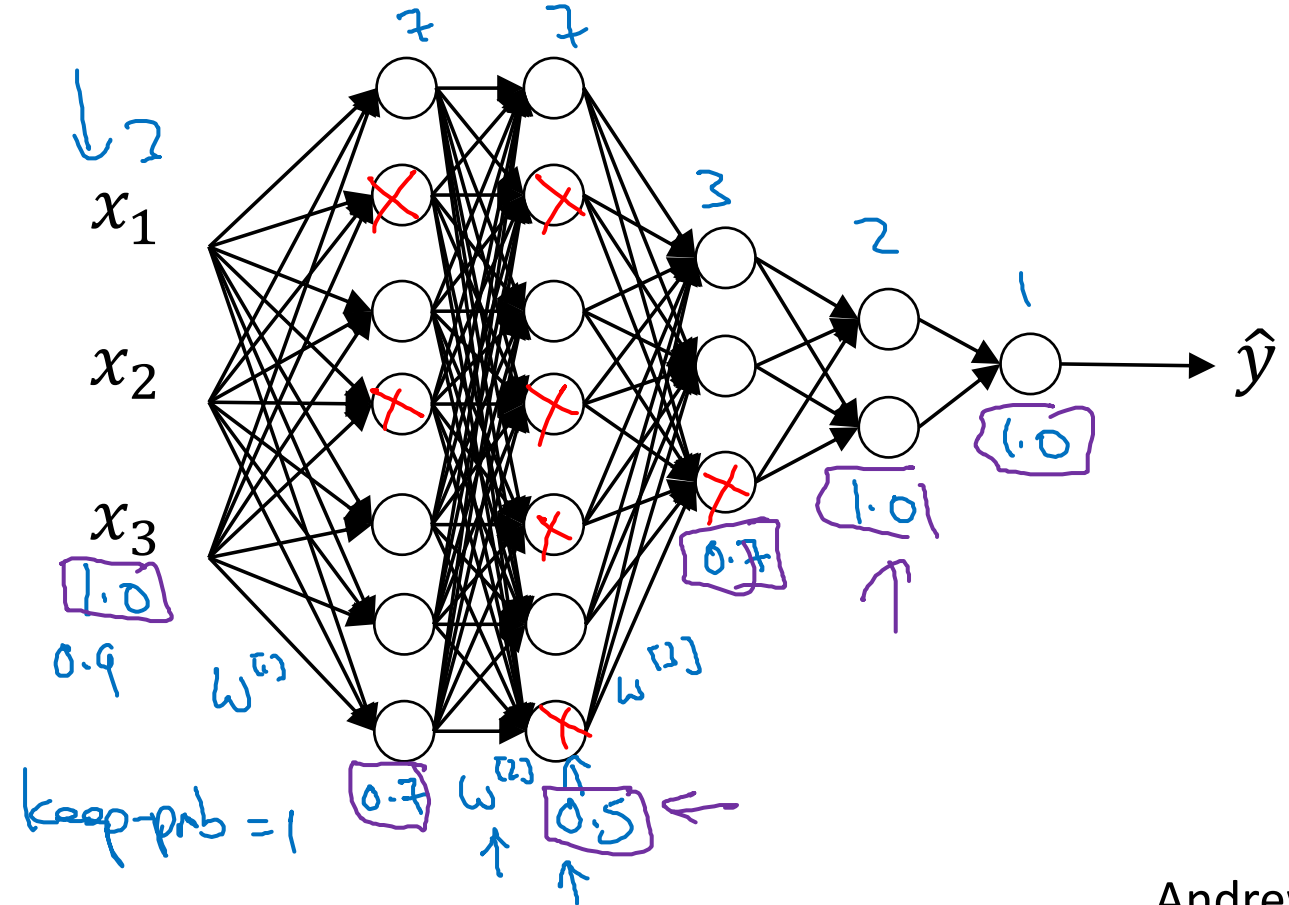
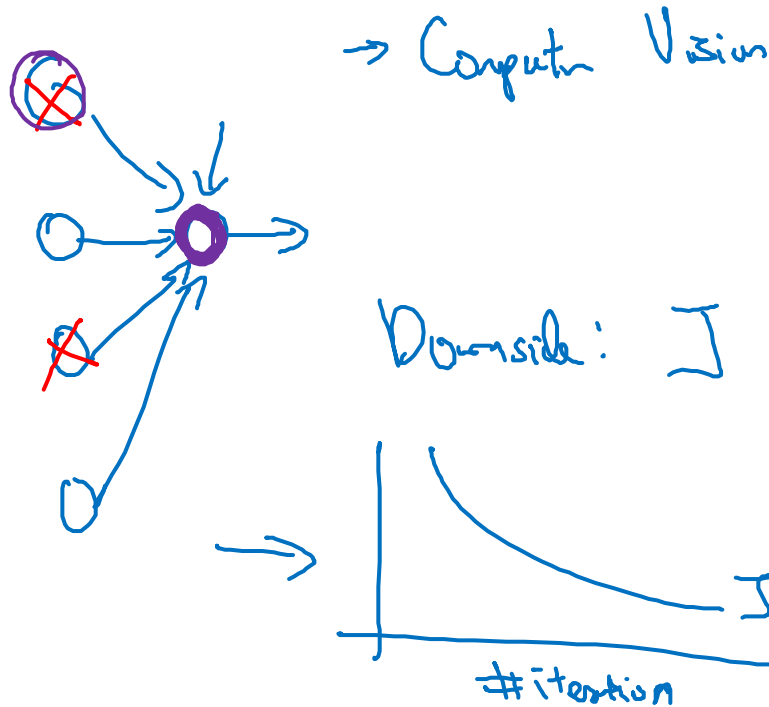
Understanding dropout

Why does drop-out work?

Jahan par overfitting ka dar lage, ya parameters jyada hon waha

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

→ Shrink weights. b_2





deeplearning.ai

Regularizing your neural network

Other regularization methods

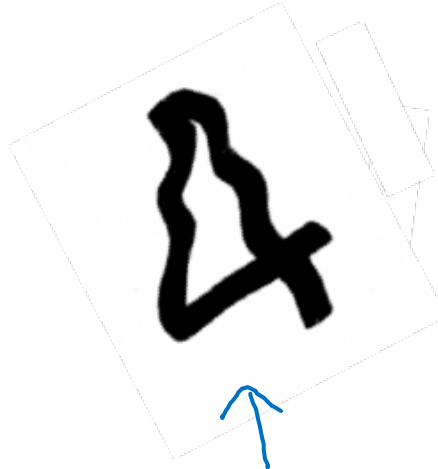
Data augmentation

Jo data hai usi ko ghuma phira ke naya data bana do. Kahin se aur dataa mil jaay

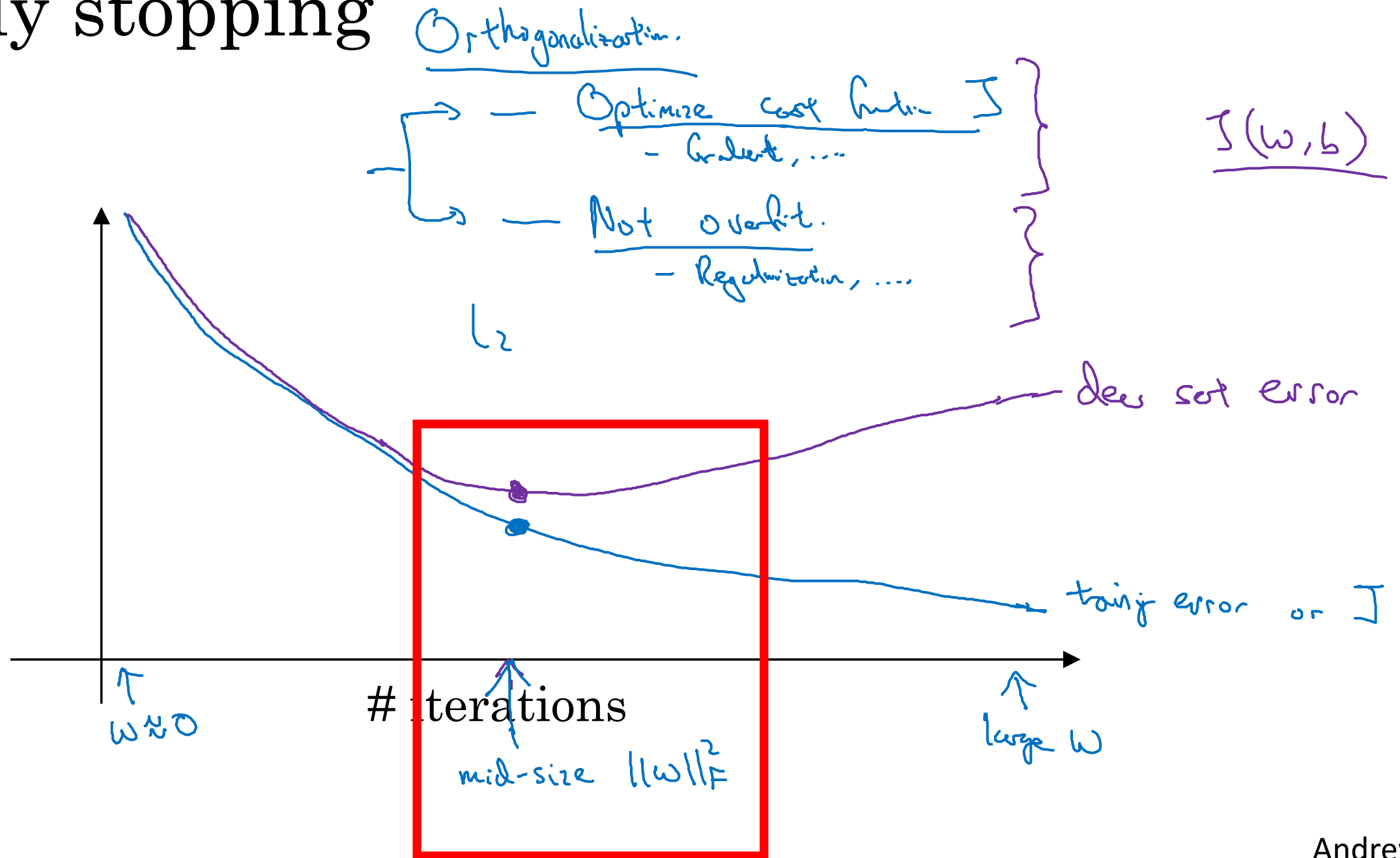


4

Augument karo



Early stopping





deeplearning.ai

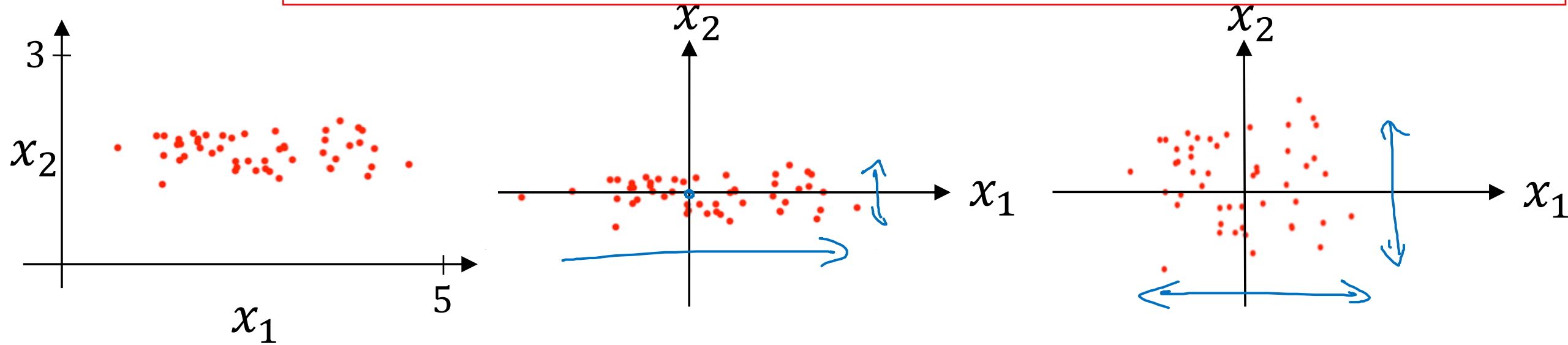
Setting up your
optimization problem

Normalizing inputs

Normalizing training sets

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

Normalizing involves two steps. The first is to subtract out means from every training set and then to normalize the variance



Subtract mean:

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

$$x := x - \mu$$

Normalize variance

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

\hookrightarrow element-wise

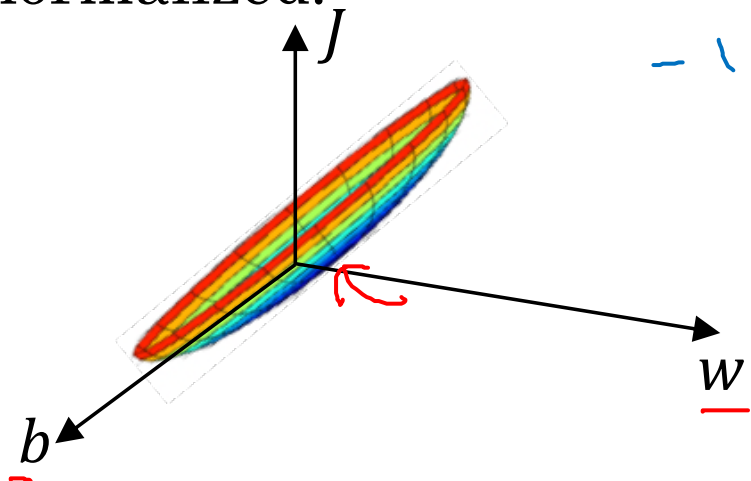
$$x / \sigma^2$$

Use same μ σ^2 to normalize test set.

Why normalize inputs?

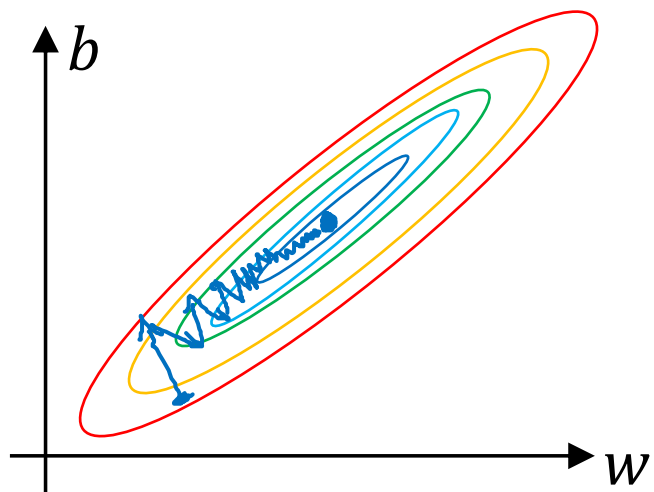
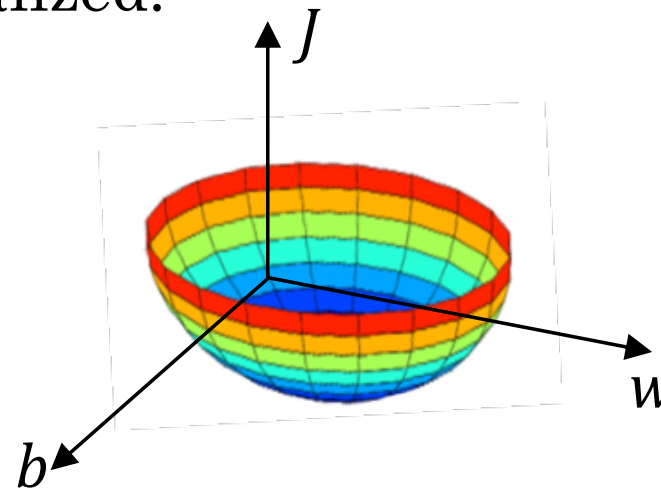
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Unnormalized:

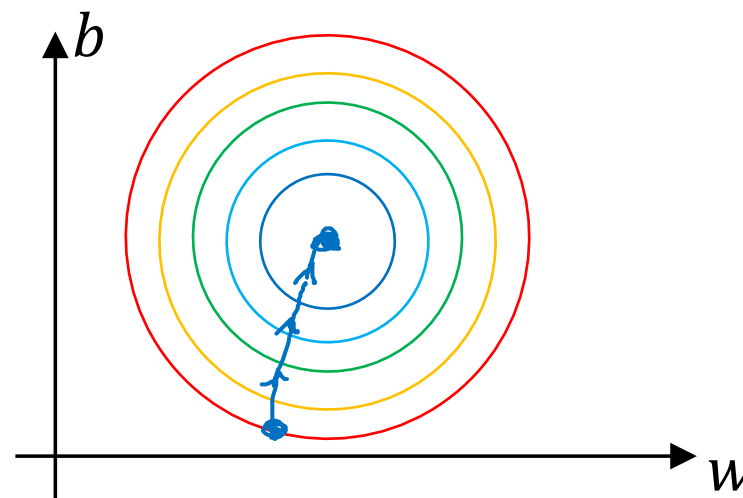


$w_1: x_1: \underline{1 \dots 1000} \leftarrow$
 $w_2: x_2: \underline{0 \dots 1} \leftarrow$
 $\quad \quad \quad -1 \dots 1$

Normalized:



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



What that means is that when you're training a very deep network your derivatives or your slopes can sometimes get either very, very big or very, very small

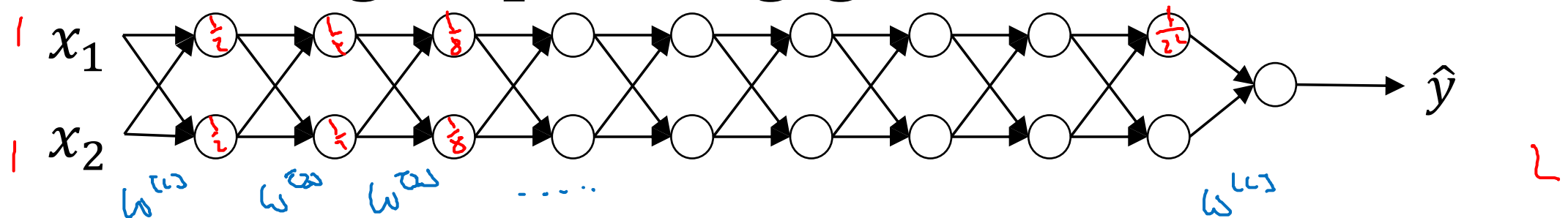


deeplearning.ai

Setting up your optimization problem

Vanishing/exploding gradients

Vanishing/exploding gradients



$g(z) = z$ $b^{(L)} = 0$

$\hat{y} = w^{(L,2)} \left(w^{(L-1,2)} w^{(L-2,2)} \dots w^{(2,2)} w^{(1,2)} x \right)$

Diagram illustrating the forward pass of the RNN. The output \hat{y} is calculated as the product of the weights $w^{(L,2)}$, $w^{(L-1,2)}$, $w^{(L-2,2)}$, ..., $w^{(2,2)}$, $w^{(1,2)}$ and the input x .

$w^{(1,2)} > I$

$w^{(2,2)} < I$ $\begin{bmatrix} 0.9 & \\ & 0.9 \end{bmatrix}$

$w^{(2,2)} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$

Diagram illustrating the weight matrix $w^{(2,2)}$ for the second layer. The matrix is a 2x2 identity matrix scaled by 1.5.

$z^{(1)} = w^{(1,2)} x$

$a^{(1)} = g(z^{(1)}) = z^{(1)}$

$a^{(2)} = g(z^{(2)}) = g(w^{(2,2)} a^{(1)})$

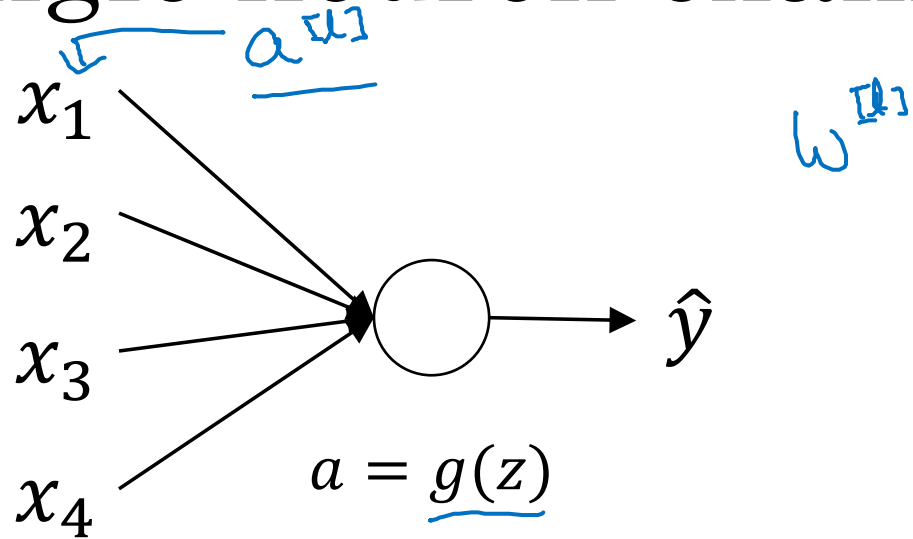
$\hat{y} = w^{(L,2)} \left[\begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} x \right]$

Diagram illustrating the forward pass of the RNN. The output \hat{y} is calculated as the product of the weights $w^{(L,2)}$ and the input x raised to the power of $L-1$.

$1.5^{L-1} x$

$0.5^{L-1} x$

Single neuron example



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

large $n \rightarrow$ Smaller w_i

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

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

ReLU $g^{[1]}(z) = \text{ReLU}(z)$

Other variants:

tanh

$$\frac{1}{n^{[1-1]}}$$

Xavier initialization

$$\frac{2}{n^{[1-1]} + n^{[1]}}$$

for other than relu, numerator is 1



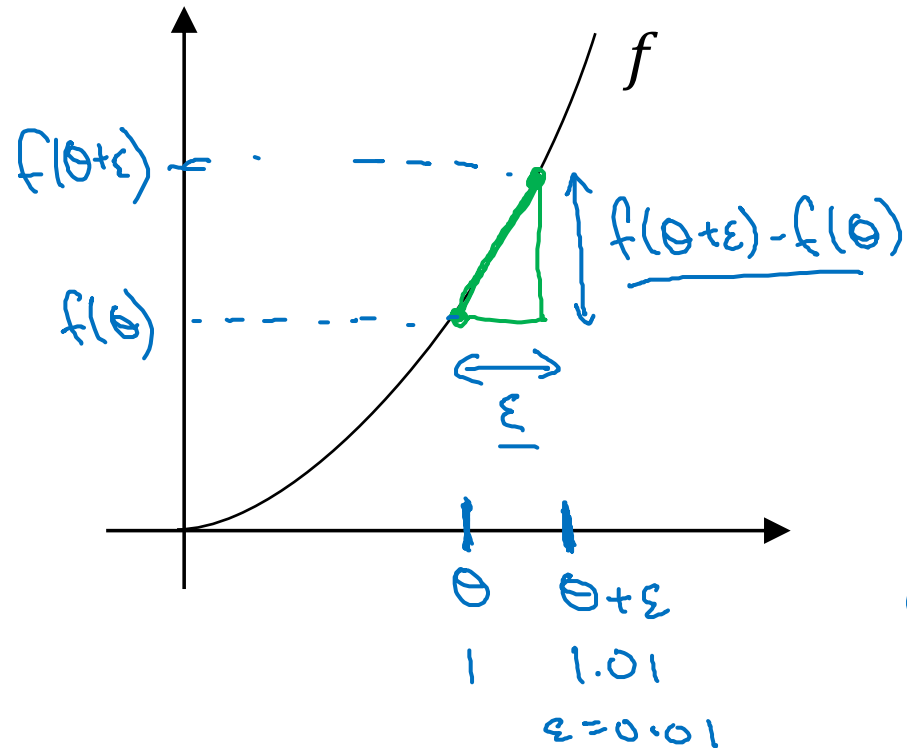
deeplearning.ai

Setting up your optimization problem

**Numerical approximation
of gradients**

Checking your derivative computation

I $f(\theta) = \theta^3$
 $\theta \in \mathbb{R}.$



$$g(\theta) = \frac{d}{d\theta} f(\theta) = f'(\theta)$$

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

$\frac{dw}{db}$

$g(\theta) = 3 \cdot (1)^2 = 3$
 when $\theta = 1$

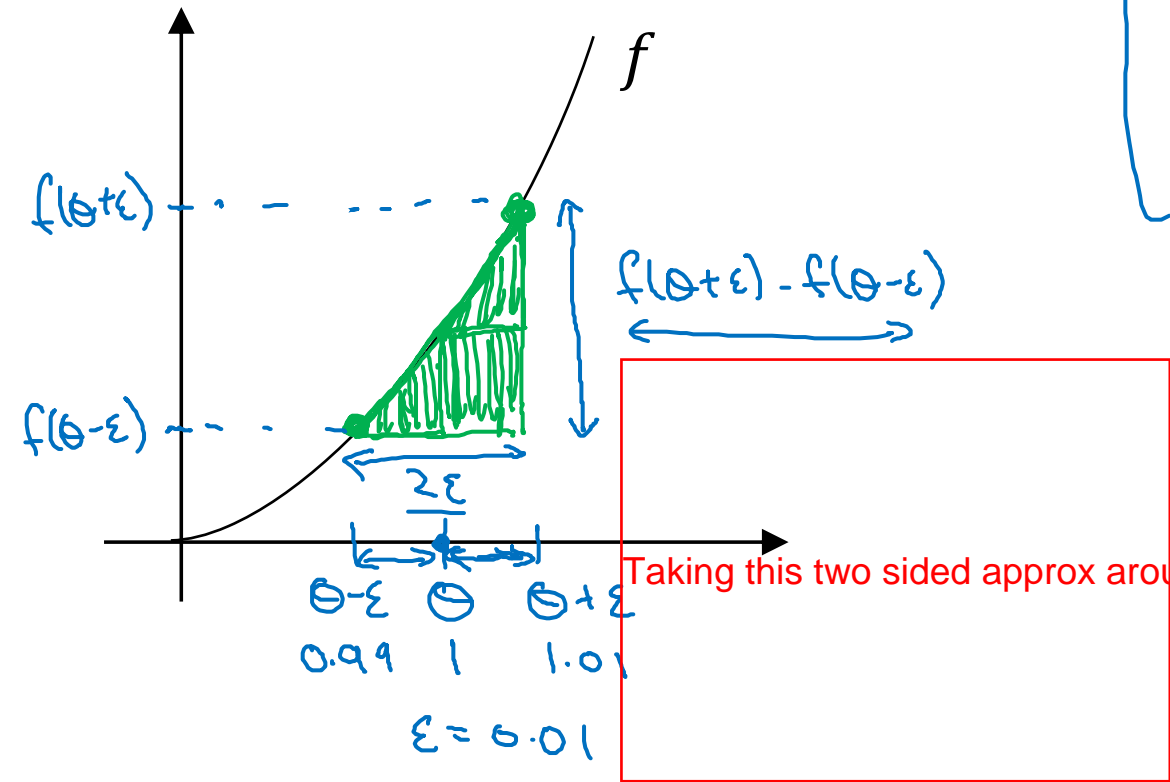
$$\frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \approx g(\theta)$$

$$\frac{(1.01)^3 - 1^3}{0.01} = \frac{1.030301 - 1}{0.01} = \frac{0.0301}{0.01} = 3.0301 \approx 3$$

Annotations: $\theta = 1$, $\theta + \epsilon = 1.01$, $\epsilon = 0.01$, 3.1 , 3.2

Checking your derivative computation

$$\underline{f(\theta) = \theta^3}$$



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

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

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

$$\text{approx error: } 0.0001$$

$$(\text{prev slide: } 3.0301, \text{ error: } 0.03)$$

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



deeplearning.ai

Setting up your optimization problem

Gradient Checking

To check correctness of gradient descent and to find bug in back propo case it is useful.

Gradient check for a neural network

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

concatenate

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

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

concatenate

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

Gradient checking (Grad check)

$$J(\theta) = J(\theta_1, \theta_2, \theta_3, \dots)$$

for each i :

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

Only change θ_i

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

Checks

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

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

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

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



deeplearning.ai

Setting up your
optimization problem

Gradient Checking
implementation notes

Gradient checking implementation notes

- Don't use in training – only to debug

$$\frac{d\theta_{\text{approx}}[\vec{i}]}{\uparrow \uparrow} \longleftrightarrow \frac{d\theta[\vec{i}]}{\uparrow}$$

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

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

- Remember regularization.

$$\underline{J(\theta)} = \frac{1}{n} \sum_i \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2n} \sum_l \|W^{[L,l]}\|_F^2$$

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

- Doesn't work with dropout.

J

$$\underline{\text{keep-prob} = 1.0}$$

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

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