# 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

# Applied ML is a highly iterative process

# layers

# hidden units

learning rates

activation functions

...

Idea

Experiment

Code
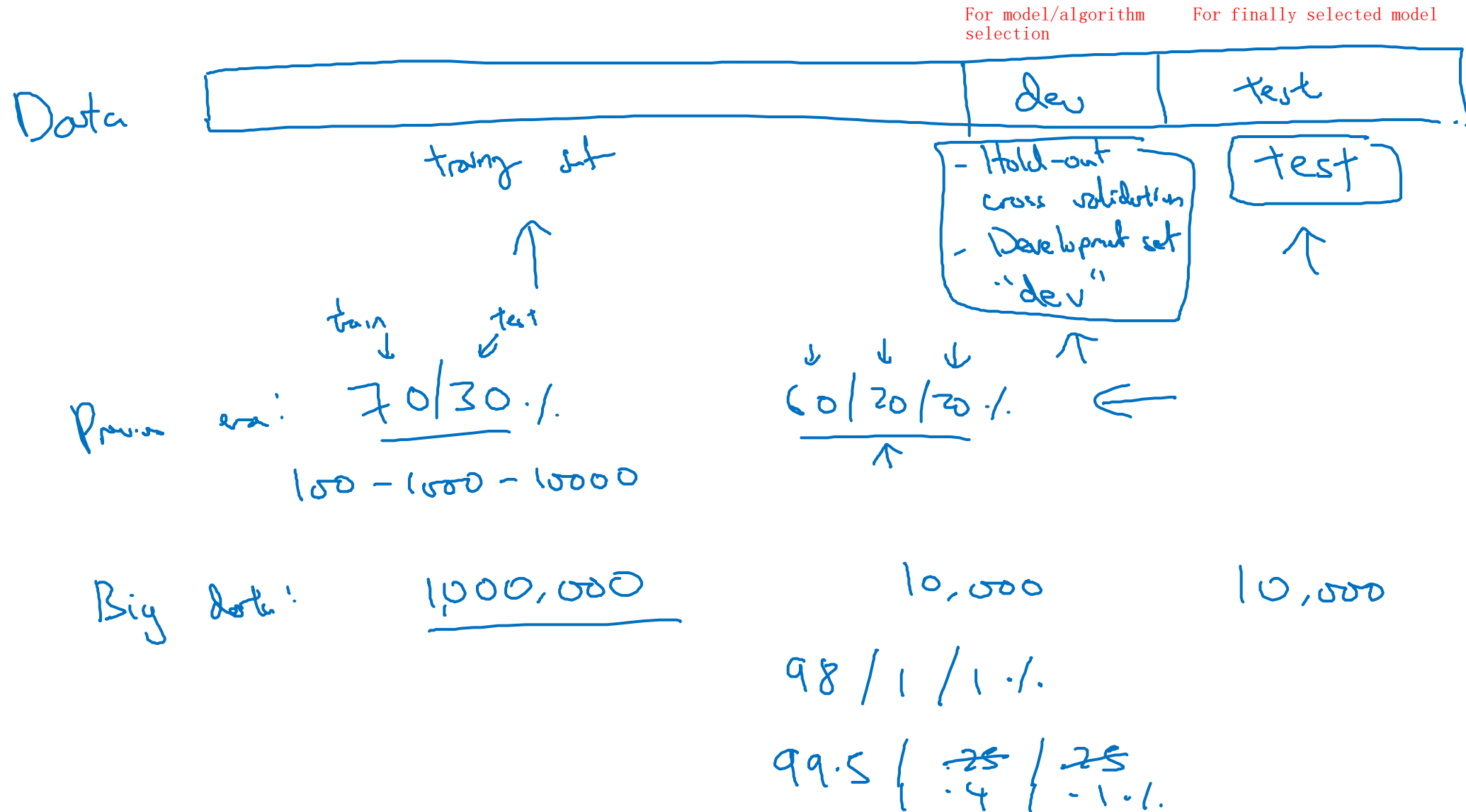
NLP, Vision, Speech, Structural data

Ads    Search    Security    logistic ....

# Train/dev/test sets

# Mismatched train/test distribution

Certs

Training set:
Cat pictures from webpages

Dev/test sets:
Cat pictures from users using your app

→ Make sure dev and test come from Same distribution.

"test"
train / dev

train / test
→ Train / dev
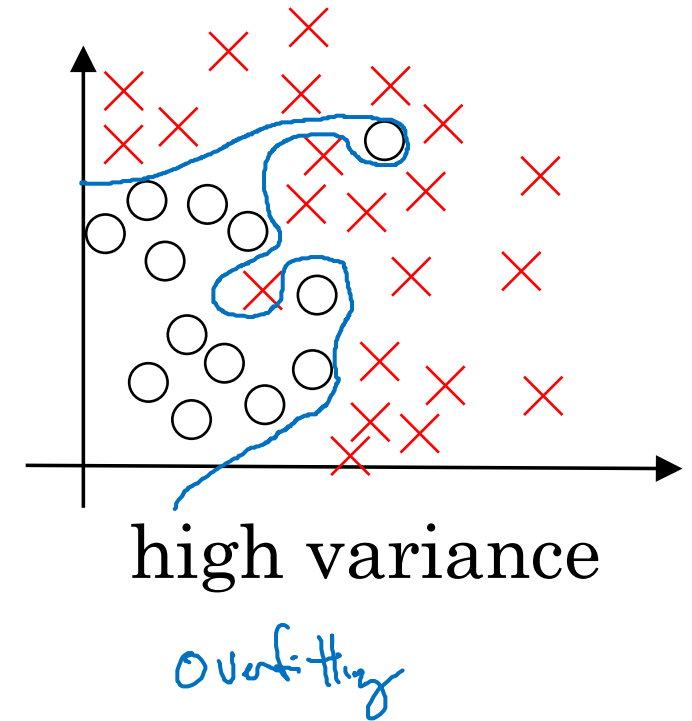
Not having a test set might be okay. (Only dev set.)

Andrew Ng
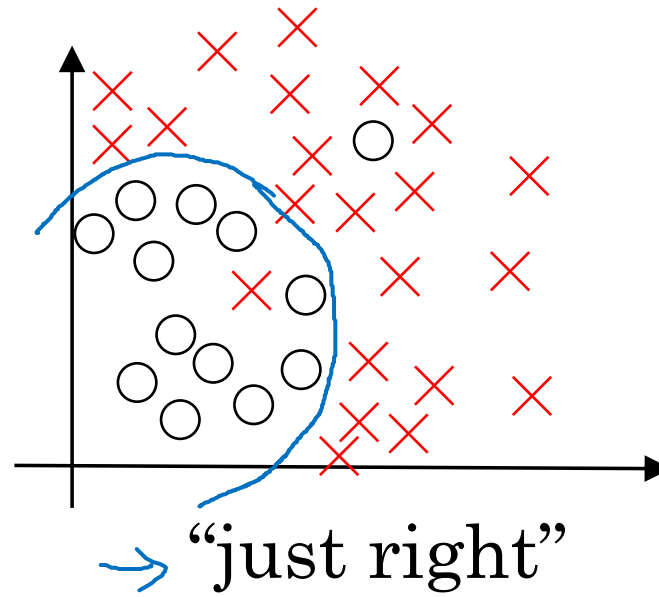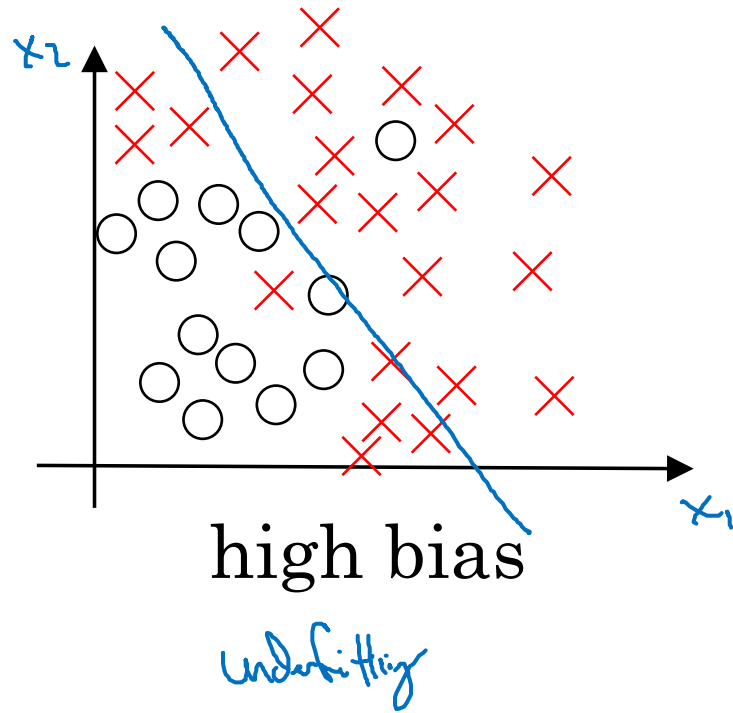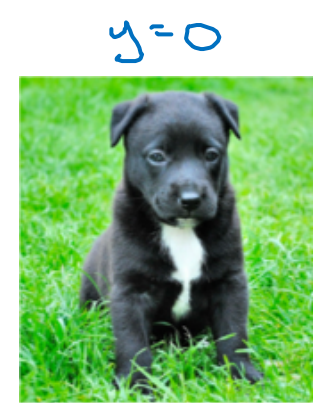
Setting up your
ML application

Bias/Variance

deeplearning.ai

# Bias and Variance



high bias

Underfitting

"just right"

high variance

Overfitting

Andrew Ng

# Bias and Variance

Cat classification

$y=1$     $y=0$



Train set error:   $1\%$    $15\%$    $15\%$    $0.5\%$

Dev set error:   $11\%$    $16\%$    $30\%$    $1\%$

high variance    high bias    high bias & high varian    low bias low variance

Human: $\approx 0\%$

Optiml (Bayes) error: $\approx 0\%$ to $15\%$    Blury images

Andrew Ng

# High bias and high variance



high bias
high varian

deeplearning.ai

Setting up your
ML application
___

Basic "recipe"
for machine learning

# Basic recipe for machine learning

High bias?
(training data performance)

Bigger network

Train longer.

(NN architecture search)

↓ N

High variance?
(dev set performance)

More data

Regularization

(NN architecture search)

↓ N

Done

Bias    Variance    tradeoff "

Andrew Ng

deeplearning.ai

# Regularizing your neural network

---

# Regularization

# Logistic regression

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

$$\lambda = \text{regularization parameter}$$
lambda     lambd

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

$$J(w,b) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}\left(\hat{y}^{(i)}, y^{(i)}\right) + \frac{\lambda}{2m} \|w\|_2^2 \quad + \frac{\lambda}{2m} b^2$$

omit

just a single number

common used

$L_2$ regularization

$$\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w \longleftarrow$$

$L_1$ regularization

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

$w$ will be sparse

compress the model

Andrew Ng

# Neural network

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

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

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

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

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

"Weight decay"

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

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

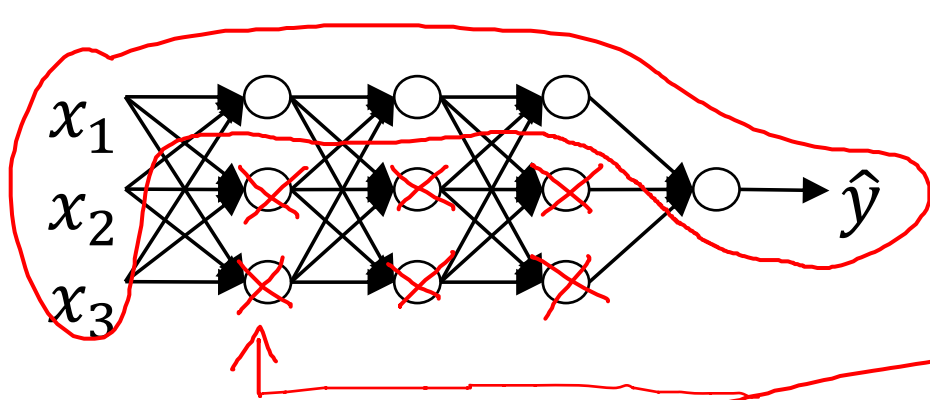$$= \underbrace{\left(1 - \frac{\alpha \lambda}{m}\right)}_{< 1} w^{[l]} - \alpha (\text{from backprop})$$

Andrew Ng

# Regularizing your neural network

---

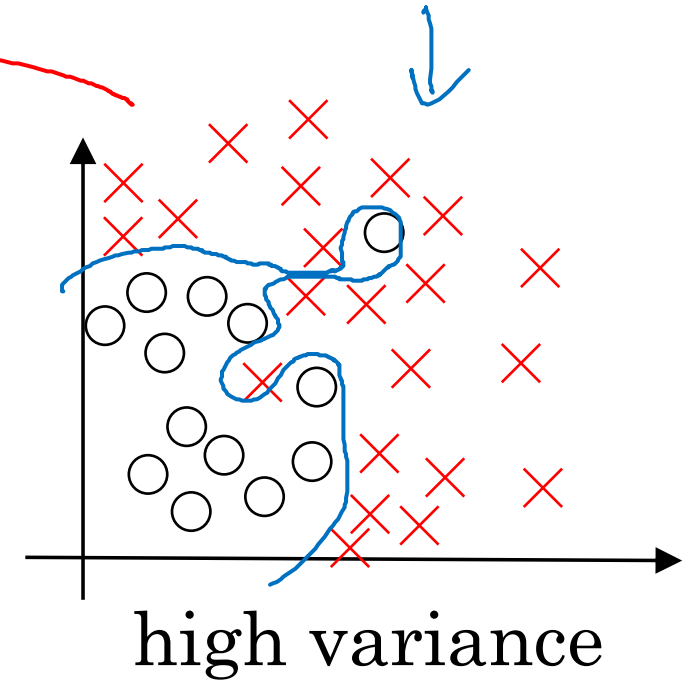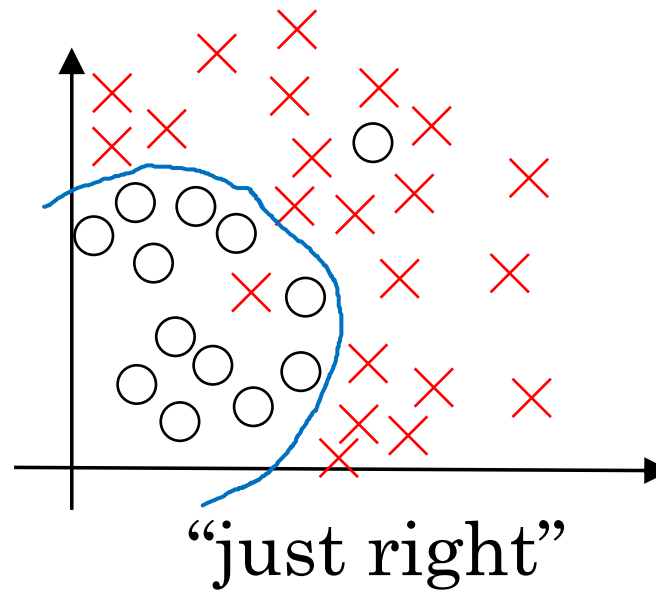# Why regularization reduces overfitting
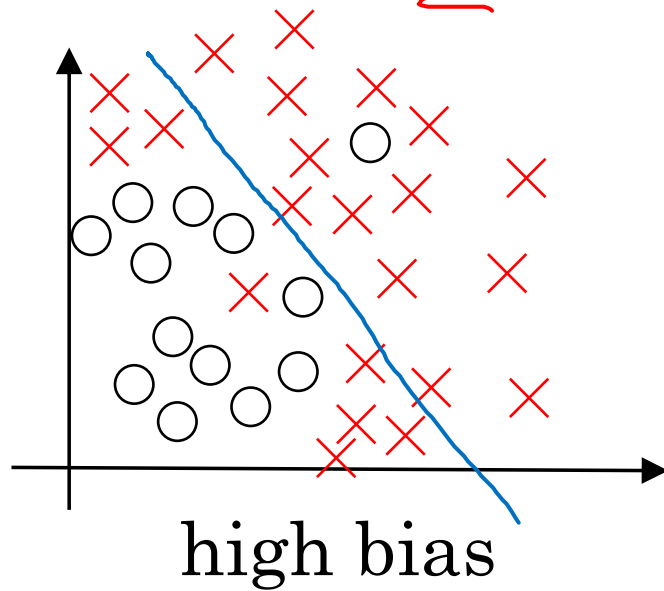
deeplearning.ai

# How does regularization prevent overfitting?
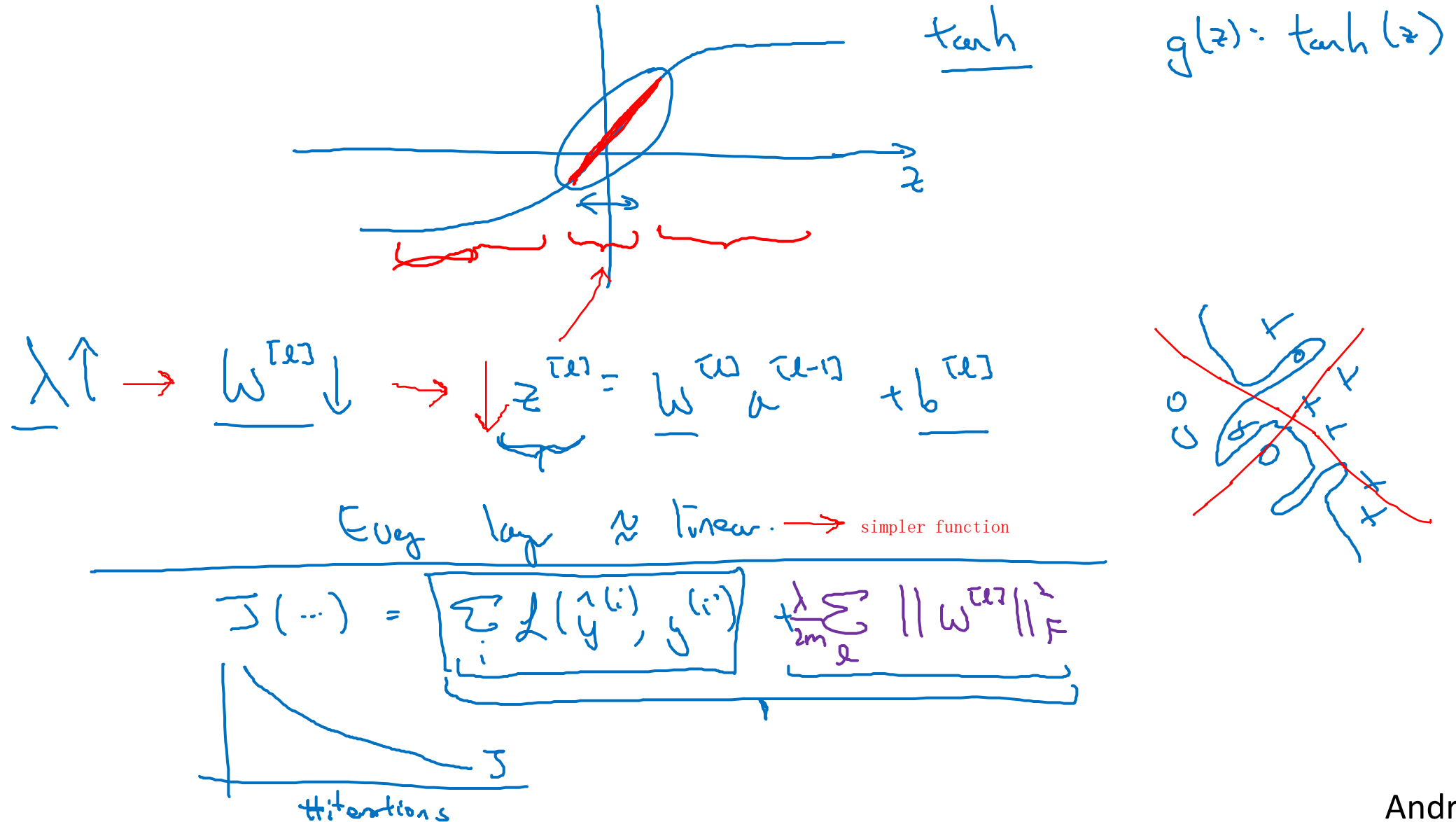
$$J(w^{[l]}, b^{[l]}) = \frac{1}{m} \sum_{i=1}^{n} \mathcal{L}(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^{L} \| w^{[l]} \|_F^2$$

penalize large weights

$\downarrow \| w^{[l]} \| \approx 0$

high bias

"just right"

high variance

Andrew Ng

# How does regularization prevent overfitting?

$tanh$

$g(z) = tanh(z)$



$\lambda \uparrow \rightarrow W^{[l]} \downarrow \rightarrow \downarrow z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$

Every layer $\approx$ linear. $\rightarrow$ simpler function

$$J(\cdots) = \sum_i \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_l \|W^{[l]}\|_F^2$$

Andrew Ng

Regularizing your neural network

Dropout regularization

deeplearning.ai

# Dropout regularization



$x_1$   $x_2$   $x_3$   $x_4$   $\hat{y}$

$x_1$   $x_2$   $x_3$   $x_4$   $\hat{y}$

0.5   0.5   0.5

Andrew Ng

# Implementing dropout ("Inverted dropout")

Illustrate with layer $l = 3$.      keep-prob = 0.8      0.2

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

a3 = np.multiply(a3, d3)      # a3 *= d3.

$\rightarrow$ $\boxed{a3 \;/=\; \cancel{0.8}\; \text{keep-prob}}$ $\leftarrow$

50 units. $\leadsto$ 10 units shut off

$z^{[4]} = w^{[4]} \cdot a^{[3]} + b^{[4]}$

$\uparrow$ reduced by 20%.

$/= 0.8$

Test

Andrew Ng

# Making predictions at test time

$$a^{[0]} = X$$

No drop out. <span style="color:red">don't want output to be random</span>

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

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

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = \dots$$

$$\downarrow$$

$$\hat{y}$$
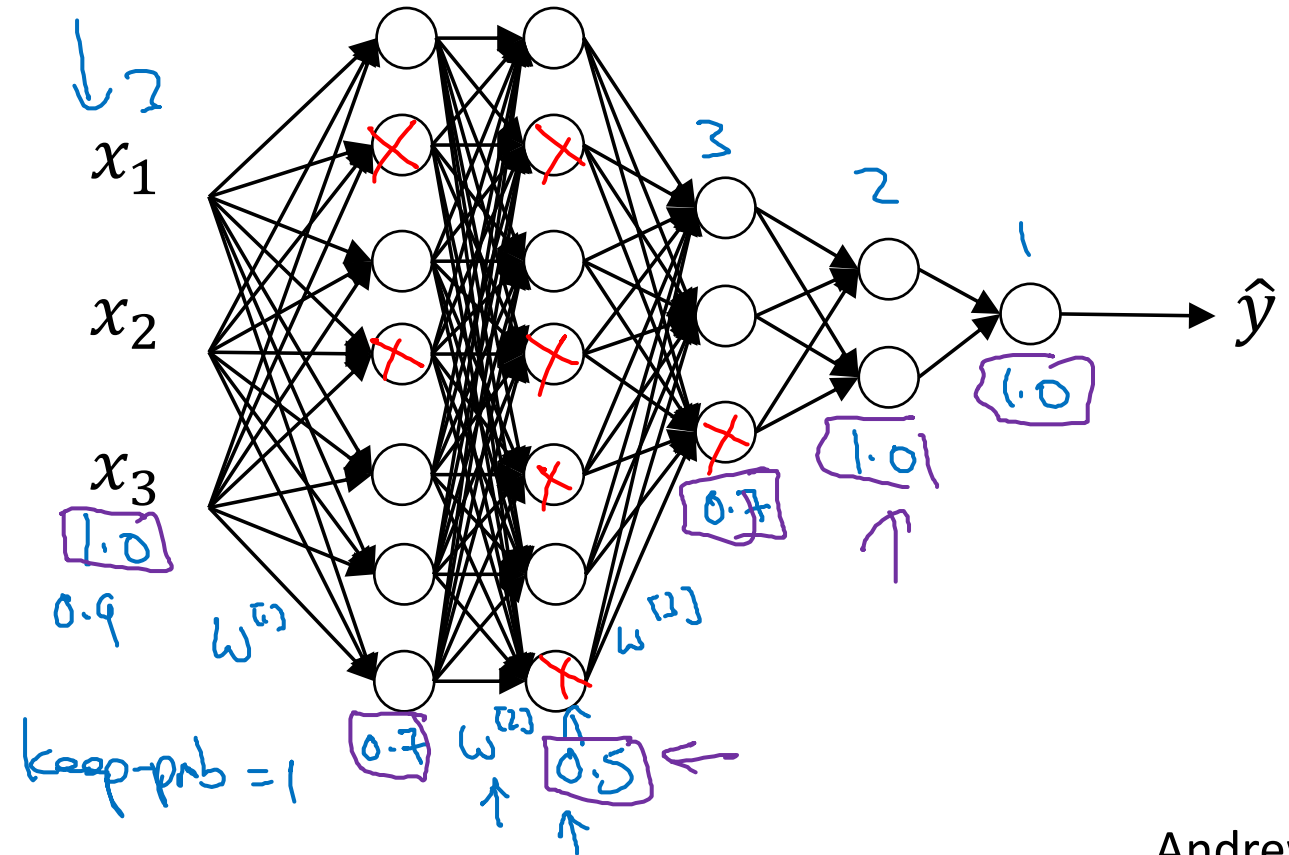
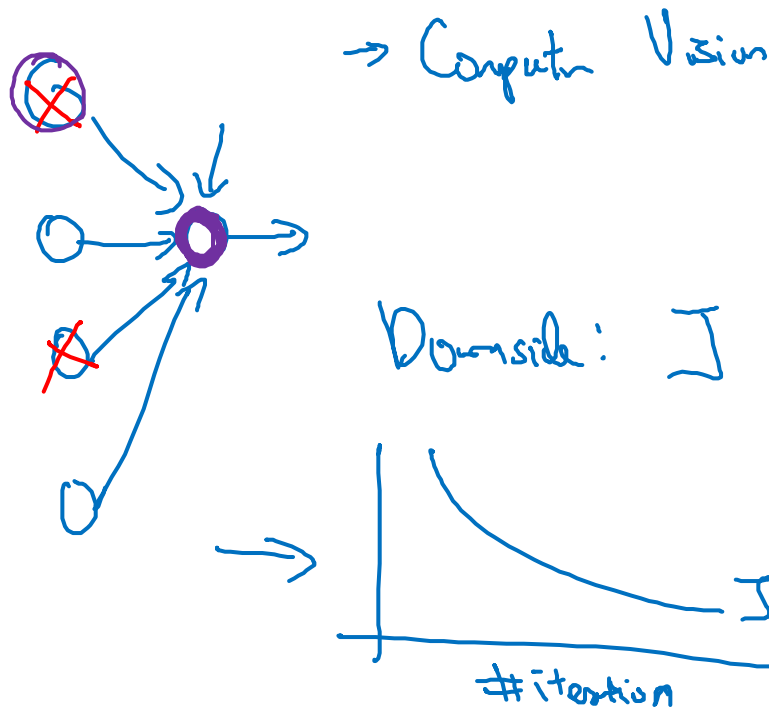$$/ = \text{keep-prob}$$

Andrew Ng

deeplearning.ai

Regularizing your
neural network

Understanding
dropout

# Why does drop-out work?

Intuition: Can't rely on any one feature, so have to spread out weights. ⟶ Shrink weights. $L_2$

norm of



Andrew Ng

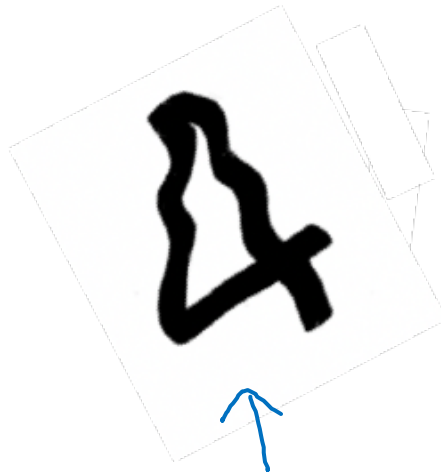# Data augmentation

# Early stopping

Orthogonalization.

Separate tasks and tools
for these two tasks,
easier to decompose and search over
(Early stopping can't)

$\to$ — Optimize cost function $J$
 — Gradient, ....

$\to$ — Not overfit.
 — Regularization, ....

$l_2$

$J(w,b)$

dev set error

training error or $J$

# iterations

$w \approx 0$

mid-size $\|w\|_F^2$

large $W$

Andrew Ng

Setting up your
optimization problem
_____

Normalizing inputs

deeplearning.ai

# 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

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

↖ element-wise

$$x /= \sigma^2$$

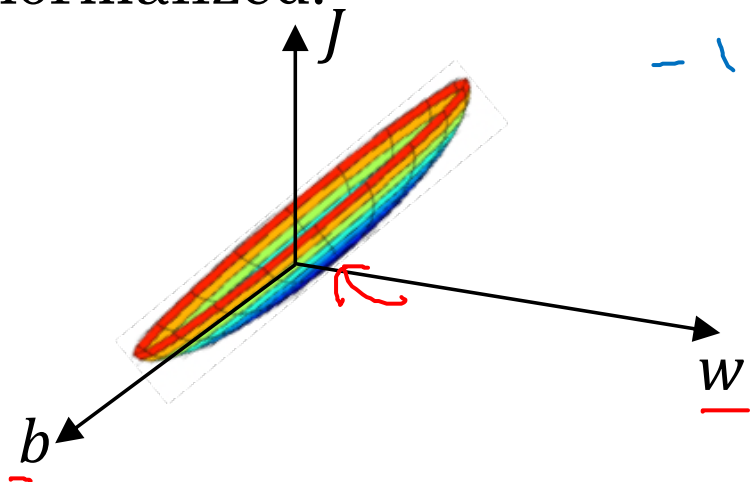Use **same** $\mu, \sigma^2$ to normalize test set.

Andrew Ng

# Why normalize inputs?

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

$w_1 \quad x_1: \quad \underline{1 \cdots 1000} \leftarrow$

$w_2 \quad x_2: \quad \underline{0 \cdots 1} \leftarrow$

$-1 \cdots 1$

Unnormalized:

Normalized:
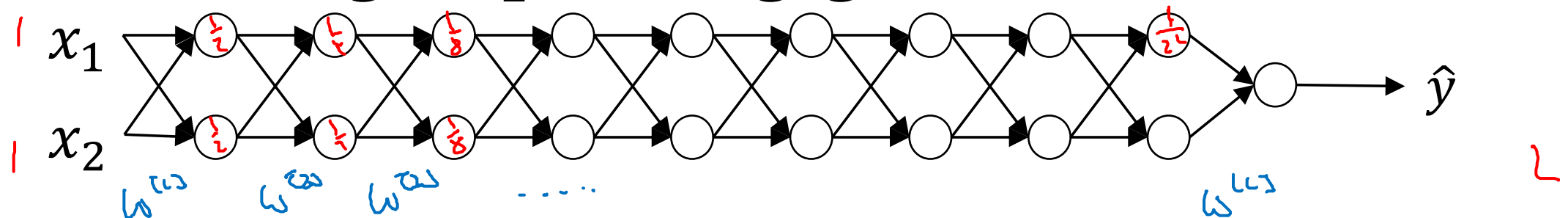
$x_1: 0 \cdots 1$

$x_2: -1 \cdots 1$

$x_3: 1 \cdots 2$



Andrew Ng

deeplearning.ai

Setting up your
optimization problem

Vanishing/exploding
gradients

# Vanishing/exploding gradients

$L = 150$



$1$   $x_1$    $\frac{1}{2}$   $\frac{1}{4}$   $\frac{1}{8}$       $\frac{1}{2^L}$

$1$   $x_2$    $\frac{1}{2}$   $\frac{1}{4}$   $\frac{1}{8}$       $\hat{y}$

$W^{[1]}$   $W^{[2]}$   $W^{[3]}$   $\cdots$   $W^{[L]}$   $2$

$g(z) = z.$    $b^{[l]} = 0.$

$a^{[3]}$    $1.5^L$

$\hat{y} = W^{[L]} \; \boxed{W^{[L-1]}} \; \boxed{W^{[L-2]}} \; \boxed{\cdots} \; W^{[3]} \; W^{[2]} \; W^{[1]} \; x$    $0.5^L$

$z^{[1]} = W^{[1]} x$

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

$W^{[l]} > I$

$W^{[l]} < I$   $\begin{bmatrix} 0.9 & \\ & 0.9 \end{bmatrix}$

$a^{[2]} = g(z^{[2]}) = g(W^{[2]} a^{[1]})$

$W^{[l]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$   0.5   0.5

$\hat{y} = W^{[L]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} x$   0.5   0.9   $1.5^{L-1} x$   $0.5^{L-1} x$

Andrew Ng

# Single neuron example

$a^{[1]}$

$x_1$

$x_2$

$x_3$

$x_4$

$\hat{y}$

$a = g(z)$

$W^{[l]}$

not too much bigger than 1 and
also not too much less than 1

$z = W_1 X_1 + W_2 X_2 + \cdots + W_n X_n + b$

Large $n$ $\rightarrow$ Smaller $W_i$

$Var(W_i) = \dfrac{1}{n} \dfrac{2}{n}$

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

ReLU

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

Other vononts:

tanh

$\dfrac{1}{n^{[l-1]}}$

Xavier initialization

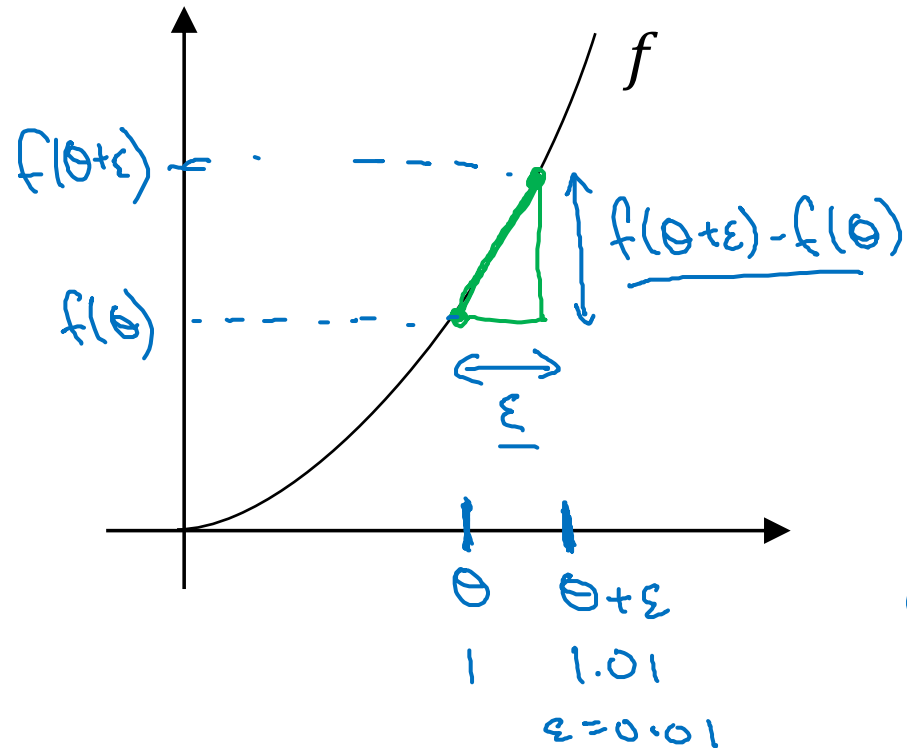$\sqrt{\dfrac{2}{n^{[l-1]} + n^{[l]}}}$

Andrew Ng

deeplearning.ai

Setting up your optimization problem

Numerical approximation of gradients

# Checking your derivative computation

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

$\theta \in \mathbb{R}.$

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

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

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

$\dfrac{dw}{db}$

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

$\dfrac{(1.01)^3 - 1^3}{0.01} = 3.0301$

$\approx 3$

$\theta = 1$

$\theta + \varepsilon = 1.01$

0.0301
3.1
3.2



$f(\theta+\varepsilon)$

$f(\theta)$

$f$

$f(\theta+\varepsilon) - f(\theta)$

$\varepsilon$

$\theta$    $\theta+\varepsilon$

1    1.01

$\varepsilon = 0.01$

Andrew Ng

# 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)

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

0.01
0.0001

0.01

Andrew Ng

Setting up your
optimization problem
_____

Gradient Checking

deeplearning.ai

# Gradient check for a neural network

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

Concatenate

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

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

Concatenate

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

# Gradient checking (Grad check)

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

for each $i$:

$$\rightarrow d\Theta_{approx}[i] = \frac{J(\theta_1, \theta_2, \ldots, \theta_i + \varepsilon, \ldots) - J(\theta_1, \theta_2, \ldots, \theta_i - \varepsilon, \ldots)}{2\varepsilon}$$

$$\approx \underline{d\Theta[i]} = \frac{\partial J}{\partial \theta_i} \qquad \Big| \qquad d\Theta_{approx} \overset{?}{\approx} d\Theta$$

$$\overset{i}{}$$

Check $\dfrac{\|d\Theta_{approx} - d\theta\|_2}{\|d\Theta_{approx}\|_2 + \|d\Theta\|_2}$

$$\rightarrow$$

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

$$\approx \boxed{10^{-7} \; - \; great!} \leftarrow$$

$$10^{-5}$$

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

Andrew Ng

deeplearning.ai

Setting up your
optimization problem

Gradient Checking
implementation notes

# Gradient checking implementation notes

- Don't use in training – only to debug

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

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

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

- Remember regularization.

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

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

- Doesn't work with dropout.

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

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

Doesn't work well when $\quad w, b \approx 0$

Andrew Ng