deeplearning.ai

Setting up your
ML application

Train/dev/test
sets

# Applied ML is a highly iterative process

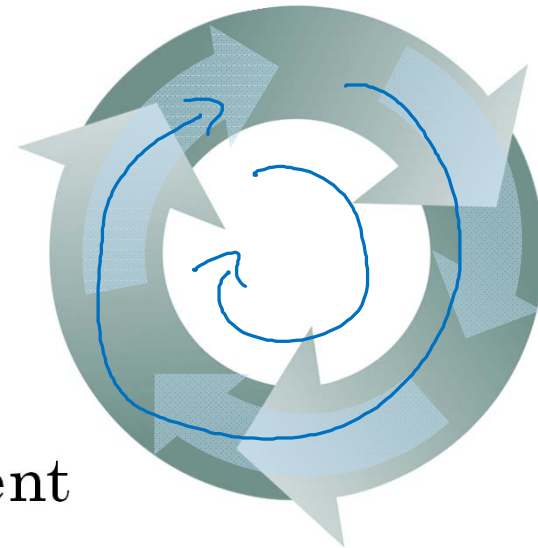# layers

# hidden units

learning rates
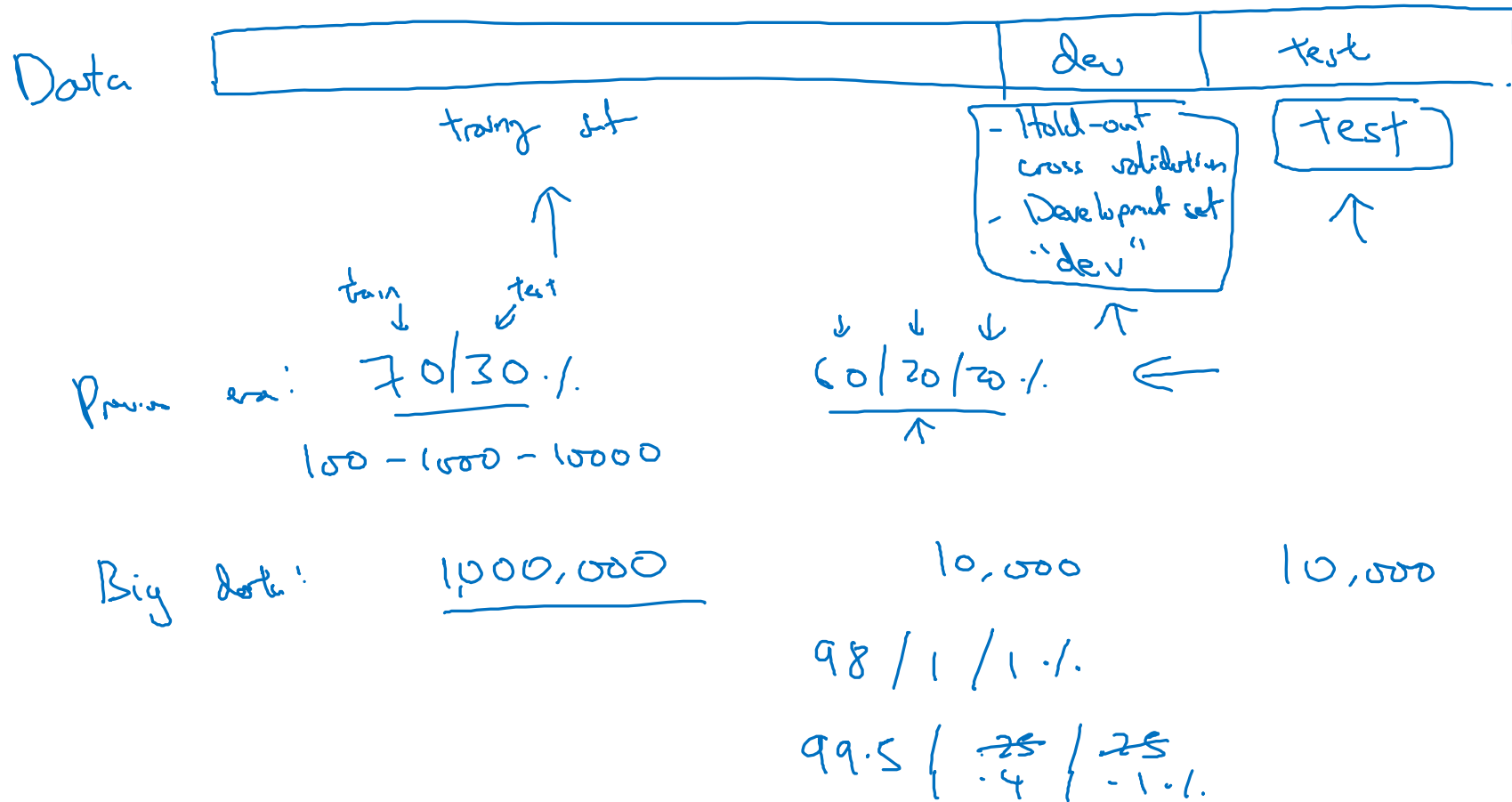
activation functions

...
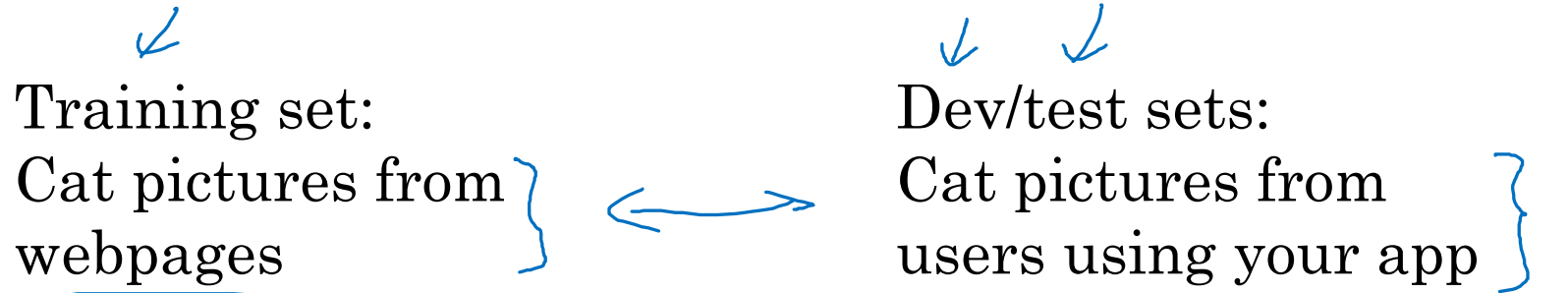


Idea

Code

Experiment

NLP, Vision, Speech, Structural data

Ads    Search    Security    logistic ....
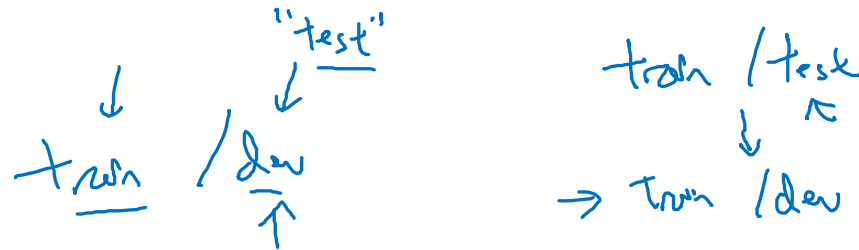
Andrew Ng

# Train/dev/test sets

Data

| | | dev | test |

training set

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

test

train    test

Previous era:   $70/30$ %.

$100 - 1000 - 10000$

$60/20/20$ %.

Big data:   $1,000,000$          $10,000$          $10,000$

$98/1/1$ %.

$99.5 / \cancel{.25}\ .4 / \cancel{.25}\ .1$ %.

Andrew Ng

# 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

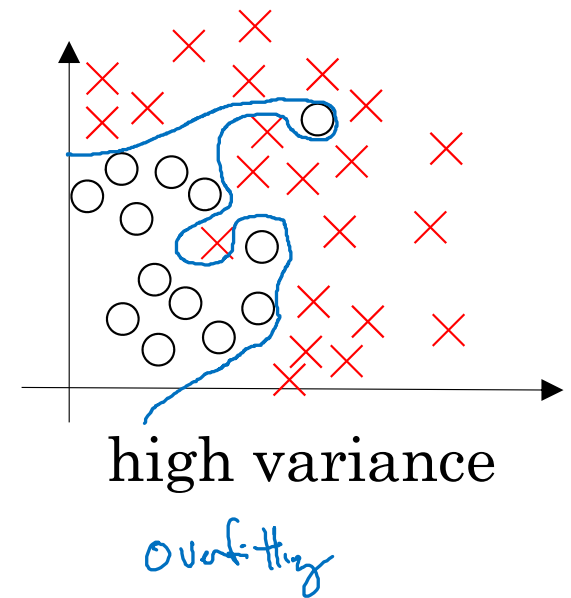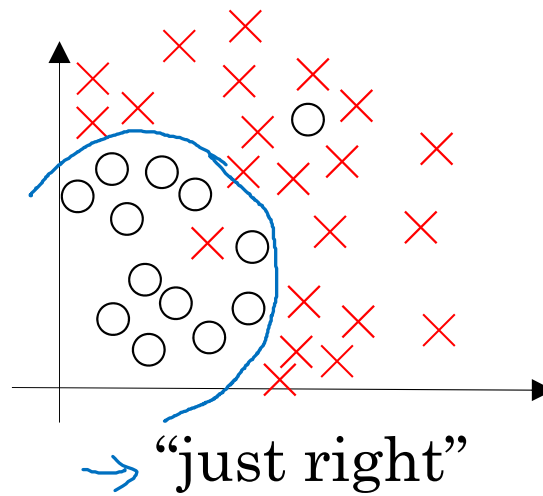Not having a test set might be okay. (Only dev set.)

Andrew Ng

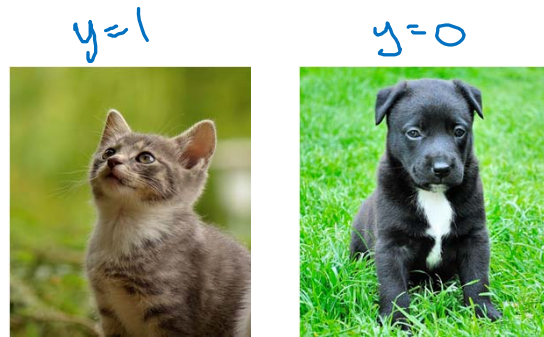deeplearning.ai

Setting up your
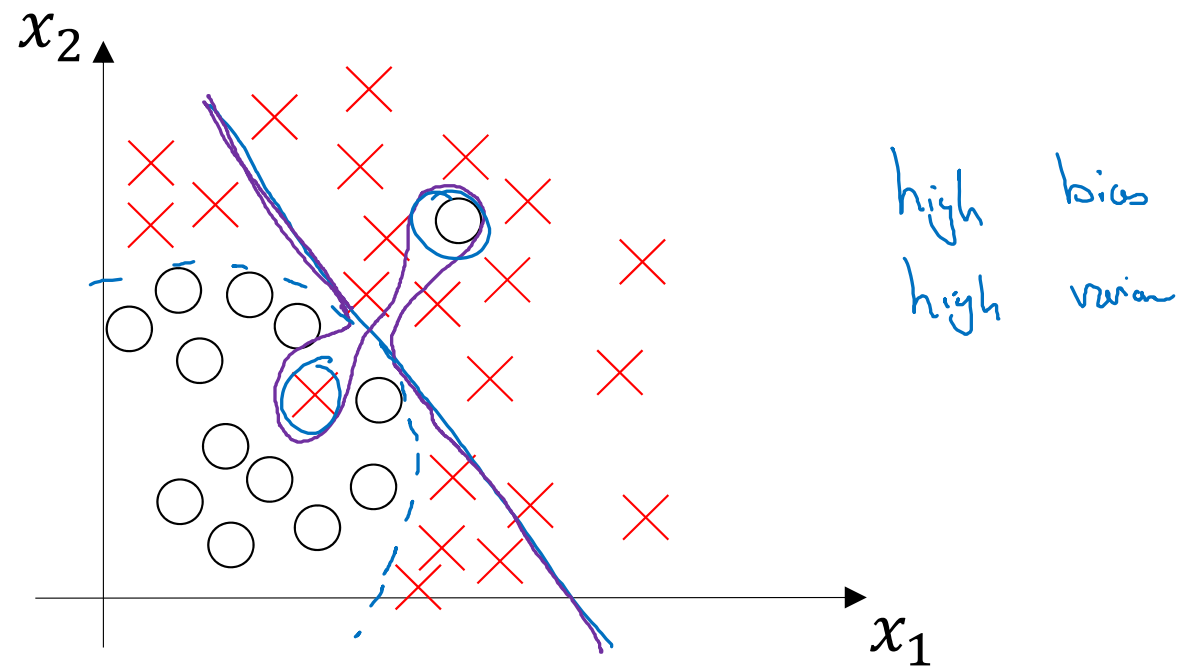ML application

Bias/Variance

# Bias and Variance



high bias

*Underfitting*

"just right"

high variance

*Overfitting*

Andrew Ng

# Bias and Variance

$y=1$        $y=0$



### Cat classification

| 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: ≈ 0%

Optimal (Bayes) error: ≈ 0 to 15%

Blurry images

Andrew Ng

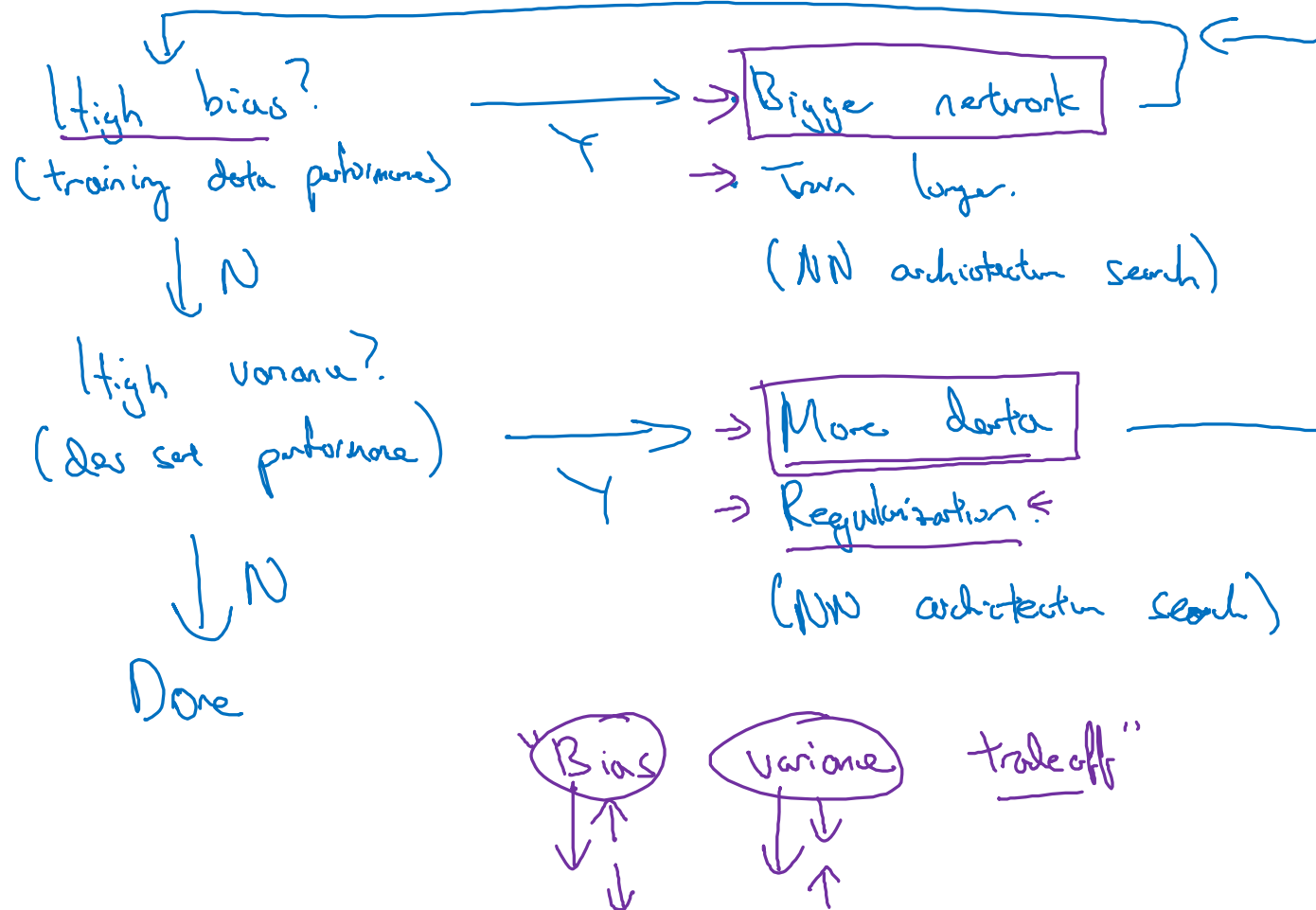# High bias and high variance



high bias
high varian

deeplearning.ai

Setting up your
ML application

Basic "recipe"
for machine learning

# Basic "recipe" for machine learning

Andrew Ng

# Basic recipe for machine learning

High bias?
(training data performance)

↓ N

High variance?
(dev set performance)

↓ N

Done

→ Bigger network
→ Train longer.
(NN architecture search)

→ More data
→ Regularization
(NN architecture search)

Bias  Variance  tradeoff"

Regularizing your
neural network

Regularization

deeplearning.ai

# Logistic regression

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

$\lambda = $ regularization parameter

lambda            lambd

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

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

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

Omit

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

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

$w$ will be sparse

Andrew Ng

# Neural network

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

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

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

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

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

"Weight decay"

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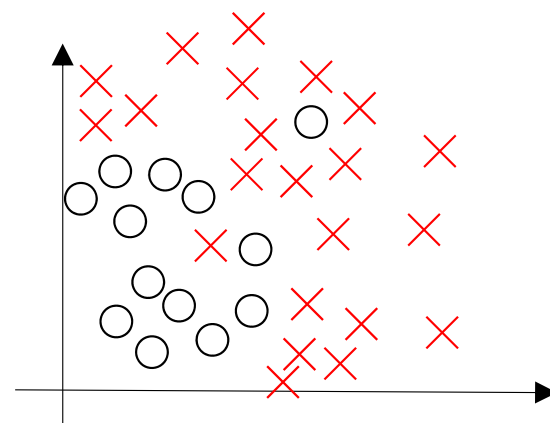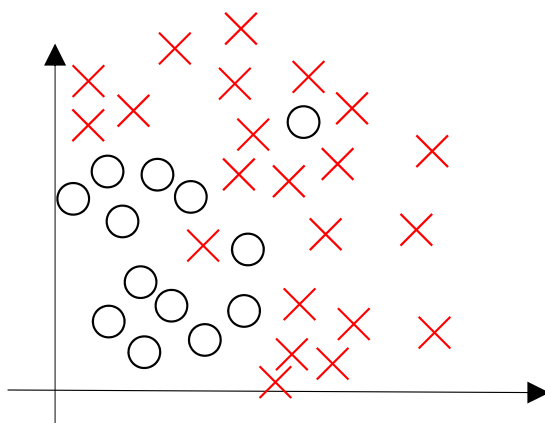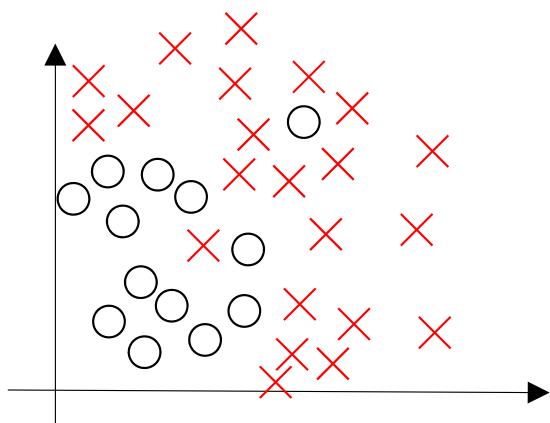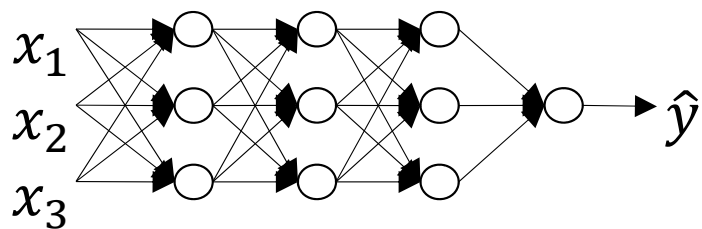
$$= W^{[l]} - \frac{\alpha \lambda}{m} W^{[l]} - \alpha (\text{from backprop})$$

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

Andrew Ng

# Neural network

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

# How does regularization prevent overfitting?

$x_1$

$x_2$

$x_3$

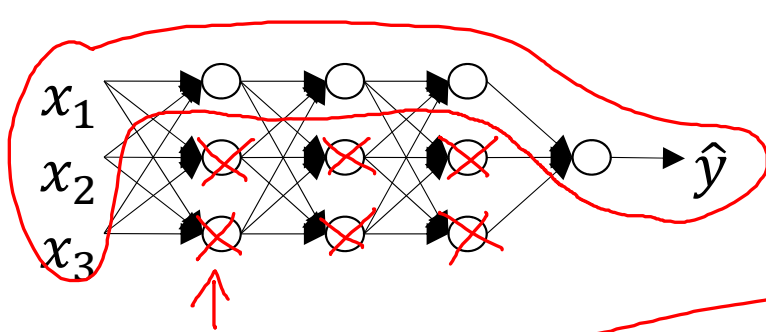$\hat{y}$

# How does regularization prevent overfitting?

deeplearning.ai

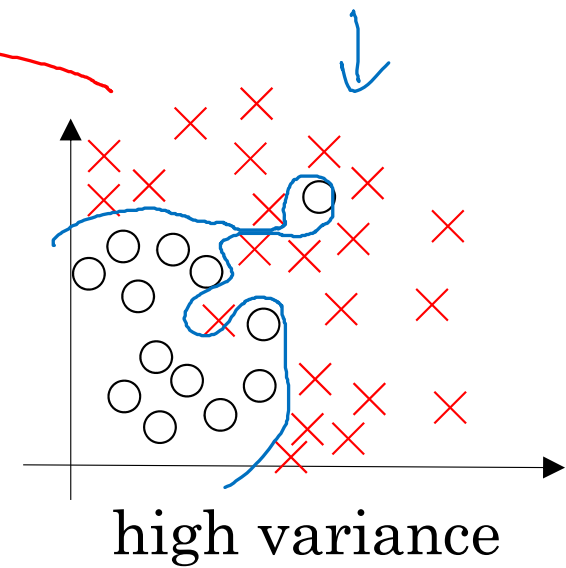Regularizing your
neural network

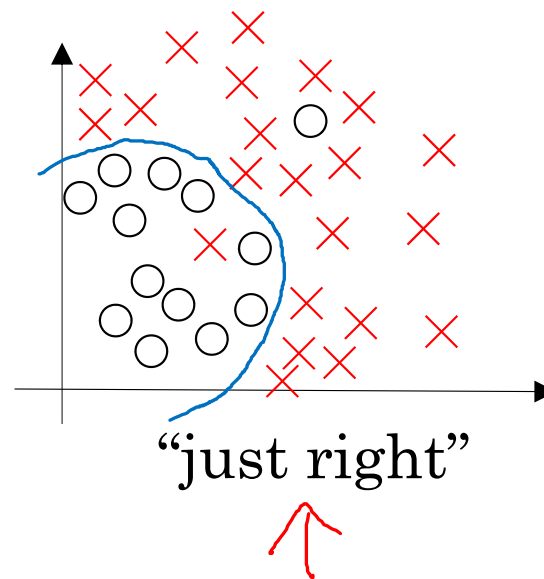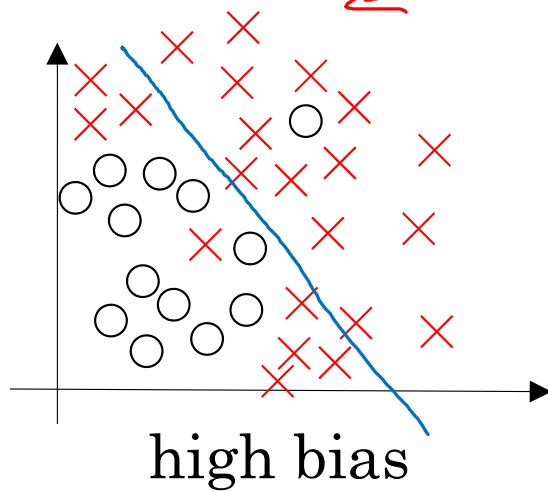Why regularization
reduces overfitting
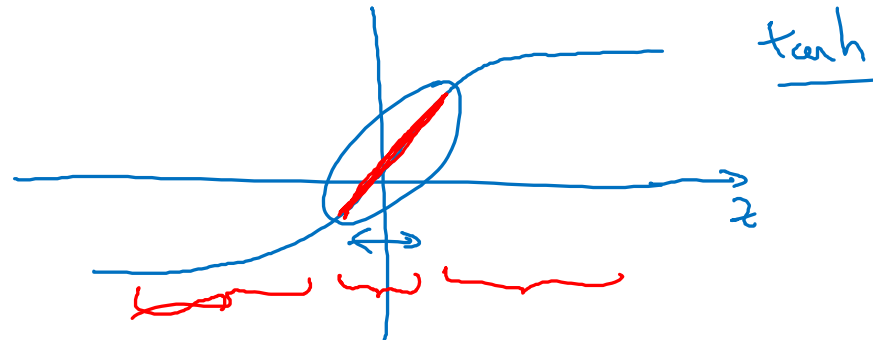
# How does regularization prevent overfitting?



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

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

high bias

"just right"

high variance

# How does regularization prevent overfitting?

$tanh$

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

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

Every layer $\approx$ linear.

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

$J$

#iterations

Andrew Ng

# Basic "recipe" for machine learning

Andrew Ng

# Basic recipe for machine learning



High bias?
(training data performance)

↓ N

High variance?
(dev set performance)

↓ N

Done

→ → Bigger network
→ Train longer.
(NN architecture search)

→ → More data
→ Regularization ←
(NN architecture search)

Bias  Variance  tradeoff"

Andrew
Ng

deeplearning.ai

Regularizing your
neural network

---

Dropout
regularization

# Dropout regularization



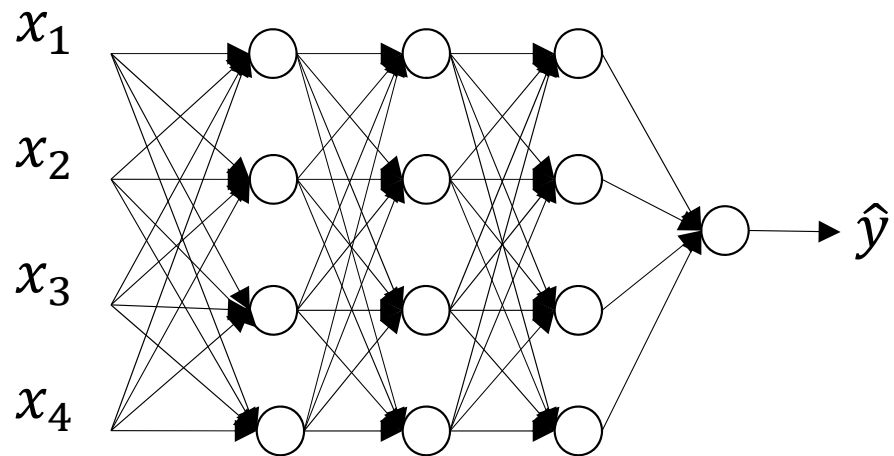0.5    0.5    0.5

# 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.

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

50 units. $\rightsquigarrow$ 10 units shut off

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

reduced by 20%.    Test

/= 0.8

Andrew Ng

# Making predictions at test time

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

__No drop out.__

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$
$$a^{[1]} = g^{[1]}(z^{[1]})$$
$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$
$$a^{[2]} = \ldots$$

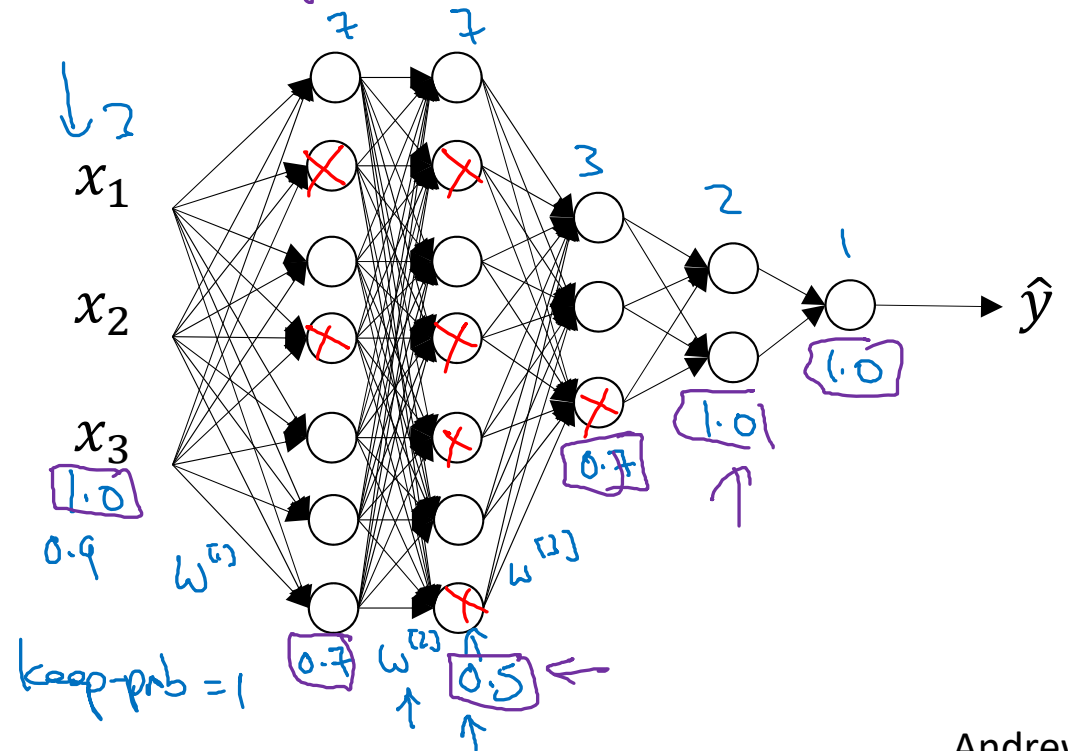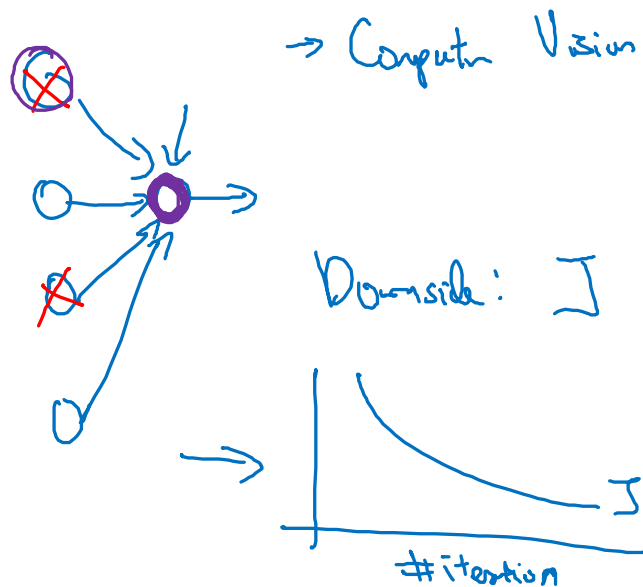$$\downarrow$$
$$\hat{y}$$

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

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. $\leadsto$ Shrink weights.



Andrew Ng

deeplearning.ai

Regularizing your
neural network

Other regularization
methods

# Data augmentation

# Early stopping

Orthogonalization.

→ — Optimize cost function $J$
   — Gradient, ....

→ — Not overfit.
   — Regularization, ....

$l_2$

$J(\omega, b)$



dev set error

train error or $J$

$\omega \approx 0$

# iterations

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

large $\omega$

Setting up your
optimization problem

Normalizing inputs

deeplearning.ai

# Normalizing training sets

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



Subtract mean:

$\mu = \dfrac{1}{m} \sum_{i=1}^{m} x^{(i)}$

$X := x - \mu$

Normalize variance

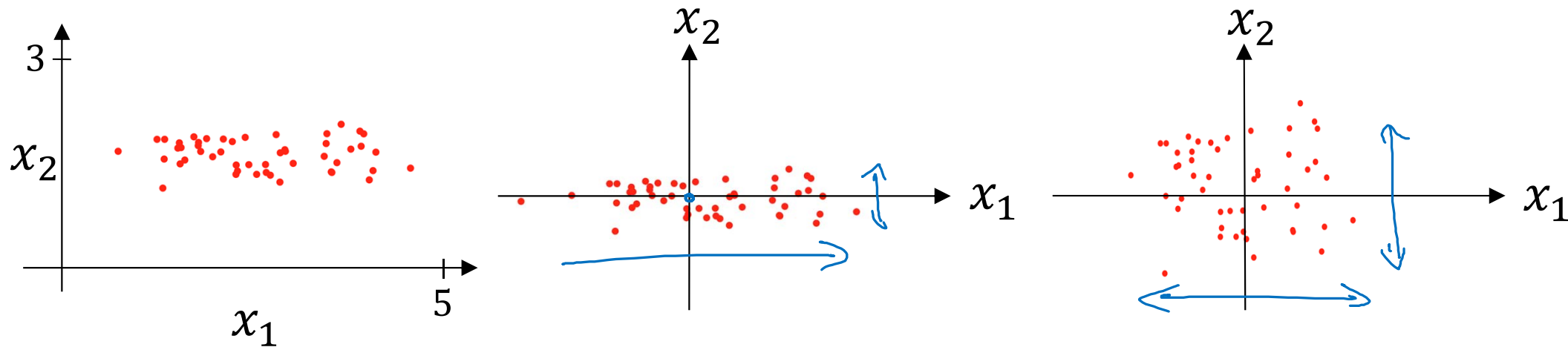$\sigma^2 = \dfrac{1}{m} \sum_{i=1}^{m} x^{(i)} ** 2$

↪ element-wise

$X /= \sigma^2$

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

Andrew Ng

# Why normalize inputs?

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

$\omega_1 \quad x_1 : \underline{1 \cdots 1000} \leftarrow$

$\omega_2 \quad x_2 : \underline{0 \cdots 1} \leftarrow$

$-1 \cdots 1$

Unnormalized:



Normalized:



$x_1 : 0 \cdots 1$

$x_2 : -1 \cdots 1$

$x_3 : 1 \cdots 2$

Andrew Ng

Setting up your
optimization problem

---

Vanishing/exploding
gradients

deeplearning.ai

# Vanishing/exploding gradients

$L = 150$

$x_1$

$x_2$



$\hat{y}$

$w^{[1]}$  $w^{[2]}$  $w^{[3]}$  ......  $w^{[L]}$

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

$\hat{y} = w^{[L]} \; w^{[L-1]} \; w^{[L-2]} \; \cdots \; w^{[3]} \; w^{[2]} \; w^{[1]} \; x$   $\to a^{[3]}$

$1.5^L$

$0.5^L$

$w^{[l]} > I$

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

$w^{[l]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$  (0.5 / 0.5)

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

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

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

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

$1.5^{L-1} x$

$0.5^{L-1} x$

Andrew Ng

# Single neuron example

$a^{[l]}$

$x_1$

$x_2$

$x_3$

$x_4$

$\hat{y}$

$W^{[l]}$

$a = g(z)$

$z = w_1 x_1 + w_2 x_2 + \cdots + w_n x_n \not{+b}$

large $n$ $\rightarrow$ Smaller $w_i$

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

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

ReLU

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

Other variants:

tanh

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

Xavier initialization

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

Andrew Ng

Setting up your
optimization problem
---
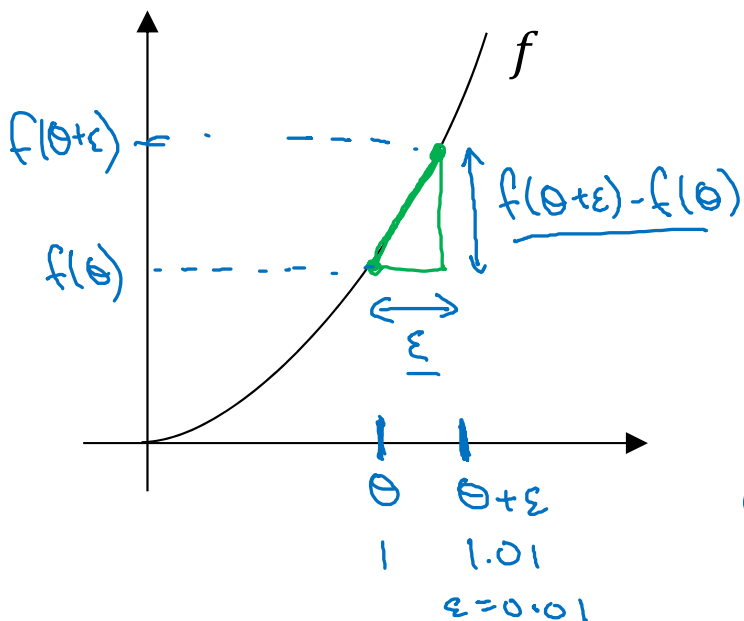Numerical approximation
of gradients

deeplearning.ai

# 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 $\theta = 1$

$\dfrac{dw}{db}$

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

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

$\theta = 1$

$\theta + \varepsilon = 1.01$



$f$

$f(\theta+\varepsilon)$

$f(\theta)$

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

$\varepsilon$

$\theta \quad \theta+\varepsilon$

$1 \quad 1.01$

$\varepsilon = 0.01$

0.0301
3 . 1
3 . 2

# Checking your derivative computation

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



$$\left[ \frac{f(\theta+\varepsilon) - f(\theta-\varepsilon)}{2\varepsilon} \approx g(\theta) \right.$$

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

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

approx error: 0.0001

(prev slide: 3.0301. error: 0.03)

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

$$0.01 \qquad\qquad\qquad\qquad\qquad 0.01$$

$$0.0001$$

Andrew Ng

Setting up your
optimization problem

Gradient Checking

deeplearning.ai

# Gradient check for a neural network

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

Concentrate

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

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

Concentrate

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

# Gradient checking (Grad check)

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

for each $i$:

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

$$\approx d\Theta[i] = \frac{\partial J}{\partial \theta_i}$$

$$d\Theta_{approx} \overset{?}{\approx} d\Theta$$

Check

$$\rightarrow \frac{\|d\Theta_{approx} - d\Theta\|_2}{\|d\Theta_{approx}\|_2 + \|d\Theta\|_2}$$

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

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

$$10^{-5}$$

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

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.

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

- Remember regularization.

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

$$d\Theta = \text{gradt of } J \text{ w.r.t. } \Theta$$

- 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