deeplearning.ai

Setting up your
ML application

Train/dev/test
sets

# Applied ML is a highly iterative process
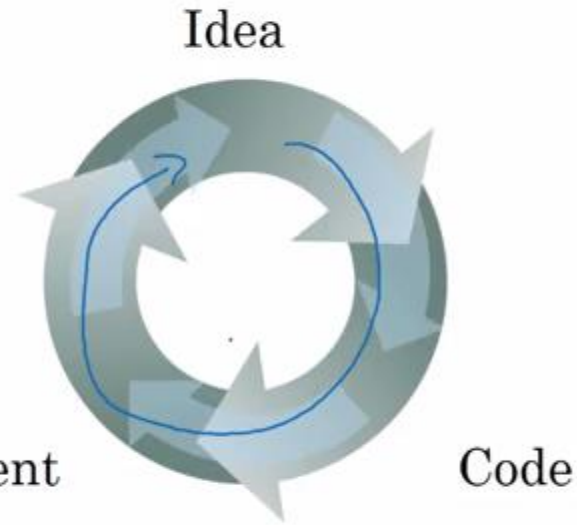
# layers

# hidden units

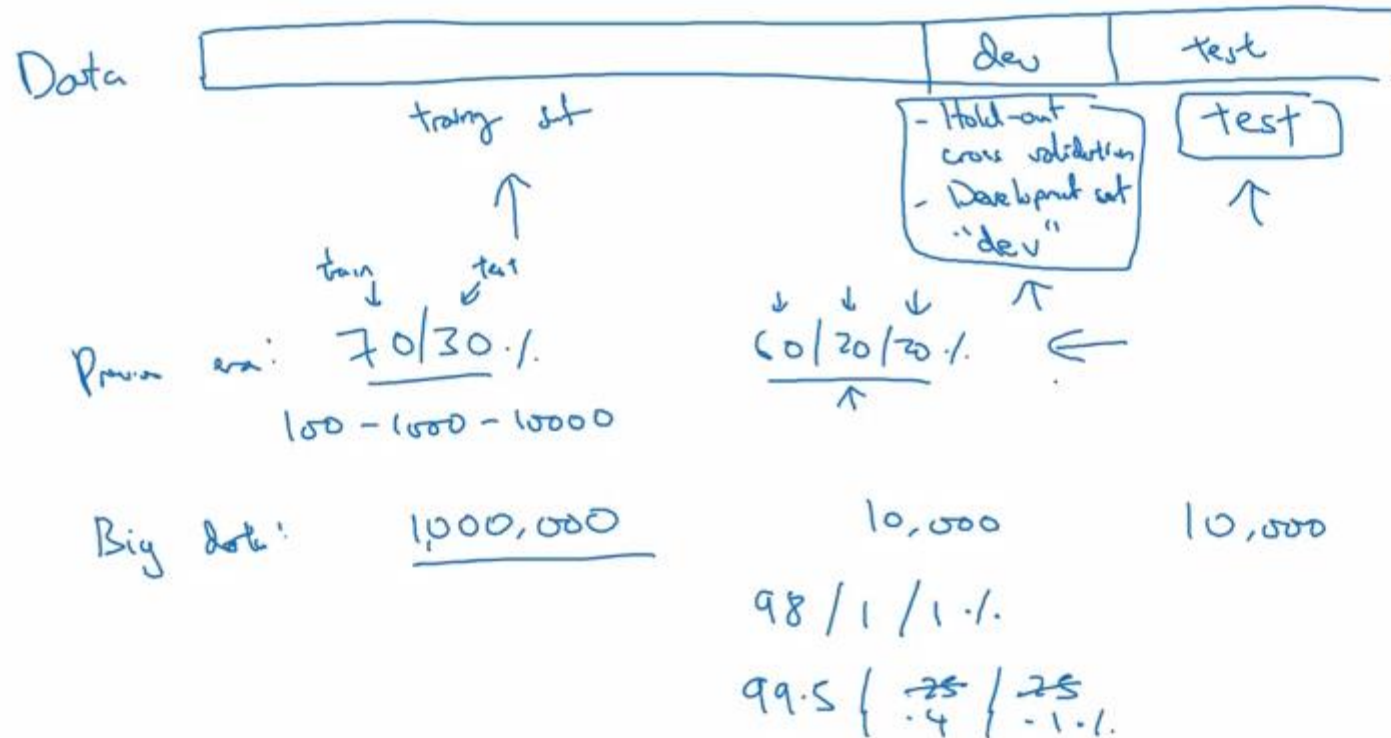learning rates

activation functions

...



Idea

Code

Experiment

NLP, Vision, Speech, Structural data

Ads     Search     Security     logistic ....

Andrew Ng

# Train/dev/test sets



Data

train set

| | dev | test |

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

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/.4/.1%

Andrew Ng

# Mismatched train/test distribution

Cats

Training set:
Cat pictures from webpages

$\longleftrightarrow$

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

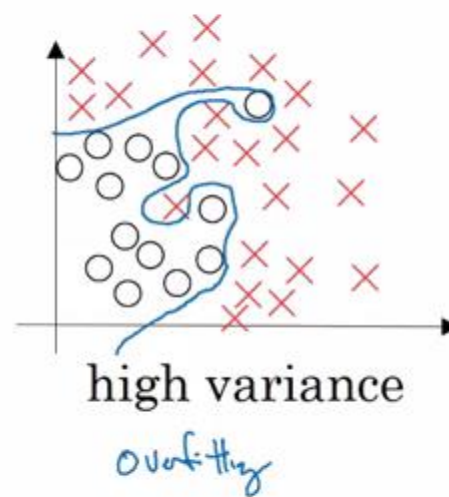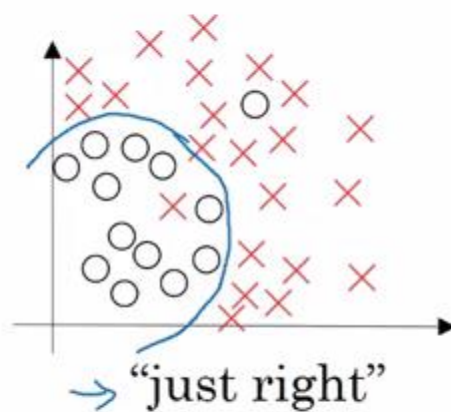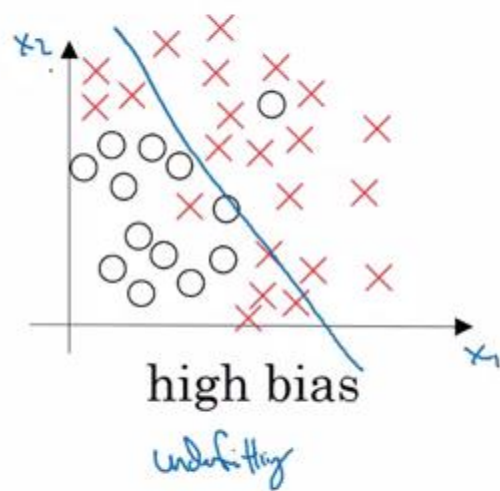Not having a **test set** might be okay. (Only **dev set**.)

# Bias and Variance



high bias — "just right" — high variance

# 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 varince ↑ |

Human: ≈ 0%

Optml (Bayes) error: ≈ 0% to 15%    Blury images

Andrew Ng

# High bias and high variance

# 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

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

$\lambda$ = regularization parameter

lambda    lambd

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

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

omit

$L_2$ regularization

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

$w$ will be sparse

$L_1$ regularization

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

Andrew Ng

# Neural network

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

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

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

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

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

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

$$\boxed{\left(1 - \frac{\alpha \lambda}{m}\right) w^{[l]}} \qquad = \boxed{w^{[l]} - \left(\frac{\alpha \lambda}{m}\right) w^{[l]}} - \alpha \, (\text{from backprop})$$
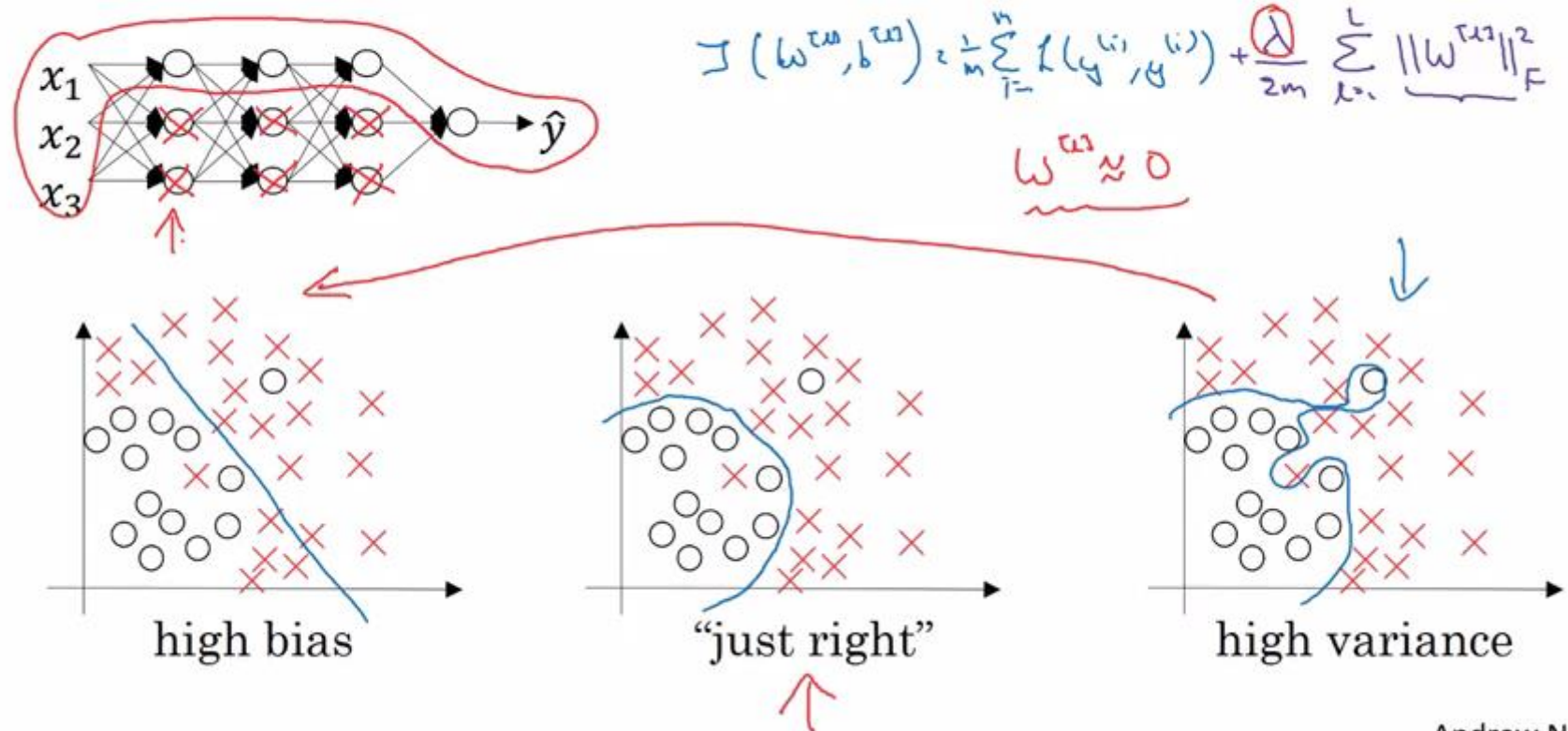
Andrew Ng

# How does regularization prevent overfitting?



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

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

high bias      "just right"      high variance

Andrew Ng

# How does regularization prevent overfitting?
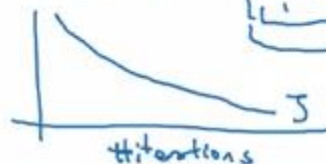


tanh

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

$\lambda \uparrow$

$W^{[l]} \downarrow$

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

Every layer $\approx$ linear.

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

J

#iterations

# Dropout regularization

# Implementing dropout ("Inverted dropout")

Illustrate with layer $l = 3$.   keep-prob $= \dfrac{0.8}{\cancel{\phantom{x}}}$          $\underline{0.2}$

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

$a3 = np.multiply(a3, d3)$          # $a3 \mathrel{*}= d3.$

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

50 units. $\rightsquigarrow$   10 units shut off

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

$J$              reduced by 20%.        Test

$\mathrel{/}= 0.8$

Andrew Ng

# Making predictions at test time

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

$$\underline{\text{No drop out.}}$$

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

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

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

$$a^{[2]} = \ldots.$$

$$\downarrow$$

$$\hat{y}$$

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

# Why does drop-out work?

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



Andrew Ng

deeplearning.ai

Regularizing your neural network

Other regularization methods

# Data augmentation
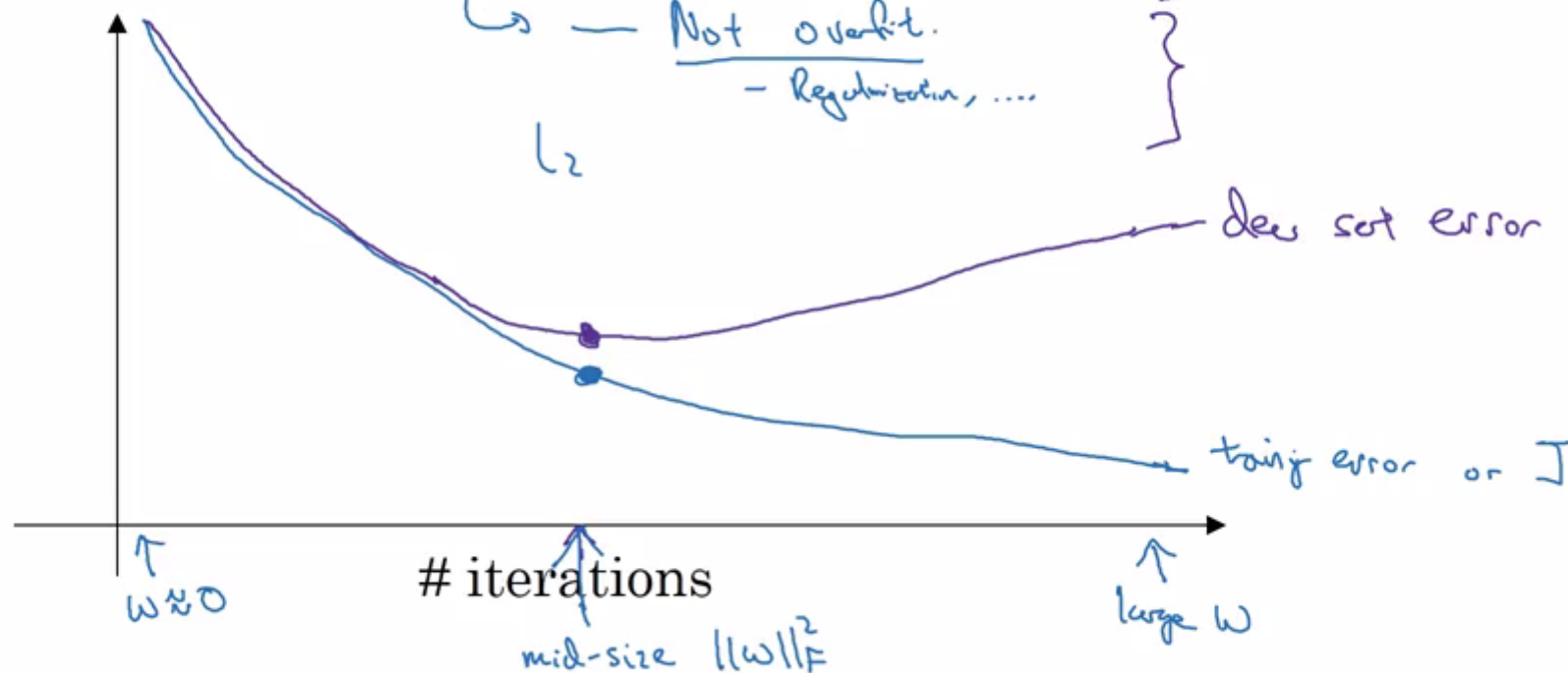
Andrew Ng

# Early stopping

Orthogonalization.

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

→ — 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 $w$

# 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

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^{m} x^{(i)} ** 2 \quad \leftarrow \text{element-wise}$$

$$x /= \sigma^2$$

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

# Why normalize inputs?

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

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

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

$-1 \cdots 1$

**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

$L = 150$

$1 \quad x_1$
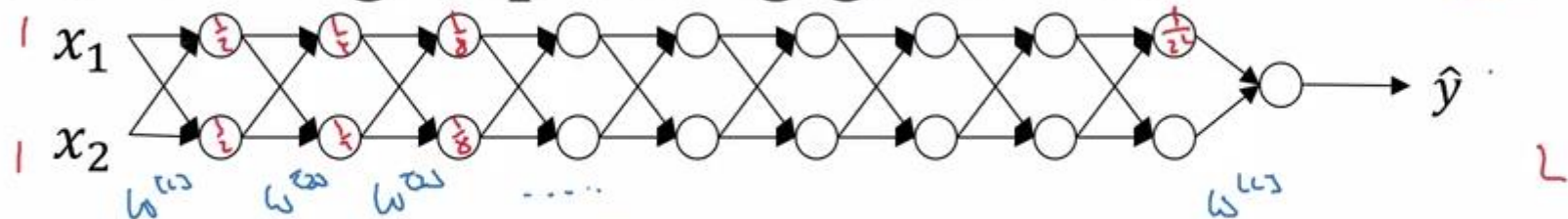$1 \quad x_2$

$W^{[1]} \quad W^{[2]} \quad W^{[2]} \quad \ldots \ldots \qquad W^{[L]}$

$\hat{y}$

$L$

$g(z) = z.$  $\qquad b^{[\ell]} = 0.$

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

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

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

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

$W^{[\ell]} > I$

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

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

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

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

$1.5^{L-1} x$

$0.5^{L-1} x$

Andrew Ng

# Single neuron example



$x_1$
$x_2$
$x_3$
$x_4$

$a^{[1]}$

$w^{[1]}$

$a = g(z)$

$\hat{y}$

$z = w_1 x_1 + w_2 x_2 + \cdots + w_n x_n \;\; \cancel{+b}$

Large $n \;\rightarrow\;$ Smaller $w_i$

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

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

ReLU

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

Other varots:

tanh

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

Xavier initialisation

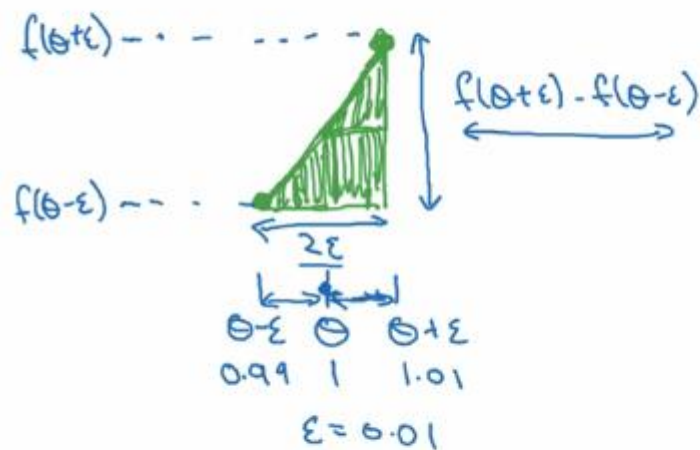$\sqrt{\dfrac{2}{n^{[l-1]} + n^{[l]}}}$

Andrew Ng

Setting up your optimization problem

Numerical approximation of gradients

deeplearning.ai

# Checking your derivative computation

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

$$\frac{f(\theta+\varepsilon) - f(\theta-\varepsilon)}{2\varepsilon} \approx g(\theta)$$



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

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

approx error: $0.0001$

(prev slide: $3.0301$.  error: $0.03$)

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

$2\varepsilon$

$\theta-\varepsilon$    $\theta$    $\theta+\varepsilon$
$0.99$      $1$      $1.01$

$\varepsilon = 0.01$

$f'(\theta) = \lim_{\varepsilon \to 0} \dfrac{f(\theta+\varepsilon) - f(\theta-\varepsilon)}{2\varepsilon}$   $O(\varepsilon^2)$

$0.01$

$0.0001$

$\dfrac{f(\theta+\varepsilon) - f(\theta)}{\varepsilon}$   error: $O(\varepsilon)$

$0.01$

Andrew Ng

# Gradient check for a neural network

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

Concatenate

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

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

concatenate

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

# Gradient checking (Grad check)

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

for each $i$:

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

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

Check $\quad \dfrac{\| d\theta_{approx} - d\theta \|_2}{\rightarrow \| d\theta_{approx} \|_2 + \| d\theta \|_2}$

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

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

$$10^{-5}$$

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

Andrew Ng

6:18 / 6:34

# Gradient checking implementation notes

- Don't use in training – only to debug

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

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

$$\underline{db_k^{[l]}} \qquad \underline{dw_k^{[l]}}$$

- Remember regularization.

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

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

- Doesn't work with dropout.    $J$    $\underline{keep\text{-}prob = 1.0}$

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

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

Andrew Ng