Setting up your
ML application

Train/dev/test
sets

deeplearning.ai

# 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

training set

dev | test

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

test

Prev era:  70/30 %

train ↓    ↘ test

60/20/20 %  ←

100 - 1000 - 10000

Big data:  1000,000      10,000        10,000

98 / 1 / 1 %

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

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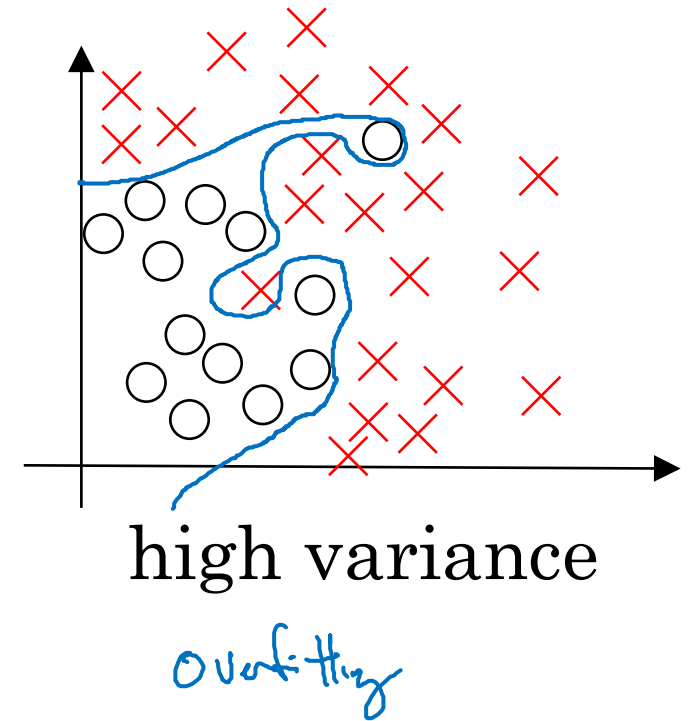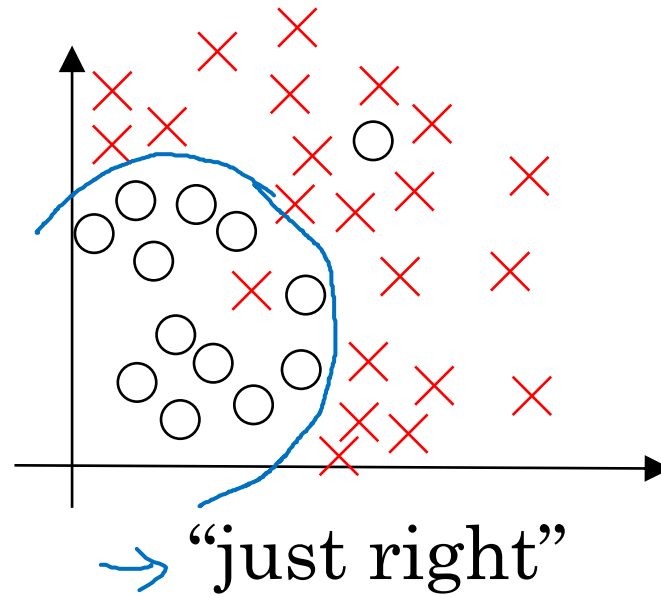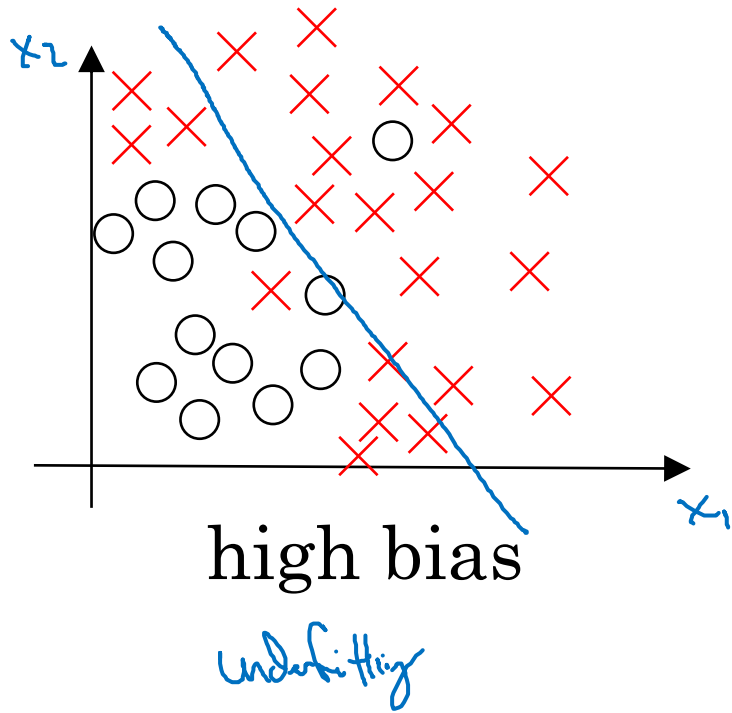
Not having a test set might be okay. (Only dev set.)
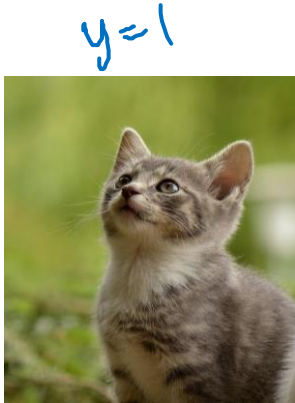
Andrew Ng

Setting up your
ML application

Bias/Variance

deeplearning.ai

# Bias and Variance



high bias           → "just right"           high variance

Underfitting                                                        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: ≈ 0%

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

Andrew Ng

# High bias and high variance



high bias
high varian

Andrew Ng

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    trade off "

Andrew Ng

Regularizing your
neural network

Regularization

deeplearning.ai

# Logistic regression

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

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

$\lambda =$ regularization parameter

lambda      lambd

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

omit

$L_2$ 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

Andrew Ng

# Neural network

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

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

$$\omega^{[l]} : (n^{[l]}, n^{[l-1]})$$

"Frobenius norm"

$$\|\cdot\|_2^2 \qquad \|\cdot\|_F^2$$

$$d\omega^{[l]} = \boxed{(\text{from backprop}) + \frac{\lambda}{m} \omega^{[l]}}$$

$$\frac{\partial J}{\partial \omega^{[l]}} = d\omega^{[l]}$$

$$\rightarrow \omega^{[l]} := \omega^{[l]} - \alpha \, d\omega^{[l]}$$

"Weight decay"

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

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

$$= \underbrace{(1 - \frac{\alpha\lambda}{m})}_{<1} \omega^{[l]} - \alpha (\text{from backprop})$$
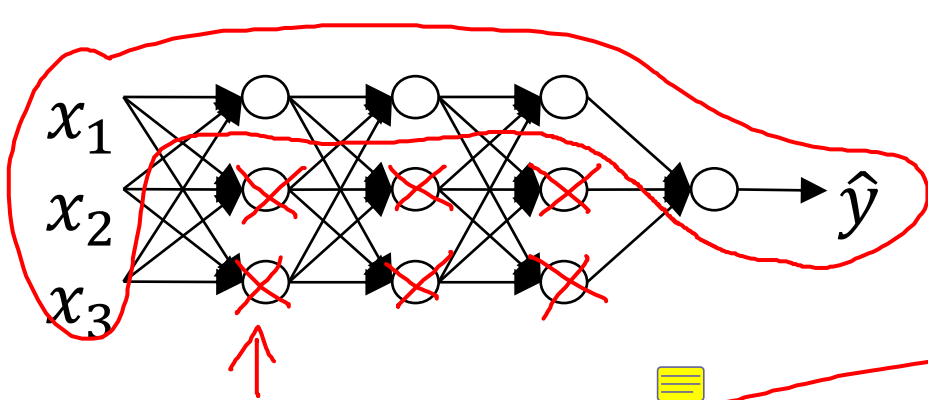
Andrew Ng

deeplearning.ai

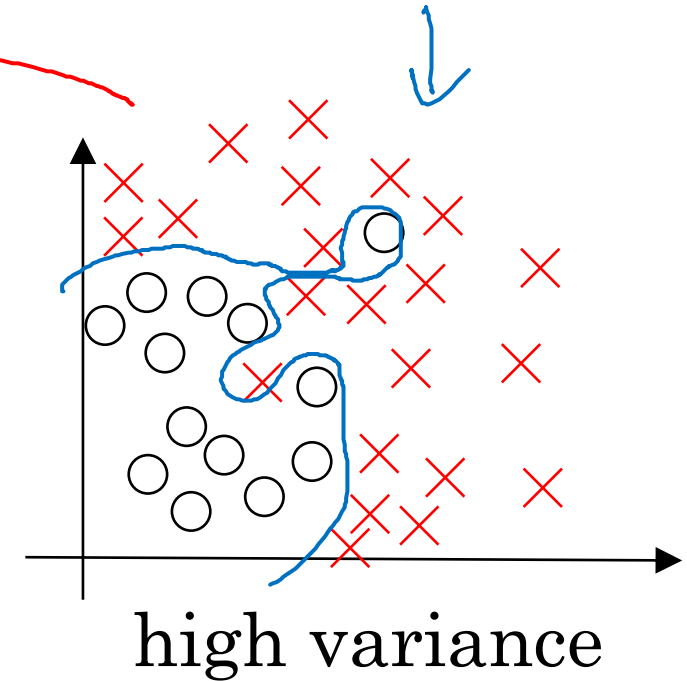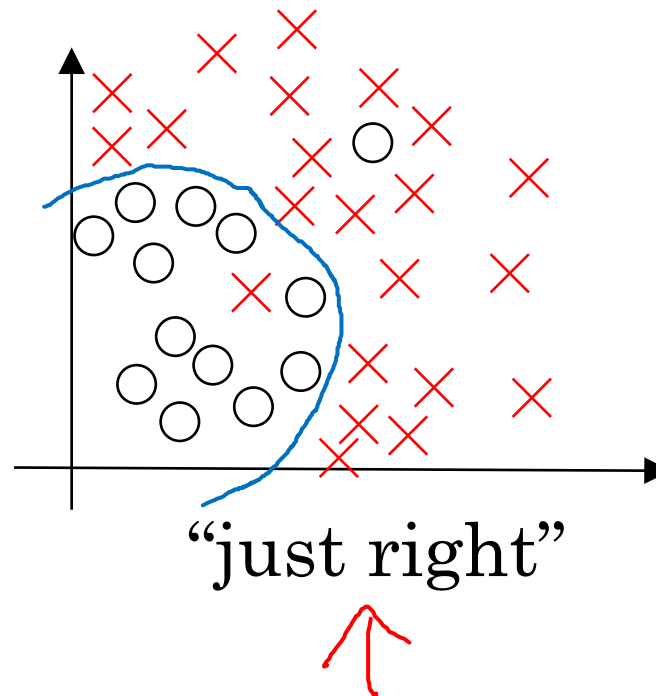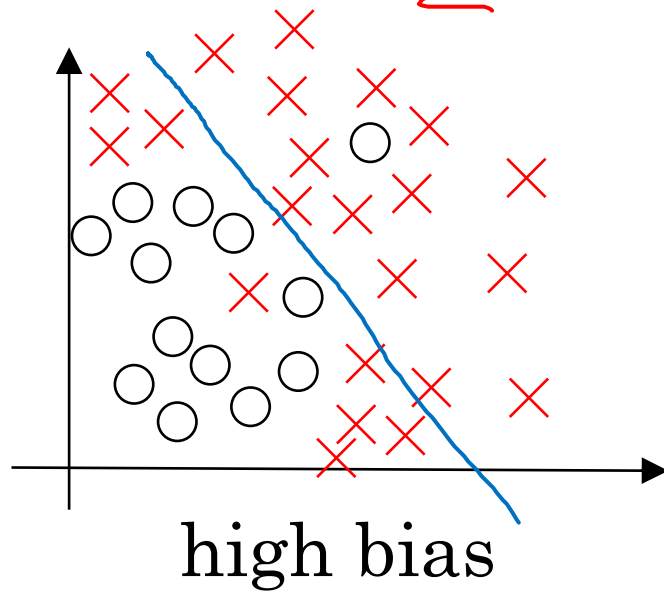# Regularizing your neural network

---

# Why regularization reduces overfitting

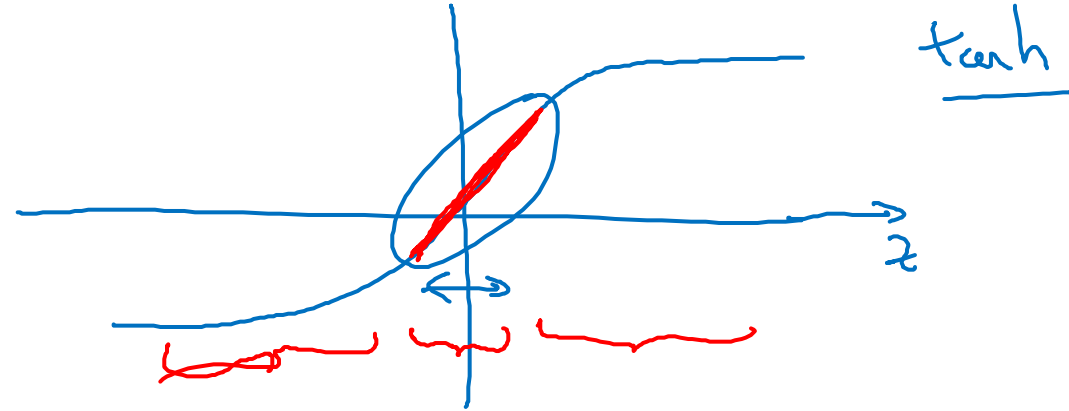# How does regularization prevent overfitting?



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

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

high bias    "just right"    high variance

Andrew Ng

# How does regularization prevent overfitting?



$tanh$

$g(z) = tanh(z)$

$z$

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

Every layer $\approx$ linear.

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

$J$

Hiterations

# 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.

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

50 units. $\rightsquigarrow$ 10 units shut off

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

reduced by 20%.    Test

$/= 0.8$

# Making predictions at test time

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

No drop out.

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

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

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

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

$$\downarrow$$

$$\hat{y}$$

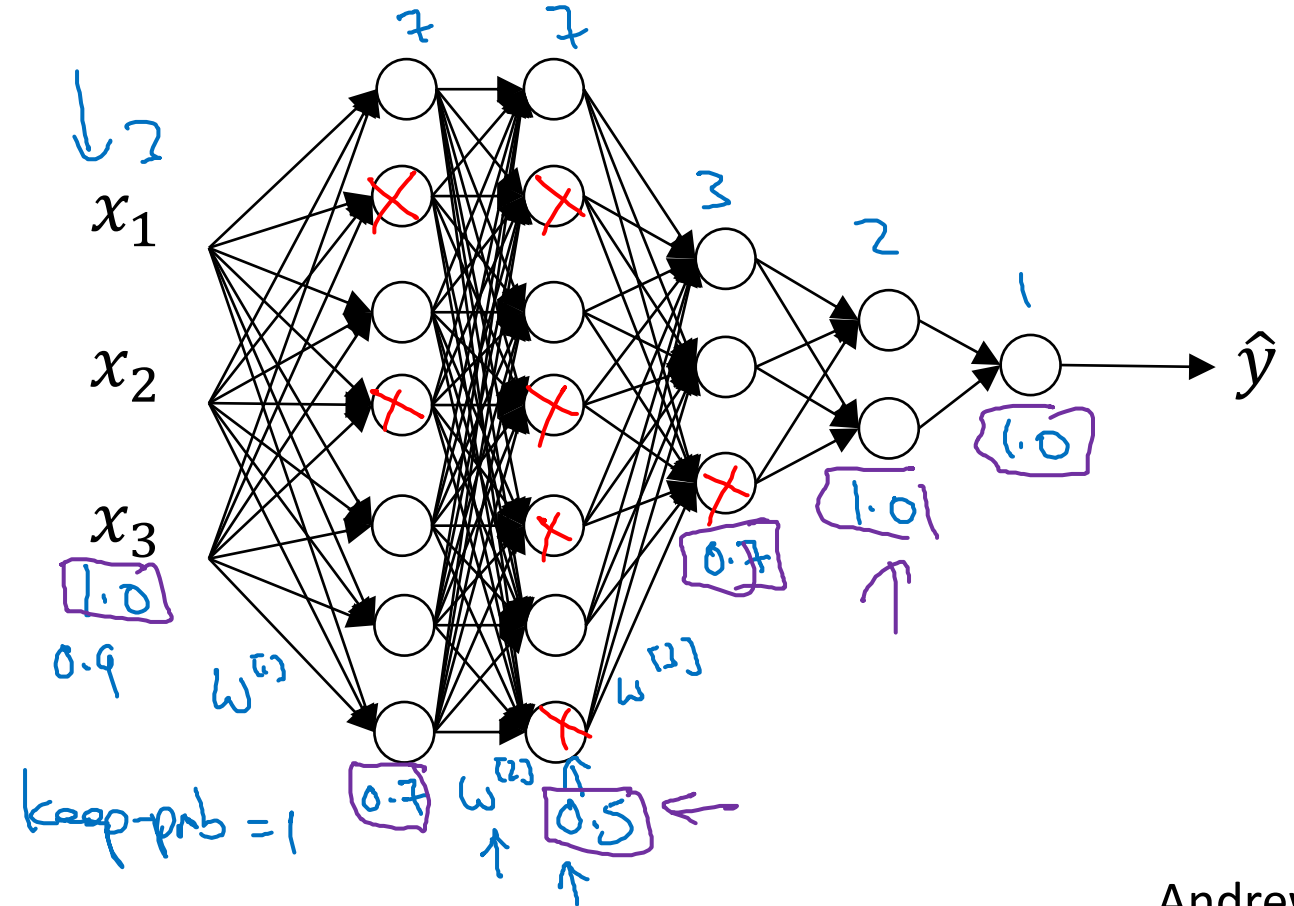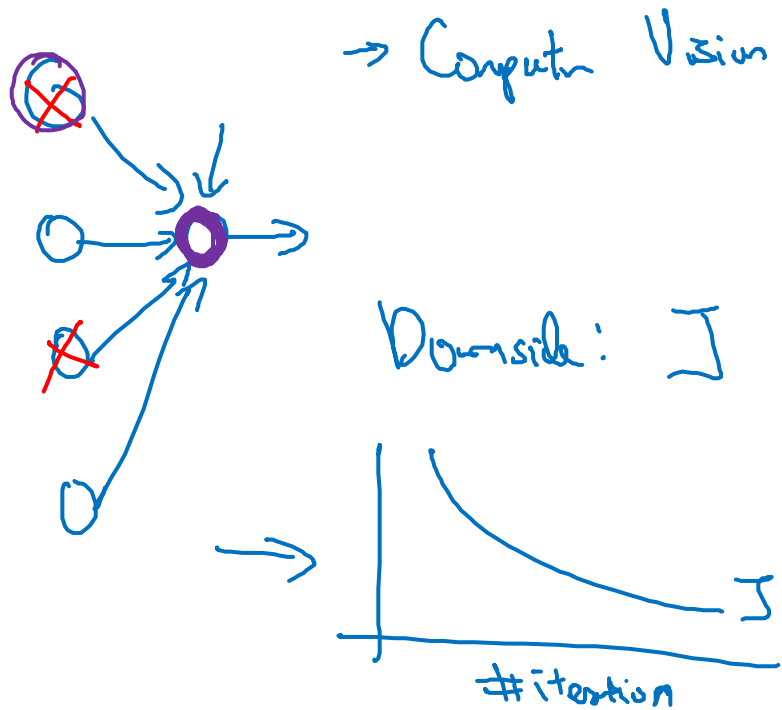$$/= \text{keep-prob}$$

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. $L_2$
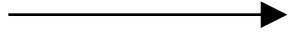


Andrew Ng

deeplearning.ai

Regularizing your
neural network
_____

Other regularization
methods

# Data augmentation

# Early stopping

Orthogonalization.

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

$\rightarrow$ — 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 $W$

Setting up your
optimization problem

Normalizing inputs

deeplearning.ai

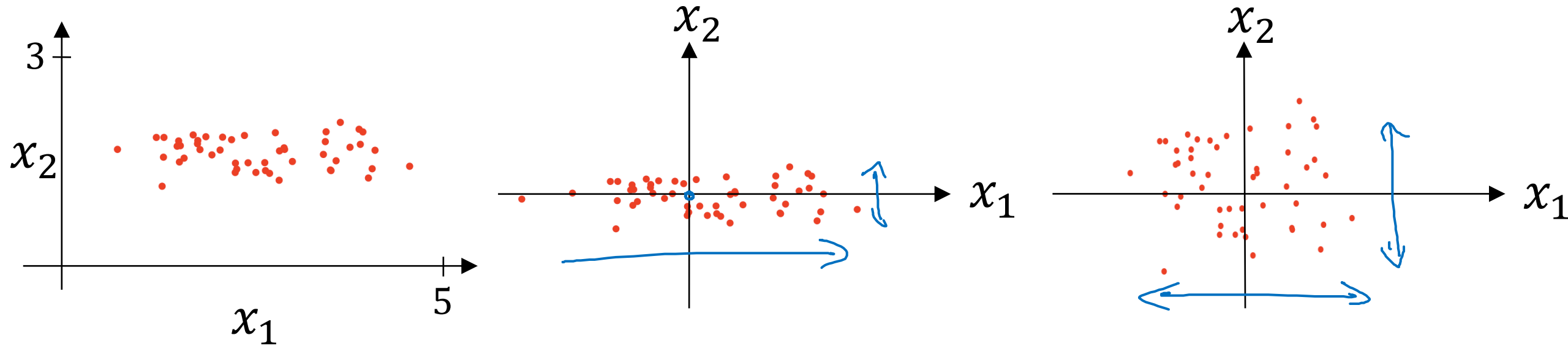# Normalizing training sets

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



Substract mean:
$$\mu = \frac{1}{m} \sum_{i=1}^{m} x^{(i)}$$

$$\boxed{x := x - \mu}$$

Normalize variance
$$\sigma^2 = \frac{1}{m} \sum_{i=1}^{m} x^{(i)} **2$$
↖ element-wise

$$\boxed{x \, / = \sigma}$$

$$\Rightarrow x = \frac{x - \mu}{\sigma}$$

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

Andrew Ng

# Why normalize inputs?

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

$w_1 \quad x_1: \quad 1 \cdots 1000$

$w_2 \quad x_2: \quad 0 \cdots 1$

$-1 \cdots 1$

**Unnormalized:**

**Normalized:**



$x_1: \quad 0 \cdots 1$

$x_2: \quad -1 \cdots 1$

$x_3: \quad 1 \cdots 2$

Andrew Ng

Setting up your
optimization problem

Vanishing/exploding
gradients

deeplearning.ai

# Vanishing/exploding gradients

$L = 150$



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

$\hat{y} = W^{[L]} \underbrace{W^{[L-1]}} \underbrace{W^{[L-2]}} \underbrace{\cdots} \boxed{W^{[3]} \boxed{W^{[2]} W^{[1]} x}} \leftarrow a^{[3]}$

$1.5^L$

$0.5^L$

$W^{[l]} > I$

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

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

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

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

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

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

$1.5^{L-1} x$

$0.5^{L-1} x$

Andrew Ng

# Single neuron example

$a^{[l]}$

$x_1$

$x_2$

$x_3$

$x_4$

$\hat{y}$

$W^{[l]}$

$a = g(z)$

$z = w_1 x_1 + w_2 x_2 + \cdots + w_n x_n$ ~~+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 variants:

tanh

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

Xavier initialization

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

Andrew Ng

Setting up your
optimization problem

Numerical approximation
of gradients

deeplearning.ai

# Checking your derivative computation

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

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

$I$

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

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

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

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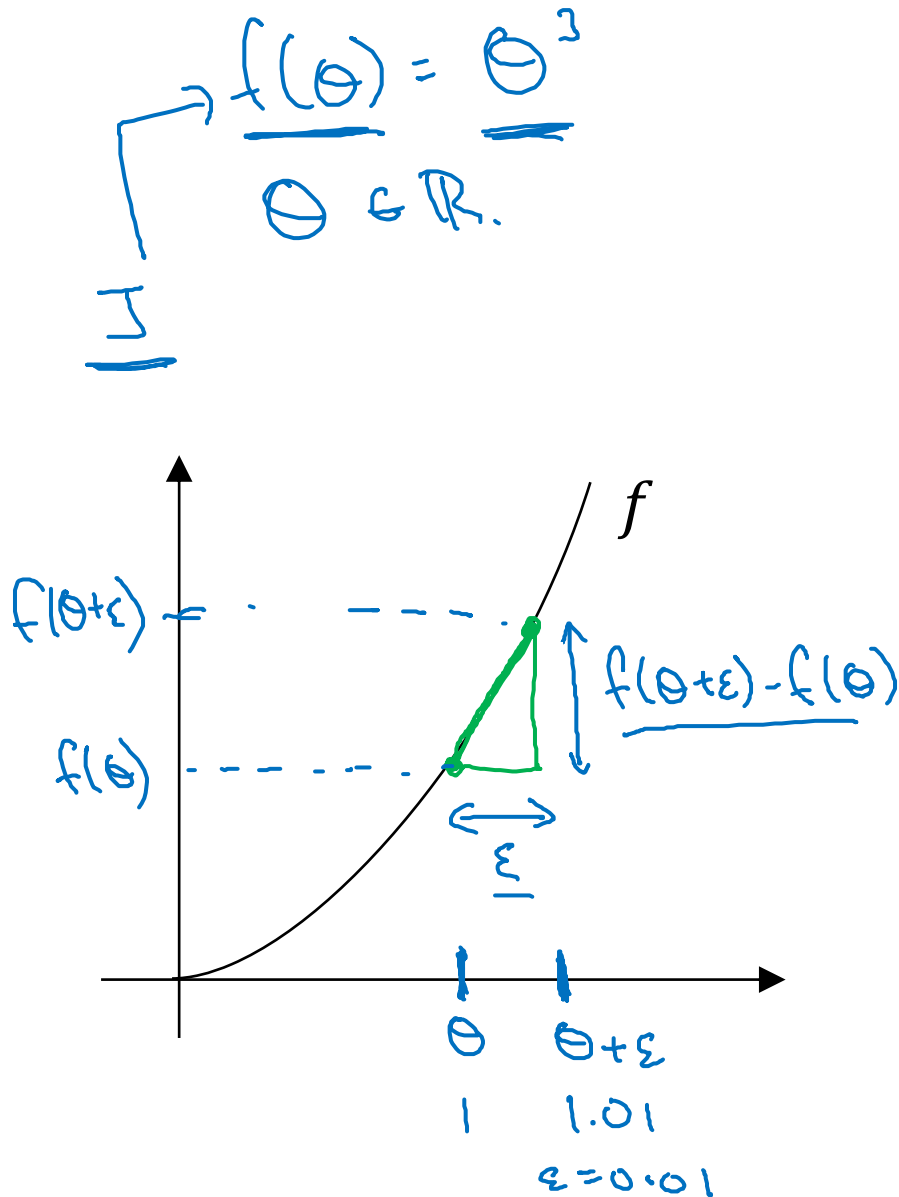
$$\frac{f(\theta + \varepsilon) - f(\theta)}{\varepsilon} \approx g(\theta)$$

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

$f$

$f(\theta + \varepsilon)$

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

$f(\theta)$

$\varepsilon$

$\theta$     $\theta + \varepsilon$

$1$     $1.01$

$\varepsilon = 0.01$

$\theta = 1$

$\theta + \varepsilon = 1.01$

$0.0301$

$3.1$
$3.2$

# 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)

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

$O(\varepsilon^2)$

0.01

0.0001

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

0.01

Andrew Ng

deeplearning.ai

Setting up your
optimization problem

Gradient Checking

# Gradient check for a neural network

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

Concatenate

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

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

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 $i$:

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

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

Check $\quad \dfrac{\| \partial\theta_{approx} - \partial\theta \|_2}{\| \partial\theta_{approx} \|_2 + \| \partial\theta \|_2} \qquad \approx \boxed{10^{-7} \; - \; great!} \leftarrow$

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

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

deeplearning.ai

Setting up your
optimization problem

Gradient Checking
implementation notes

# Gradient checking implementation notes

- Don't use in training – only to debug

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

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

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

- Remember regularization.

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

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

- 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