# Applied ML is a highly iterative process
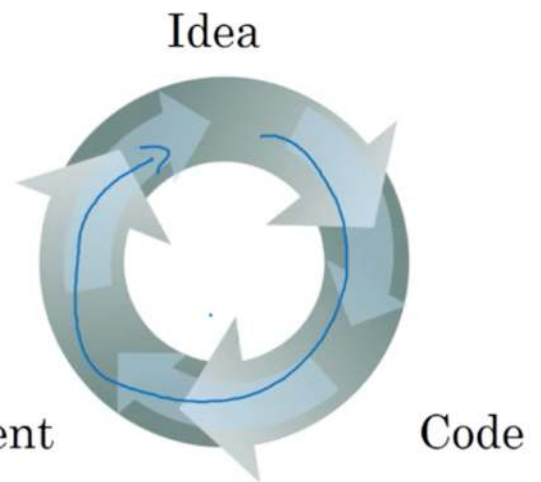
# layers

# hidden units

learning rates

activation functions
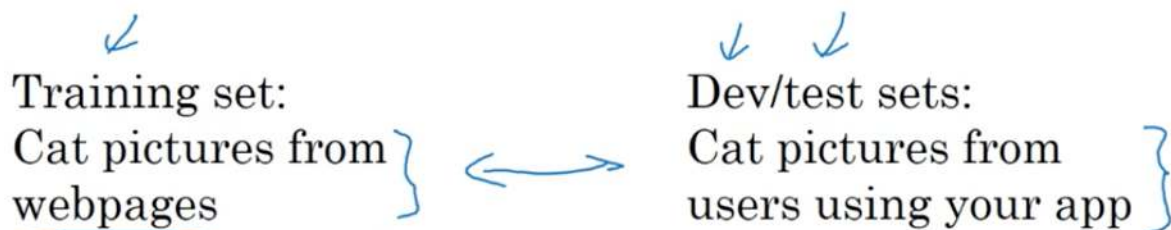
...

NLP, Vision, Speech, Structural Data

Idea

Experiment

Code

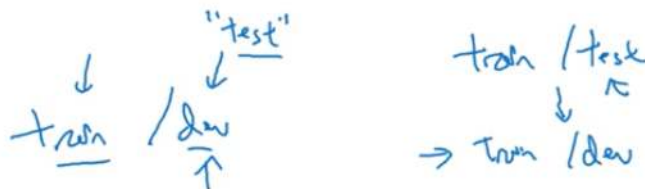to hopefully find a good choice
of network for your application.

Andrew Ng

# Train/dev/test sets

Data 

training set

$$\uparrow$$

train    test

Previous era: $70/30\%$

$100 - 1000 - 10000$

Big data: $1,000,000$

dev    test

test

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

$60/20/20\%$ $\Leftarrow$

$10,000$           $10,000$

$98/1/1\%$

$99.5 | \frac{.75}{.4} | \frac{.75}{-.1}\%$

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
↓ ↖
→ Trom / dev

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

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 varan | low bias low variance |
| | ↑ | ↑↑ | | ↑ |

Human: ≈ 0%

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

Andrew Ng

6:47 / 8:46

# High bias and high variance

# Basic recipe for machine learning

High bias?
(training data performance)

$\downarrow$ N

High variance?
(dev set performance)

$\downarrow$ N

Done

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

→ More data
→ Regularization
(NN architecture search)

Bias ⟷ Variance "tradeoff"

Andrew Ng

# Logistic regression

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

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

$$\text{lambda} \qquad \text{lambd}$$

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

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

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

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

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

$$\underline{\text{"Frobenius norm"}} \qquad \|\cdot\|_2^2 \qquad \|\cdot\|_F^2$$

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

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

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

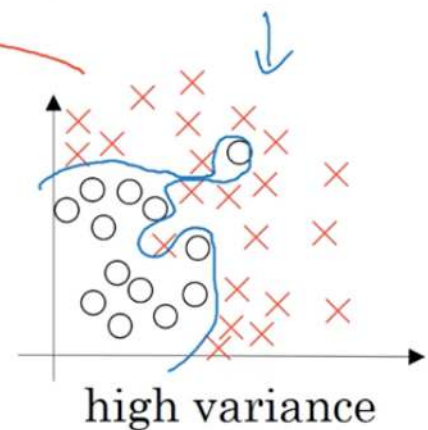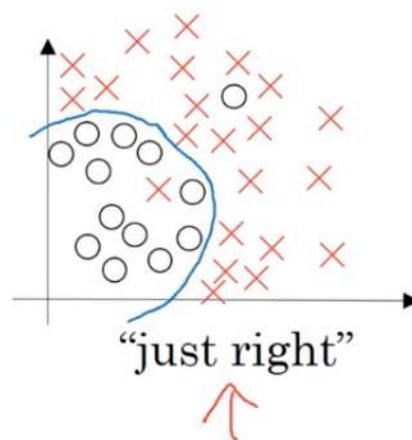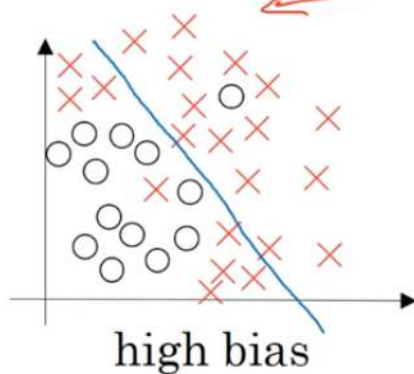$$\boxed{(1 - \frac{\alpha \lambda}{m}) w^{[l]}} = w^{[l]} - \boxed{\left(\frac{\alpha \lambda}{m}\right) w^{[l]}} - \alpha (\text{from backprop})$$
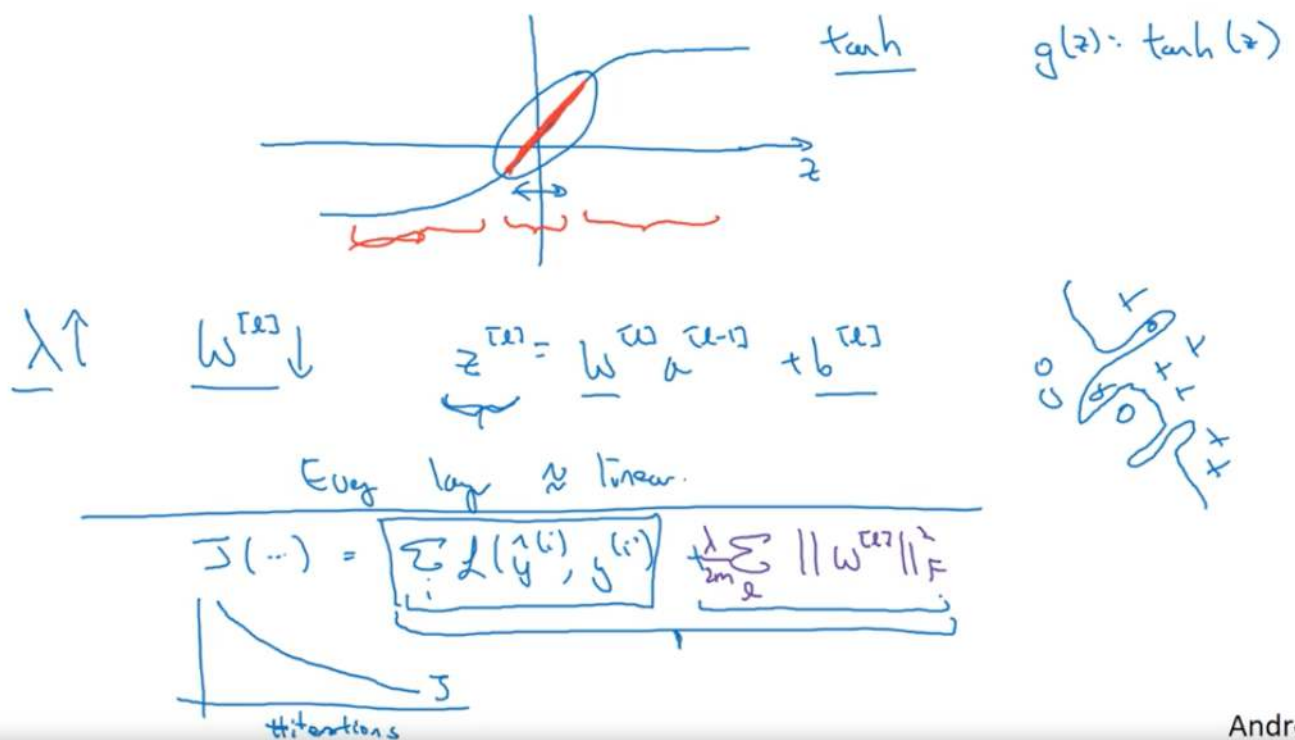
# How does regularization prevent overfitting?



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

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

high bias      "just right"      high variance

# How does regularization prevent overfitting?



$\text{tanh}$

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

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

Every layer $\approx$ linear.

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

#iterations   J

Andrew Ng

# Dropout regularization



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

$\rightarrow \hat{y}$

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

$\rightarrow \hat{y}$

0.5    0.5    0.5

Andrew Ng

# Implementing dropout ("Inverted dropout")

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

$\rightarrow$ $\boxed{d3}$ = 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

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

$\underset{\text{reduced by } \underline{20\%.}}{\underline{a}}$          $\underline{\text{Test}}$

$/= 0.8$

Andrew Ng

# Making predictions at test time

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

No drop out.

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

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

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

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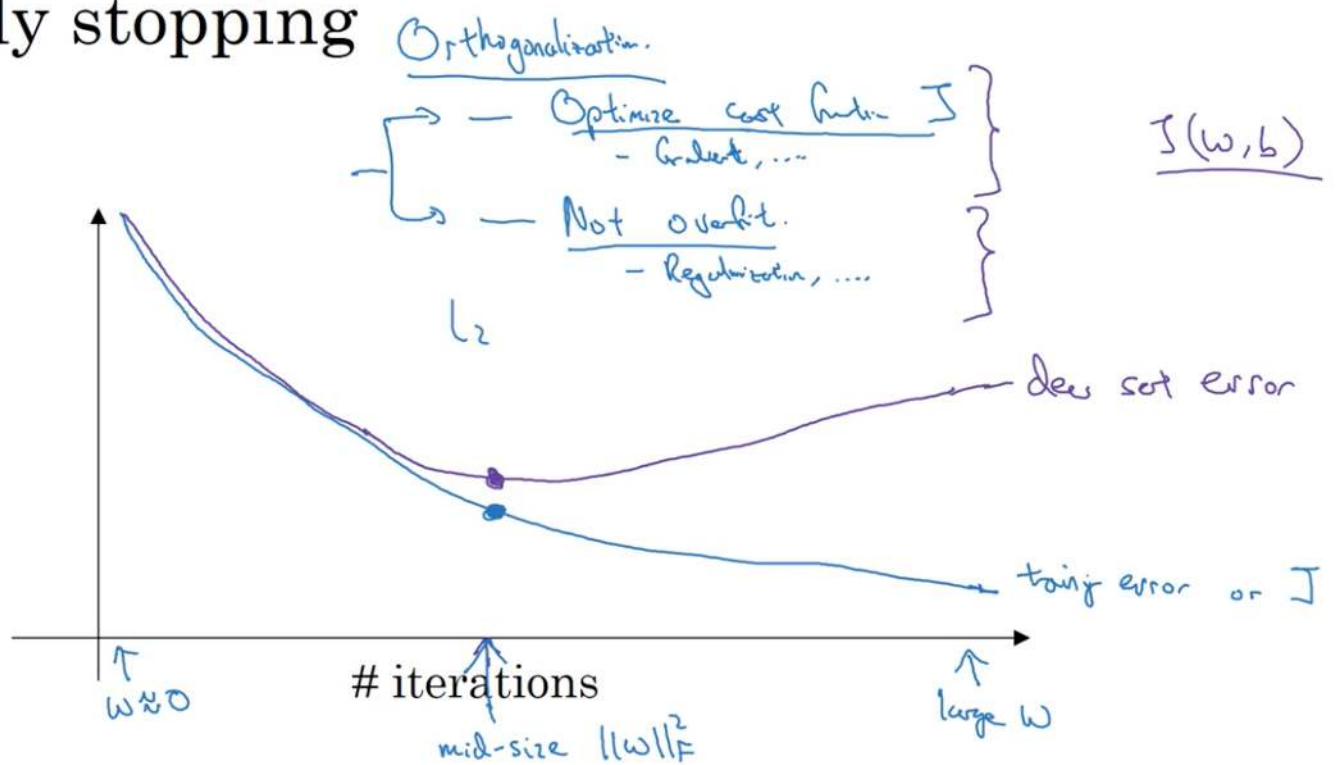
$$\downarrow$$

$$\hat{y}$$

$/ = $ keep-prob

Andrew Ng

# Why does drop-out work?

Intuition: <u>Can't rely on any one feature</u>, so have to
<u>spread out weights.</u> $\leadsto$ Shrink weights.



Andrew Ng

# Early stopping

Orthogonalization.

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

↳ — Not overfit.
    — Regularization, ....

$l_2$

$J(w, b)$



dev set error

training error or $J$
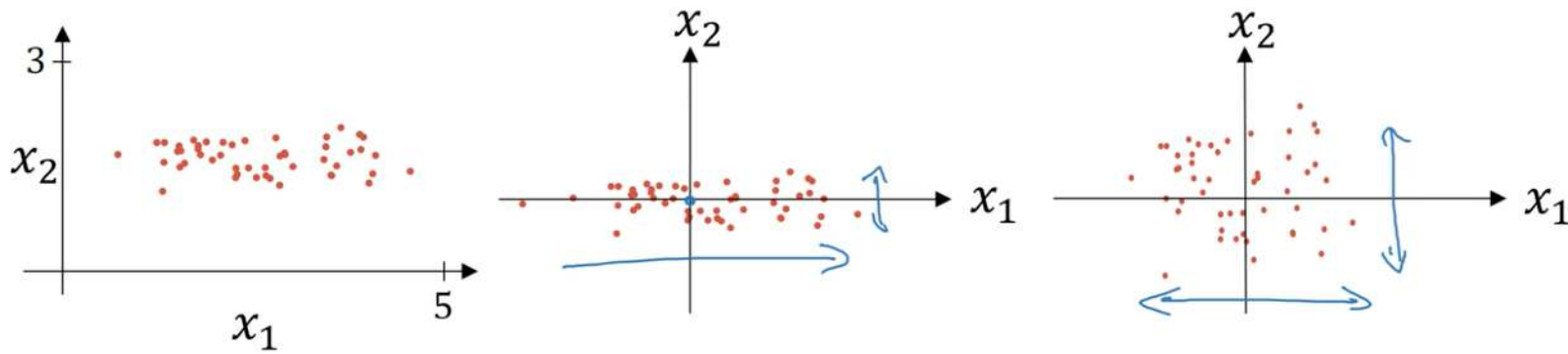
# iterations

$w \approx 0$

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

large $W$

# Normalizing training sets

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

Substract 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

$$x /= \sigma^2$$
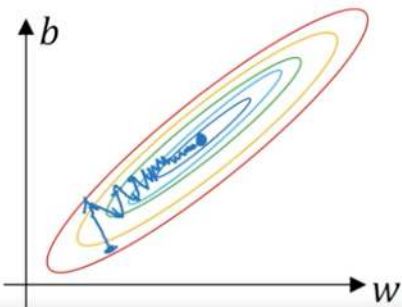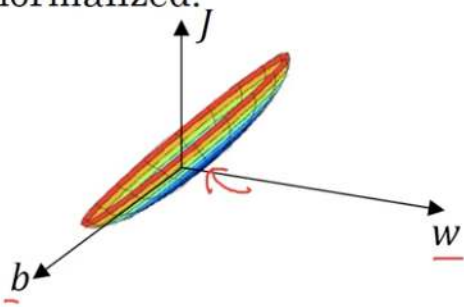
Use same $\mu$, $\sigma^2$ to normalize test set.

Andrew Ng

2:11 / 5:30

# Why normalize inputs?

$w_1 \quad x_1 : 1 \cdots 1000$
$w_2 \quad x_2 : 0 \cdots 1$

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

Unnormalized:

Normalized:



Andrew Ng

# Single neuron example

$a^{[1]}$

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

$w^{[1]}$

$\hat{y}$

$a = g(z)$

$z = w_1 x_1 + w_2 x_2 + \cdots + w_n x_n + b$

large $n$ → Smaller $w_i$

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

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

ReLU

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

Other variants:

tanh

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

Xavier initialization ↑

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

↑

# Vanishing/exploding gradients



$x_1$

$x_2$

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

$\hat{y}$

$L$

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

$\hat{y} = w^{[L]} \, w^{[L-1]} \, w^{[L-2]} \, (\cdots) \, w^{[3]} \, w^{[2]} \, w^{[1]} \, x$    $a^{[3]}$

$1.5^L$

$0.5^L$

$w^{[l]} > I$

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

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

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

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

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

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

$1.5^{L-1} x$

$0.5^{L-1} x$

Andrew Ng

# Gradient checking (Grad check)

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

for each $i$:

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

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

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

$\rightarrow$

$\varepsilon = 10^{-7}$

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

$10^{-5}$

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

Andrew Ng