

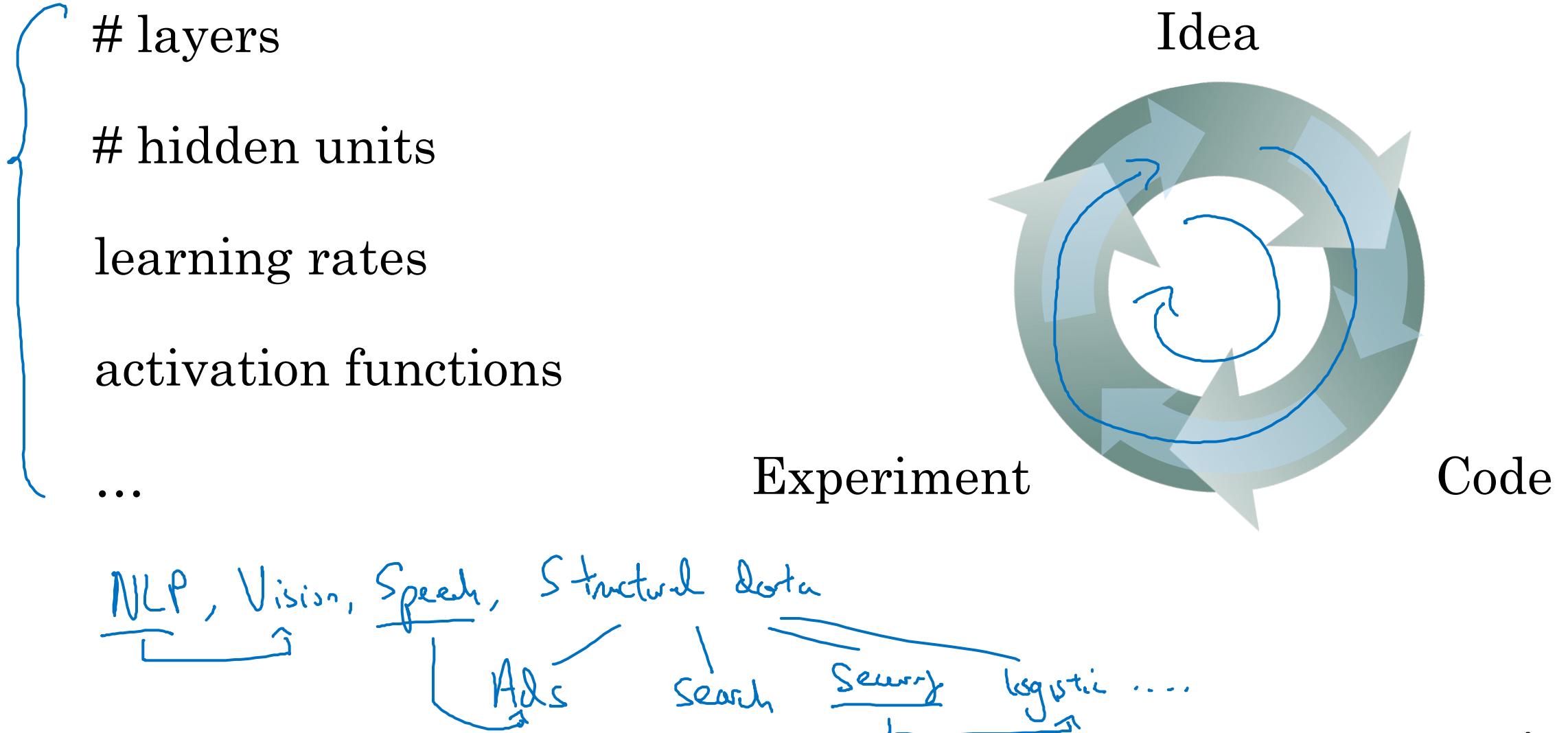


deeplearning.ai

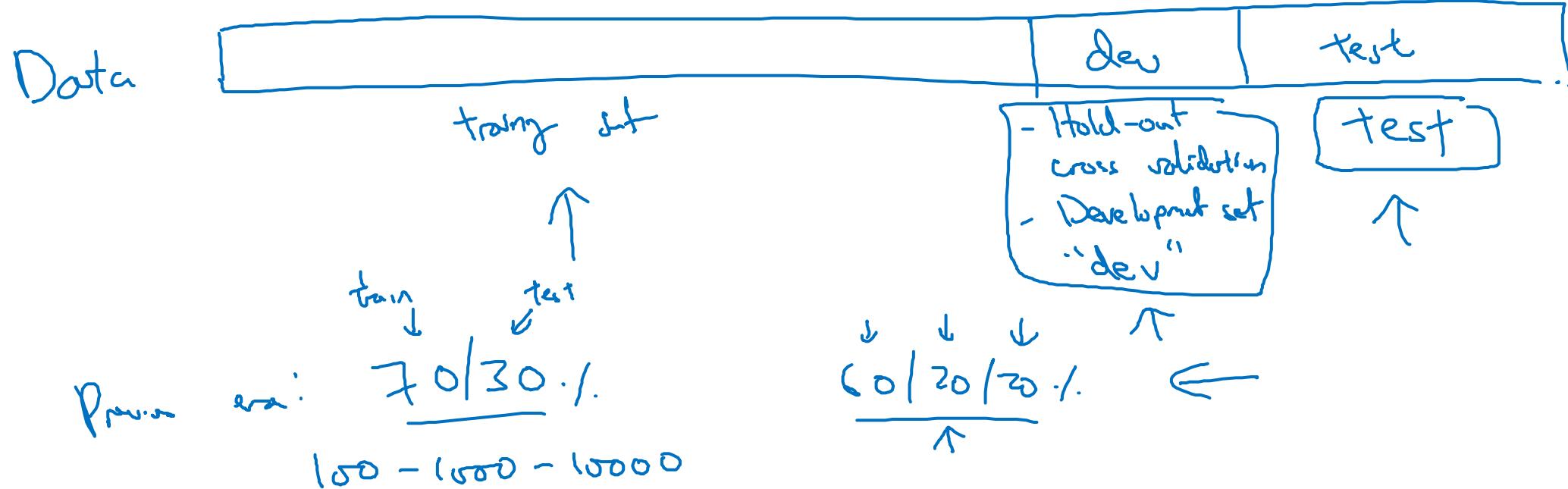
Setting up your
ML application

Train/dev/test
sets

Applied ML is a highly iterative process



Train/dev/test sets



Big data! 1,000,000

10,000 10,000
98 / 1 / 1%.
99.5 { 25 / 25
· 4 { - 1 - 1.

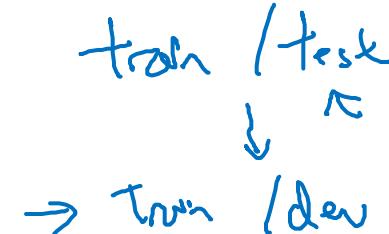
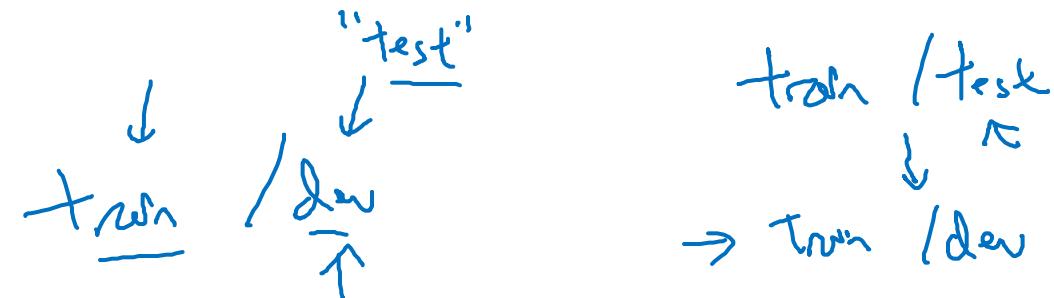
Mismatched train/test distribution

Conts

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.



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

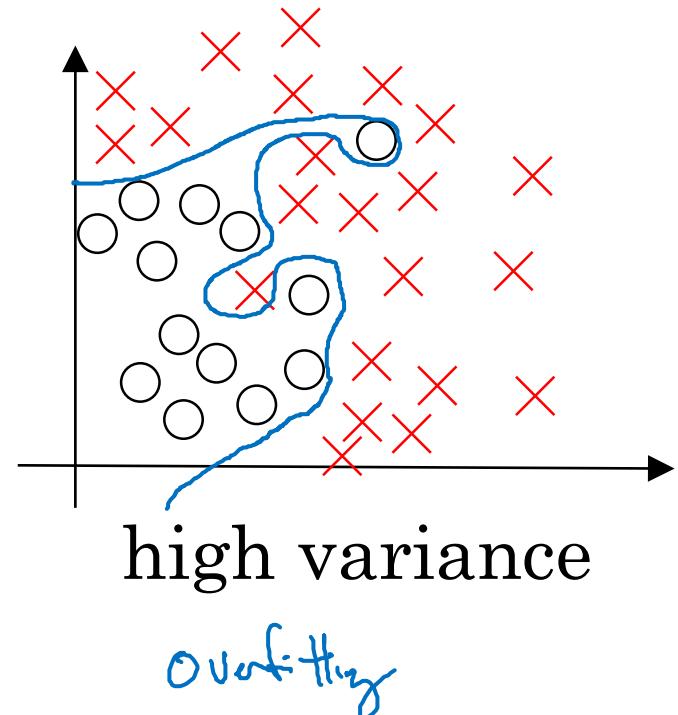
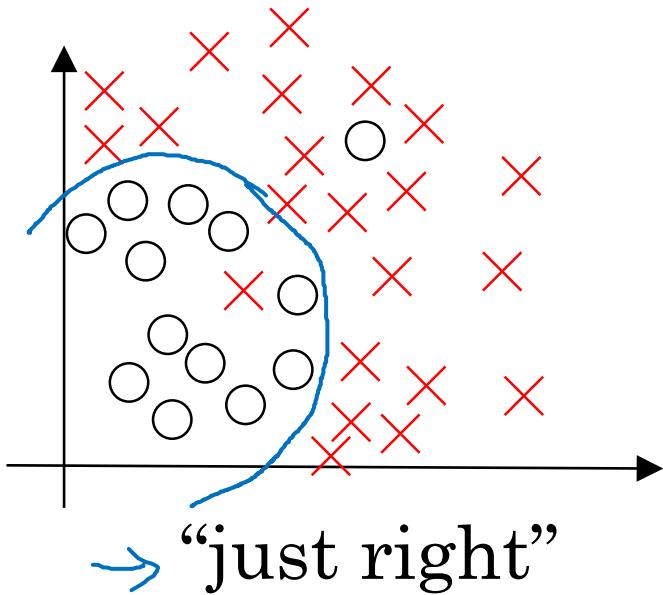
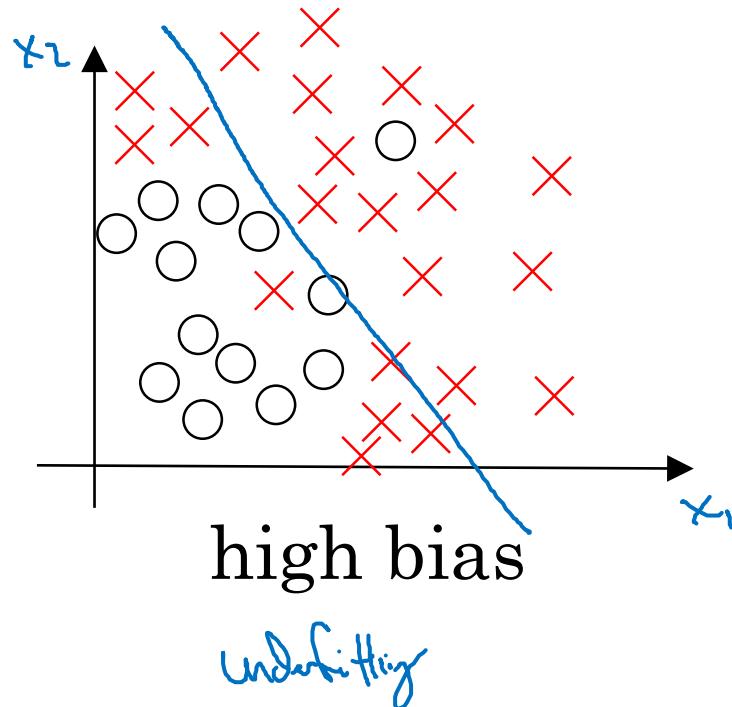


deeplearning.ai

Setting up your
ML application

Bias/Variance

Bias and Variance



Bias and Variance

Cat classification



$$y=1$$



$$y = 0$$

Train set error:

Dev set error:

Herran : $\approx 0\%$

Optimal (Bayes) error: ~~> 0 to~~ 15%

10

10/2

high variance

1

15% ↗

16% 

high bias

1

15.1.

30%

high bias
& high varian

G.S.I.

11

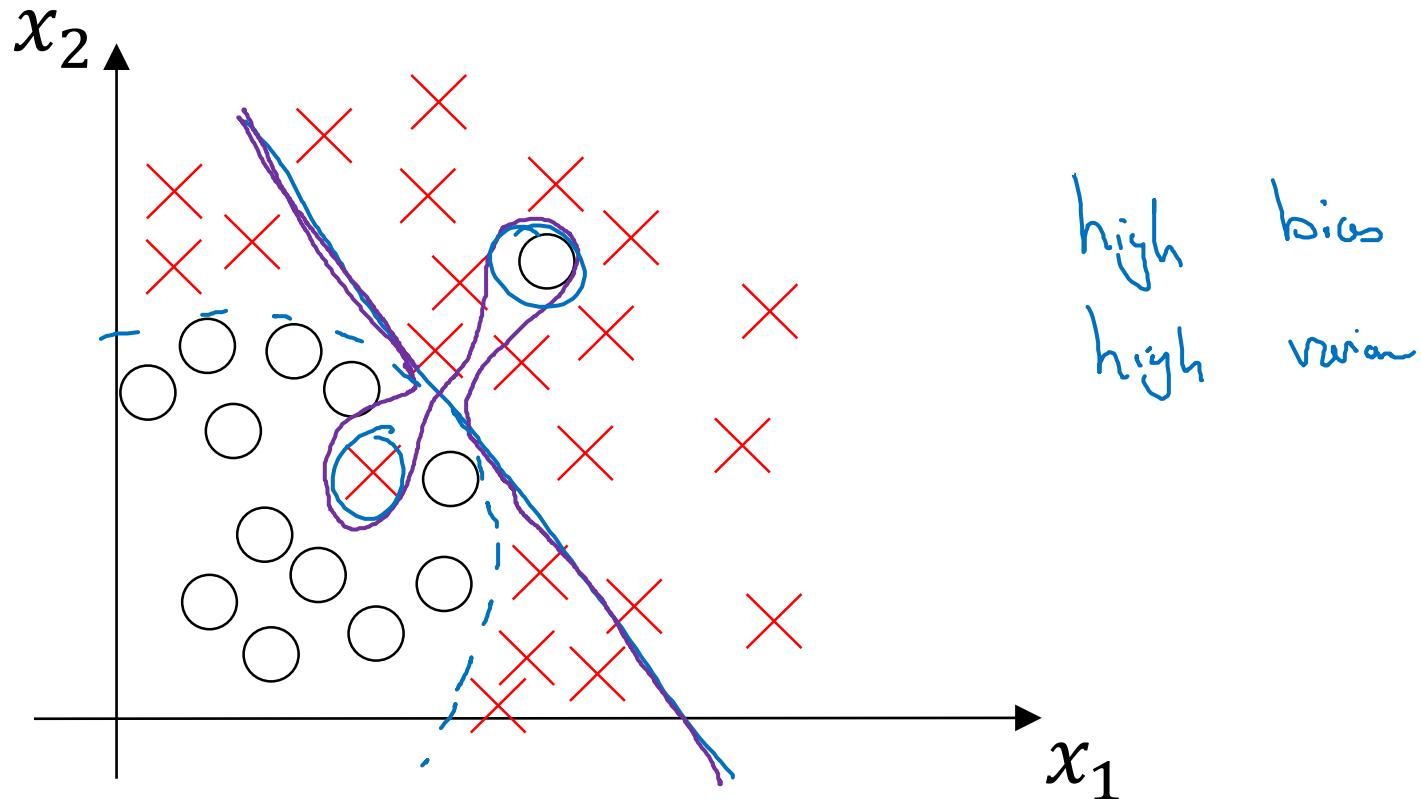
low bias

low variance

1

Berry images

High bias and high variance





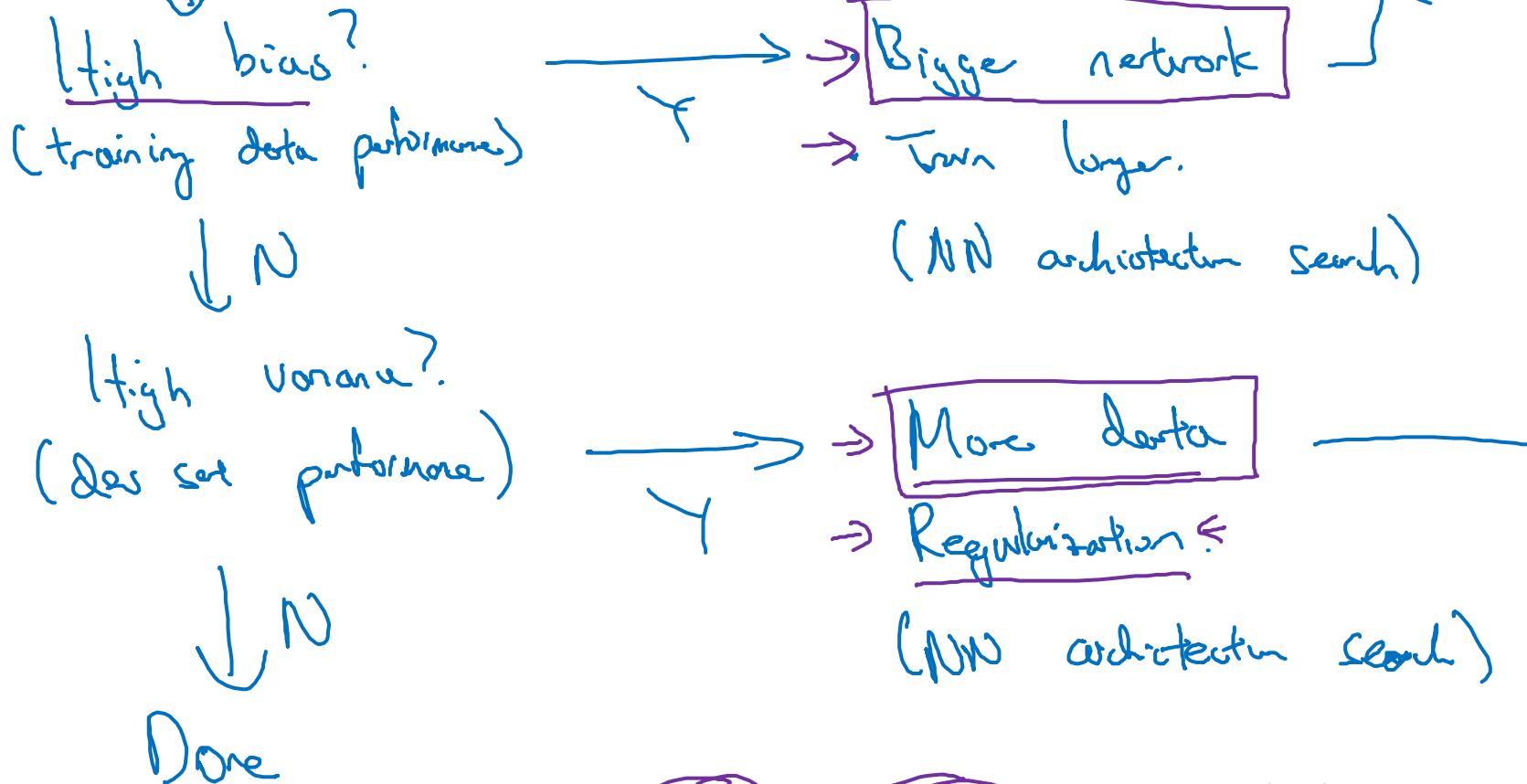
deeplearning.ai

Setting up your ML application

Basic “recipe” for machine learning

Basic “recipe” for machine learning

Basic recipe for machine learning





deeplearning.ai

Regularizing your
neural network

Regularization

Logistic regression

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

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

λ = regularization parameter
lambda lambd

$$J(w, b) = \underbrace{\frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)})}_{\text{L1 regularization}} + \frac{\lambda}{2m} \|w\|_2^2$$

$$+ \cancel{\frac{\lambda}{2m} \|w\|_2^2}$$

omit

L₂ regularization

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

L₁ regularization

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

w will be sparse

Neural network

$$\rightarrow J(w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}) = \underbrace{\frac{1}{m} \sum_{i=1}^m f(y^{(i)}, \hat{y}^{(i)})}_{n^{(1)} \times n^{(L-1)}} + \underbrace{\frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2}_{\text{regularization}}$$

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

$w^{(l)}: (n^{(l)}, n^{(l-1)})$

"Frobenius norm"

$$\|\cdot\|_2^2$$

$$\|\cdot\|_F^2$$

$$dW^{(l)} = \boxed{(\text{from backprop}) + \frac{\lambda}{m} w^{(l)}}$$

$$\rightarrow w^{(l)} := w^{(l)} - \alpha dW^{(l)}$$

$$\frac{\partial J}{\partial w^{(l)}} = dw^{(l)}$$

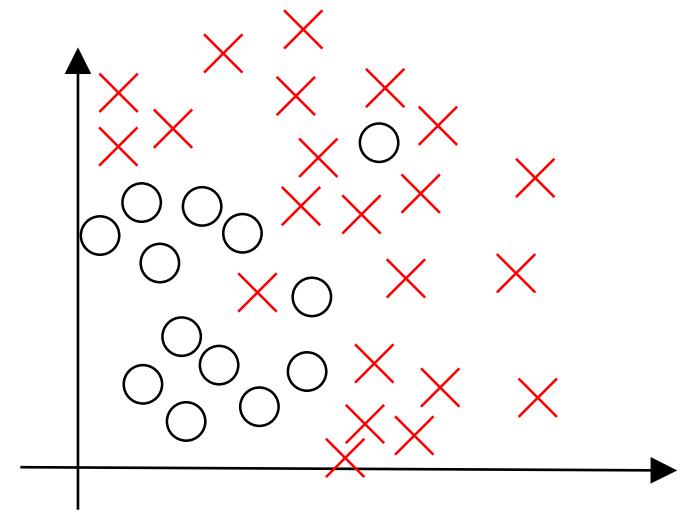
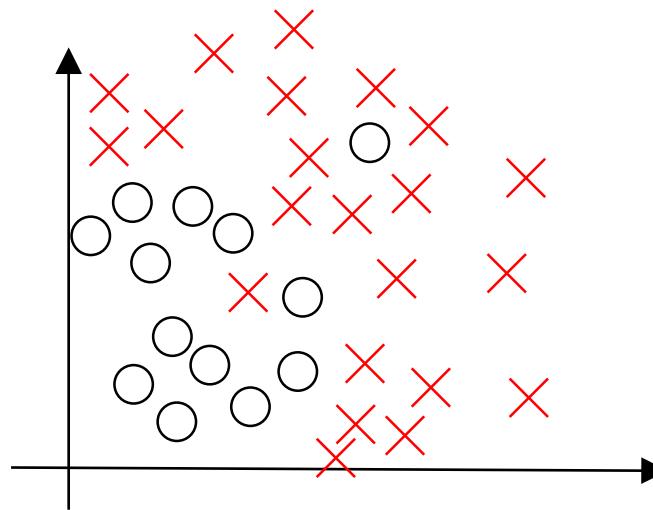
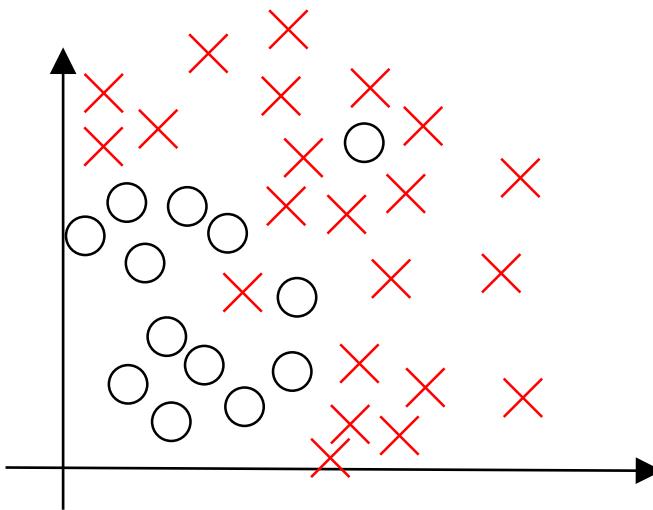
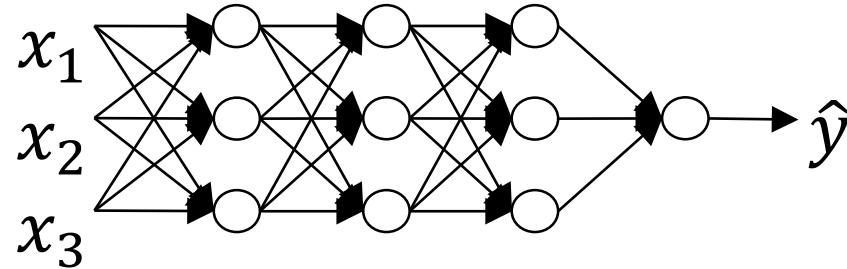
"Weight decay"

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

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

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

How does regularization prevent overfitting?



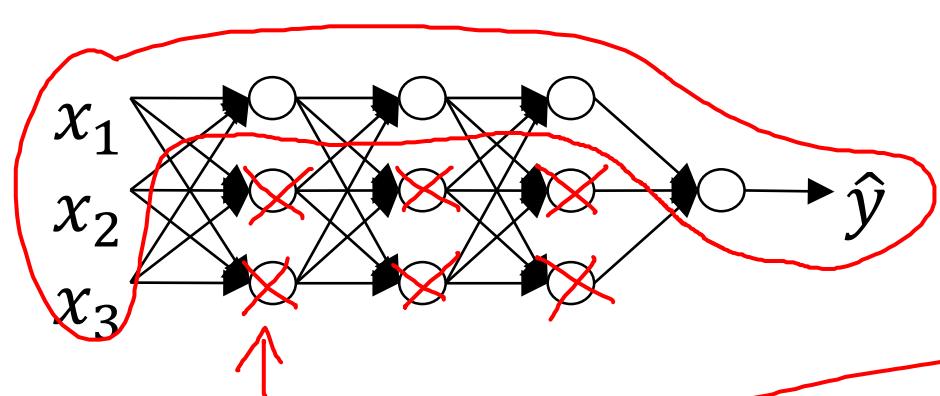


deeplearning.ai

Regularizing your neural network

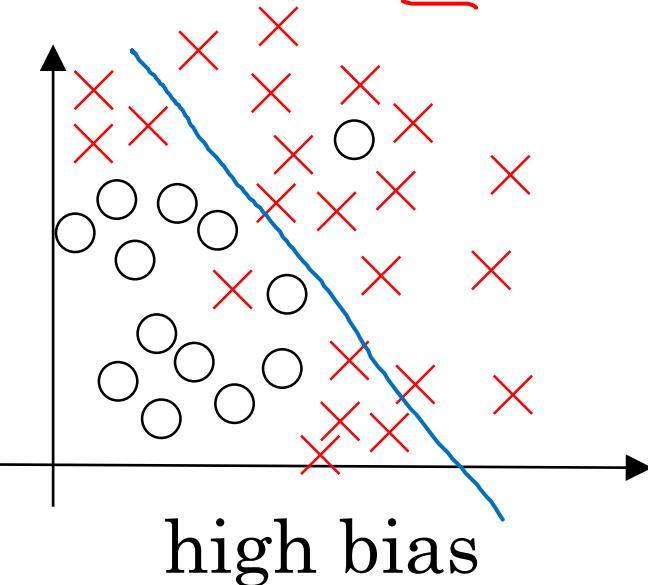
Why regularization reduces overfitting

How does regularization prevent overfitting?

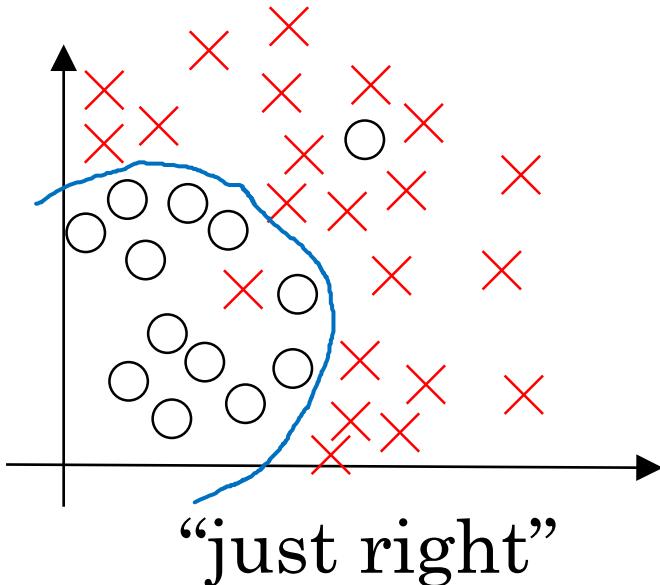


$$J(\boldsymbol{\omega}^{(u)}, \boldsymbol{b}^{(u)}) = \frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|\boldsymbol{w}^{(l)}\|_F^2$$

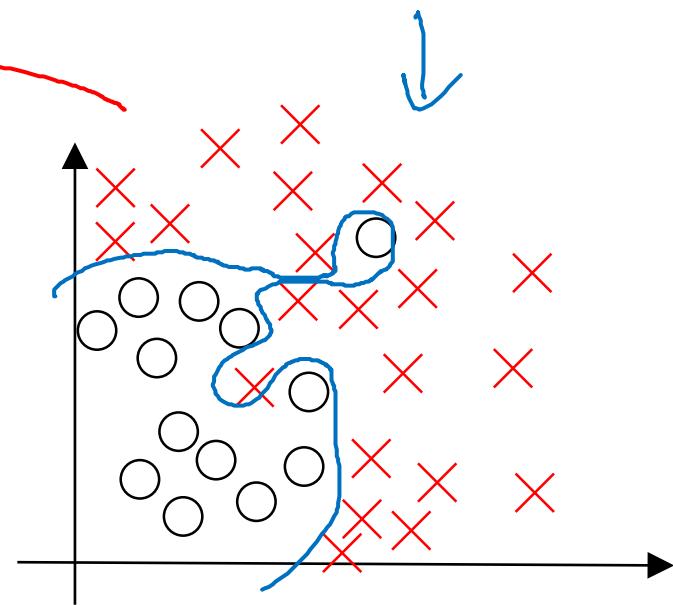
$\boldsymbol{w}^{(u)} \approx 0$



high bias

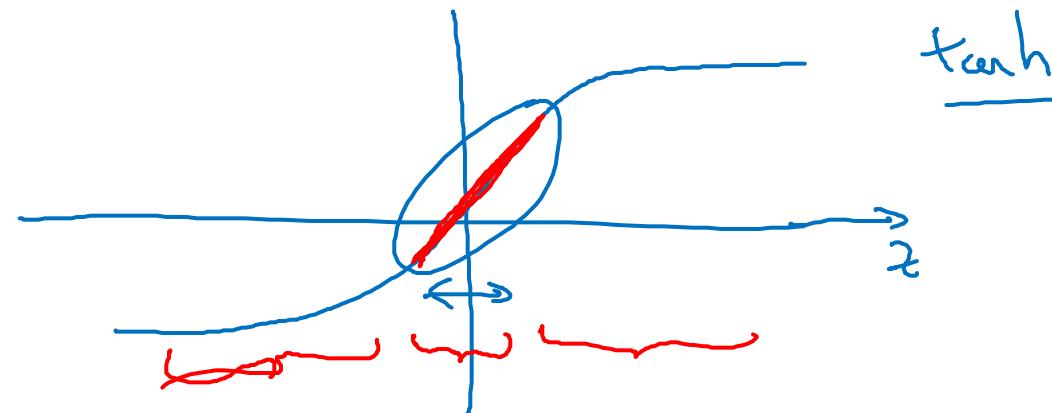


“just right”



high variance

How does regularization prevent overfitting?



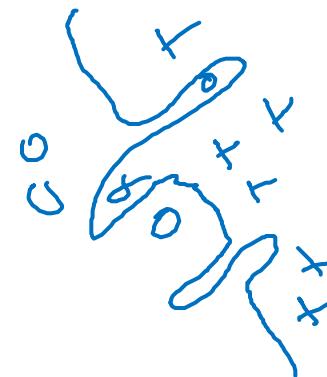
$$\lambda \uparrow$$

$$\underline{w^{[l]}} \downarrow$$

$$z^{[l]} = \underline{w^{[l]}} \underline{a^{[l-1]}} + b^{[l]}$$

Every layer \approx linear.

$$J(\dots) = \boxed{\sum_i L(\hat{y}^{(i)}, y^{(i)})} + \lambda \sum_l \|\underline{w^{[l]}}\|_F^2$$



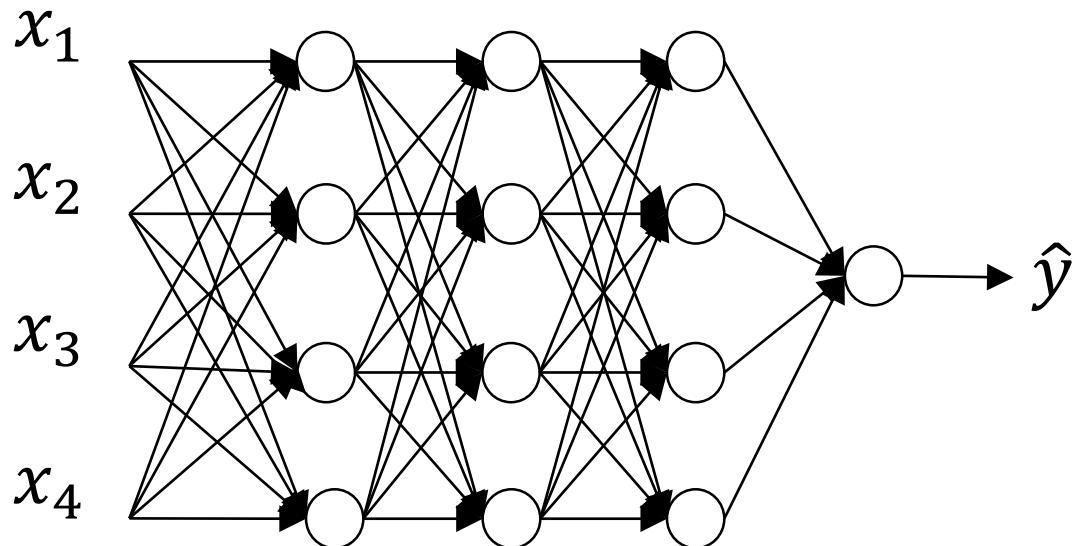


deeplearning.ai

Regularizing your
neural network

Dropout
regularization

Dropout regularization



\uparrow
0.5 \uparrow
0.5 \uparrow
0.5

Implementing dropout (“Inverted dropout”)

Illustrate with layer $\underline{a^3}$. $\text{keep-prob} = \frac{0.8}{x}$ $\underline{0.2}$

$\rightarrow \underline{d^3} = \underline{\text{np.random.rand}(a^3.\text{shape}[0], a^3.\text{shape}[1]) < \text{keep-prob}}$

$\underline{a^3} = \text{np.multiply}(a^3, d^3)$ $\# a^3 * d^3$.

$\rightarrow \underline{a^3 / \cancel{= 0.8} \text{ keep-prob}} \leftarrow$

50 units. \rightsquigarrow 10 units shut off

$$z^{(4)} = w^{(4)} \cdot \underline{\frac{a^{(3)}}{x}} + b^{(4)}$$

\cancel{x} reduced by $\underline{20\%}$.

Test

$$\underline{1 = 0.8}$$

Making predictions at test time

$$a^{(0)} = X$$

No drop out.

$$\uparrow z^{(1)} = w^{(1)} \underline{a^{(0)}} + b^{(1)}$$

$$a^{(1)} = g^{(1)}(\underline{z^{(1)}})$$

$$z^{(2)} = w^{(2)} \underline{a^{(1)}} + b^{(2)}$$

$$a^{(2)} = \dots$$

$$\downarrow \hat{y}$$

λ = keep-prob



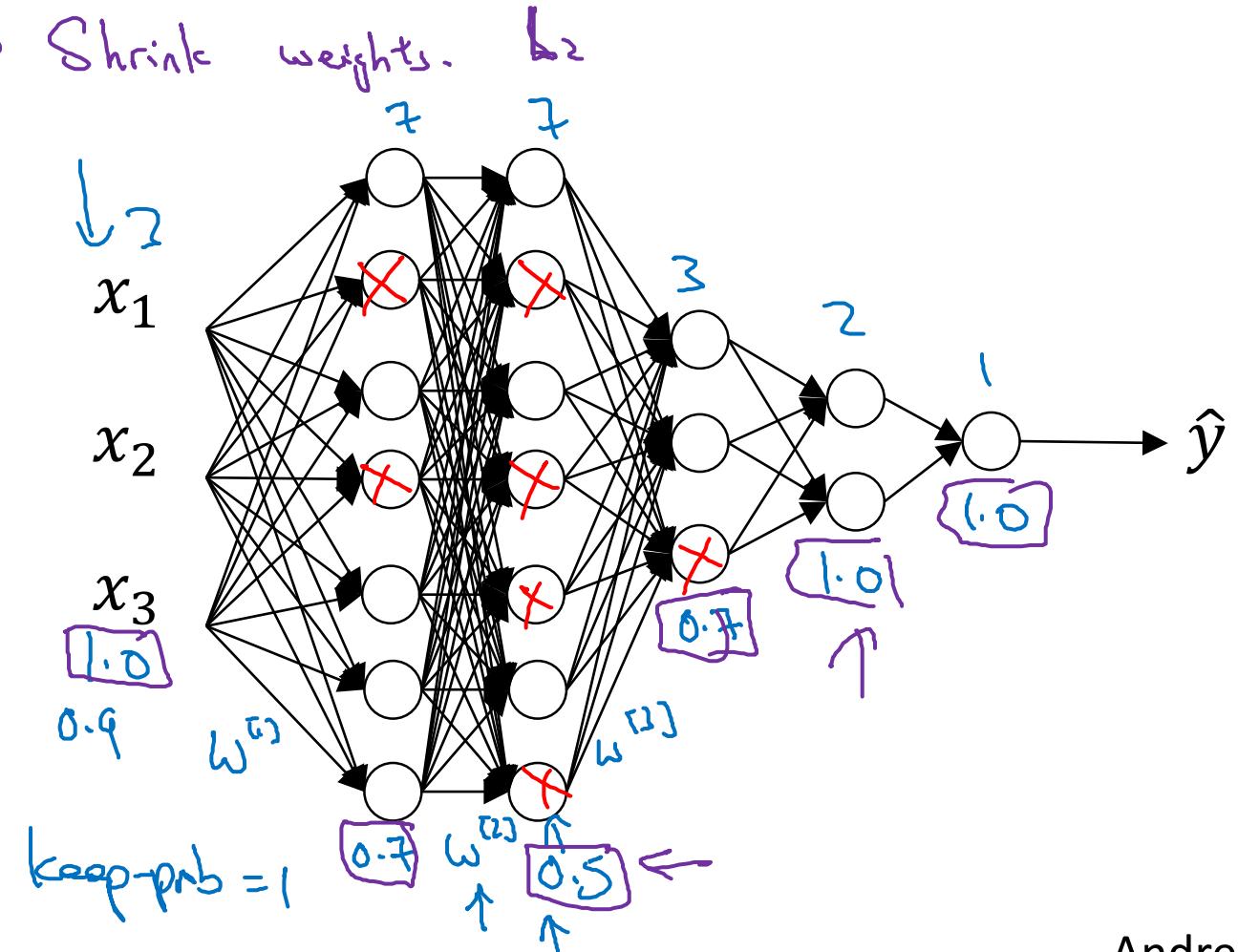
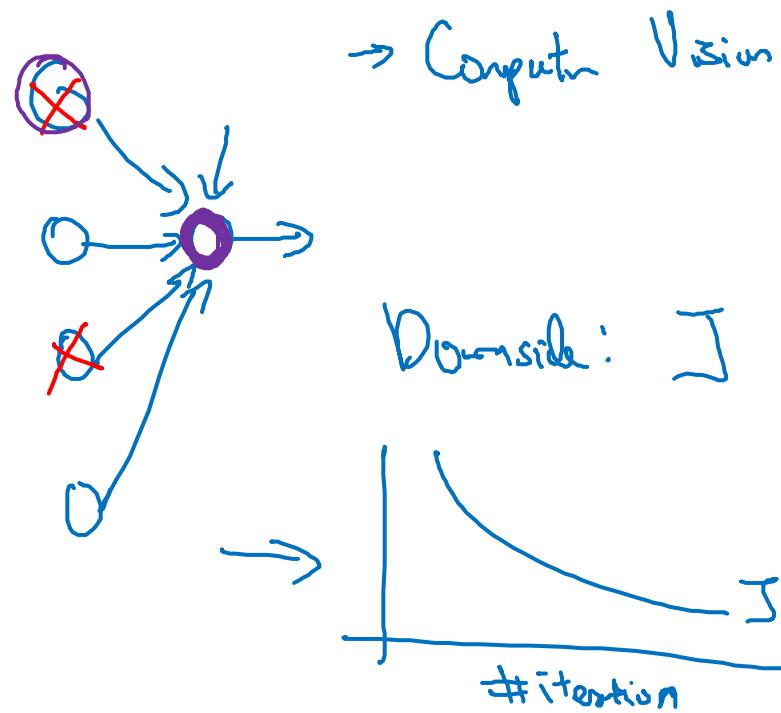
deeplearning.ai

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. \rightsquigarrow Shrink weights. b_2





deeplearning.ai

Regularizing your neural network

Other regularization methods

Data augmentation



4

A large black digit '4' centered on the page.

4

A black silhouette of the digit '4' on a white background.

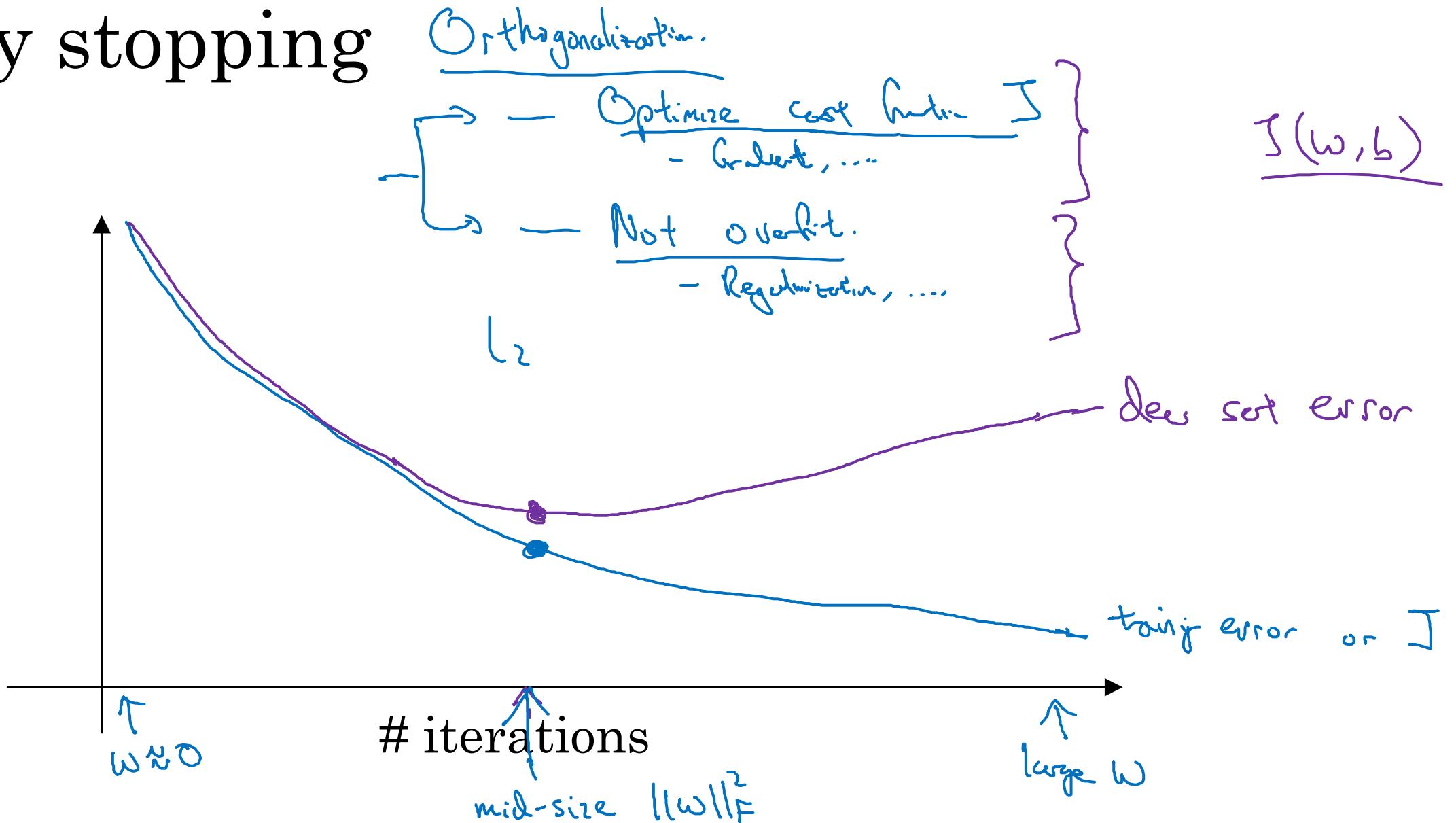
4

A black silhouette of the digit '4' on a white background.

4

A black silhouette of the digit '4' on a white background.

Early stopping





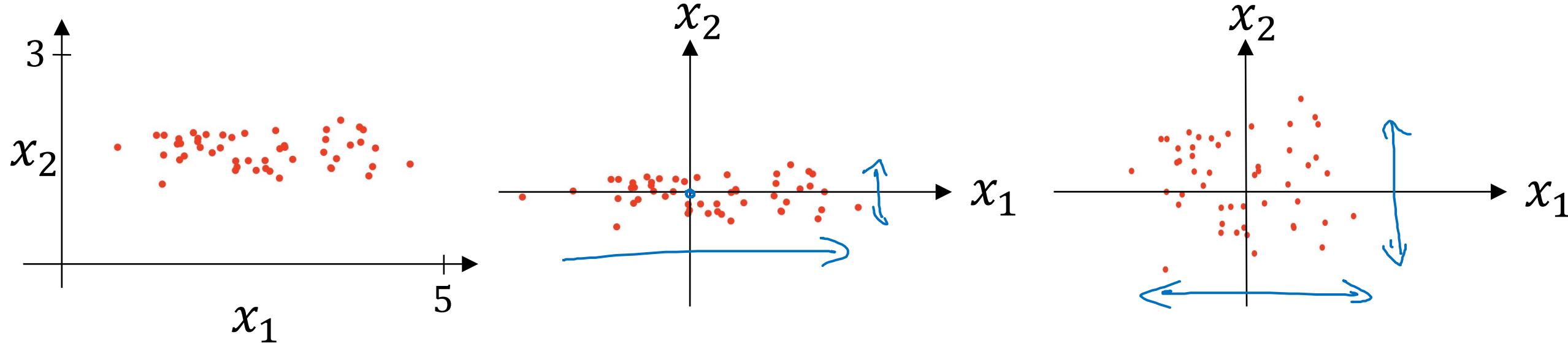
deeplearning.ai

Setting up your
optimization problem

Normalizing inputs

Normalizing training sets

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



Subtract mean:

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

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

Normalize variance

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

~ element-wise

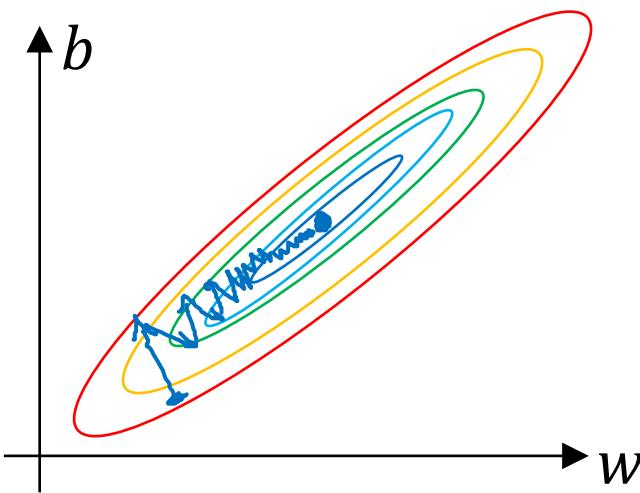
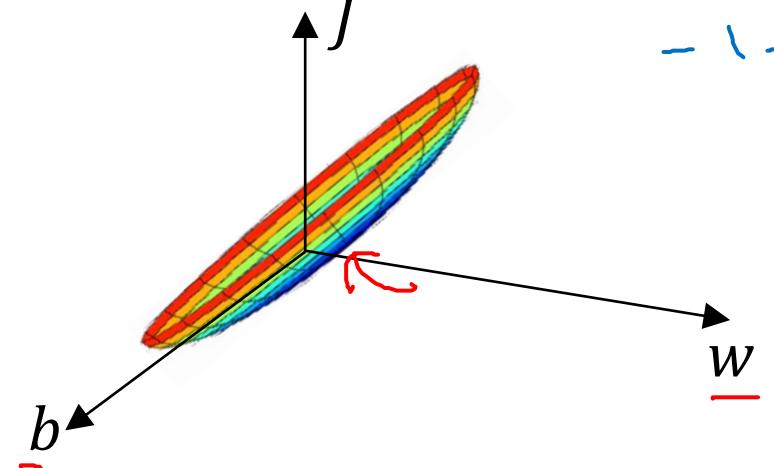
$$\underline{x / = \sigma^{-2}}$$

Use same μ, σ^2 to normalize test set.

Why normalize inputs?

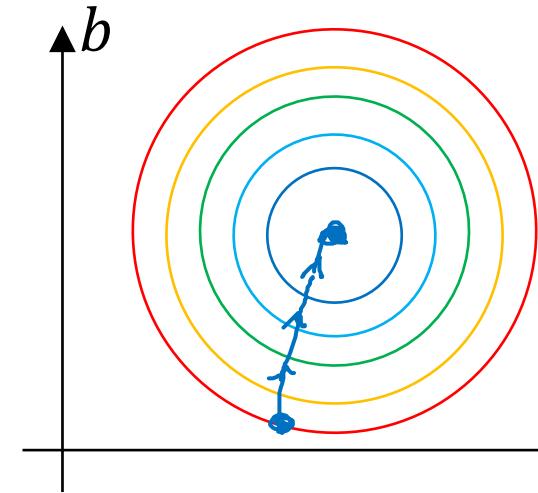
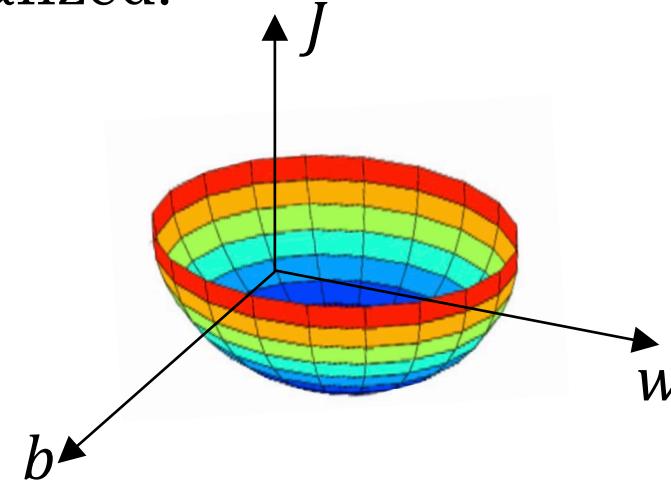
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Unnormalized:
 ω_1 $x_1: \underline{1...100}$ ←
 ω_2 $x_2: \underline{0...1}$ ←
- ... -



$x_1: 0...1$
 $x_2: -1...1$
 $x_3: 1...2$

Normalized:





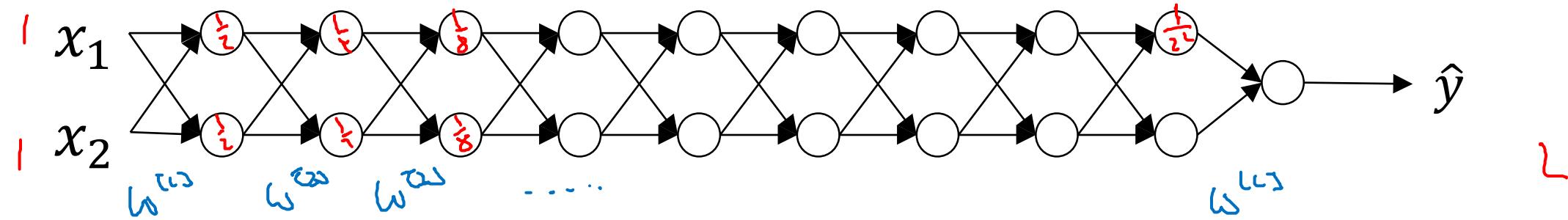
deeplearning.ai

Setting up your
optimization problem

Vanishing/exploding
gradients

Vanishing/exploding gradients

$L=150$



$$\underline{g(z) = z} \quad b^{[L]} = 0$$



$$w^{[1]} > I$$

$$w^{[2]} < I \quad \begin{bmatrix} 0.9 & \\ & 0.9 \end{bmatrix}$$

$$w^{[L]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 6.5 \end{bmatrix}$$

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

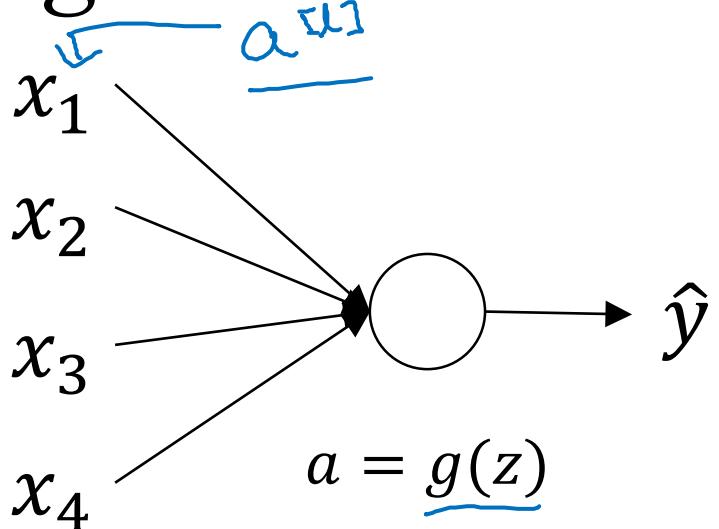
$$a^{[1]} = g(z^{[1]}) = z^{[1]}$$

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

$$\hat{y} = w^{[1]} \begin{bmatrix} 0.5 & 0 \\ 0 & 6.5 \end{bmatrix}^{L-1} x$$

$$1.5^{L-1} x \\ 6.5^{L-1} x$$

Single neuron example

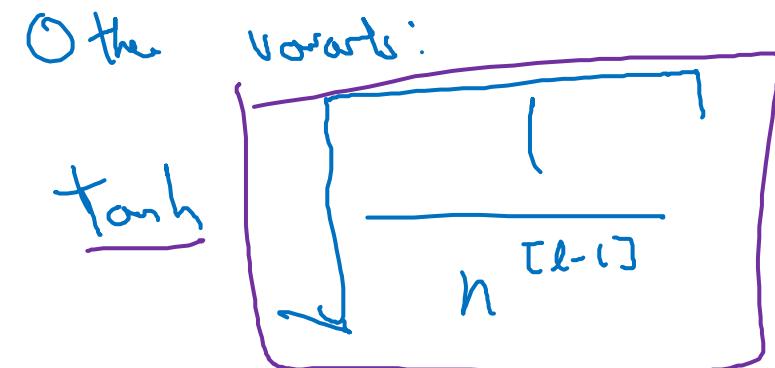


$$z = \underline{w_1}x_1 + \underline{w_2}x_2 + \cdots + \underline{w_n}x_n \quad \cancel{\text{if } n \text{ is large}}$$

Large $n \rightarrow$ Smaller w_i

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

$$\underline{w^{[l]}} = \text{np.random.randn}(\text{shape}) * \sqrt{\frac{2}{n^{[l-1]}}} \quad \text{ReLU}(z) = \text{ReLU}(z)$$



Xavier initialization ↑

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

↑

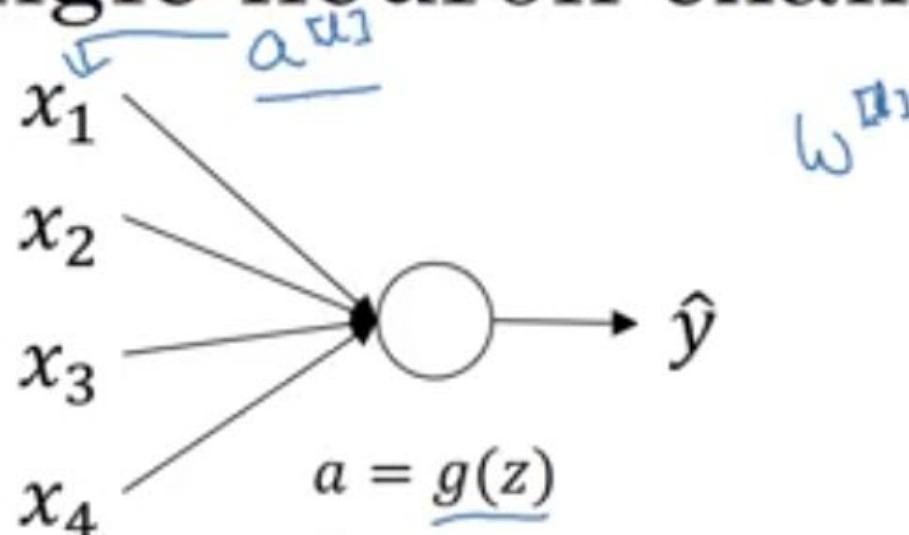


deeplearning.ai

Setting up your
optimization problem

Weight initialization
for deep networks

Single neuron example



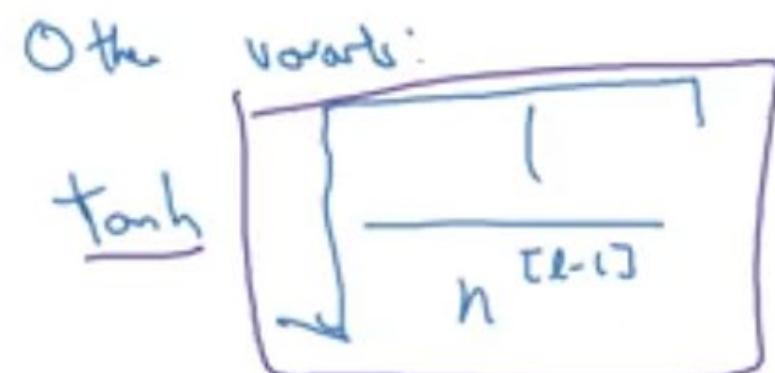
$$z = \underline{w_1}x_1 + \underline{w_2}x_2 + \cdots + \underline{w_n}x_n \quad \cancel{\text{X}}$$

$\text{Large } n \rightarrow \text{Smaller } w_i$

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

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

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



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



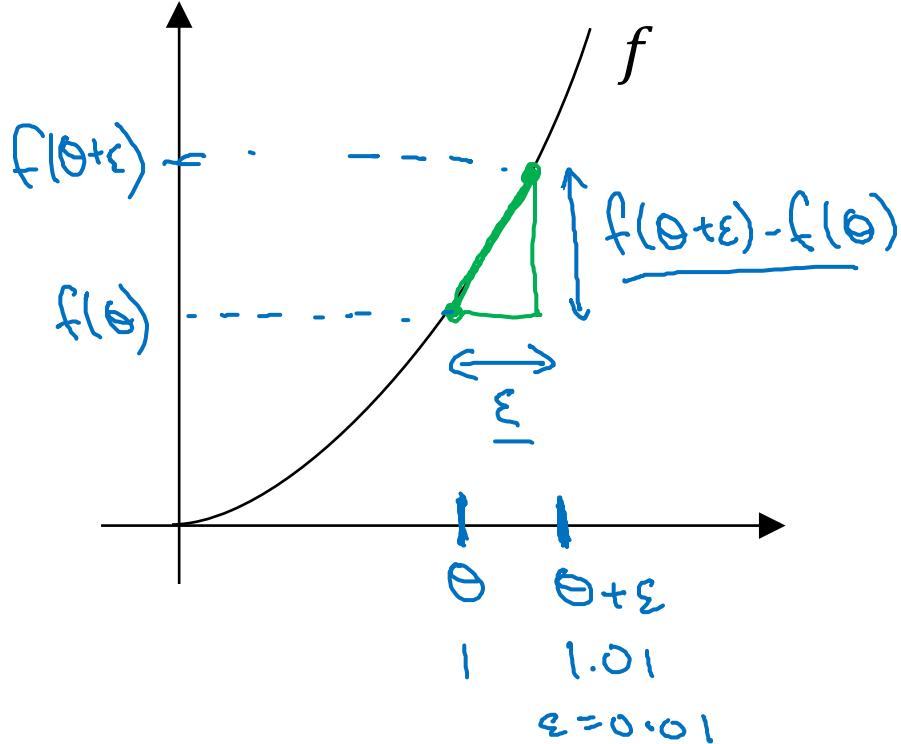
deeplearning.ai

Setting up your optimization problem

Numerical approximation of gradients

Checking your derivative computation

$$\begin{aligned} f(\theta) &= \underline{\theta^3} \\ \theta &\in \mathbb{R}. \\ \text{I} \end{aligned}$$



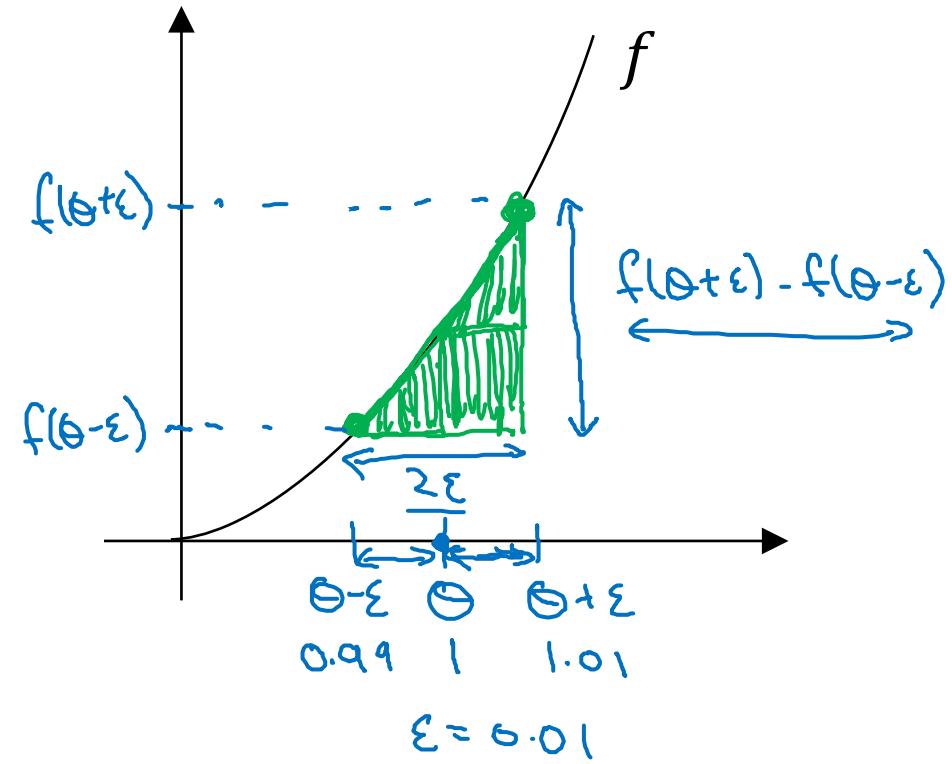
$$\begin{aligned} g(\theta) &= \frac{d}{d\theta} f(\theta) = f'(\theta) \\ g(\theta) &= 3\theta^2. \\ g(1) &= 3 \cdot (1)^2 = 3 \\ \text{when } \theta &= 1 \\ \frac{dw}{db} \end{aligned}$$

$$\begin{aligned} \frac{f(\theta+\epsilon) - f(\theta)}{\epsilon} &\approx g(\theta) \\ \frac{(1.01)^3 - 1^3}{0.01} &= 3.0301 \\ \frac{3.0301}{0.0301} &\approx 3 \end{aligned}$$

$$\begin{aligned} \theta &= 1 \\ \theta + \epsilon &= 1.01 \end{aligned}$$

Checking your derivative computation

$$\underline{f(\theta) = \theta^3}$$



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

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

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

approx error: 0.0001

(prev slide: 3.0301. error: 0.03)

$\left\{ f'(\theta) = \lim_{\varepsilon \rightarrow 0} \frac{f(\theta + \varepsilon) - f(\theta - \varepsilon)}{2\varepsilon} \right.$	$\frac{\mathcal{O}(\varepsilon^2)}{0.01} = \underline{0.0001}$	$\left \frac{f(\theta + \varepsilon) - f(\theta)}{\varepsilon} \right $	$\text{error: } \mathcal{O}(\varepsilon) = \underline{0.01}$
--	--	--	--



deeplearning.ai

Setting up your
optimization problem

Gradient Checking

Gradient check for a neural network

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

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

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

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

Gradient checking (Grad check)

$$J(\theta) = J(\theta_0, \theta_1, \theta_2, \dots)$$

for each i :

$$\rightarrow \underline{d\theta_{\text{approx}}[i]} = \frac{J(\theta_0, \theta_1, \dots, \theta_i + \varepsilon, \dots) - J(\theta_0, \theta_1, \dots, \theta_i - \varepsilon, \dots)}{\varepsilon}$$

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

$$d\theta_{\text{approx}} \stackrel{?}{\approx} d\theta$$

Check

$$\rightarrow \frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2}$$

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

$$\approx \boxed{10^{-7} - \text{great!}} \leftarrow 10^{-5}$$

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



deeplearning.ai

Setting up your optimization problem

Gradient Checking implementation notes

Gradient checking implementation notes

- Don't use in training – only to debug

$$\frac{\partial \theta_{\text{approx}}^{[i]}}{\uparrow} \longleftrightarrow \frac{\partial \theta^{[i]}}{\uparrow}$$

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

$$\frac{\partial b^{[l]}}{\uparrow} \quad \frac{\partial w^{[l]}}{\uparrow}$$

- Remember regularization.

$$J(\theta) = \frac{1}{m} \sum_i f(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_l \|w^{(l)}\|_F^2$$

$\frac{\partial \theta}{\theta} = \text{gradt of } J \text{ wrt. } \theta$

- Doesn't work with dropout.

J

keep-prob = 1.0

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

w, b no



deeplearning.ai

Optimization Algorithms

Mini-batch gradient descent

Batch vs. mini-batch gradient descent

X, Y

$X^{\{t\}}, Y^{\{t\}}$

Vectorization allows you to efficiently compute on m examples.

$$X = [x^{(1)} \ x^{(2)} \ x^{(3)} \ \dots \ x^{(500)} \ | \ x^{(1001)} \ \dots \ x^{(2000)} \ | \ \dots \ | \ \dots \ x^{(m)}]$$

$\underbrace{x^{(1)} \ \dots \ x^{(1000)}}_{(n_x, m)}$ $\underbrace{x^{(2)} \ \dots \ x^{(1000)}}_{(n_x, 1000)}$ \dots $\underbrace{x^{(5,000)} \ \dots \ x^{(m)}}_{(n_x, 1000)}$

$$Y = [y^{(1)} \ y^{(2)} \ y^{(3)} \ \dots \ y^{(1000)} \ | \ y^{(1001)} \ \dots \ y^{(2000)} \ | \ \dots \ | \ \dots \ y^{(m)}]$$

$\underbrace{y^{(1)} \ \dots \ y^{(1000)}}_{(1, m)}$ $\underbrace{y^{(2)} \ \dots \ y^{(1000)}}_{(1, 1000)}$ \dots $\underbrace{y^{(5,000)} \ \dots \ y^{(m)}}_{(1, 1000)}$

What if $m = 5,000,000$?

5,000 mini-batches of 1,000 each

Mini-batch t : $X^{\{t\}}, Y^{\{t\}}$

$x^{(i)}$
 $z^{[l]}$
 $X^{\{t\}}, Y^{\{t\}}$

Mini-batch gradient descent

repeat {
for $t = 1, \dots, 5000$ {

Forward prop on $X^{\{t\}}$.

$$Z^{(l)} = W^{(l)} X^{\{t\}} + b^{(l)}$$

$$A^{(l)} = g^{(l)}(Z^{(l)})$$

:

$$A^{(L)} = g^{(L)}(Z^{(L)})$$

Compute cost $J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^L f(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_l \|W^{(l)}\|_F^2$.

Backprop to compute gradients wrt $J^{\{t\}}$ (using $(X^{\{t\}}, Y^{\{t\}})$)

$$W^{(l)} := W^{(l)} - \alpha \nabla W^{(l)}, \quad b^{(l)} := b^{(l)} - \alpha \nabla b^{(l)}$$

3 } 3 }

"1 epoch"
└ pass through training set.

1 step of gradient descent
using $\frac{X^{\{t+1\}}}{Y^{\{t+1\}}}$
(as if $t=5000$)

X, Y



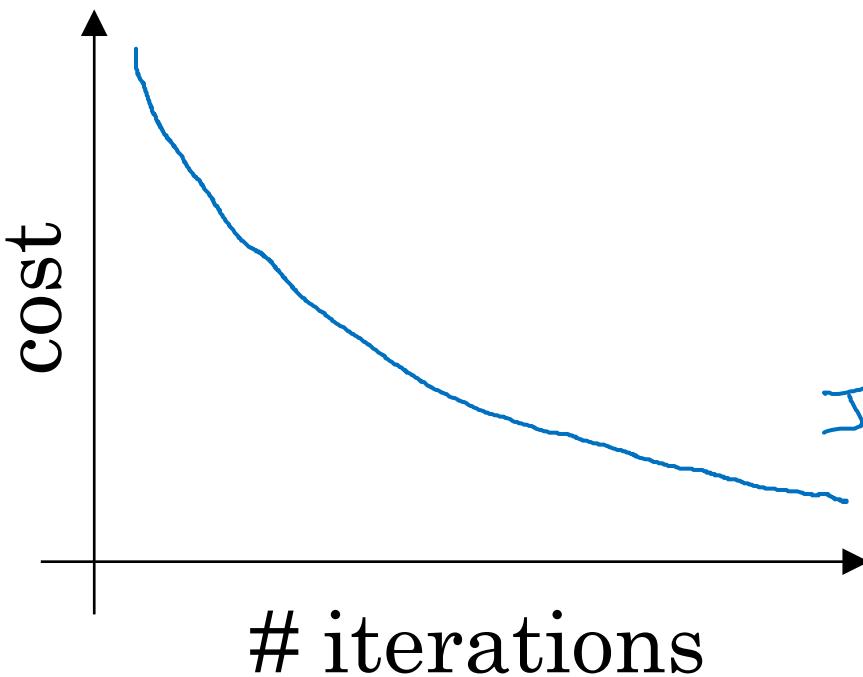
deeplearning.ai

Optimization Algorithms

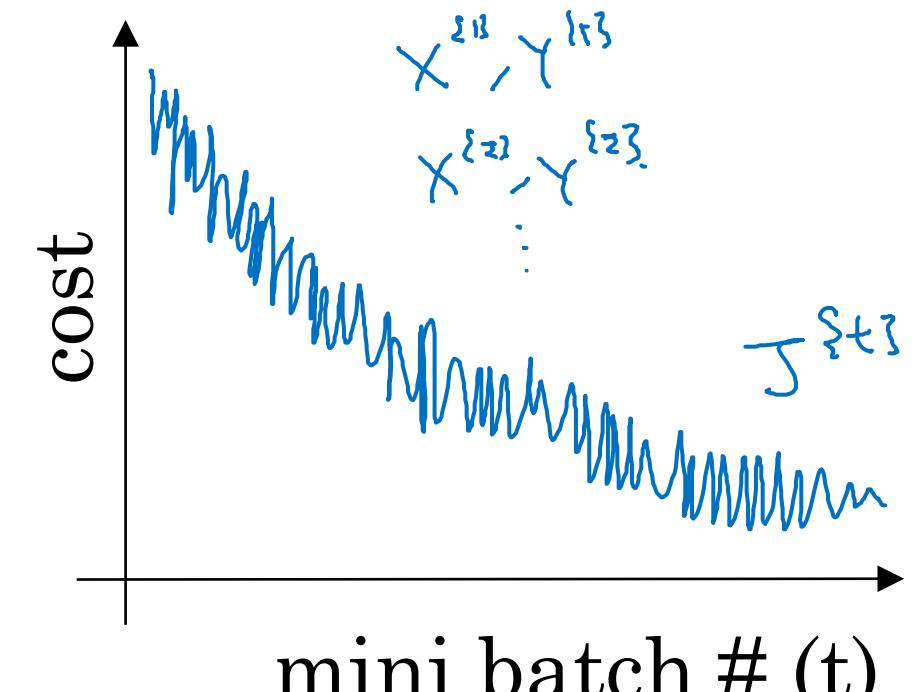
Understanding mini-batch gradient descent

Training with mini batch gradient descent

Batch gradient descent



Mini-batch gradient descent



Plot J^{st} computed using $X^{\{st\}}, Y^{\{st\}}$

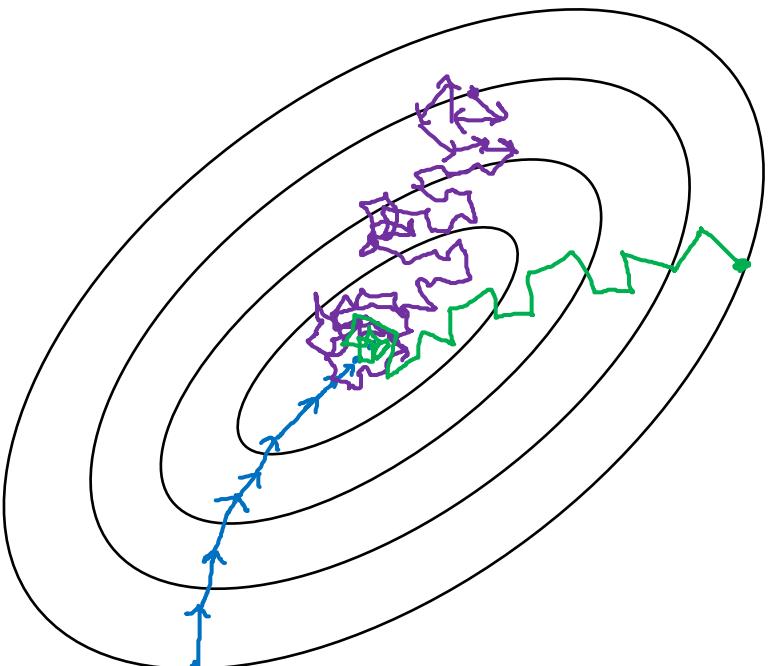
Choosing your mini-batch size

→ If mini-batch size = m : Batch gradient descent.

$$(X^{\{1\}}, Y^{\{1\}}) = (X, Y)$$

→ If mini-batch size = 1 : Stochastic gradient descent. Every example is its own
 $(X^{(1)}, Y^{(1)}) = (x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})$ mini-batch.

In practice: Somehw in-between 1 and m



Stochastic
gradient
descent

{
Use speedup
from vectorization

In-between
(mini-batch size
not too big/small)

{
Fastest learning:

- Vectorization.
($n \approx 1000$)
- Make passes without
processing entire training set.

Batch
gradient descent
(mini-batch size = m)

{
Two long
per iteration

Choosing your mini-batch size

If small training set : Use batch gradient descent.
 $(m \leq 2000)$

Typical mini-batch sizes:

$$\rightarrow 64, 128, 256, 512 \quad \frac{1024}{2^{10}}$$

$2^6 \quad 2^7 \quad 2^8 \quad 2^9$



Make sure mini-batch fits in CPU/GPU memory.

$$X^{\{t\}}, Y^{\{t\}}$$



deeplearning.ai

Optimization Algorithms

Exponentially weighted averages

Temperature in London

$$\theta_1 = 40^{\circ}\text{F} \quad 4^{\circ}\text{C} \quad \leftarrow$$

$$\theta_2 = 49^{\circ}\text{F} \quad 9^{\circ}\text{C}$$

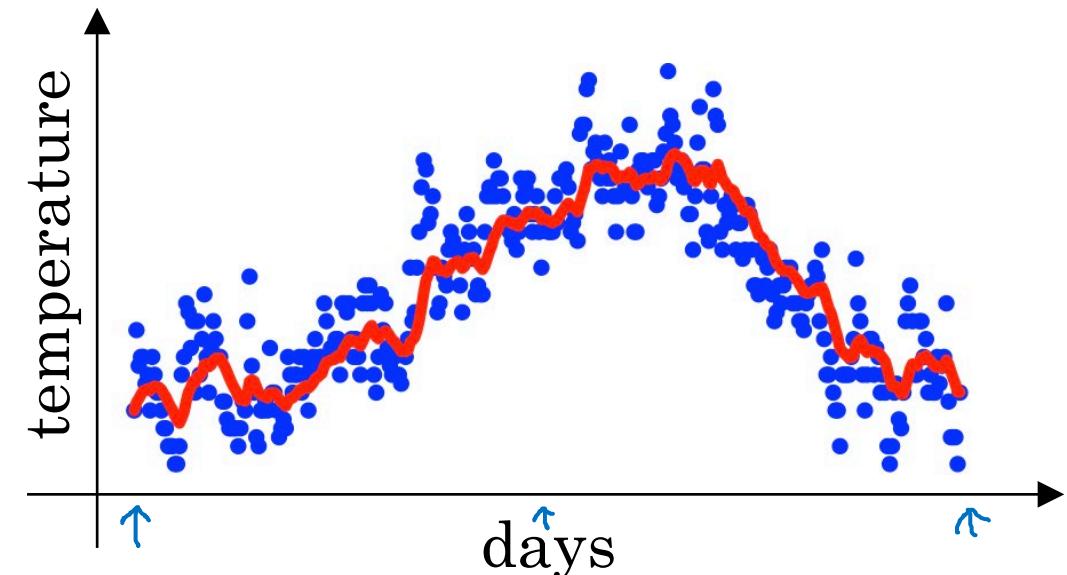
$$\theta_3 = 45^{\circ}\text{F} \quad \vdots$$

⋮

$$\theta_{180} = 60^{\circ}\text{F} \quad 15^{\circ}\text{C}$$

$$\theta_{181} = 56^{\circ}\text{F} \quad \vdots$$

⋮



$$V_0 = 0$$

$$V_1 = 0.9 V_0 + 0.1 \theta_1$$

$$V_2 = 0.9 V_1 + 0.1 \theta_2$$

$$V_3 = 0.9 V_2 + 0.1 \theta_3$$

⋮

$$V_t = 0.9 V_{t-1} + 0.1 \theta_t$$

Exponentially weighted averages ^{Moving}

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t \leftarrow$$

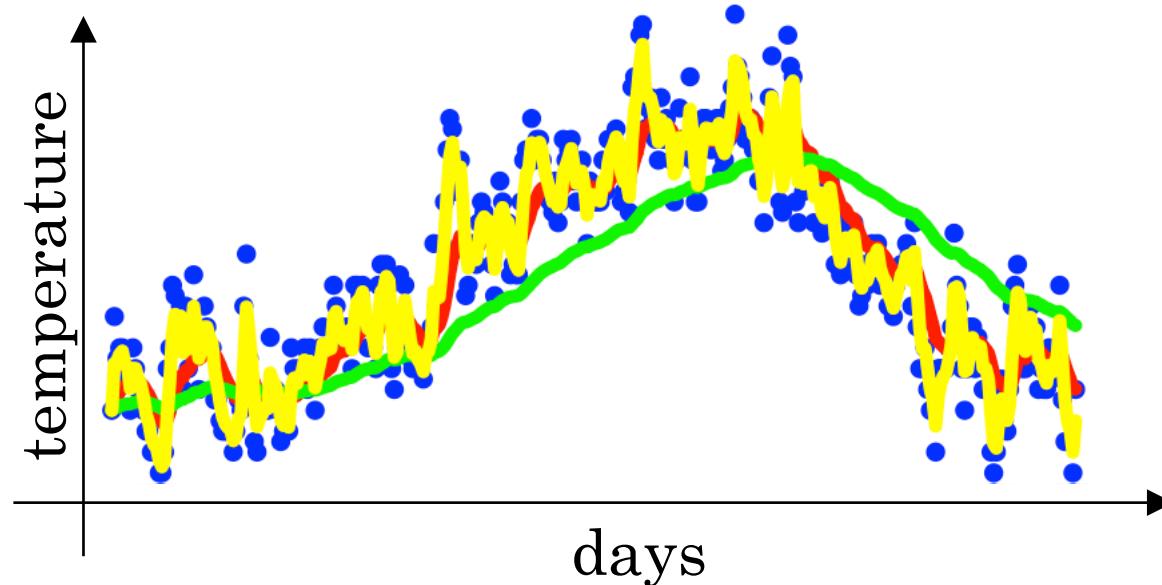
$\beta = 0.9$: ≈ 10 days' temperatur.

$\beta = 0.98$: ≈ 50 days

$\beta = 0.5$: ≈ 2 days

V_t is approximately
average over
 $\rightarrow \approx \frac{1}{1-\beta}$ days'
temperature.

$$\frac{1}{1-0.98} = 50$$





deeplearning.ai

Optimization Algorithms

Understanding exponentially weighted averages

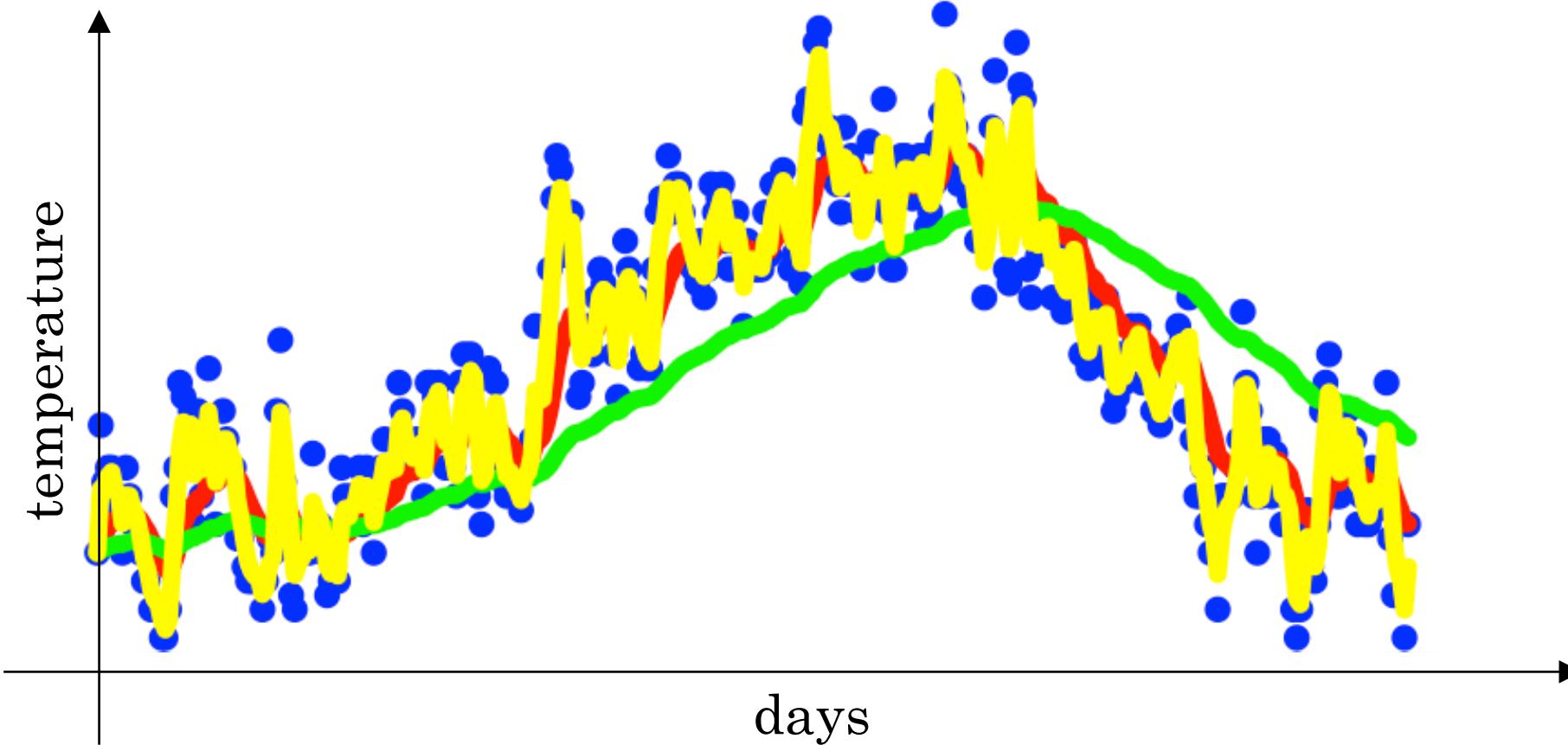
Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

$$\beta = 0.9$$

$$0.98$$

$$0.5$$



Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

$$v_{100} = 0.9v_{99} + 0.1\theta_{100}$$

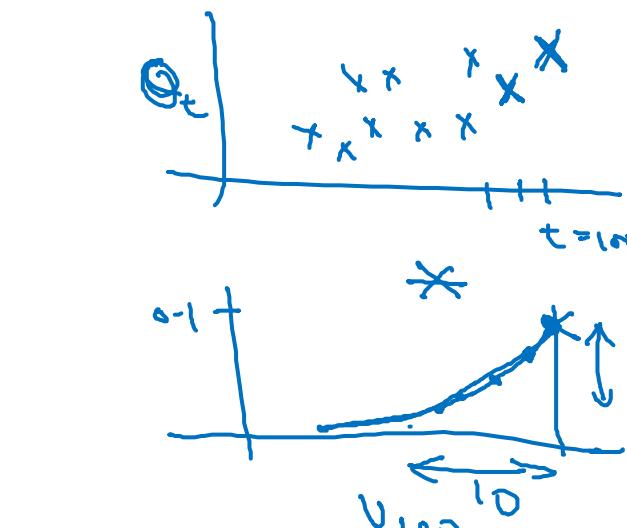
$$v_{99} = 0.9v_{98} + 0.1\theta_{99}$$

$$v_{98} = 0.9v_{97} + 0.1\theta_{98}$$

...

$$\begin{aligned} \underline{v_{100}} &= 0.1 \underline{\theta_{100}} + 0.9 \cancel{(0.1 \underline{\theta_{99}})} + 0.9 \cancel{(0.1 \underline{\theta_{98}})} \\ &= 0.1 \underline{\theta_{100}} + \underline{0.1 \times 0.9 \cdot \theta_{99}} + \underline{0.1 (0.9)^2 \theta_{98}} + \underline{0.1 (0.9)^3 \theta_{97}} + \underline{0.1 (0.9)^4 \theta_{96}} + \dots \end{aligned}$$

$$\underline{0.9^{10}} \approx \underline{0.35} \approx \frac{1}{e}$$



$$\frac{1}{1-\beta}$$

$$\Sigma = 1 - \beta$$

$$0.1 \underline{\theta_{98}} + 0.9 \underline{v_{97}}$$

$$\frac{(1-\epsilon)^{1/\epsilon}}{0.9} = \frac{1}{e}$$

0.98?

$$\epsilon = 0.02 \rightarrow \underline{0.98^{50}} \approx \frac{1}{e}$$

Implementing exponentially weighted averages

$$v_0 = 0$$

$$v_1 = \beta v_0 + (1 - \beta) \theta_1$$

$$v_2 = \beta v_1 + (1 - \beta) \theta_2$$

$$v_3 = \beta v_2 + (1 - \beta) \theta_3$$

...

$$\left| \begin{array}{l} v_0 := 0 \\ v_0 := \beta v + (1-\beta) \theta_1 \\ v_0 := \beta v + (1-\beta) \theta_2 \\ \vdots \\ \hline \end{array} \right. \rightarrow v_0 = 0$$

Repeat {

Get next θ_t

$$v_0 := \beta v_0 + (1-\beta) \theta_t \quad \leftarrow \}$$

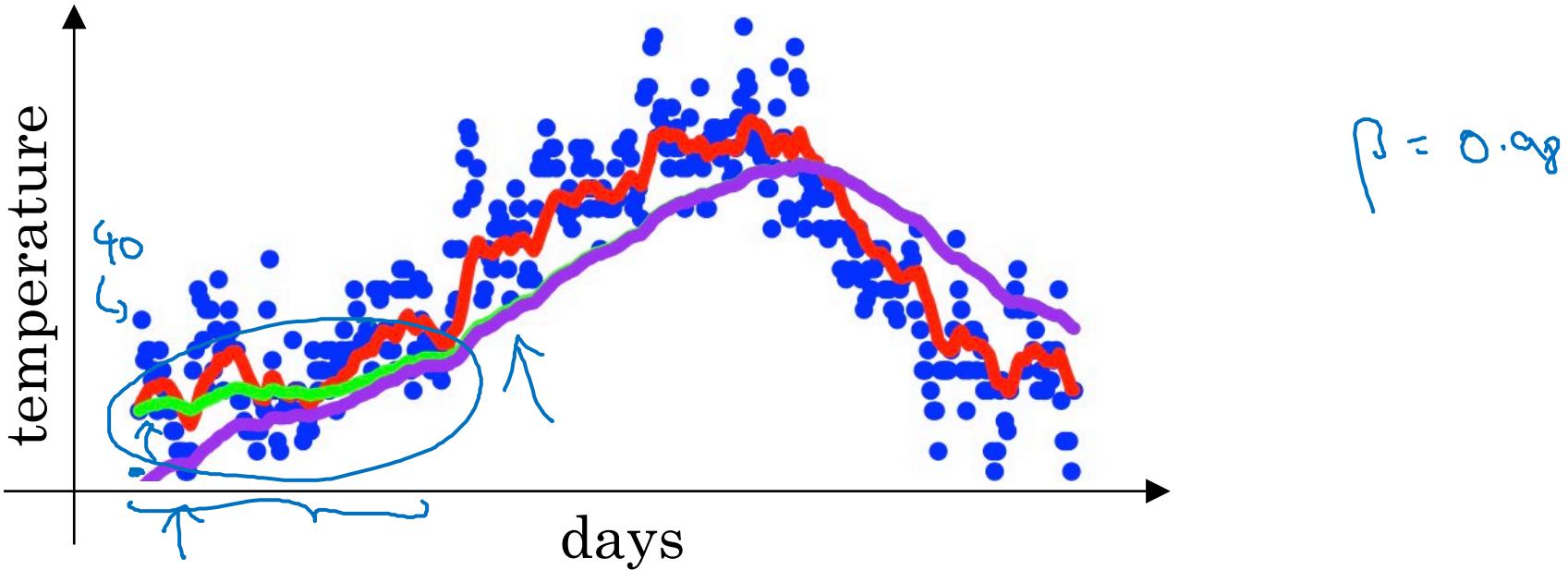


deeplearning.ai

Optimization Algorithms

Bias correction
in exponentially
weighted average

Bias correction



$$\rightarrow v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

$$v_0 = 0$$

$$v_1 = \cancel{0.98v_0} + \underline{0.02\theta_1}$$

$$\begin{aligned} v_2 &= 0.98 v_1 + 0.02 \theta_2 \\ &= 0.98 \times 0.02 \times \theta_1 + 0.02 \theta_2 \\ &= \underline{0.0196\theta_1} + \underline{0.02\theta_2} \end{aligned}$$

$$\frac{v_t}{1 - \beta^t}$$

$$t=2: 1 - \beta^t = 1 - (0.98)^2 = 0.0396$$

$$\frac{v_2}{0.0396} =$$

$$\frac{\underline{0.0196\theta_1} + \underline{0.02\theta_2}}{0.0396}$$

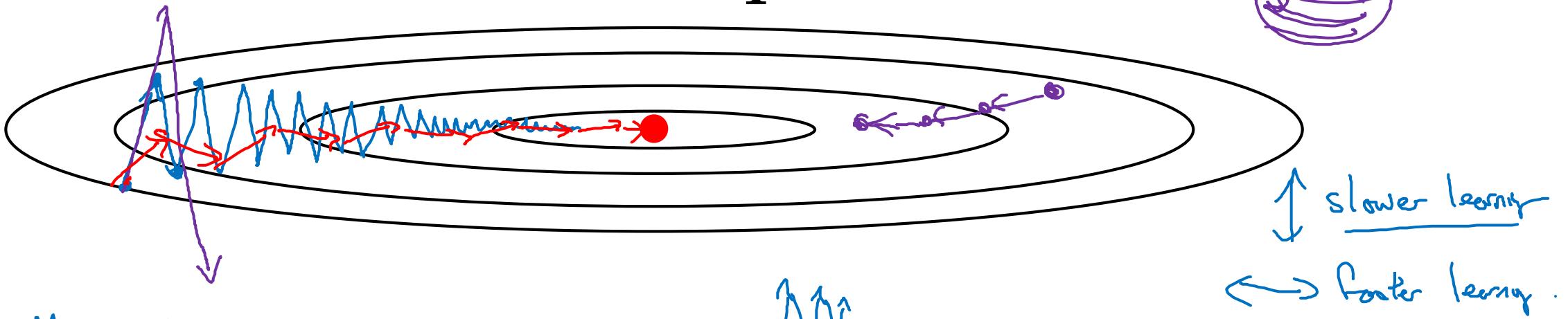


deeplearning.ai

Optimization Algorithms

Gradient descent with momentum

Gradient descent example



Momentum:

On iteration t :

Compute $\Delta w, \Delta b$ on current mini-batch.

$$v_{dw} = \beta v_{dw} + (1-\beta) \Delta w$$

$$v_{db} = \beta v_{db} + (1-\beta) \Delta b$$

Friction ↑ velocity

$$w := w - \alpha v_{dw}$$

$$"v_\theta = \beta v_\theta + (1-\beta) \theta_t"$$



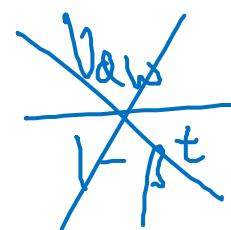
Implementation details

$$v_{dw} = 0, v_{db} = 0$$

On iteration t :

Compute dW, db on the current mini-batch

$$\begin{aligned} \rightarrow v_{dw} &= \beta v_{dw} + (1 - \beta) dW \\ \rightarrow v_{db} &= \beta v_{db} + (1 - \beta) db \end{aligned} \quad \left| \begin{array}{l} v_{dw} = \beta v_{dw} + dW \leftarrow \\ v_{db} = \beta v_{db} + db \end{array} \right.$$
$$W = W - \underbrace{\alpha v_{dw}}, b = \underbrace{b - \alpha v_{db}}$$



Hyperparameters: α, β

$$\beta = 0.9$$

average over last ≈ 10 gradients

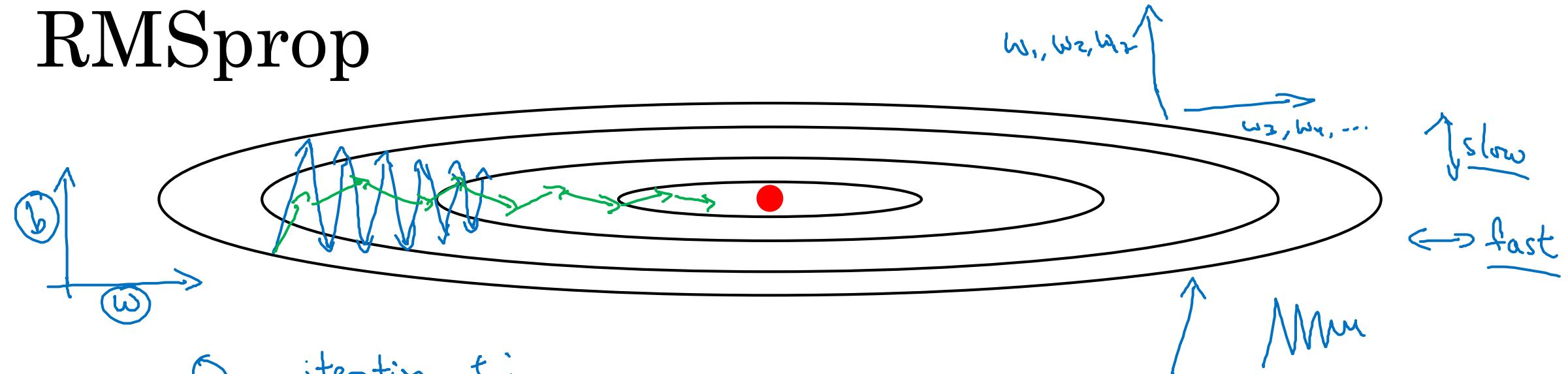


deeplearning.ai

Optimization Algorithms

RMSprop

RMSprop



On iteration t :

Compute $d\mathbf{w}, d\mathbf{b}$ on current mini-batch

$$\underline{S_{d\mathbf{w}}} = \beta_2 \underline{S_{d\mathbf{w}}} + (1-\beta_2) \underline{d\mathbf{w}^2} \leftarrow \begin{matrix} \text{element-wise} \\ \text{small} \end{matrix}$$

$$\rightarrow \underline{S_{d\mathbf{b}}} = \beta_2 \underline{S_{d\mathbf{b}}} + (1-\beta_2) \underline{d\mathbf{b}^2} \leftarrow \text{large}$$

$$\mathbf{w} := \mathbf{w} - \frac{\alpha}{\sqrt{\underline{S_{d\mathbf{w}}} + \epsilon}} \underline{d\mathbf{w}}$$

$$\mathbf{b} := \mathbf{b} - \frac{\alpha}{\sqrt{\underline{S_{d\mathbf{b}}} + \epsilon}} \underline{d\mathbf{b}}$$

$$\epsilon = 10^{-8}$$



deeplearning.ai

Optimization Algorithms

Adam optimization algorithm

Adam optimization algorithm

$$V_{dw} = 0, S_{dw} = 0. \quad V_{db} = 0, S_{db} = 0$$

On iteration t :

Compute $\delta w, \delta b$ using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) \delta w, \quad V_{db} = \beta_1 V_{db} + (1 - \beta_1) \delta b \quad \leftarrow \text{"moment"} \beta_1$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) \delta w^2, \quad S_{db} = \beta_2 S_{db} + (1 - \beta_2) \delta b \quad \leftarrow \text{"RMSprop"} \beta_2$$

$$\text{yhat} = \text{np.array}([.9, 0.2, 0.1, .4, .9])$$

$$V_{dw}^{\text{corrected}} = V_{dw} / (1 - \beta_1^t), \quad V_{db}^{\text{corrected}} = V_{db} / (1 - \beta_1^t)$$

$$S_{dw}^{\text{corrected}} = S_{dw} / (1 - \beta_2^t), \quad S_{db}^{\text{corrected}} = S_{db} / (1 - \beta_2^t)$$

$$w := w - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}}} + \epsilon}$$

$$b := b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}}} + \epsilon}$$

Hyperparameters choice:

- α : needs to be tune
- β_1 : 0.9 $\rightarrow (\underline{dw})$
- β_2 : 0.999 $\rightarrow (\underline{dw^2})$
- ϵ : 10^{-8}

Adam: Adaptive moment estimation



Adam Coates



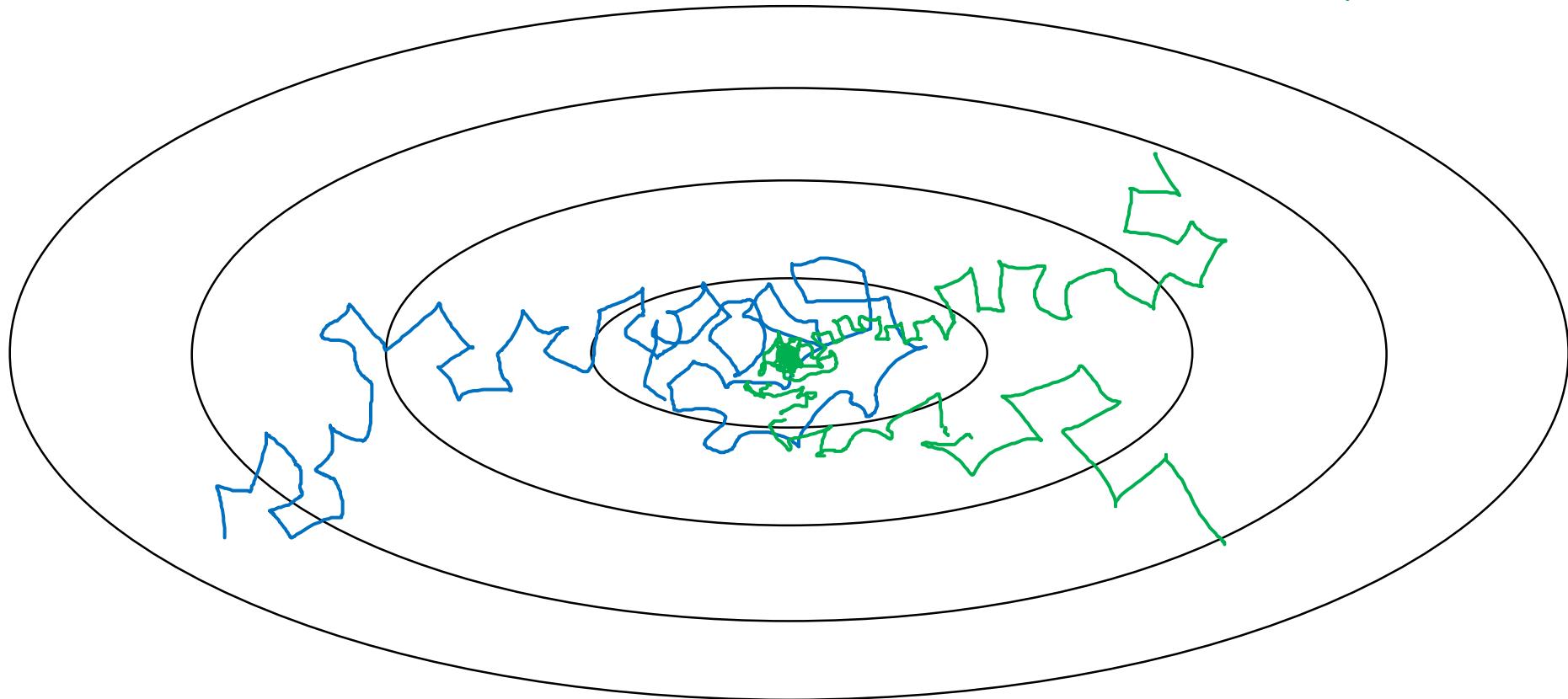
deeplearning.ai

Optimization Algorithms

Learning rate decay

Learning rate decay

Slowly reduce λ

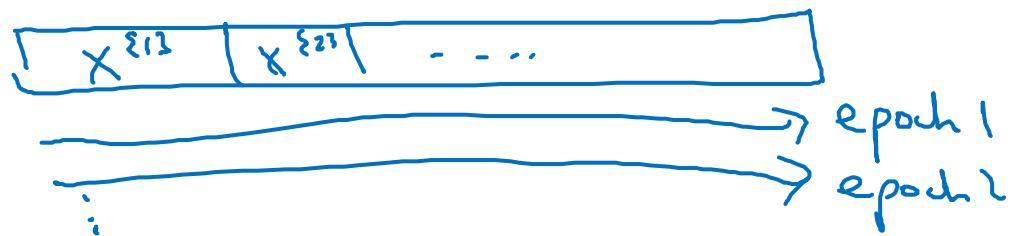


Learning rate decay

1 epoch = 1 pass through data.

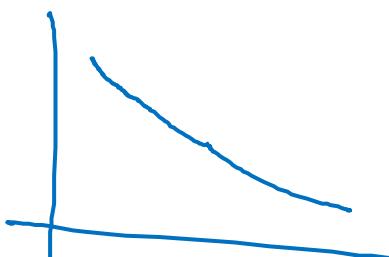
$$\alpha = \frac{\alpha_0}{1 + \text{decay-rate} * \text{epoch-num}}$$

Epoch	α
1	0.1
2	0.67
3	0.5
4	0.4
:	:

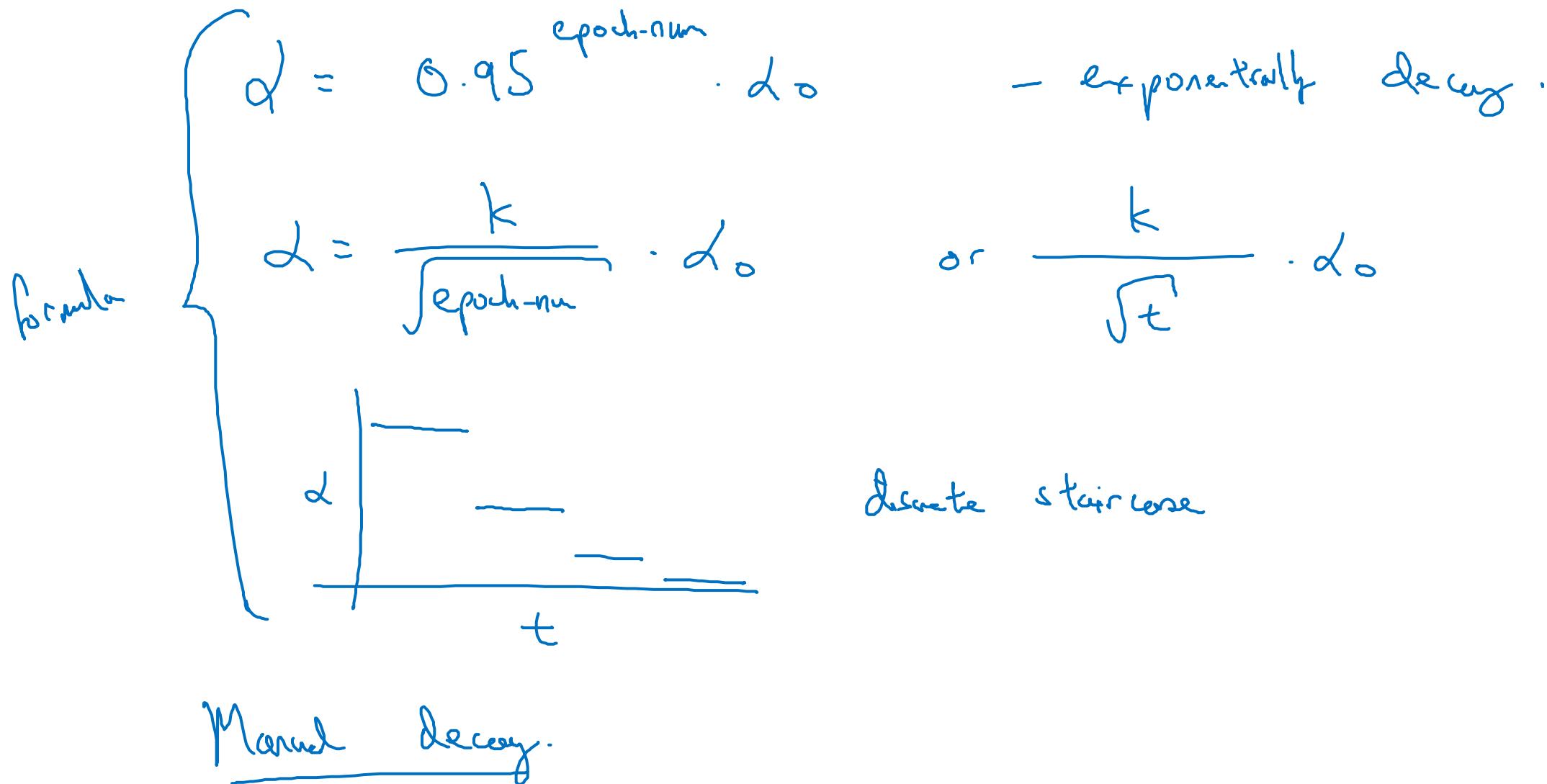


$$\alpha_0 = 0.2$$

$$\text{decay-rate} = 1$$



Other learning rate decay methods



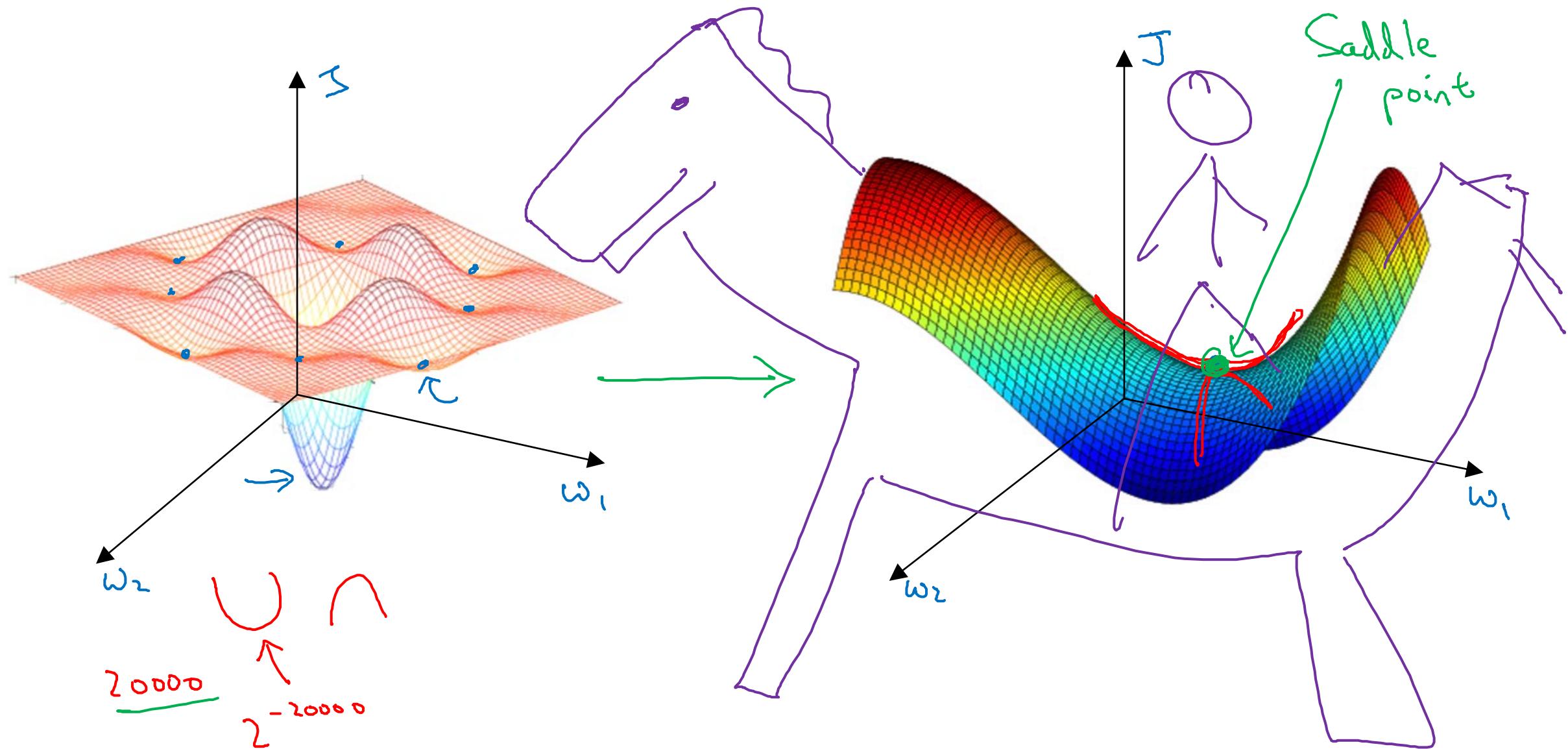


deeplearning.ai

Optimization Algorithms

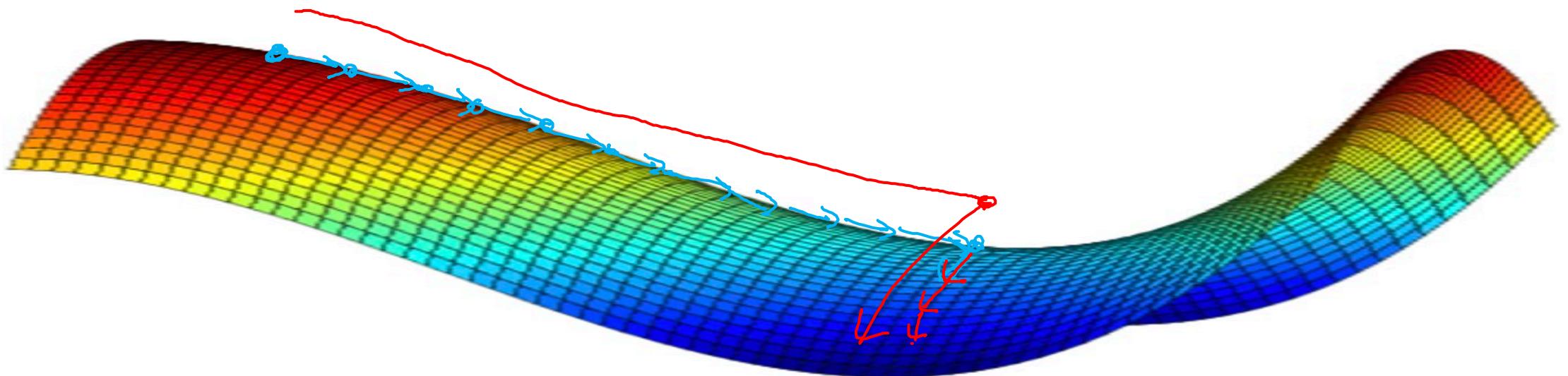
The problem of local optima

Local optima in neural networks



Andrew Ng

Problem of plateaus



- Unlikely to get stuck in a bad local optima
- Plateaus can make learning slow



deeplearning.ai

Hyperparameter tuning

Tuning process

Hyperparameters

$$\begin{array}{c} \xrightarrow{\quad} \boxed{\alpha} \\ \boxed{\beta}^{NO \cdot q} \end{array}$$

$$\beta_1, \beta_2, \epsilon_{0.9}, \epsilon_{0.99}, 10^8$$

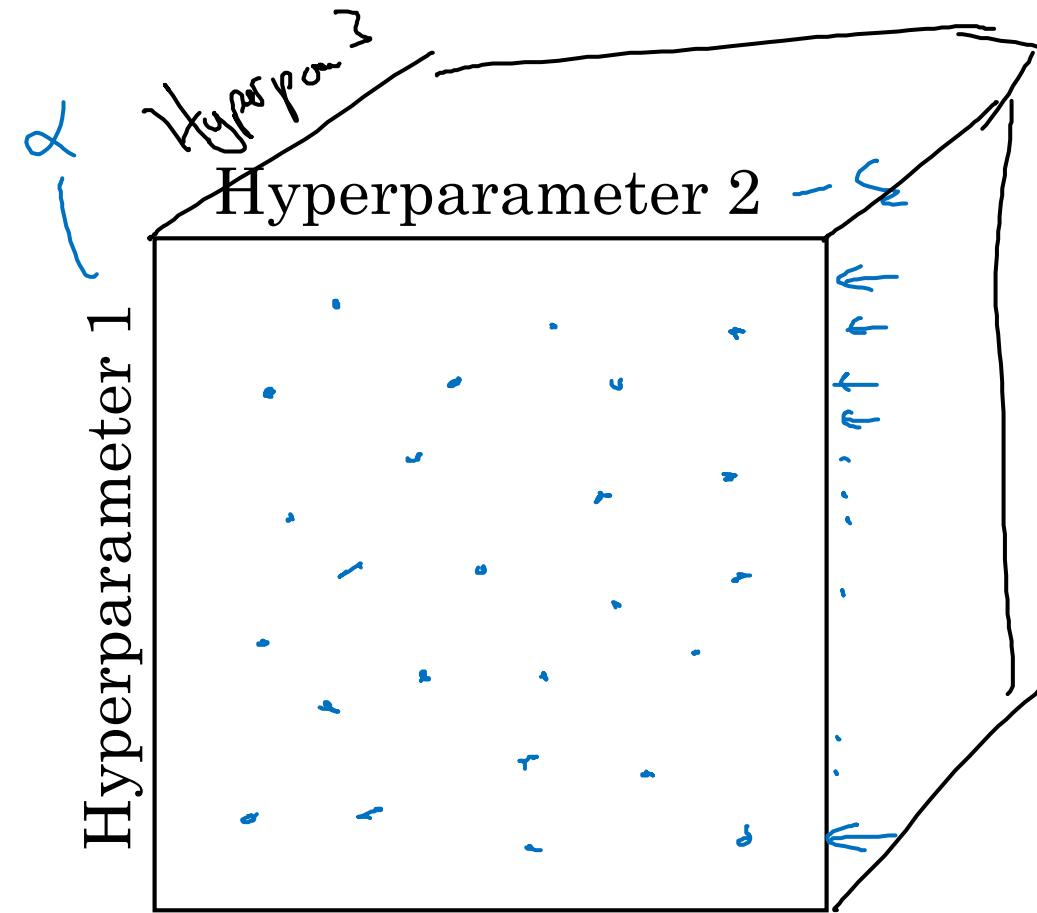
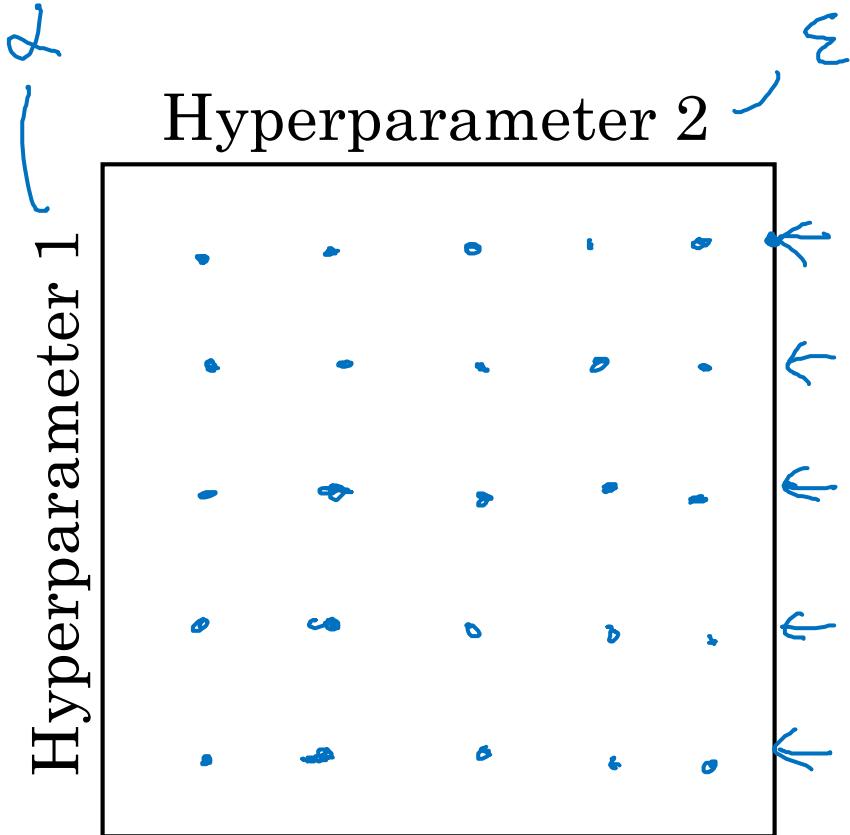
layers

hidden units

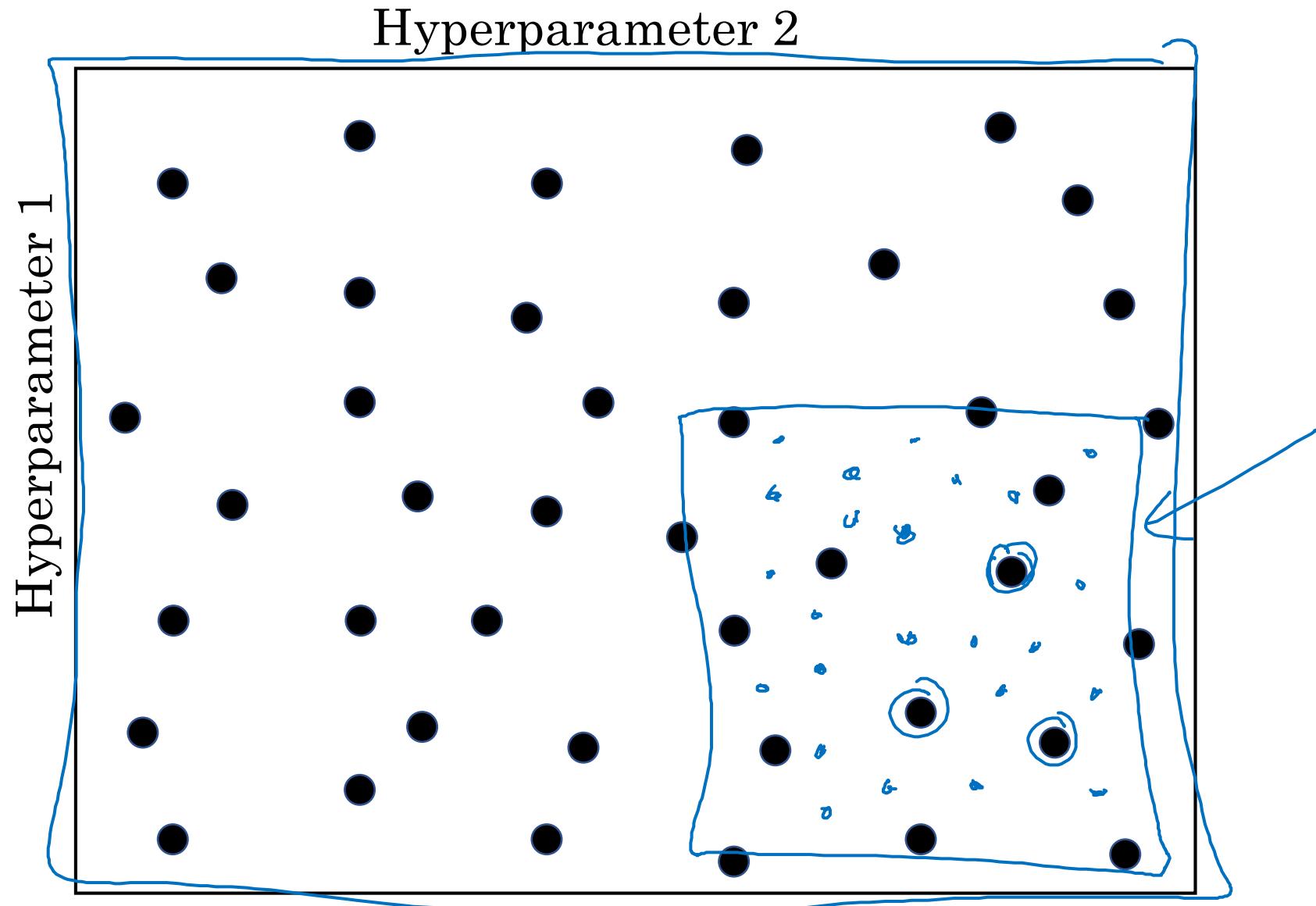
learning rate decay

mini-batch size

Try random values: Don't use a grid



Coarse to fine





deeplearning.ai

Hyperparameter tuning

Using an appropriate
scale to pick
hyperparameters

Picking hyperparameters at random

→ $n^{[l]} = 50, \dots, 100$

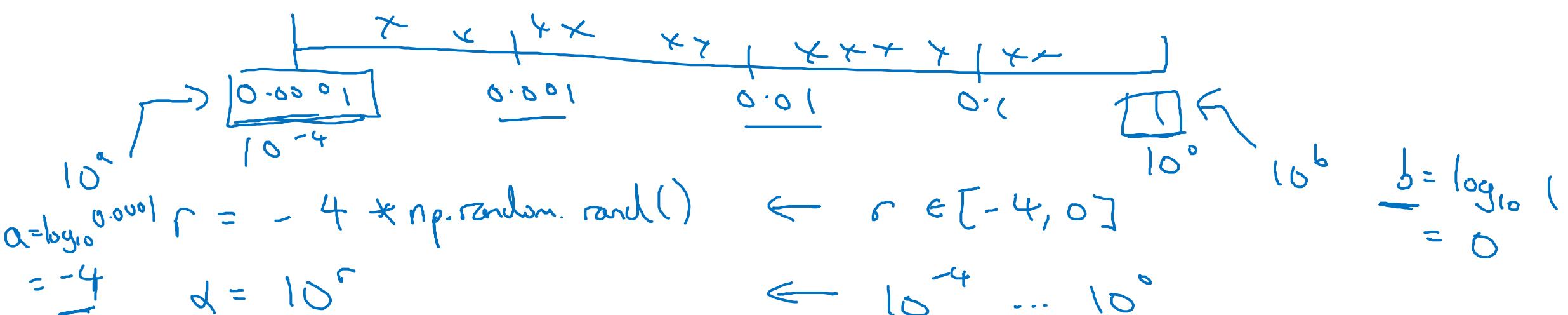
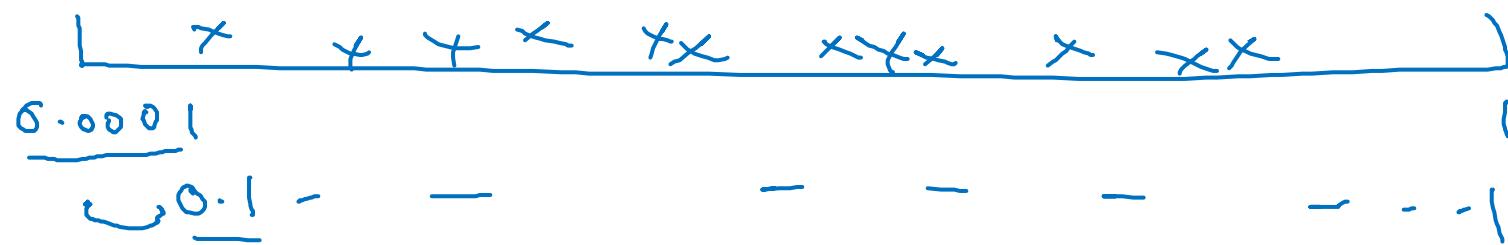


→ #layers $L : 2 - 4$

2, 3, 4

Appropriate scale for hyperparameters

$$\lambda = 0.0001, \dots, 1$$



$$10^a \dots 10^b$$

$$\frac{r \in [a, b]}{[-4, 0]}$$

$$\lambda = 10^r$$

Hyperparameters for exponentially weighted averages

$$\beta = 0.9 \dots 0.999$$

\downarrow \downarrow
 10 1000

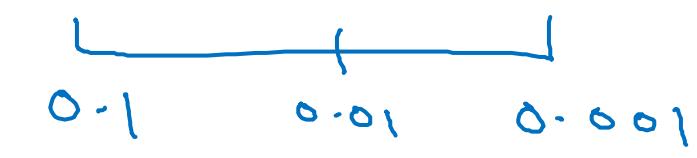
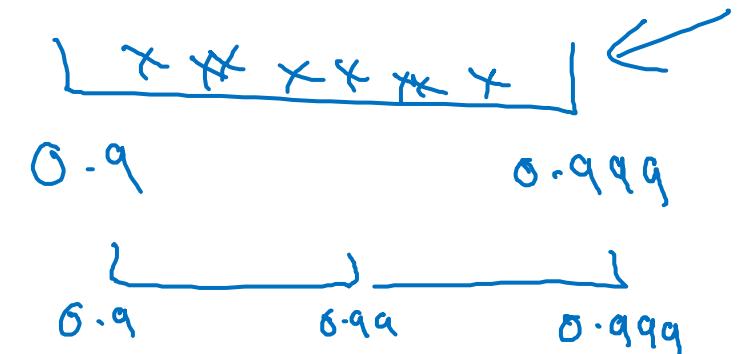
$$1-\beta = 0.1 \dots 0.001$$

$$\beta: 0.900 \rightarrow 0.9005 \quad \} \sim 10$$

$$\beta: 0.999 \rightarrow 0.9995$$

~ 1000 ~ 2000

$$\frac{1}{1-\beta}$$



$$\frac{10^{-1}}{1-\beta} \quad \frac{10^{-3}}{\beta}$$

$r \in [-3, -1]$

$$1-\beta = 10^r$$

$$\beta = 1 - 10^r$$

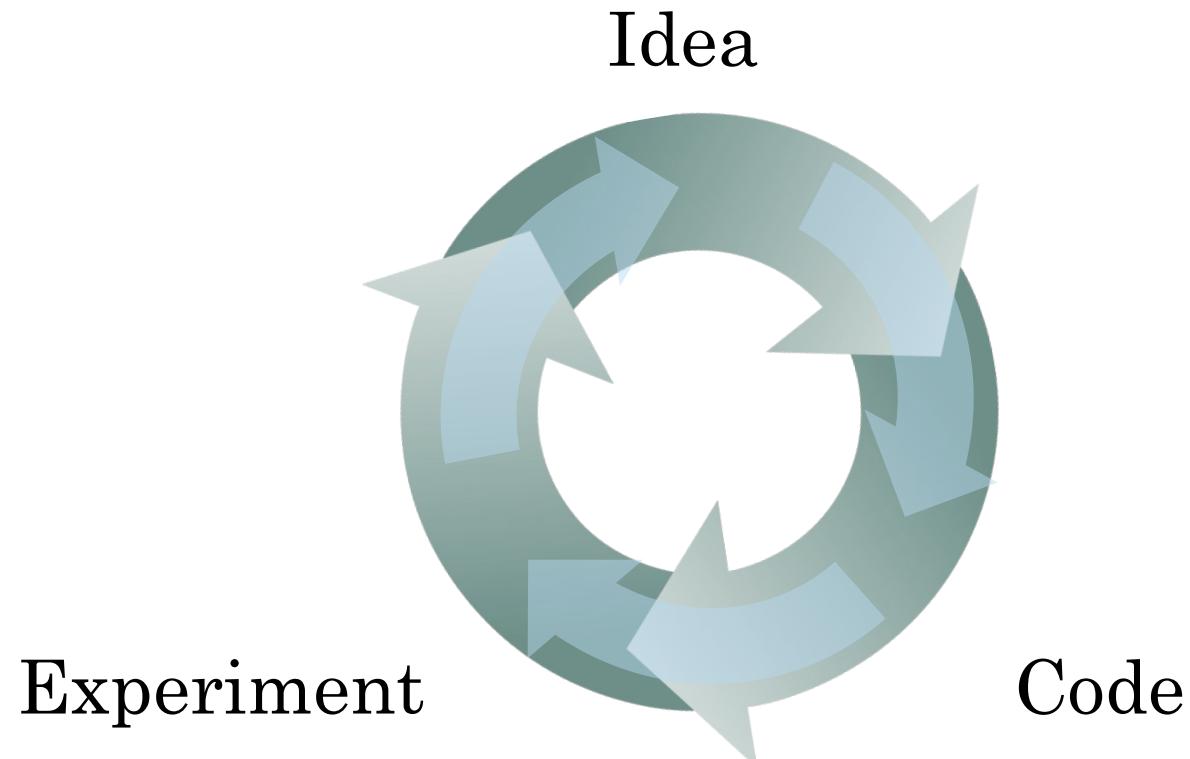


deeplearning.ai

Hyperparameters tuning

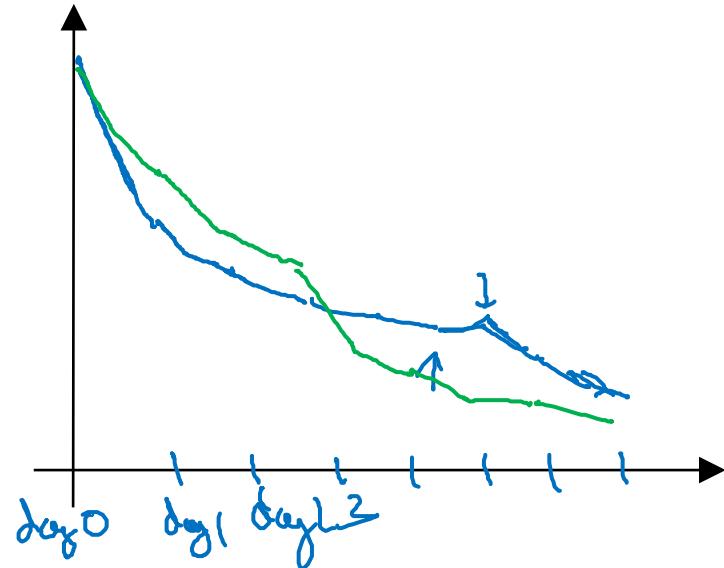
Hyperparameters tuning in practice: Pandas vs. Caviar

Re-test hyperparameters occasionally



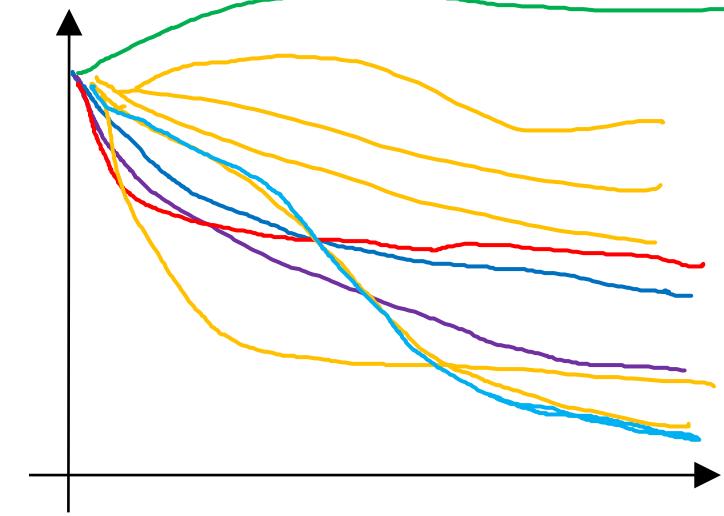
- NLP, Vision, Speech,
Ads, logistics,
- Intuitions do get stale.
Re-evaluate occasionally.

Babysitting one model



Panda ↪

Training many models in parallel



Caviar ↪

Andrew Ng

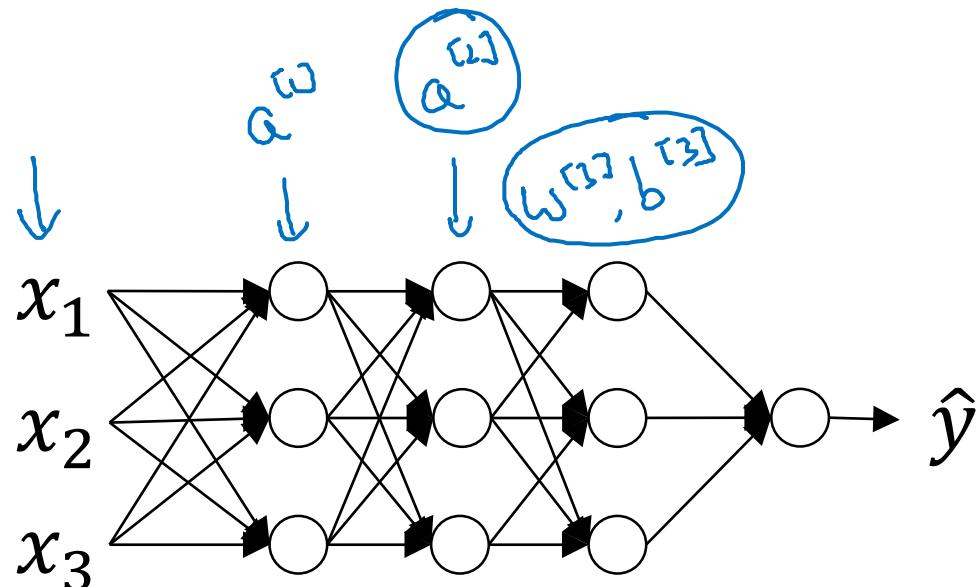
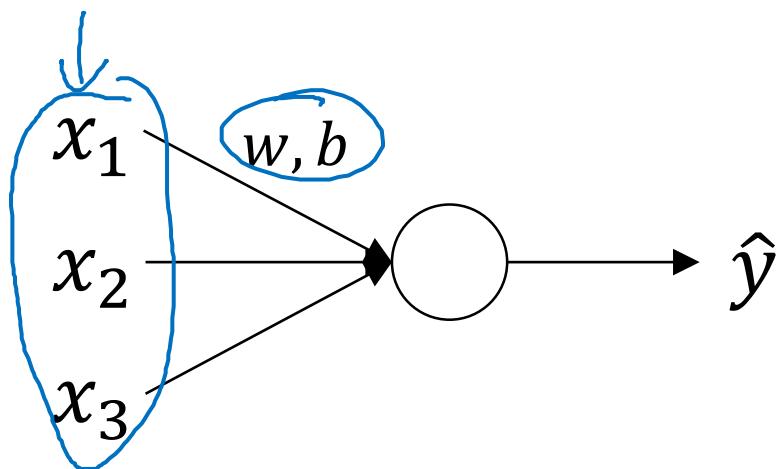


deeplearning.ai

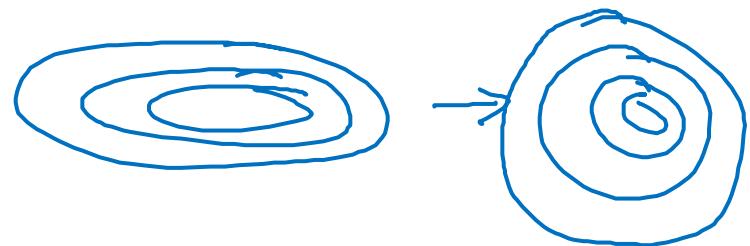
Batch Normalization

Normalizing activations in a network

Normalizing inputs to speed up learning



$$\mu = \frac{1}{m} \sum_i x^{(i)}$$
$$X = X - \mu$$
$$\sigma^2 = \frac{1}{m} \sum_i (x^{(i)} - \mu)^2 \quad \text{← element-wise}$$
$$X = X / \sigma^2$$



Can we normalize $\frac{a^{[2]}}{w^{[2]}, b^{[2]}}$ so
as to train $w^{[2]}, b^{[2]}$ faster

Normalize $\frac{z^{[2]}}{\uparrow}$

Implementing Batch Norm

Given some intermediate values in NN

$$\mu = \frac{1}{m} \sum_i z^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2$$

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\hat{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

Use $\hat{z}^{(i)}$ instead of $z^{(i)}$.

If $\gamma = \sqrt{\sigma^2 + \epsilon}$ ←
then $\hat{z}^{(i)} = z^{(i)}$ ←

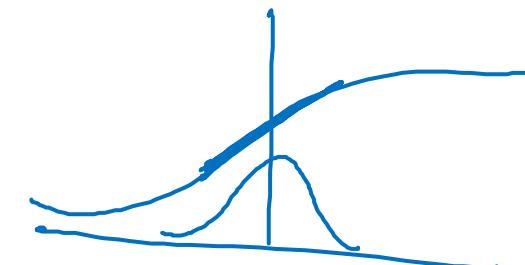
learnable parameters
of model.

$$z^{(1)}, \dots, z^{(m)}$$

$$z^{[l]}(:)$$

$$x \leftarrow$$

$$z^{(i)} \leftarrow$$



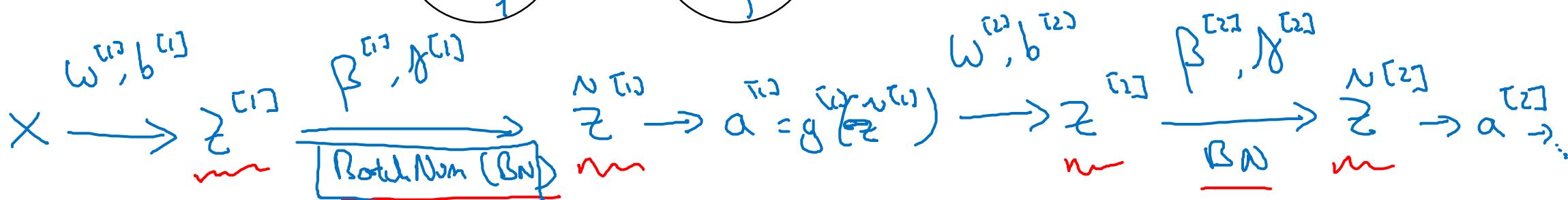
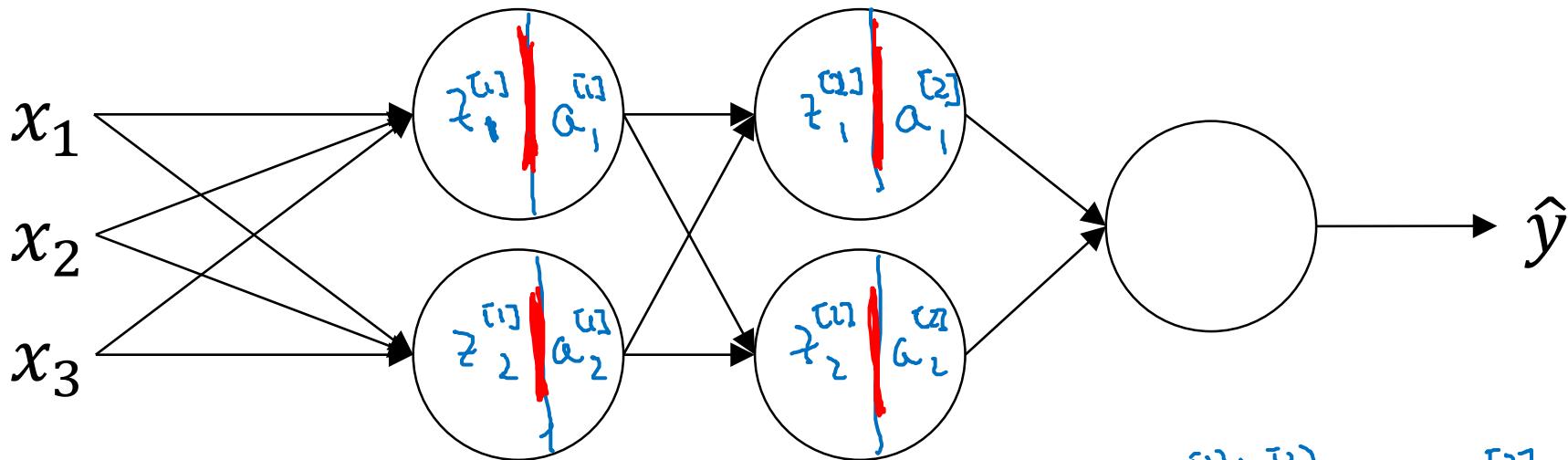


deeplearning.ai

Batch Normalization

Fitting Batch Norm
into a neural network

Adding Batch Norm to a network



Parameters:

$$\left\{ w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}, \right.$$

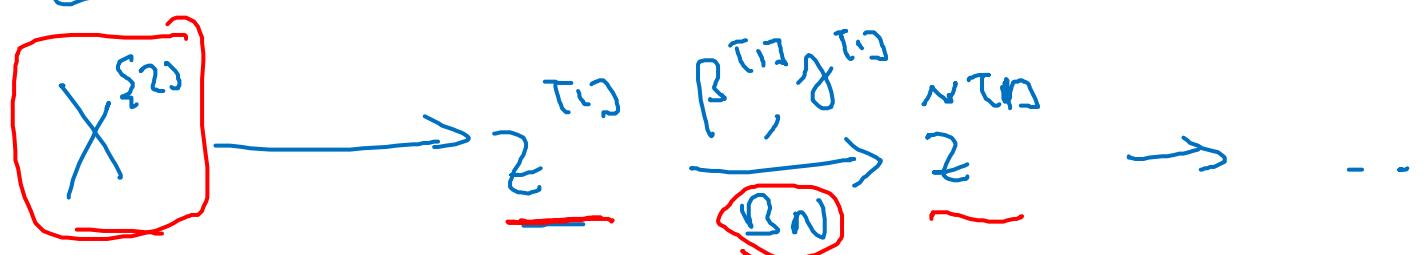
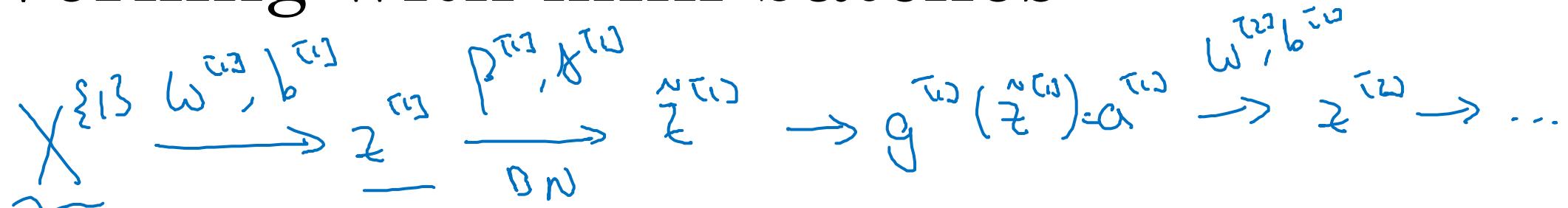
$$\left. \rightarrow \beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \dots, \beta^{[L]}, \gamma^{[L]} \right\}$$

$$\rightarrow \beta$$

$$\beta = \bar{\beta} - d \delta \beta^{[L]}$$

`tf.nn.batch_normalization` ←

Working with mini-batches



$X^2 \xrightarrow{\dots}$

Parameters: $\{W^{[l]}, \cancel{b^{[l]}}, \beta^{[l]}, \gamma^{[l]}\}$.

$$z^{[l]} \\ (n^{[l]}, 1)$$

$$\begin{matrix} | & | & | \\ (n^{[l]}, 1) & (n^{[l]}, 1) & (n^{[l]}, 1) \end{matrix}$$

$$\begin{aligned} \rightarrow z^{[l]} &= W^{[l]} a^{[l-1]} + \cancel{b^{[l]}} \\ z^{[l]} &= W^{[l]} a^{[l-1]} \\ z^{[l]}_{norm} &= \gamma^{[l]} z^{[l]}_{norm} \\ \rightarrow z^{[l]} &= \gamma^{[l]} z^{[l]}_{norm} \end{aligned}$$

$T\beta^{[l]}$

Andrew Ng

Implementing gradient descent

for $t = 1 \dots \text{num MiniBatches}$
Compute forward prop on $X^{[t]}$.

In each hidden layer, use BN to replace $\underline{z}^{[l]}$ with $\hat{\underline{z}}^{[l]}$.

Use backprop to compute $\underline{dw}^{[l]}$, ~~$\underline{db}^{[l]}$~~ , $\underline{d\beta}^{[l]}$, $\underline{dg}^{[l]}$

Update parameters $\left. \begin{array}{l} w^{[l]} := w^{[l]} - \alpha \underline{dw}^{[l]} \\ \beta^{[l]} := \beta^{[l]} - \alpha \underline{d\beta}^{[l]} \\ g^{[l]} := \dots \end{array} \right\} \leftarrow$

Works w/ momentum, RMSprop, Adam.

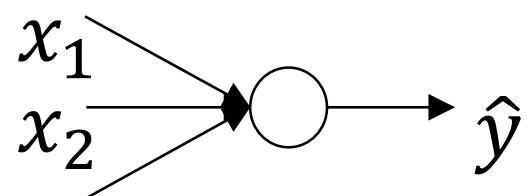


deeplearning.ai

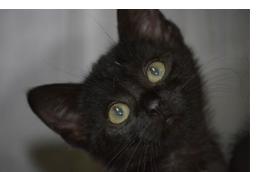
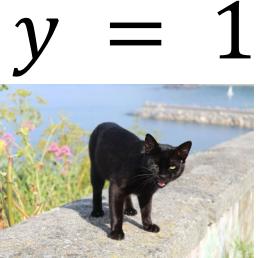
Batch Normalization

Why does
Batch Norm work?

Learning on shifting input distribution



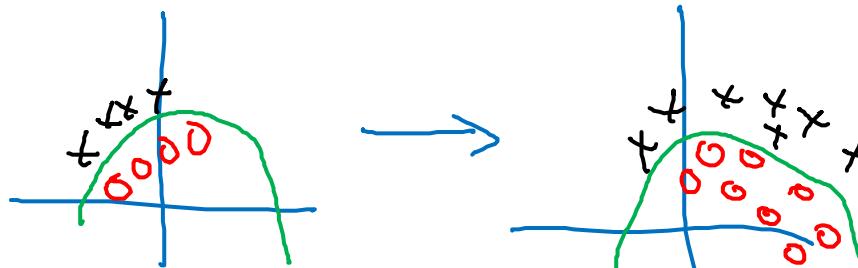
Cat



Non-Cat



$y = 1$



$y = 1$



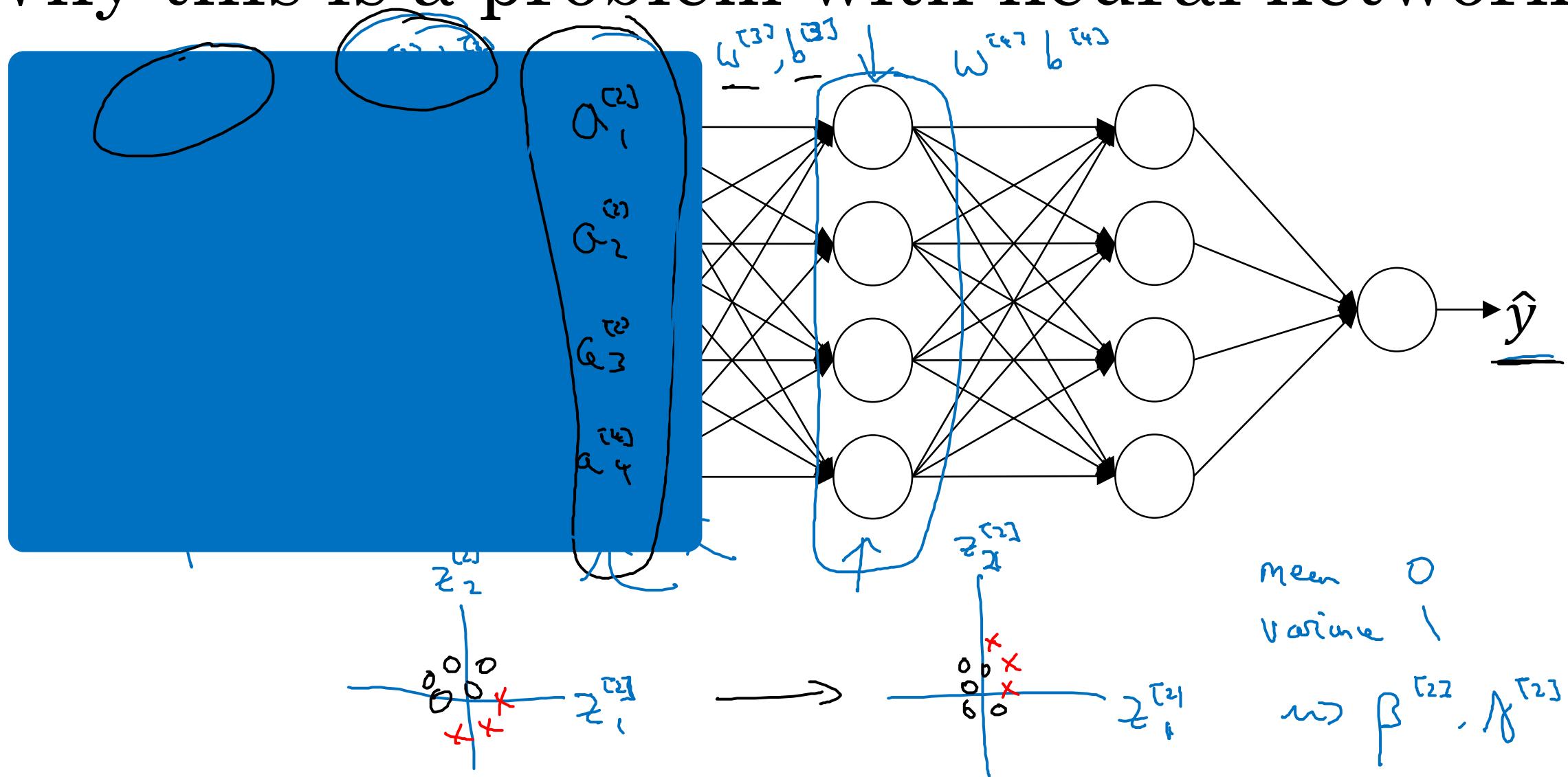
$y = 0$



"Covariate shift"

$X \rightarrow Y$

Why this is a problem with neural networks?



Batch Norm as regularization

X

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.
 $\xrightarrow{\hat{z}^{[l]}}$ μ, σ^2 $\{z^{[l]}\}$
 $\underline{64}, \underline{128}$
- This adds some noise to the values $z^{[l]}$ within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations.
 μ, σ^2
- This has a slight regularization effect.

mini-batch : $\underline{64} \longrightarrow \underline{512}$



deeplearning.ai

Batch Normalization

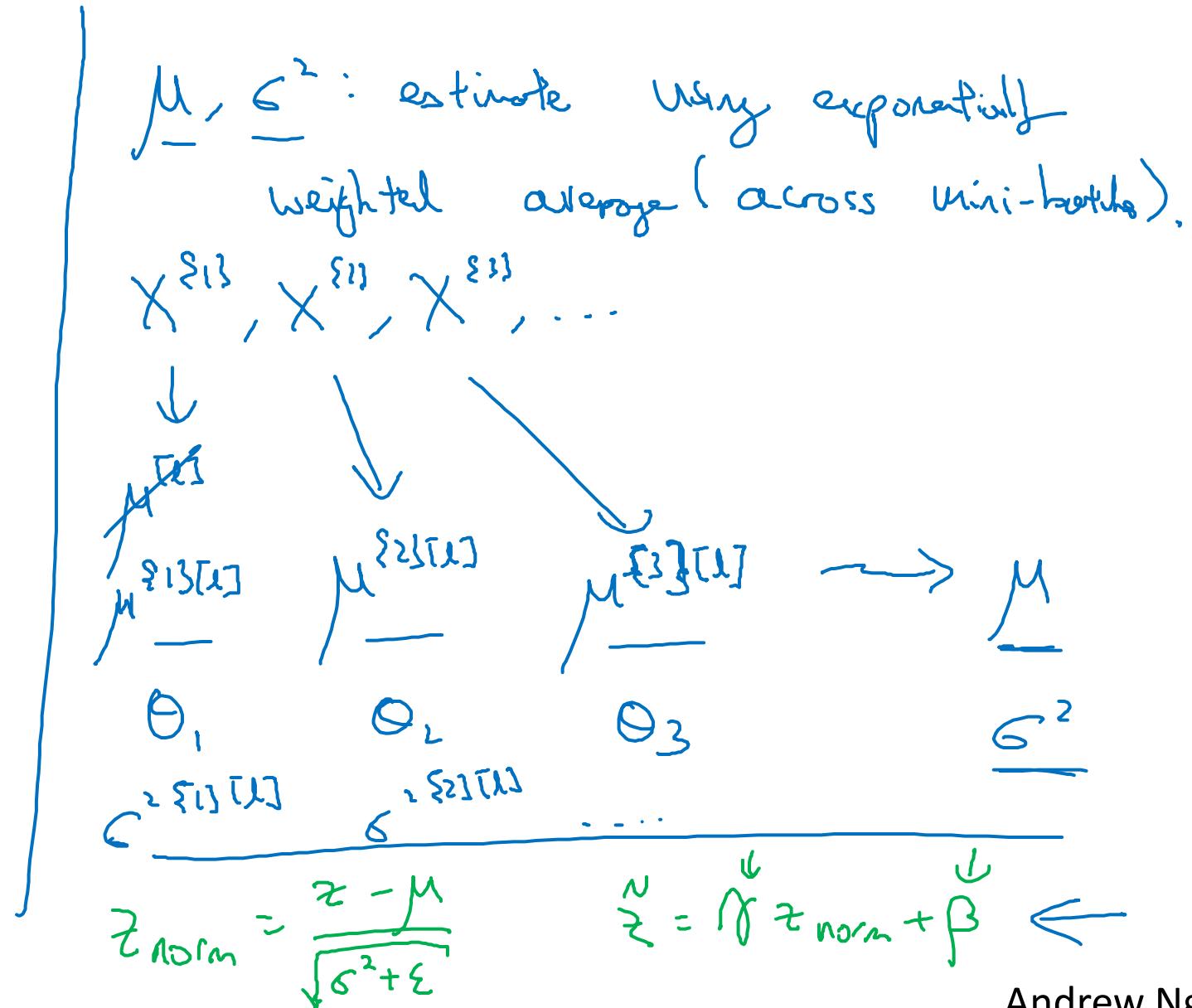
Batch Norm at test time

Batch Norm at test time

$$\mu = \frac{1}{m} \sum_i z^{(i)}$$
$$\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2$$

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$





deeplearning.ai

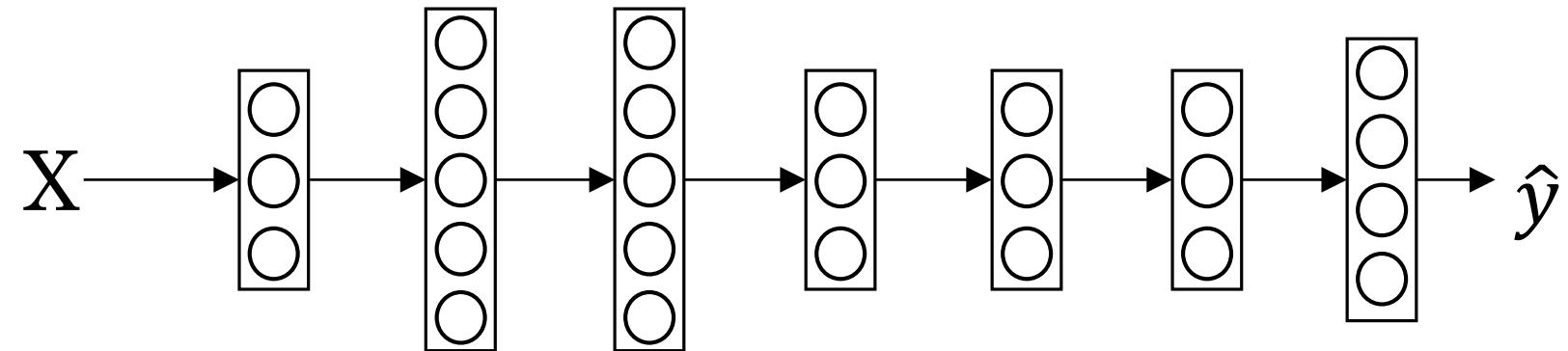
Multi-class
classification

Softmax regression

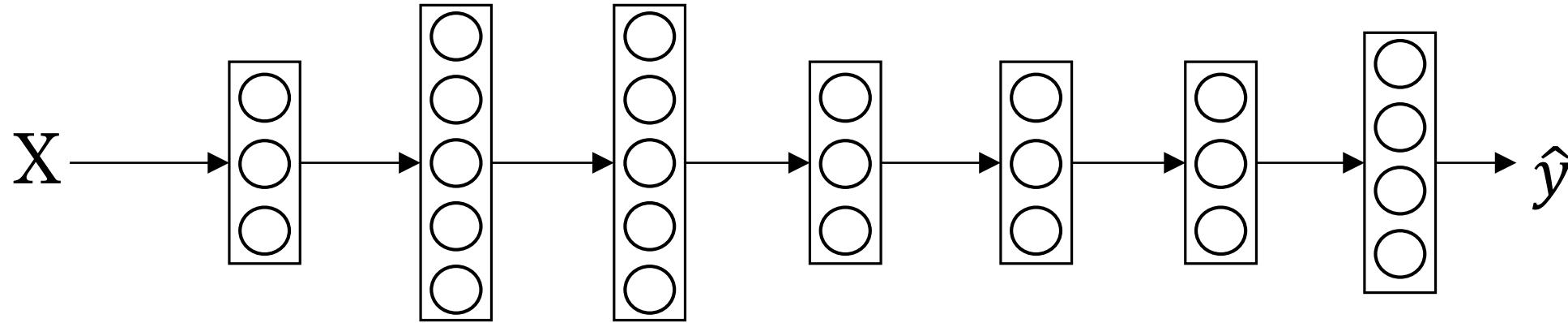
Recognizing cats, dogs, and baby chicks



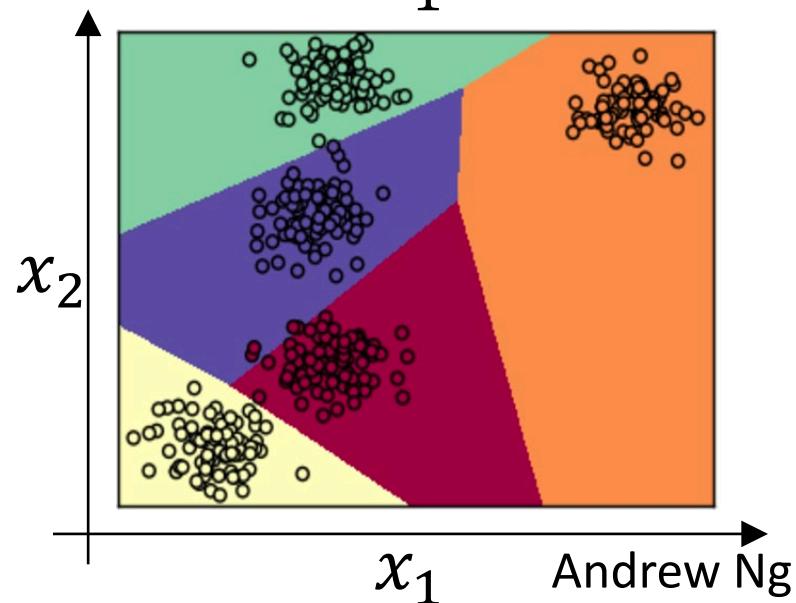
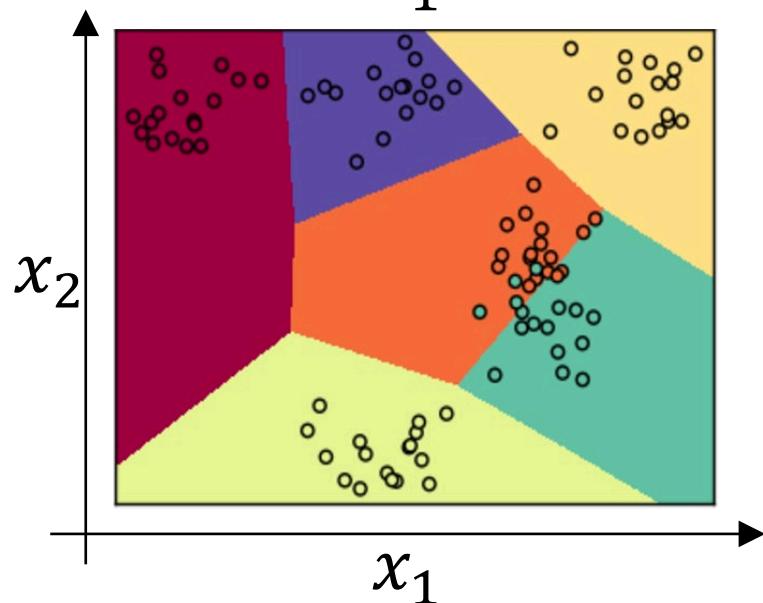
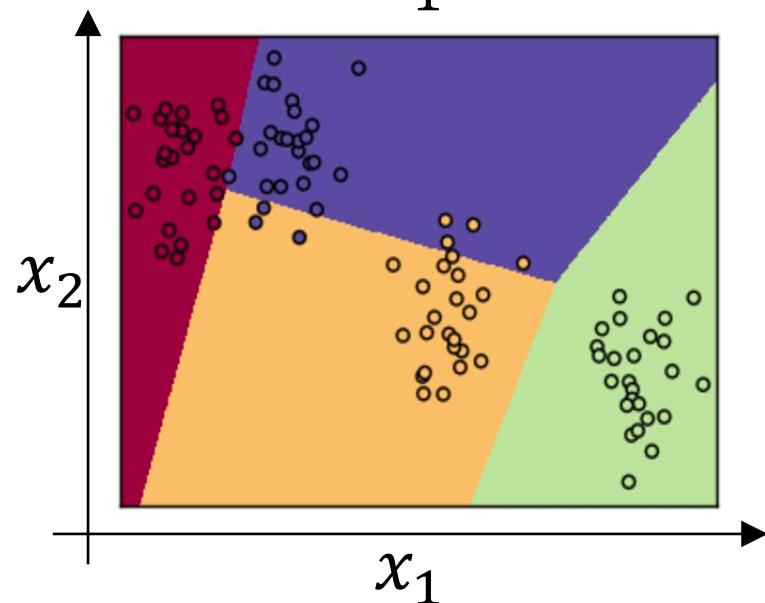
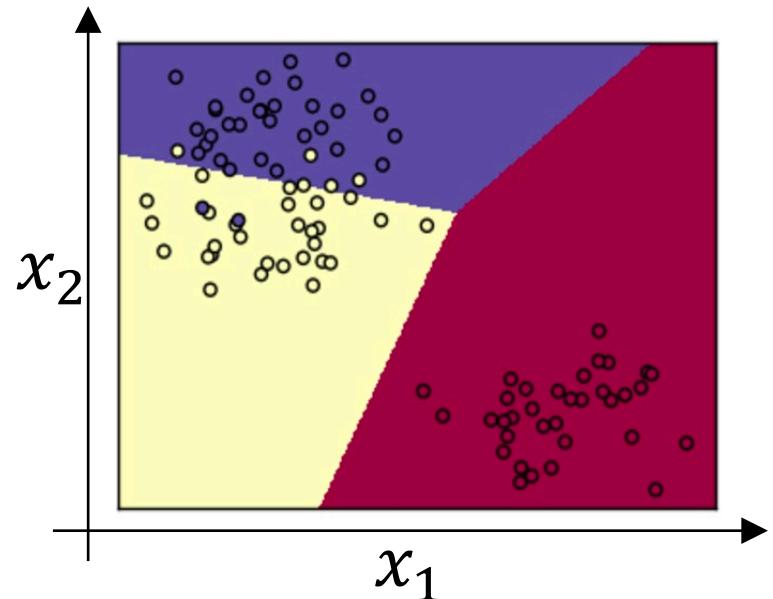
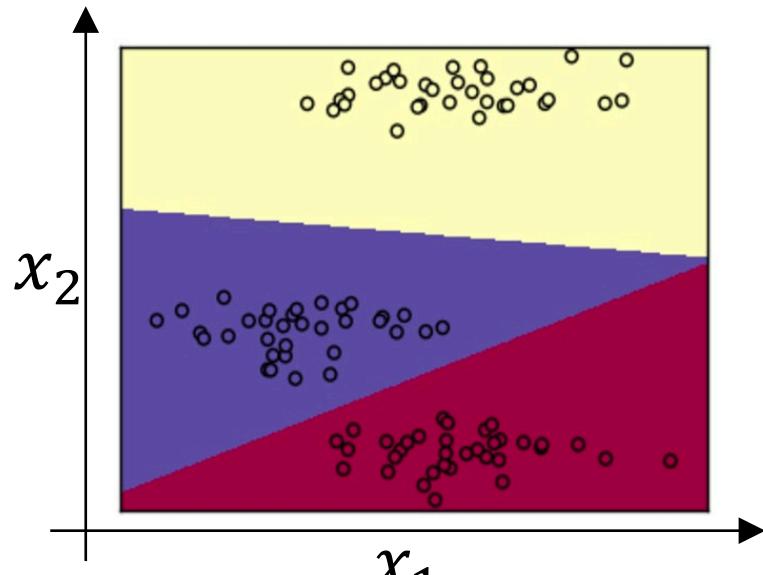
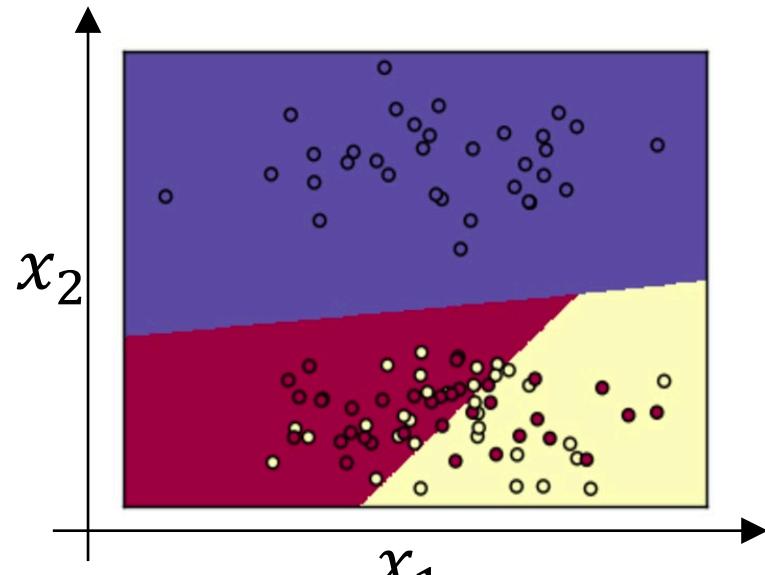
3 1 2 0 3 2 0 1



Softmax layer



Softmax examples





deeplearning.ai

Programming Frameworks

Deep Learning frameworks

Deep learning frameworks

- Caffe/Caffe2
- CNTK
- DL4J
- Keras
- Lasagne
- mxnet
- PaddlePaddle
- TensorFlow
- Theano
- Torch

Choosing deep learning frameworks

- Ease of programming (development and deployment)
 - Running speed
- - Truly open (open source with good governance)



deeplearning.ai

Programming Frameworks

TensorFlow

Motivating problem

$$J(w) = \frac{1}{2} \left[w^2 - 10w + 25 \right]$$

\uparrow

$(w-5)^2$

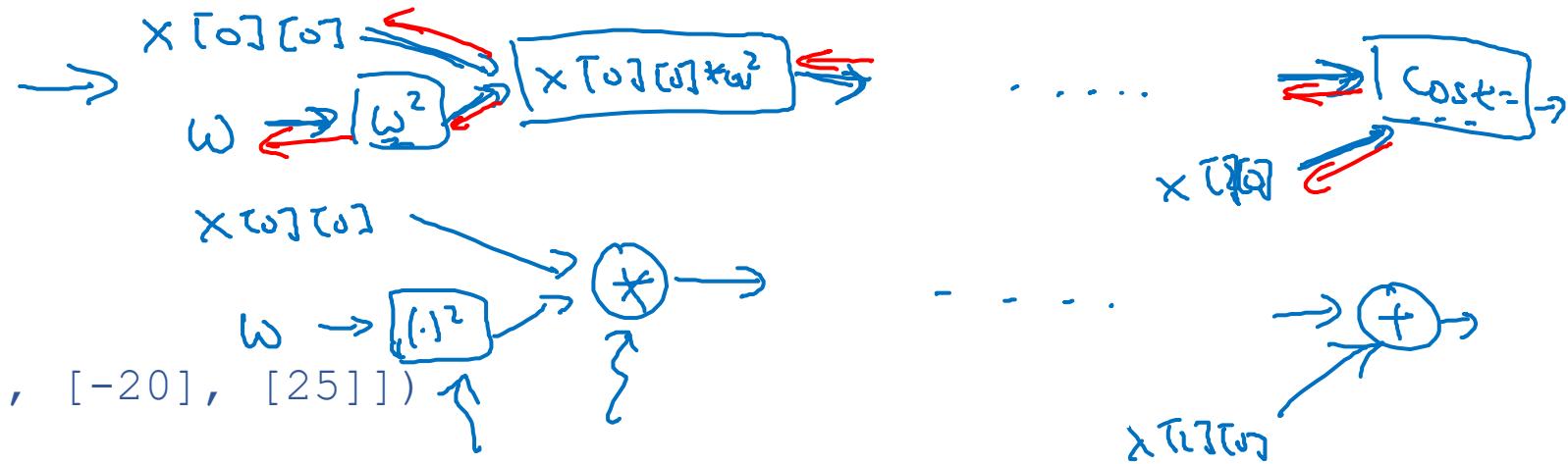
$w = 5$

$J(w, b)$

$\uparrow \uparrow$

Code example

```
import numpy as np  
import tensorflow as tf  
  
coefficients = np.array([[1], [-20], [25]])  
  
w = tf.Variable([0], dtype=tf.float32)  
x = tf.placeholder(tf.float32, [3,1])  
  
cost = x[0][0]*w**2 + x[1][0]*w + x[2][0] # (w-5)**2  
train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)  
init = tf.global_variables_initializer()  
session = tf.Session()  
session.run(init)  
print(session.run(w))  
  
for i in range(1000):  
    session.run(train, feed_dict={x:coefficients})  
print(session.run(w))
```



```
with tf.Session() as session:  
    session.run(init)  
    print(session.run(w))
```