deeplearning.ai

Setting up your
ML application

Train/dev/test
sets

# Applied ML is a highly iterative process
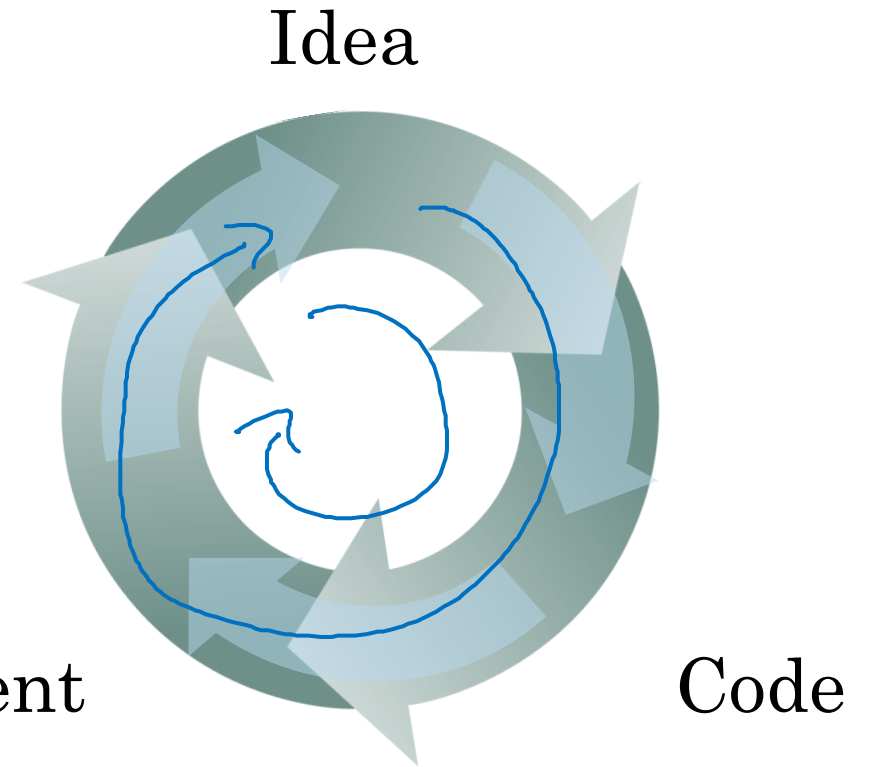
# layers

# hidden units

learning rates
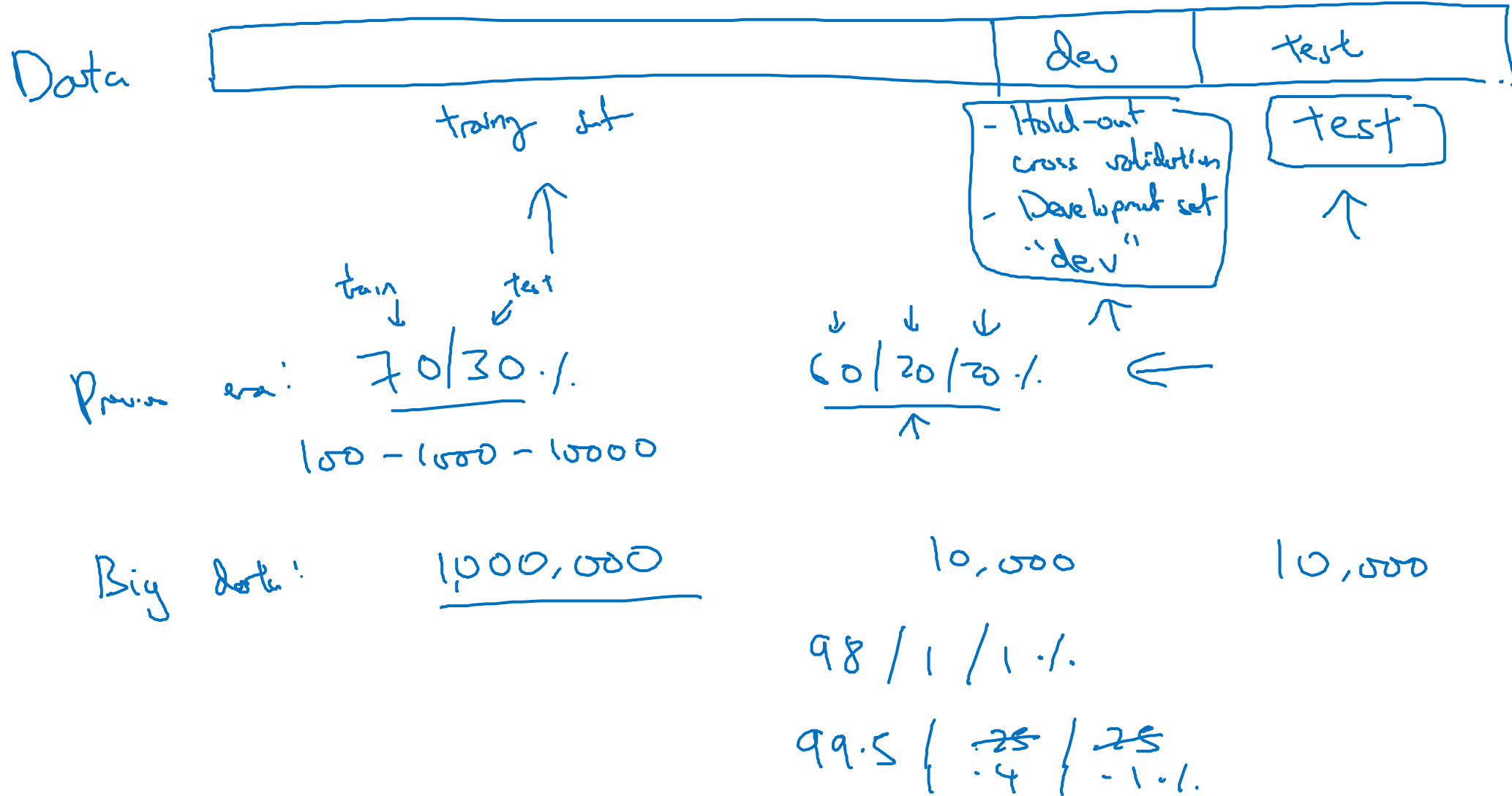
activation functions

...

Idea

Experiment

Code

NLP, Vision, Speech, Structural data

Ads    Search    Security    logistic ....

# Train/dev/test sets

# 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
↓ ↖
→ Trn / dev

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

Setting up your
ML application

Bias/Variance

deeplearning.ai

# Bias and Variance



high bias

*Underfitting*

$\rightarrow$ "just right"

high variance

*Overfitting*

# 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%

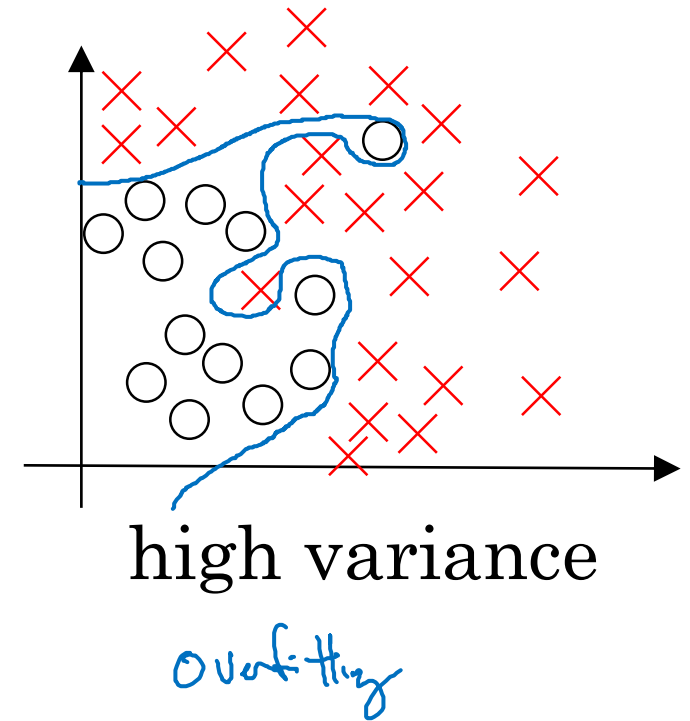Optimal (Bayes) error: ≈ 0% to 15%   Blury images

# 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

High bias?
(training data performance)

$\longrightarrow$ Bigger network

$\longrightarrow$ Train longer.

(NN architecture search)

↓ N

High variance?
(dev set performance)

$\longrightarrow$ More data

$\longrightarrow$ Regularization

(NN architecture search)

↓ N

Done

"Bias Variance tradeoff"

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

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

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

Omit

$L_2$ regularization

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

$L_1$ regularization

$w$ will be sparse

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

# Neural network

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

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

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

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

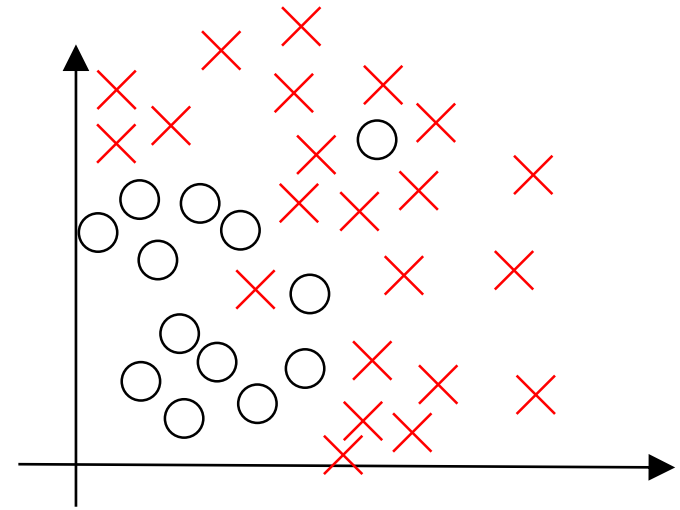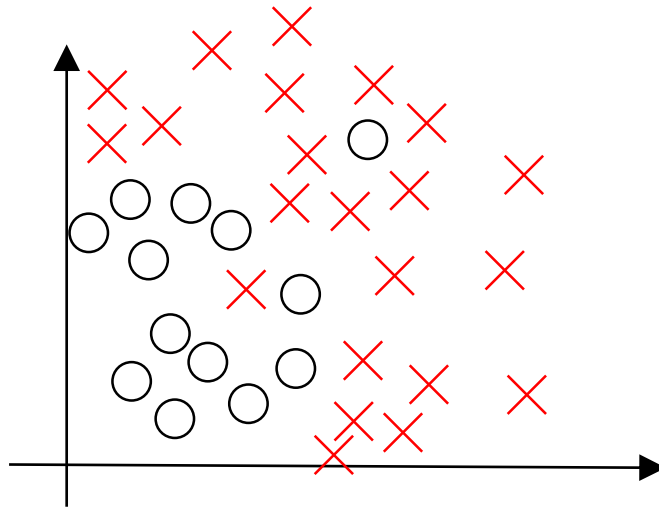$$\rightarrow W^{[l]} := W^{[l]} - \alpha \, dW^{[l]}$$

"Weight decay"

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

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

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

# How does regularization prevent overfitting?

# How does regularization prevent overfitting?

deeplearning.ai

Regularizing your
neural network

Why regularization
reduces overfitting

# How does regularization prevent overfitting?



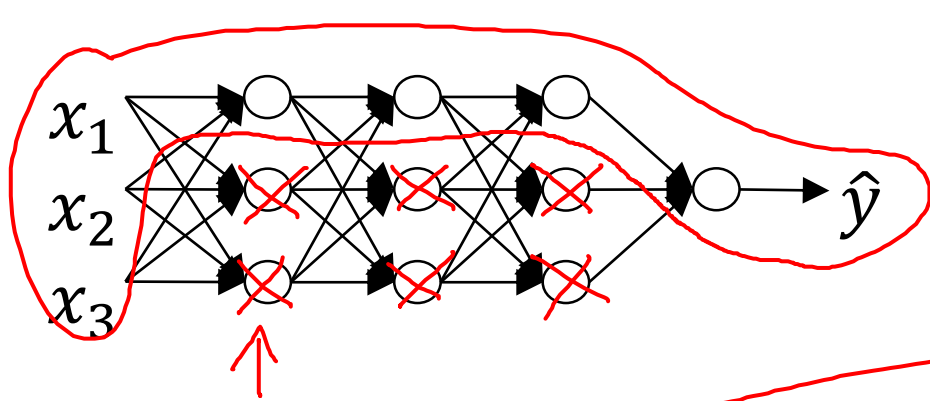$$J\left(\omega^{[l]}, b^{[l]}\right) = \frac{1}{m} \sum_{i=1}^{n} \mathcal{L}\left(y^{(i)}, \hat{y}^{(i)}\right) + \frac{\lambda}{2m} \sum_{l=1}^{L} \|\omega^{[l]}\|_F^2$$

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

high bias        "just right"        high variance

# How does regularization prevent overfitting?



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

$\tanh$

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

Every layer $\approx$ linear.

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

deeplearning.ai

Regularizing your
neural network

Dropout
regularization

# Dropout regularization

# 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. $\leadsto$ 10 units shut off

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

reduced by 20%.    Test

$\mathrel{/}= 0.8$

# Making predictions at test time

$a^{[0]} = X$

## No drop out.

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

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

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

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

$$\downarrow$$

$$\hat{y}$$

$/ = $ 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. ⟿ Shrink weights. $l_2$

deeplearning.ai

Regularizing your
neural network

---

Other regularization
methods

# Data augmentation

# Early stopping

Orthogonalization.

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

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

$J(w,b)$

$l_2$



dev set error

traing error or $J$

# iterations

$w \approx 0$

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

large $W$

deeplearning.ai

Setting up your optimization problem

Normalizing inputs

# Normalizing training sets

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



Subtract mean:

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

$$x := x - \mu$$

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.

# Why normalize inputs?

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

Unnormalized:

Normalized:

$w2 \quad x_2 : 0 \cdots 1 \leftarrow$

$-1 \cdots 1$

$x_1 : 0 \cdots 1$

$x_2 : -1 \cdots 1$

$x_3 : 1 \cdots 2$

Setting up your optimization problem

---

Vanishing/exploding gradients

deeplearning.ai

# Vanishing/exploding gradients

$L = 150$



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

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

$1.5^L$

$0.5^L$

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

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

$$W^{[l]} > I$$

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

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

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

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

$1.5^{L-1} x$

$0.5^{L-1} x$

# Single neuron example



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

$a^{[l]}$

$w^{[l]}$

$\hat{y}$

$a = g(z)$

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

Large $n \rightarrow$ Smaller $w_i$

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

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

ReLU

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

Other variants:

tanh

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

Xavier initialization

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

deeplearning.ai

Setting up your
optimization problem

---

Numerical approximation
of gradients

# 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 $\theta = 1$

$\dfrac{dw}{db}$

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

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

$\approx 3$



$f$

$f(\theta + \varepsilon)$

$f(\theta)$

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

$\varepsilon$

$\theta$   $\theta + \varepsilon$

$1$   $1.01$

$\varepsilon = 0.01$

$\theta = 1$

$\theta + \varepsilon = 1.01$

$0.0301$

$3.1$

$3.2$

# Checking your derivative computation

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



$f(\theta+\varepsilon)$

$f(\theta-\varepsilon)$

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

$2\varepsilon$

$\theta-\varepsilon \quad \theta \quad \theta+\varepsilon$

$0.99 \quad 1 \quad 1.01$

$\varepsilon = 0.01$

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

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

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

approx error: $0.0001$

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

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

$0.01$
$0.0001$

$0.01$

Setting up your
optimization problem
_____

Gradient Checking

deeplearning.ai

# Gradient check for a neural network

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

Concentrate

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

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

concentrate

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

# Gradient checking (Grad check)

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

for each $i$:

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

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

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

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

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

$$10^{-5}$$

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

Setting up your
optimization problem

Gradient Checking
implementation notes

deeplearning.ai

# Gradient checking implementation notes

- Don't use in training – only to debug

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

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

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

- Remember regularization.

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

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

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

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

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

Optimization
Algorithms

Mini-batch
gradient descent

deeplearning.ai

# Batch vs. mini-batch gradient descent

$X, Y$

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

Vectorization allows you to efficiently compute on $m$ examples.

$$X = [\, x^{(1)} \; x^{(2)} \; x^{(3)} \; \cdots \; x^{(1000)} \mid x^{(1001)} \; \cdots \; x^{(2000)} \mid \cdots \mid \cdots \; x^{(m)} \,]$$

$(n_x, m)$

$\underbrace{\qquad}_{X^{\{1\}} \; (n_x, 1000)}$  $\underbrace{\qquad}_{X^{\{2\}} \; (n_x, 1000)}$ $\cdots$ $\underbrace{\qquad}_{X^{\{5000\}} \; (n_x, 1000)}$

$$Y = [\, y^{(1)} \; y^{(2)} \; y^{(3)} \; \cdots \; y^{(1000)} \mid y^{(1001)} \; \cdots \; y^{(2000)} \mid \cdots \mid \cdots \; y^{(m)} \,]$$

$(1, m)$

$\underbrace{\qquad}_{Y^{\{1\}} \; (1, 1000)}$  $\underbrace{\qquad}_{Y^{\{2\}} \; (1, 1000)}$ $\cdots$ $\underbrace{\qquad}_{Y^{\{5000\}} \; (1, 1000)}$

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

5,000 mini-batches of 1,000 each

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

$x^{(i)}$

$z^{[l]}$

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

# Mini-batch gradient descent

repeat $\{$

for $t = 1, \ldots, 5000$ $\{$

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

$$Z^{[1]} = W^{[1]} X^{\{t\}} + b^{[1]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

$$\vdots$$

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

$\left.\begin{array}{c} \\ \\ \\ \\ \end{array}\right\}$ Vectorized implementation (1000 examples)

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

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

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

$$W^{[l]} := W^{[l]} - \alpha \, dW^{[l]}, \quad b^{[l]} := b^{[l]} - \alpha \, db^{[l]}$$

$\}$

$\}$

"1 epoch"
$\hookleftarrow$ pass through training set.

1 step of gradient descent
using $X^{\{t\}}, Y^{\{t\}}$.
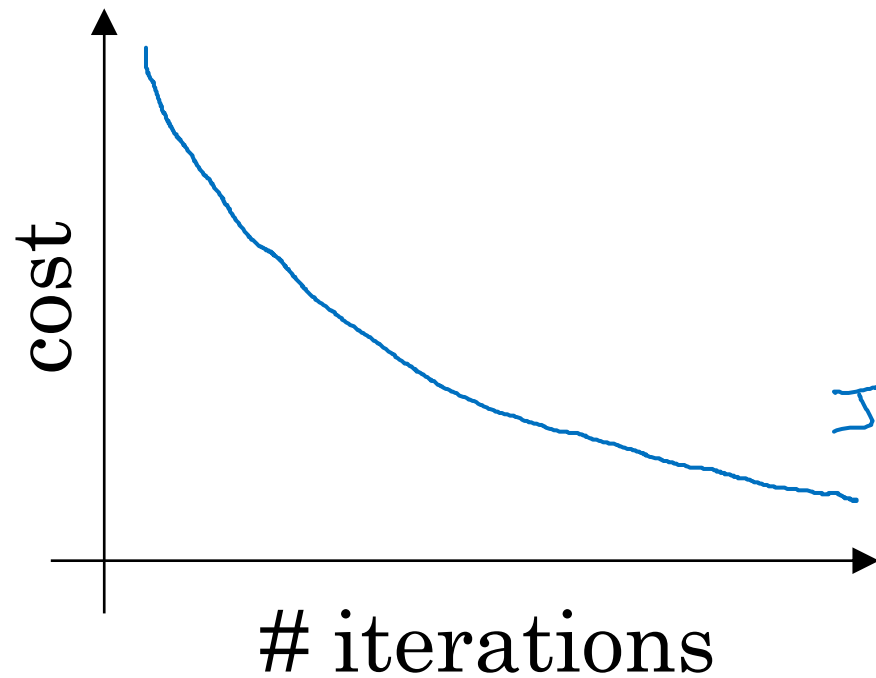(as if $m = 1000$)

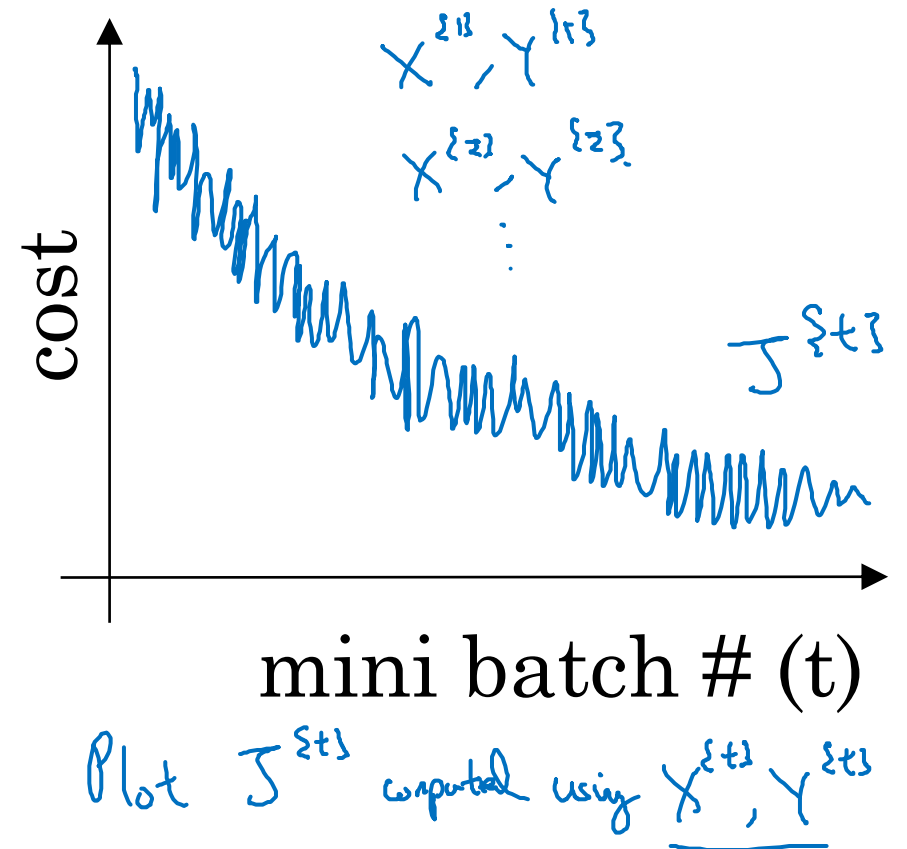$\underline{X, Y}$

deeplearning.ai

Optimization
Algorithms

Understanding
mini-batch
gradient descent

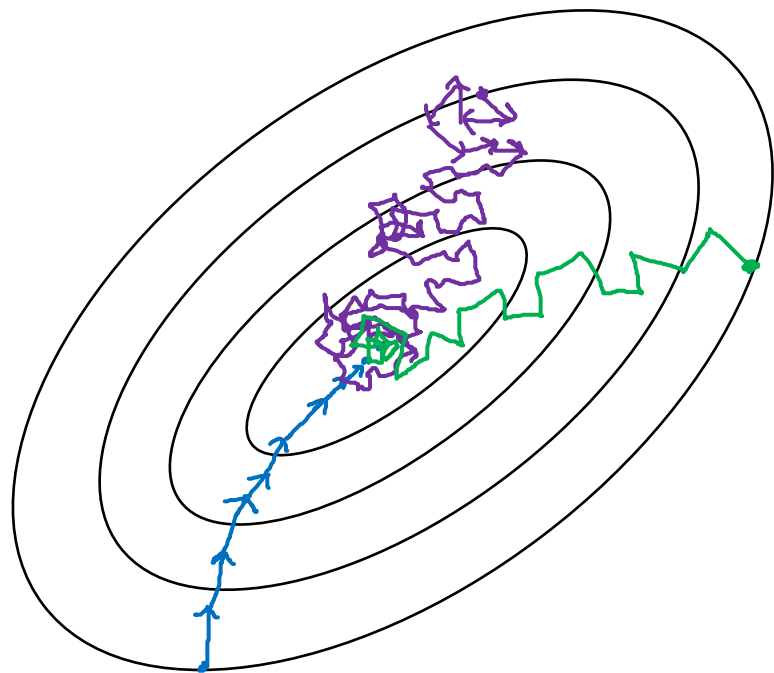# Training with mini batch gradient descent

### Batch gradient descent



cost

J

# iterations

### Mini-batch gradient descent



cost

$X^{\{1\}}, Y^{\{1\}}$

$X^{\{2\}}, Y^{\{2\}}$

$J^{\{t\}}$

mini batch # (t)

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

# 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 it own
$(X^{\{1\}}, Y^{\{1\}}) = (x^{(1)}, y^{(1)}) \dots (x^{(2)}, y^{(1)})$ mini-batch.

In practice: Somewhere in-between $\underline{1}$ and $\underline{m}$



Stochastic
gradient
descent

$\}$

Lose speedup
from vectorization

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

$\}$

Fastest learning.

● Vectorization.
($\sim 1000$)

● Make progress without
processing entire training set.

Batch
gradient descent
(mini-batch size = m)

$\}$

Too long
per iteration

# Choosing your mini-batch size

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

Typical mini-batch sizes:

$$\longrightarrow \quad \underbrace{64 \quad , \quad 128, \quad 256, \quad 512}_{} \qquad\qquad \frac{1024}{2^{10}}$$
$$\quad\quad\; 2^6 \qquad\quad 2^7 \qquad 2^8 \qquad\quad 2^9$$

Make sure mini-batch fit in CPU/GPU memory.
  $X^{\{t\}}, Y^{\{t\}}$

# Optimization Algorithms

## Exponentially weighted averages

deeplearning.ai

# Temperature in London

$\theta_1 = 40°F$   4°C ←

$\theta_2 = 49°F$   9°C

$\theta_3 = 45°F$   ⋮

⋮

$\theta_{180} = 60°F$   15°C

$\theta_{181} = 56°F$   ⋮

⋮



$V_0 = 0$

$V_1 = 0.9 V_0 + 0.1 \theta_1$

$V_2 = 0.9 V_1 + 0.1 \theta_2$

$V_3 = 0.9 V_2 + 0.1 \theta_3$

⋮

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

# Exponentially weighted averages (moving)
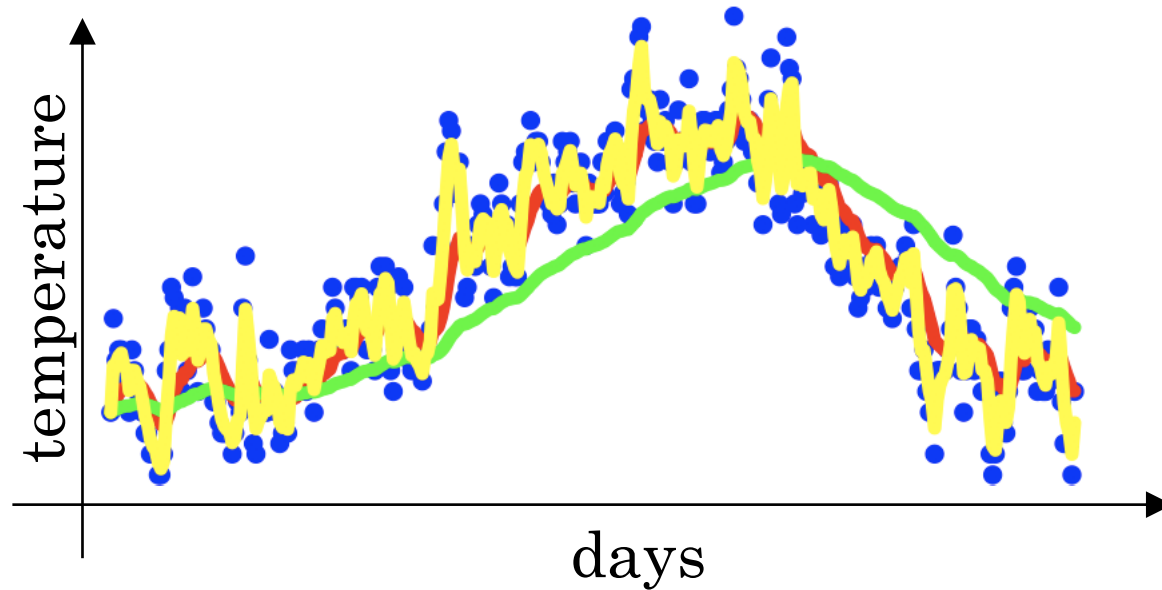
$$V_t = \beta V_{t-1} + (1-\beta) \theta_t \leftarrow$$

$\beta = 0.9$ : $\approx 10$ days' temperature.

$\beta = 0.98$ : $\approx 50$ days

$\beta = 0.5$ : $\approx 2$ days

$V_t$ as approximately average over

$\rightarrow \approx \frac{1}{1-\beta}$ days' temperature.

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



temperature

days

deeplearning.ai

Optimization
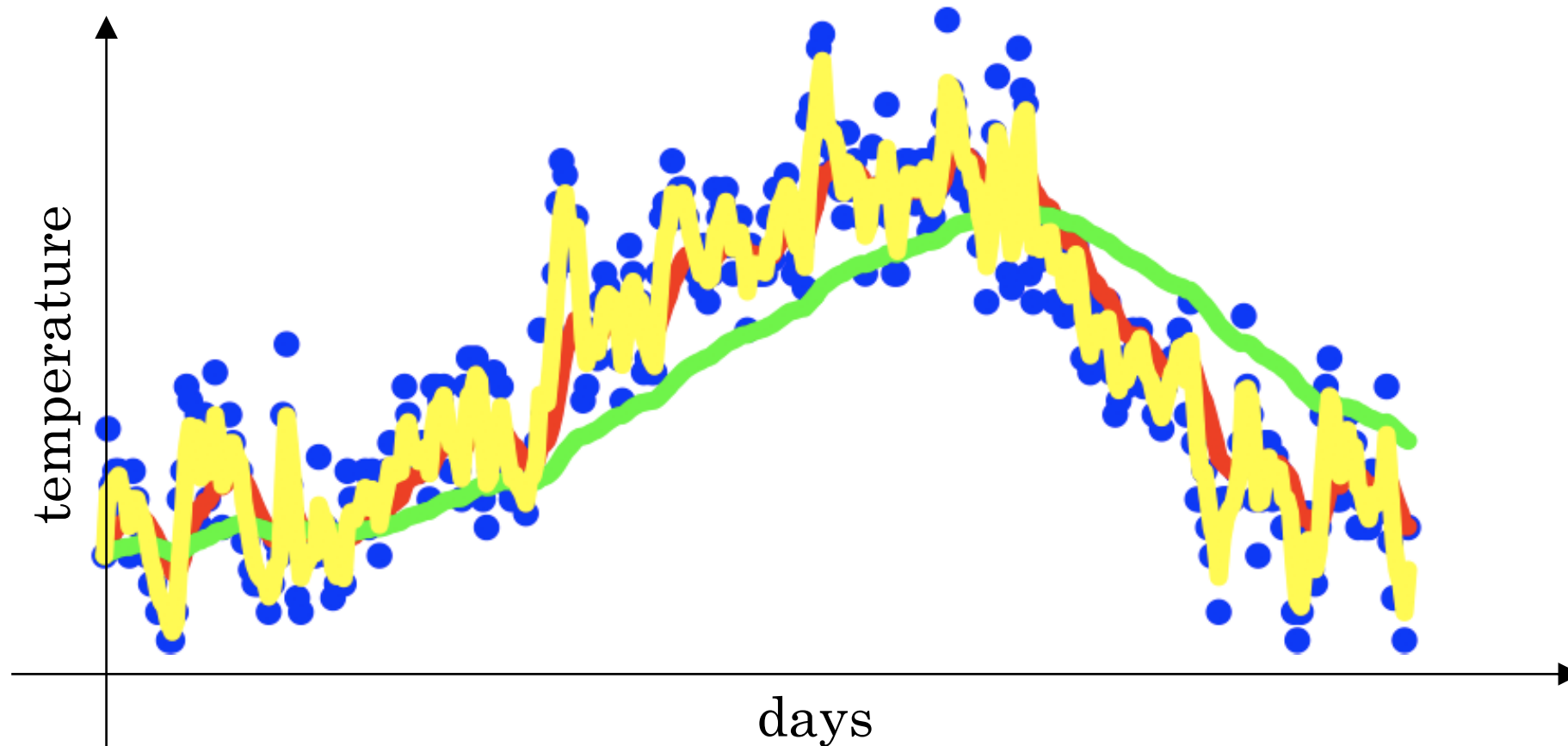Algorithms

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

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

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

$$v_{99} = 0.9v_{98} + 0.1\theta_{99}$$
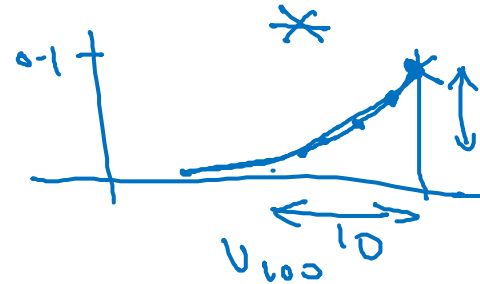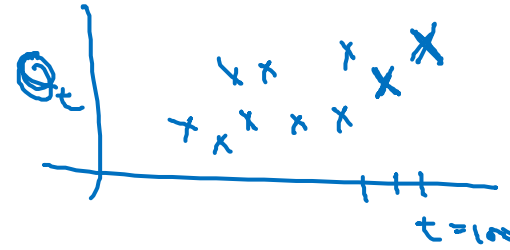
$$v_{98} = 0.9v_{97} + 0.1\theta_{98}$$

...



$\approx \dfrac{1}{1-\beta}$

$\varepsilon = 1-\beta$

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

$V_{100} = 0.1\,\theta_{100} + 0.9\,\cancel{v_{99}}(0.1\,\theta_{99} + 0.9\,\cancel{v_{98}})$

$= 0.1\,\theta_{100} + 0.1 \times 0.9 \cdot \theta_{99} + 0.1(0.9)^2\theta_{98} + 0.1(0.9)^3\theta_{97} + 0.1(0.9)^4\theta_{96}$
$+ \cdots$

$0.9^{10} \approx 0.35 \approx \dfrac{1}{e}$

$(1-\varepsilon)^{1/\varepsilon} = \dfrac{1}{e}$

$0.9$

$0.98\,?$

$\varepsilon = 0.02 \rightarrow 0.98^{50} \approx \dfrac{1}{e}$

# Implementing exponentially weighted averages

$v_0 = 0$

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

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

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

$\ldots$

$V_\theta := 0$

$V_\theta := \beta v + (1 - \beta)\, \theta_1$

$V_\theta := \beta v + (1 - \beta)\, \theta_2$

$\vdots$

$\rightarrow V_\theta = 0$

Repeat $\{$

    Get next $\theta_t$

    $V_\theta := \beta V_\theta + (1 - \beta)\, \theta_t \quad \leftarrow$

$\}$

# Optimization Algorithms

---

# Bias correction in exponentially weighted average

deeplearning.ai

# Bias correction



$\beta = 0.98$

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

$v_0 = 0$

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

$v_2 = 0.98 v_1 + 0.02\theta_2$

$\quad = 0.98 \times 0.02 \times \theta_1 + 0.02\theta_2$

$\quad = 0.0196\theta_1 + 0.02\theta_2$

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

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

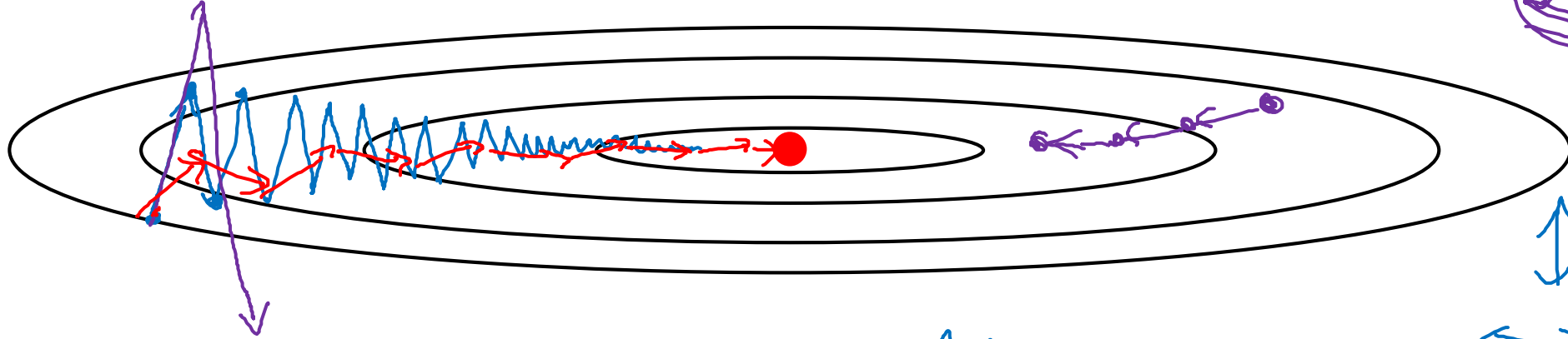$\dfrac{v_2}{0.0396} = \dfrac{0.0196\theta_1 + 0.02\theta_2}{0.0396}$

deeplearning.ai

Optimization
Algorithms

Gradient descent
with momentum

# Gradient descent example



slower learning

faster learning.

Momentum:

On iteration $t$:

Compute $dW, db$ on current mini-batch.

$$V_{dW} = \beta V_{dW} + (1-\beta) dW$$

$$V_{db} = \beta V_{db} + (1-\beta) db$$

friction ⟶ velocity

$$W = W - \alpha V_{dW}, \quad b = b - \alpha V_{db}$$

acceleration

$$" V_{\theta} = \beta V_{\theta} + (1-\beta) \theta_t "$$

deeplearning.ai

# Implementation details

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

On iteration $t$:

    Compute $dW, db$ on the current mini-batch

$$\rightarrow v_{dW} = \beta v_{dW} + (1 - \beta)dW \}$$

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

$$v_{dW} = \beta v_{dW} + \quad dW \leftarrow$$

$$W = W - \alpha v_{dW}, \quad b = b - \alpha v_{db}$$

$$\frac{v_{dw}}{1 - \beta^t}$$

Hyperparameters: $\alpha, \beta$          $\beta = 0.9$

average over last $\approx 10$ gradients

Optimization
Algorithms

---

RMSprop

deeplearning.ai

# RMSprop



$W_1, W_2, W_3$

$W_3, W_4, \ldots$

slow

fast

On iteration $t$:

Compute $dW, db$ on current mini-batch

$S_{dW} = \beta_2 S_{dW} + (1-\beta_2) dW^2$ ← element-wise ← small

$\rightarrow S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2$ ← large

$W := W - \alpha \dfrac{dW}{\sqrt{S_{dW} + \varepsilon}}$ ←      $b := b - \alpha \dfrac{db}{\sqrt{S_{db} + \varepsilon}}$ ←

$\varepsilon = 10^{-8}$

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 $dW, db$ using current mini-batch

    $V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dW \ , \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) db \quad \Leftarrow \text{"momentum" } \beta_1$

    $S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dW^2 \ , \quad S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2 \quad \Leftarrow \text{"RMSprop" } \beta_2$

yhat = np.array([.9, 0.2, 0.1, .4, .9])

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

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

$W := W - \alpha \dfrac{V_{dw}^{corrected}}{\sqrt{S_{dw}^{corrected}} + \varepsilon} \qquad\qquad b := b - \alpha \dfrac{V_{db}^{corrected}}{\sqrt{S_{db}^{corrected}} + \varepsilon}$

# Hyperparameters choice:

$\rightarrow \alpha$ : needs to be tune

$\rightarrow \beta_1$ : 0.9 $\longrightarrow (\underline{dw})$

$\rightarrow \beta_2$ : 0.999 $\longrightarrow (\underline{dw^2})$

$\rightarrow \varepsilon$ : $10^{-8}$

Adam : Adaptive moment estimation



Adam Coates

Optimization
Algorithms

Learning rate
decay

deeplearning.ai
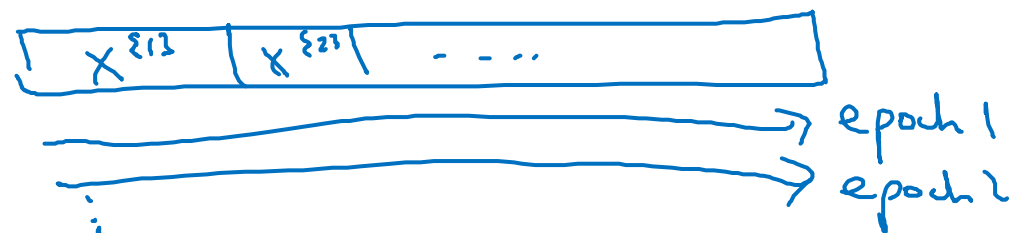
# Learning rate decay

Slowly reduce $\alpha$

# Learning rate decay

1 epoch = 1 pass through data.

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

$X^{\{1\}}$ | $X^{\{2\}}$ | - - - -

→ epoch 1
→ epoch 2
⋮

$\alpha_0 = 0.2$

$\text{decay-rate} = 1$

| Epoch | $\alpha$ |
|-------|------|
| 1 | 0.1 |
| 2 | 0.67 |
| 3 | 0.5 |
| 4 | 0.4 |
| ⋮ | ⋮ |

# Other learning rate decay methods

$$\alpha = 0.95^{epoch\text{-}num} \cdot \alpha_0 \qquad - \text{exponentially decay.}$$

$$\alpha = \frac{k}{\sqrt{epoch\text{-}num}} \cdot \alpha_0 \qquad or \qquad \frac{k}{\sqrt{t}} \cdot \alpha_0$$

discrete staircase



Manual decay.

# Optimization Algorithms

---

# The problem of local optima

# Local optima in neural networks

# Problem of plateaus



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

Hyperparameter tuning

Tuning process

deeplearning.ai

# Hyperparameters

$\rightarrow$ $\boxed{\alpha}$

$\boxed{\beta}$ No.9

$\beta_1$, $\beta_2$, $\varepsilon$
$0.9$  $0.999$  $10^{-8}$

$\boxed{\text{\# layers}}$

$\boxed{\text{\# hidden units}}$

$\boxed{\text{learning rate decay}}$

$\boxed{\text{Mini-batch size}}$

# Try random values: Don't use a grid

# Coarse to fine

# Hyperparameter tuning

deeplearning.ai

---

# Using an appropriate scale to pick hyperparameters

# Picking hyperparameters at random

$\longrightarrow n^{[\ell]} = 50, \ldots, 100$

$$\left[ \underset{50}{\underline{\quad \times \quad \times \times \times \quad \times \times \quad \times \quad \times \quad \times \quad}} \right]_{100}$$

$\longrightarrow$ #layers $\quad L: \quad 2 - 4$

$$2, 3, 4$$

# Appropriate scale for hyperparameters

$\alpha = 0.0001 \ldots, 1$



$0.0001$

$0.1$ — — — — — — — 1

$0.0001$  $0.001$  $0.01$  $0.1$  $10^0$  $10^6$

$10^{-4}$

$10^a$

$a = \log_{10} 0.0001$  $r = -4 * np.random.rand()$  $\leftarrow$  $r \in [-4, 0]$  $b = \log_{10} 1$

$= -4$  $\alpha = 10^r$  $\leftarrow$  $10^{-4} \ldots 10^0$  $= 0$

$10^a \ldots 10^b$  $\dfrac{r \in [a, b]}{[-4, 0]}$  $\alpha = 10^r$

# Hyperparameters for exponentially weighted averages

$\beta = 0.9 \quad \cdots \quad 0.999$

$\downarrow \qquad\qquad\qquad \downarrow$

$10 \qquad\qquad\qquad 1000$

$1-\beta = 0.1 \quad \cdots \quad 0.001$

---

$\beta: \quad 0.9000 \rightarrow 0.9005 \quad \} \sim 10$

$\beta: \quad 0.999 \rightarrow 0.9995$

$\qquad \sim 1000 \qquad\qquad \sim 2000$

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

$\underset{\uparrow}{\overline{\phantom{xxxx}}} \; \underset{x}{\times} \; x \; x \; \underset{x}{\times} \; x \quad \leftarrow$

$0.9 \qquad\qquad\qquad\qquad 0.999$

$0.9 \qquad\qquad 0.99 \qquad\qquad 0.999$

$0.1 \qquad 0.01 \qquad 0.001$

$\underline{10^{-1}} \qquad\qquad\qquad\qquad \overline{10^{-3}}$
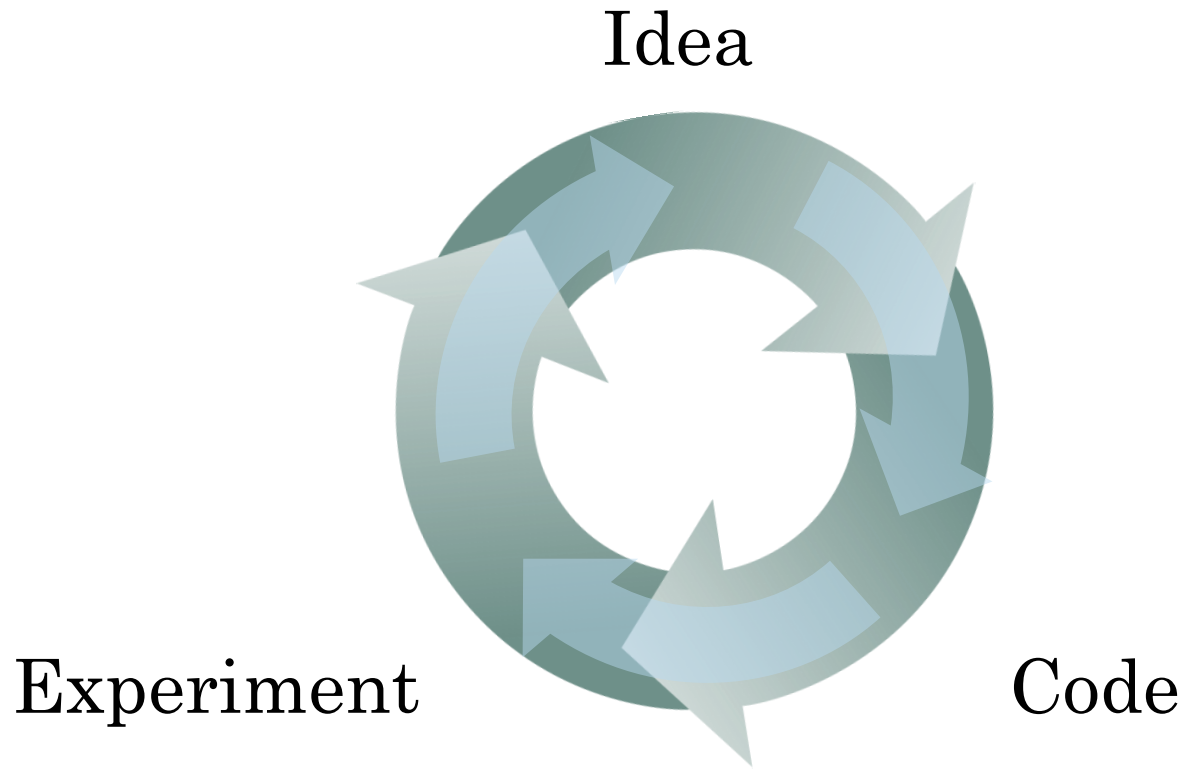
$r \in [-3, -1]$

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

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

deeplearning.ai

Hyperparameters
tuning

---

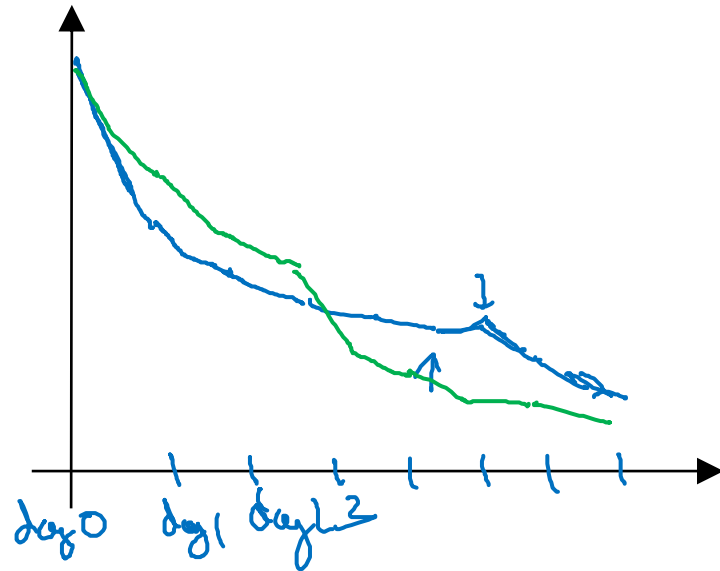Hyperparameters
tuning in practice:
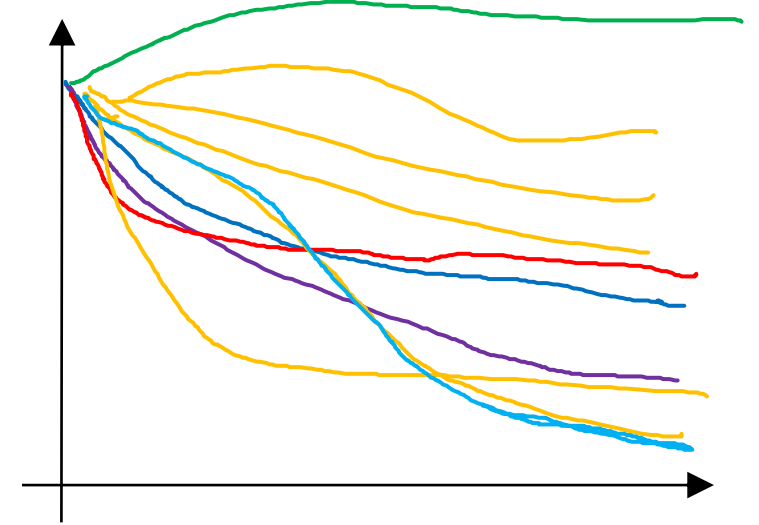Pandas vs. Caviar

# Re-test hyperparameters occasionally

Idea

Experiment

Code

- NLP, Vision, Speech, Ads, logistics, ....

- Intuitions do get stale. Re-evaluate occasionally.

# Babysitting one model


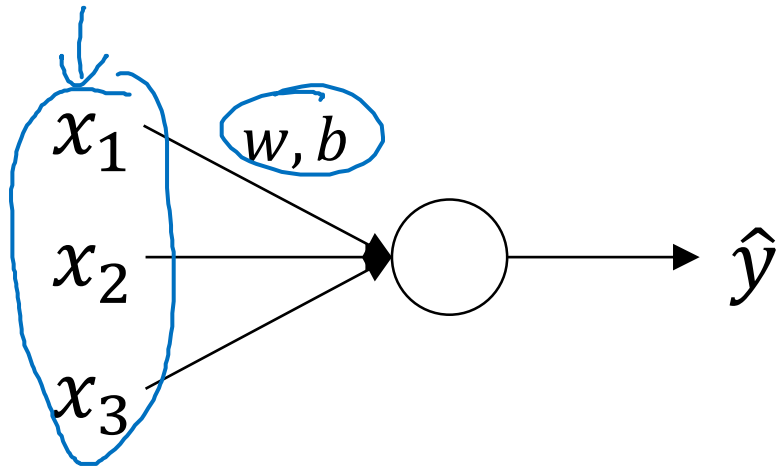
Panda

# Training many models in parallel



Caviar

Batch
Normalization
_____

Normalizing activations
in a network

deeplearning.ai

# Normalizing inputs to speed up learning

$x_1$

$w, b$
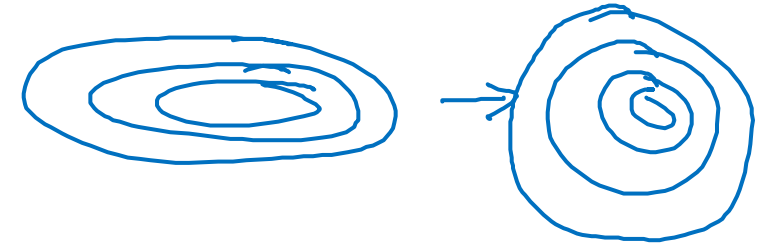
$x_2$

$x_3$

$\hat{y}$

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

$$X = X - \mu$$

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

$$X = X / \sigma^2$$

$a^{[1]}$

$a^{[2]}$

$W^{[3]}, b^{[3]}$

$x_1$

$x_2$

$x_3$

$\hat{y}$

Can we normalize $a^{[2]}$ so as to train $W^{[3]}, b^{[3]}$ faster

Normalize $z^{[2]}$ ↑

# Implementing Batch Norm

Given some intermediate values in NN     $\downarrow \quad \downarrow$
$$z^{(1)}, \ldots, z^{(m)}$$

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

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

$$z_{norm}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \varepsilon}} \quad \Leftarrow$$
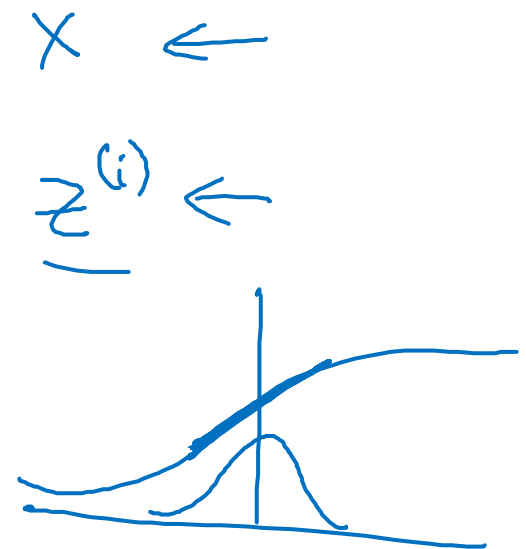
$$\tilde{z}^{(i)} = \gamma \, z_{norm}^{(i)} + \beta$$

learnable parameters of model.

If

$$\gamma = \sqrt{\sigma^2 + \varepsilon} \quad \Leftarrow$$

$$\beta = \mu \quad \Leftarrow$$

then $\tilde{z}^{(i)} = z^{(i)}$

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

$X \Leftarrow$

$z^{(i)} \Leftarrow$
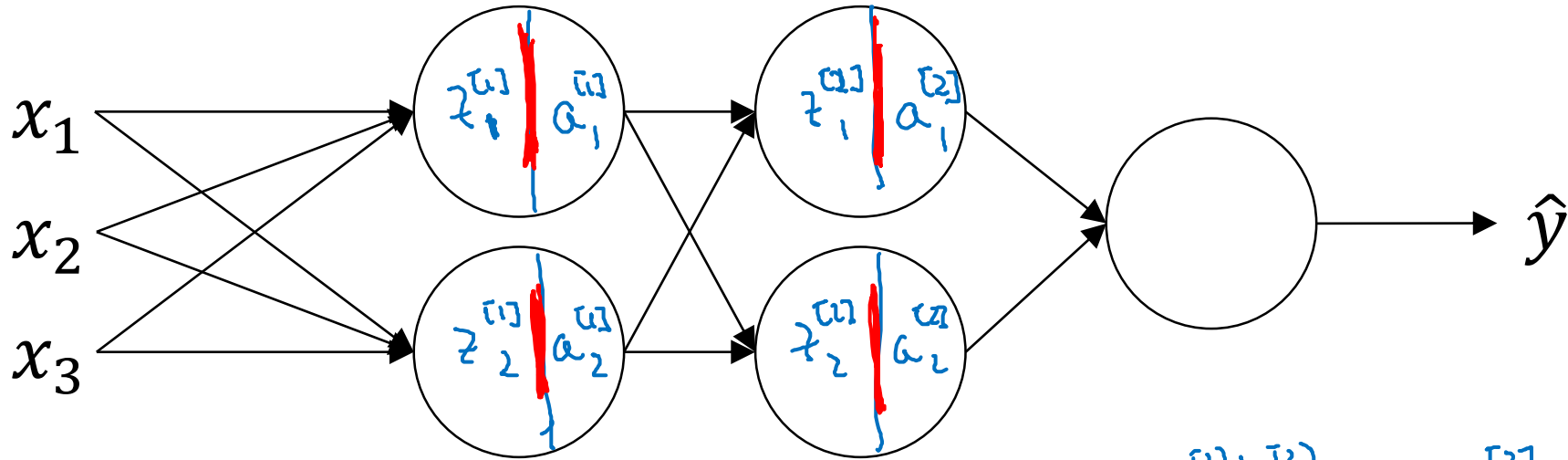
Use $\tilde{z}^{[l](i)}$ instead of $z^{[l](i)}$.

Batch
Normalization

Fitting Batch Norm
into a neural network

deeplearning.ai

# Adding Batch Norm to a network



$$X \xrightarrow{W^{[1]}, b^{[1]}} Z^{[1]} \xrightarrow[\text{Batch Norm (BN)}]{\beta^{[1]}, \gamma^{[1]}} \tilde{Z}^{[1]} \rightarrow a^{[1]} = g^{[1]}(\tilde{Z}^{[1]}) \xrightarrow{W^{[2]}, b^{[2]}} Z^{[2]} \xrightarrow[\text{BN}]{\beta^{[2]}, \gamma^{[2]}} \tilde{Z}^{[2]} \rightarrow a^{[2]} \rightarrow \cdots$$

Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \ldots, W^{[L]}, b^{[L]},$
$\rightarrow \beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \ldots, \beta^{[L]}, \gamma^{[L]}$

$d\beta^{[L]}$

$\beta^{[L]} = \beta^{[L]} - \alpha \, d\beta^{[L]}$

$\rightarrow \beta$

tf.nn.batch-normalization $\Leftarrow$

# Working with mini-batches

$$X^{\{1\}} \xrightarrow{W^{[1]}, b^{[1]}} Z^{[1]} \xrightarrow[BN]{\beta^{[1]}, \gamma^{[1]}} \tilde{Z}^{[1]} \rightarrow g^{[1]}(\tilde{Z}^{[1]}) = a^{[1]} \xrightarrow{W^{[2]}, b^{[2]}} Z^{[2]} \rightarrow \dots$$

$$\boxed{X^{\{2\}}} \rightarrow Z^{[1]} \xrightarrow[\boxed{BN}]{\beta^{[1]}, \gamma^{[1]}} \tilde{Z}^{[1]} \rightarrow \dots$$

$$X^{\{3\}} \rightarrow \dots$$

---

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

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

$Z^{[l]}$

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

$$\rightarrow Z^{[l]} = W^{[l]} a^{[l-1]} + \cancel{\boxed{b^{[l]}}}$$

$$Z^{[l]} = W^{[l]} a^{[l-1]}$$

$$Z^{[l]}_{norm}$$

$$\rightarrow \tilde{Z}^{[l]} = \gamma^{[l]} Z^{[l]}_{norm} + \boxed{\beta^{[l]}} \leftarrow$$

# Implementing gradient descent

for $t = 1 \ldots$ num MiniBatches

   Compute forward prop on $X^{\{t\}}$.

     In each hidden layer, use BN to repar $Z^{[l]}$ with $\tilde{Z}^{[l]}$.

   Use backprop to compute $dW^{[l]}$, ~~$db^{[l]}$~~, $d\beta^{[l]}$, $d\gamma^{[l]}$

   Update params
$$W^{[l]} := W^{[l]} - \alpha \, dW^{[l]}$$
$$\beta^{[l]} := \beta^{[l]} - \alpha \, d\beta^{[l]}$$
$$\gamma^{[l]} := \ldots$$
$\Big\}\ \leftarrow$

Works w/ momentum, RMSprop, Adam.

Batch
Normalization

deeplearning.ai

Why does
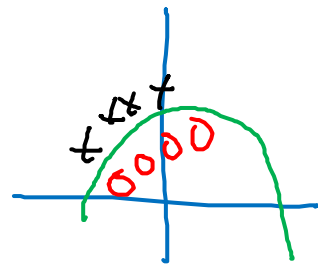Batch Norm work?

# Learning on shifting input distribution



$x_1$
$x_2$
$x_3$

$\hat{y}$

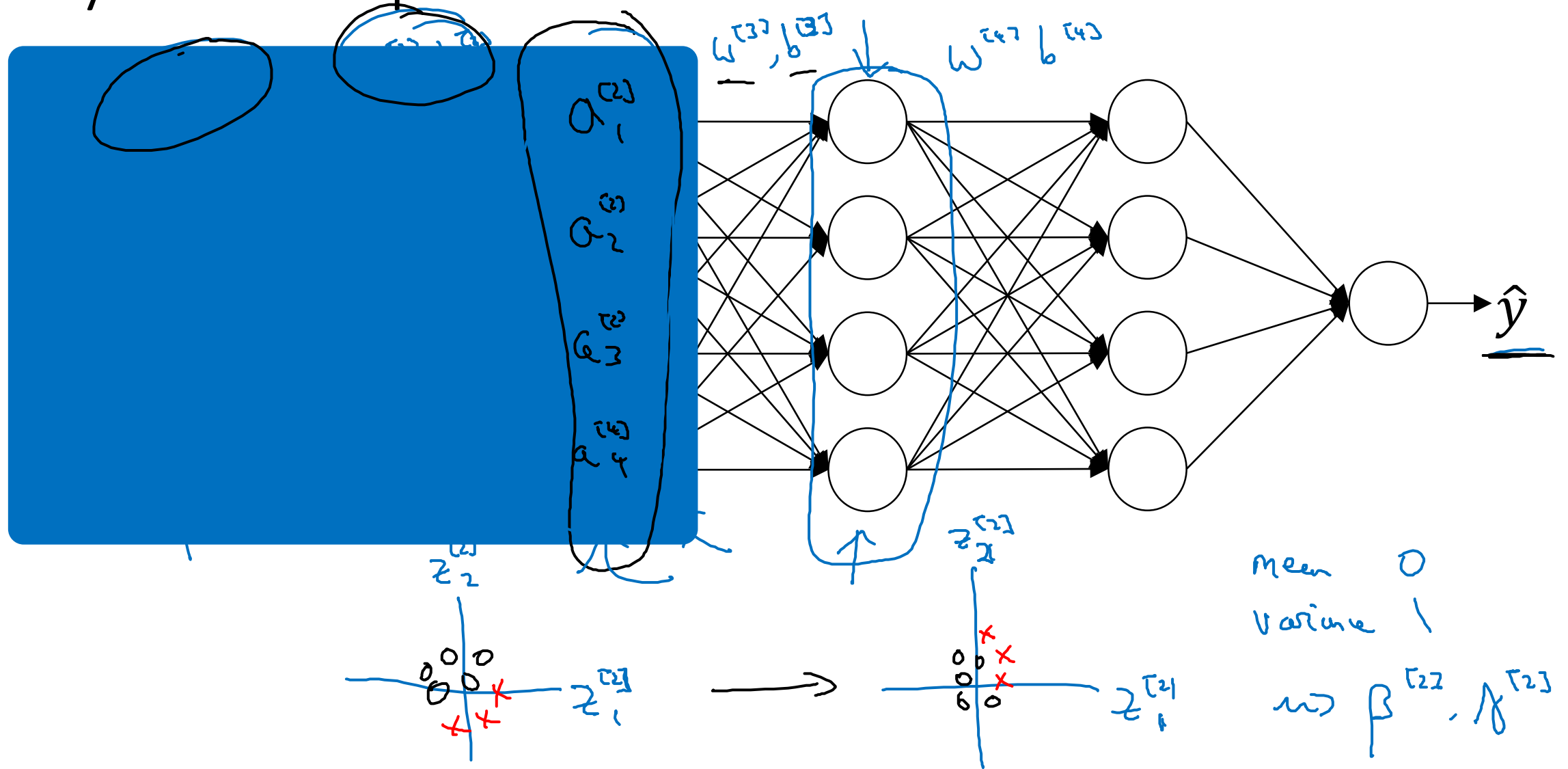Cat $\quad$ Non-Cat

$y = 1 \quad\quad y = 0$

$y = 1 \quad\quad y = 0$

"Covariate shift"

$x \longrightarrow y$

# Why this is a problem with neural networks?

# Batch Norm as regularization

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.

- 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.

- This has a slight regularization effect.

$$\rightarrow \tilde{z}^{[l]} \qquad \{64, 128\} \qquad z^{[l]}$$

$$X \qquad X^{\{t\}}$$

$$\mu, \sigma^2$$

$$\text{Mini-batch}: \underline{64} \longrightarrow \underline{512}$$

Batch Normalization

deeplearning.ai

Batch Norm at test time

# Batch Norm at test time

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

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

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

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

$\mu, \sigma^2$: estimate using exponentially weighted average (across mini-batches).

$X^{\{1\}}, X^{\{2\}}, X^{\{3\}}, \ldots$

$\mu^{\{1\}[l]} \qquad \mu^{\{2\}[l]} \qquad \mu^{\{3\}[l]} \longrightarrow \mu$

$\theta_1 \qquad\qquad \theta_2 \qquad\qquad \theta_3 \qquad\qquad\qquad \sigma^2$

$\sigma^{2\{1\}[l]} \qquad \sigma^{2\{2\}[l]} \qquad \ldots$

$$z_{\text{norm}} = \frac{z - \mu}{\sqrt{\sigma^2 + \varepsilon}} \qquad\qquad \tilde{z} = \gamma z_{\text{norm}} + \beta$$

Multi-class
classification

Softmax regression

deeplearning.ai

# Recognizing cats, dogs, and baby chicks



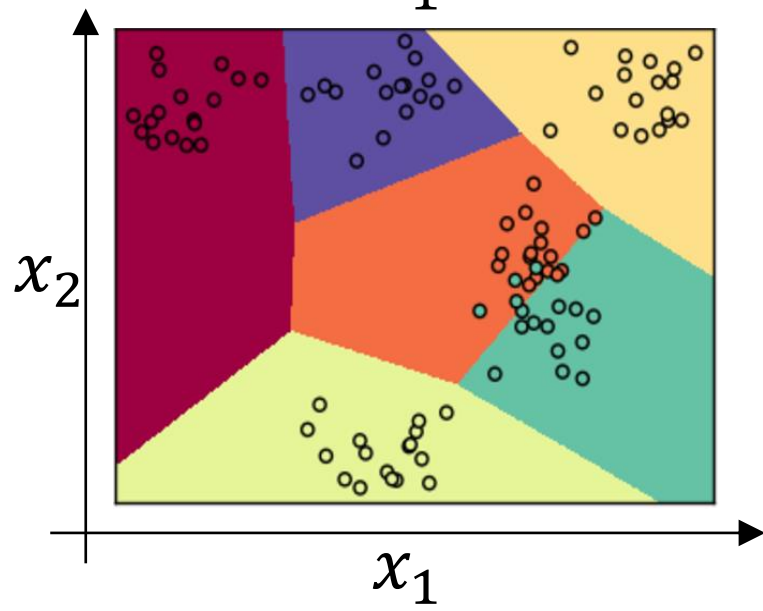3      1      2      0      3      2      0      1

$X \rightarrow \hat{y}$

# Softmax layer

$$X \longrightarrow \hat{y}$$
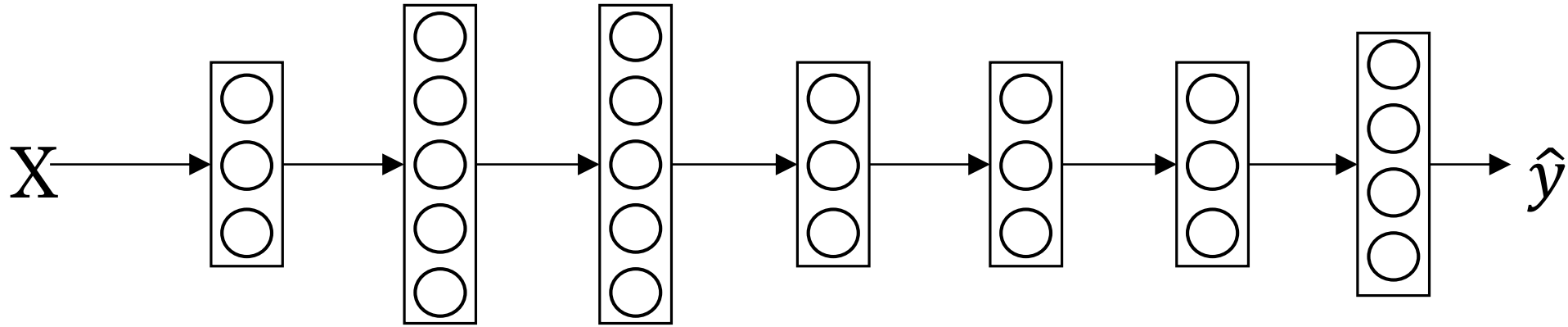
# Softmax examples

Multi-class classification

Trying a softmax classifier

deeplearning.ai

# Understanding softmax

# Loss function

# Summary of softmax classifier

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)

Programming
Frameworks

TensorFlow

deeplearning.ai

# Motivating problem

$$J(\omega) = \boxed{\omega^2 - 10\omega + 25}$$

(cost)

$(\omega - 5)^2$

$\omega = 5$

$$J(W, b)$$
$\uparrow \quad \uparrow$

# Code example

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

for i in range(1000):
    session.run(train, feed_dict={x:coefficients})
print(session.run(w))
```