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.

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

$(n_x, m)$

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

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

$(1, m)$

$\underbrace{\qquad\qquad}_{Y^{\{1\}} \ (1, 1000)}$  $\underbrace{\qquad\qquad}_{Y^{\{2\}} \ (1, 1000)}$  $\underbrace{\qquad}_{Y^{\{5,000\}} \ (1, 1000).}$

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

5,000 mini-batches of 1,000 each

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

$x^{(i)}$

$z^{[l]}$

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

Andrew Ng

# Mini-batch gradient descent

repeat {

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

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

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

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

$$\vdots$$

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

} Vectorized implementation (1000 examples)

$X, Y$

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

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"
    └ pass through training set.

Andrew Ng

# Optimization Algorithms

---

# Understanding mini-batch gradient descent

deeplearning.ai

# Training with mini batch gradient descent

### Batch gradient descent



cost

# iterations

J

### Mini-batch gradient descent



cost

mini batch # (t)

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

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

$\vdots$

$J^{\{t\}}$

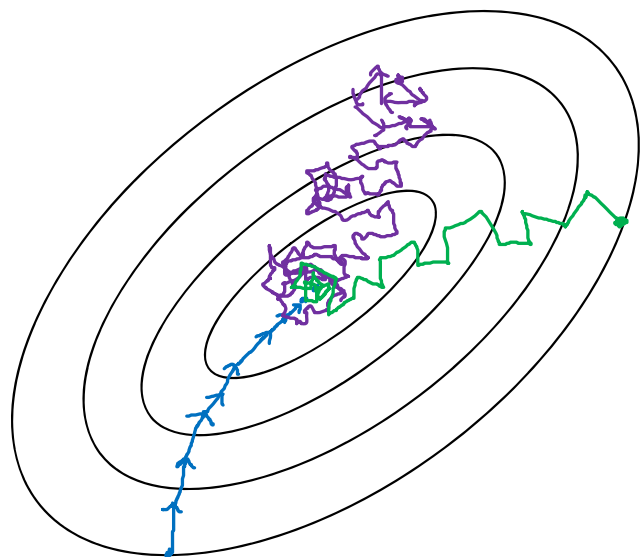Plot $J^{\{t\}}$ computed using $\underline{X^{\{t\}}, Y^{\{t\}}}$

# 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
  $(X^{\{1\}}, Y^{\{1\}}) = (x^{(1)}, y^{(1)}) \ldots (x^{(m)}, y^{(m)})$ 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.
  (~1000)
• Make progress without
  processing entire training set.

Batch
gradient descent
(mini-batch size = m)

↓

Too long
per iteration

Andrew Ng

# Choosing your mini-batch size

If small tray set : Use batch gradient descent.

$\quad\quad$ $(m \leq 2000)$

Typical mini-batch sizes :

$\longrightarrow$ $\underset{2^6}{64}$ , $\underset{2^7}{128}$, $\underset{2^8}{256}$, $\underset{2^9}{512}$ $\quad\quad\quad\quad\quad$ $\dfrac{1024}{2^{10}}$

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

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

deeplearning.ai

Optimization
Algorithms

Exponentially
weighted averages

# Temperature in London

$\theta_1 = 40°F$  4°C ←

$\theta_2 = 49°F$  9°C

$\theta_3 = 45°F$

$\vdots$

$\theta_{180} = 60°F$  15°C

$\theta_{181} = 56°F$

$\vdots$

temperature

↑

days

↑        ↑        ↑

$V_0 = 0$

$V_1 = 0.9 V_0 + 0.1 \theta_1$

$V_2 = 0.9 V_1 + 0.1 \theta_2$

$V_3 = 0.9 V_2 + 0.1 \theta_3$

$\vdots$

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

Andrew Ng

# Exponentially weighted averages

_moving_

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

$\beta = 0.9 \quad : \quad \approx 10 \text{ days' temperature}$

$\beta = 0.98 \quad : \quad \approx 50 \text{ days}$

$\beta = 0.5 \quad : \quad \approx 2 \text{ days}$

$V_t$ as approximately average over

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

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



temperature / days

deeplearning.ai

# Optimization Algorithms

## Understanding exponentially weighted averages

# Exponentially weighted averages

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

$\beta = 0.9 \qquad 0.98 \qquad 0.5$



Andrew Ng

# Exponentially weighted averages

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

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

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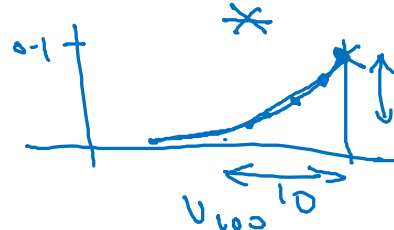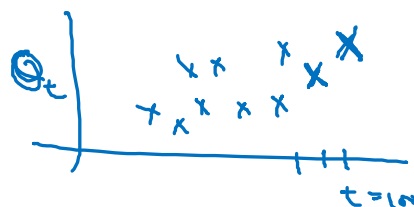
$$v_{98} = 0.9v_{97} + 0.1\theta_{98}$$

...

$\theta_t$

$t=100$

$0.1$

$v_{100}$  $10$

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

$\varepsilon = 1-\beta$

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

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

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

$+ \ldots$

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

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

$0.9$

$0.98\,^?$

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

Andrew Ng

# Implementing exponentially weighted averages

$v_0 = 0$

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

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

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

...

$V_\theta := 0$

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

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

$\vdots$

$\rightarrow V_\theta = 0$

Repeat {

$\quad$ Get next $\theta_t$

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

}

Andrew Ng

deeplearning.ai

# Optimization Algorithms

---

## Bias correction in exponentially weighted average

# Bias correction



temperature / days

$\beta = 0.98$

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

$v_0 = 0$

$v_1 = 0.98 v_0 + 0.02 \theta_1$

$v_2 = 0.98 v_1 + 0.02 \theta_2$

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

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

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

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

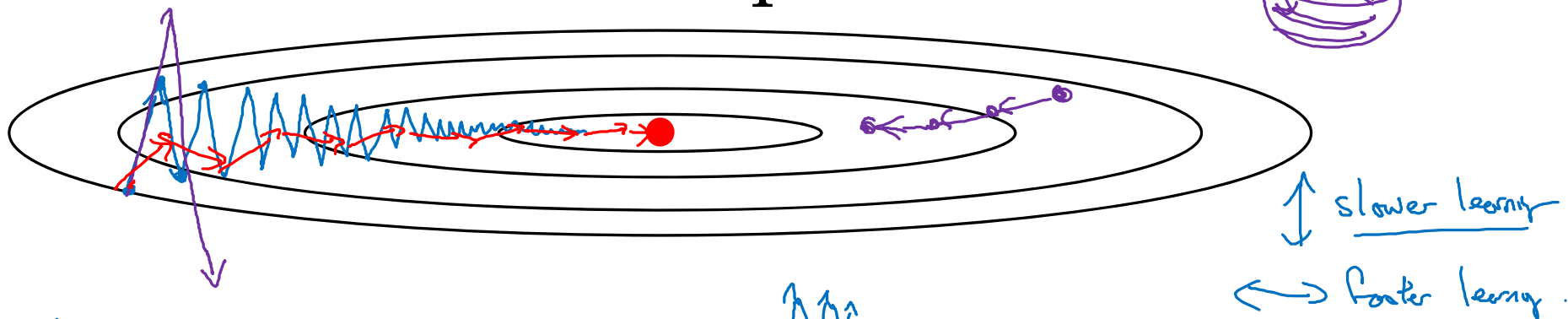$\dfrac{v_2}{0.0396} = \dfrac{0.0196\theta_1 + 0.02\theta_2}{0.0396}$

Andrew Ng

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 $\quad$ velocity

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

acceleration

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

deeplearning.ai

Andrew Ng

# Implementation details

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

On iteration $t$:

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

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

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

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

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

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

average over last $\approx 10$ gradients

Andrew Ng

Optimization
Algorithms

RMSprop

deeplearning.ai

# RMSprop



$b$    $w$

$w_1, w_2, w_2$

$w_3, w_4, \dots$

$\uparrow$ slow

$\leftrightarrow$ fast

On iteration $t$:

Compute $dW, db$ on current mini-batch

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

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

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

$\varepsilon = 10^{-8}$

Andrew Ng

Optimization
Algorithms

Adam optimization
algorithm

deeplearning.ai

# Adam optimization algorithm

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

On iteration $t$:

Compute $dW, db$ using current mini-batch

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

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

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

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

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

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

Andrew Ng

# 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

Andrew Ng

Optimization Algorithms

Learning rate decay
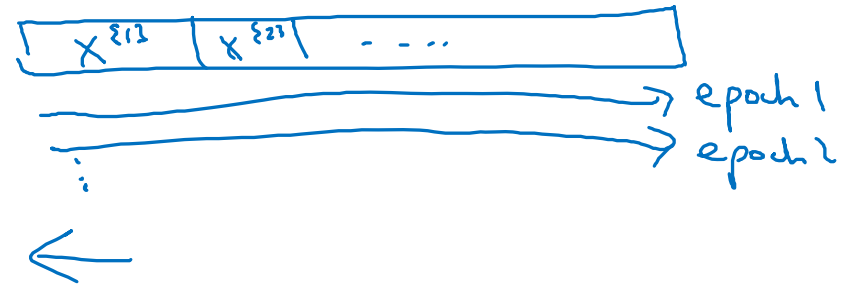
deeplearning.ai

# Learning rate decay



Slowly reduce $\alpha$

Andrew Ng

# Learning rate decay

1 epoch = 1 pass through data.

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

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

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

⟶ epoch 1
⟶ epoch 2
⋮

$\alpha_0 = 0.2$

decay-rate = 1

# Other learning rate decay methods

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

formula $\Big\{$

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

discrete staircase

$\alpha$ | graph of decreasing staircase vs $t$
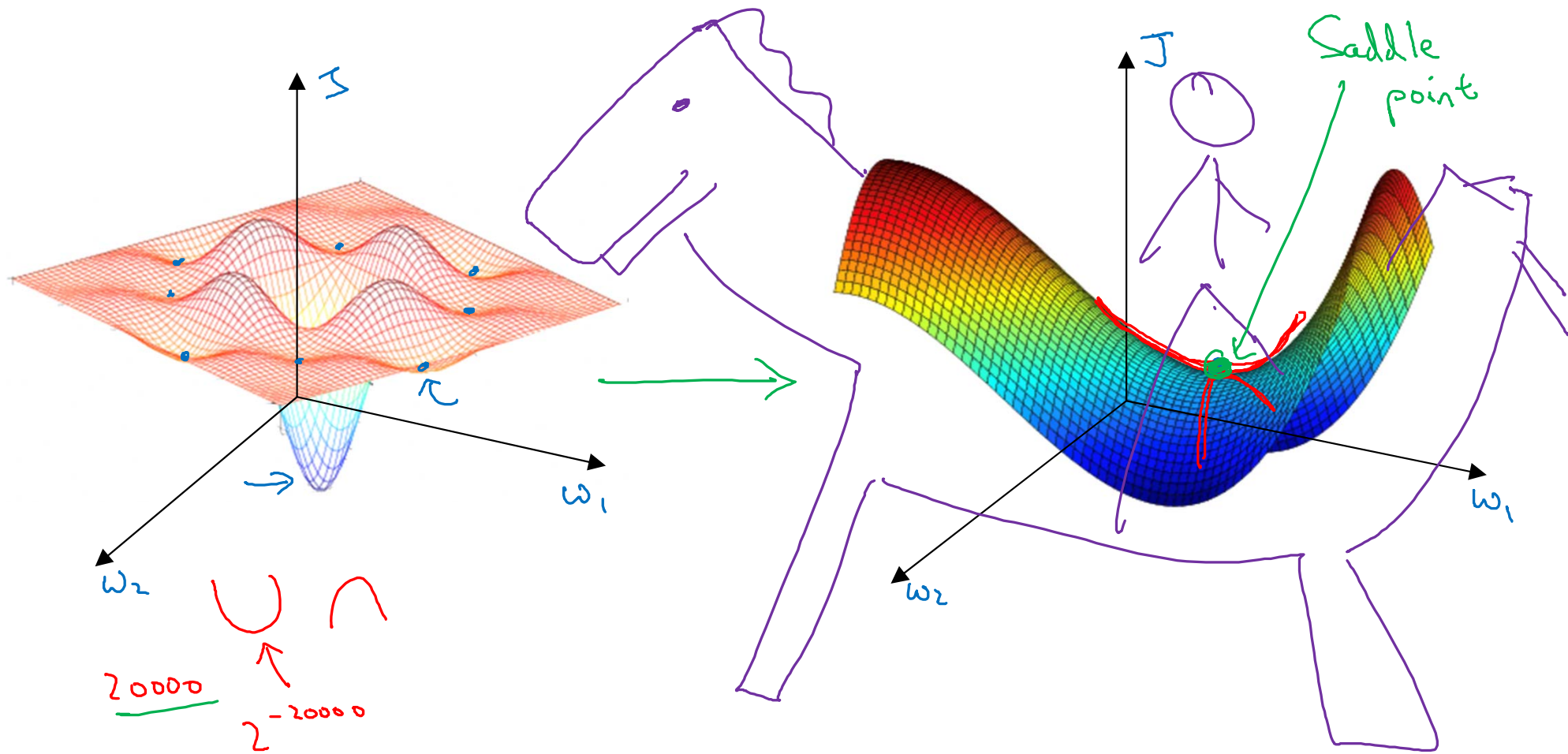
Manual decay.

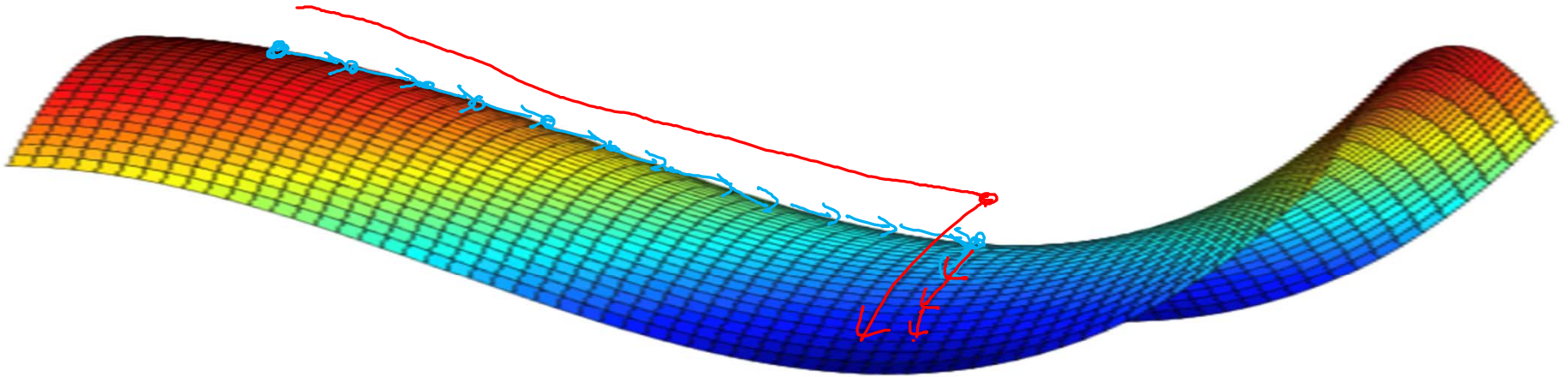Andrew Ng

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