

l.



deeplearning.ai

# Optimization Algorithms

## Mini-batch gradient descent

실용적이고 반복적인 과정 → 빠르게 훈련!  
→ 좋은 최적화 알고리즘!

### Batch vs. mini-batch gradient descent

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

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(1000)} & | & x^{(1001)} & \dots & x^{(2000)} & | & \dots & | & \dots & x^{(m)} \end{bmatrix}$$

$(n_x, m)$        $X^{\{1\}}$   $(n_x, 1000)$        $X^{\{2\}}$   $(n_x, 1000)$        $X^{\{5,000\}}$   $(n_x, 1000)$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(1000)} & | & y^{(1001)} & \dots & y^{(2000)} & | & \dots & | & \dots & y^{(m)} \end{bmatrix}$$

$(1, m)$        $Y^{\{1\}}$   $(1, 1000)$        $Y^{\{2\}}$   $(1, 1000)$        $Y^{\{5,000\}}$   $(1, 1000)$

What if  $m = 5,000,000$ ? 벡터화 여전히 노련  
5,000 mini-batches of 1,000 each (정확히 1,000씩 끊어서)  
Mini-batch  $t$ :  $X^{\{t\}}, Y^{\{t\}}$   
작은 훈련세트

$x^{(i)}$   $i$  번째 훈련샘플  
 $z^{[L]}$  " 신경망  
 $X^{\{t\}}, Y^{\{t\}}$

Andrew Ng

### Mini-batch gradient descent

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

① 정방향 전파 Forward prop on  $X^{\{t\}}$   $\leftarrow$  추가!

$$\begin{aligned} z^{[L]} &= W^{[L]} X^{\{t\}} + b^{[L]} \\ A^{[L]} &= g^{[L]}(z^{[L]}) \\ &\vdots \\ A^{[2]} &= g^{[2]}(z^{[2]}) \end{aligned}$$

Vectorized implementation (1000 examples)  
for  $X^{\{t\}}, Y^{\{t\}}$

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

③ Backprop to compute gradients w.r.t  $J^{\{t\}}$  (minimize  $J(X^{\{t\}}, Y^{\{t\}})$ )  
 $W^{[L]} = W^{[L]} - \alpha dW^{[L]}, b^{[L]} = b^{[L]} - \alpha db^{[L]}$

"1 epoch" single pass through training set.

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

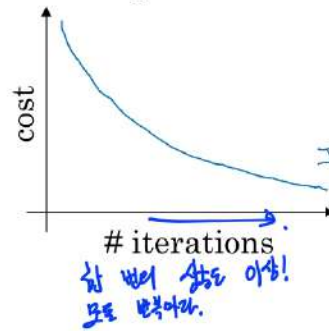
$X, Y$

정기화항

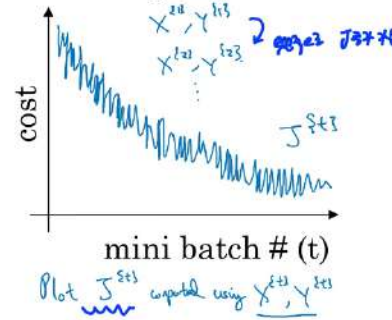
Andrew Ng



Batch gradient descent



Mini-batch gradient descent

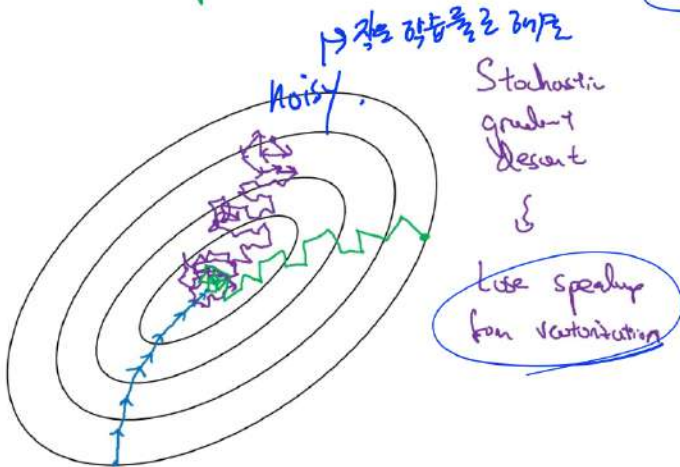


Andrew Ng

정정하기.

## Choosing your mini-batch size

- If mini-batch size  $\in m$ : Batch gradient descent.  $(X^{t+1}, Y^{t+1}) = (X, Y)$ .  $\# \text{ of updates} = \text{entire training set}.$
- If mini-batch size  $\in 1$ : Stochastic gradient descent. Every example is its own mini-batch.  $(X^{t+1}, Y^{t+1}) = (x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})$  mini-batch.
- In practice: Search in-between 1 and  $m$ .



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

Faster learning.

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

Batch gradient descent (mini-batch size =  $m$ )  
 많은 반복이 필요함.  
 Too long per iteration

Andrew Ng

## Choosing your mini-batch size

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

Typical mini-batch sizes:

→ 64, 128, 256, 512, 1024

$2^6, 2^7, 2^8, 2^9, 2^{10}$

일반적 4배, 2배 비율.

⊕ Make sure mini-batch fits in CPU/GPU memory.  $X^{t+1}, Y^{t+1}$

Andrew Ng





deeplearning.ai

## Optimization Algorithms

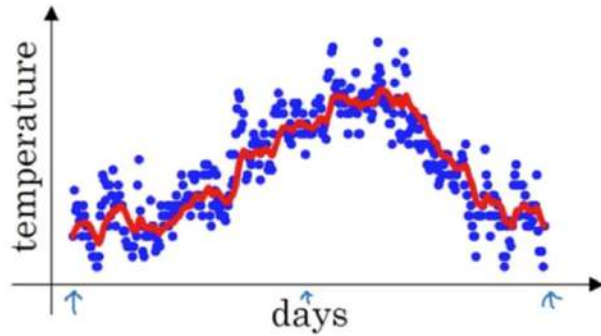
경사하강법 보다 빠른  
최적화 알고리즘

## Exponentially weighted averages

2/5가중이동평균

### Temperature in London

$\theta_1 = 40^\circ\text{F}$   $4^\circ\text{C} \leftarrow$   
 $\theta_2 = 49^\circ\text{F}$   $9^\circ\text{C}$   
 $\theta_3 = 45^\circ\text{F}$   $\vdots$   
 $\vdots$   
 $\theta_{180} = 60^\circ\text{F}$   $15^\circ\text{C}$   
 $\theta_{181} = 56^\circ\text{F}$   $\vdots$   
 $\vdots$



$$\begin{aligned}
 v_0 &= 0 \\
 v_1 &= 0.9 v_0 + 0.1 \theta_1 \\
 v_2 &= 0.9 v_1 + 0.1 \theta_2 \\
 v_3 &= 0.9 v_2 + 0.1 \theta_3 \\
 &\vdots \\
 v_t &= 0.9 v_{t-1} + 0.1 \theta_t
 \end{aligned}$$

Andrew Ng

### Exponentially weighted averages

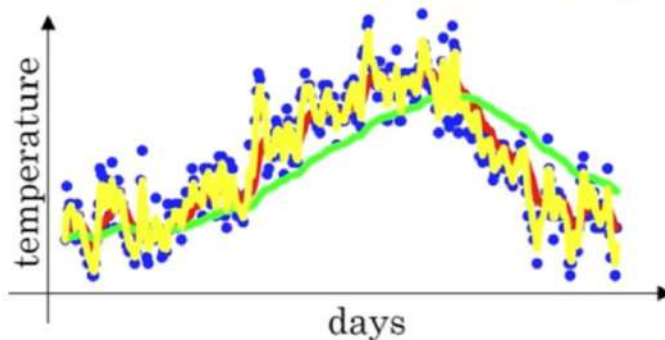
moving

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

$\beta = 0.9$  :  $\approx 10$  days' temperature red.  
 $\beta = 0.98$  :  $\approx 50$  days green  
 $\beta = 0.5$  :  $\approx 2$  days yellow  
 (Note:  $\beta$ 가 클수록 과거에 더 많은 가중치를 준다.)  
 (Note:  $\beta$ 가 작을수록 미래에 더 많은 가중치를 준다.)

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

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



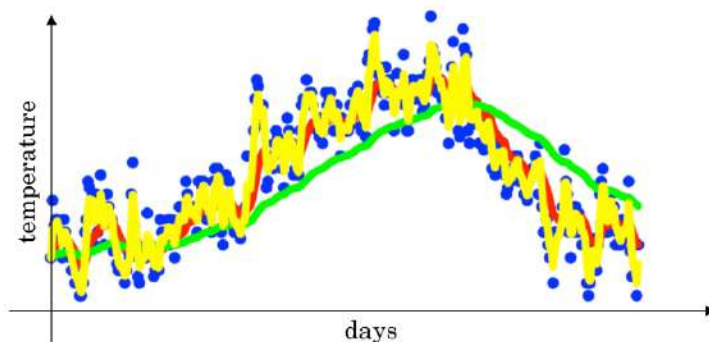
Andrew Ng



## Exponentially weighted averages

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

$$\beta = 0.9 \text{ or } 0.98 \text{ or } 0.5$$



Andrew Ng

## Exponentially weighted averages

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

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

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

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

Handwritten derivation of the exponentially weighted average formula:

$$v_{100} = 0.1\theta_{100} + 0.9(0.1\theta_{99} + 0.9(0.1\theta_{98} + 0.9(0.1\theta_{97} + 0.9(0.1\theta_{96} + \dots)))$$

Handwritten notes and diagrams:

- Diagram showing a sequence of points  $\theta_t$  and a line representing the average.
- Equation:  $\frac{1}{1-\beta} \approx \frac{1}{1-0.9} = 10$
- Equation:  $\sum = 1-\beta$
- Equation:  $0.1\theta_{99} + 0.9v_{99}$
- Equation:  $0.1\theta_{100} + 0.9(0.1\theta_{99} + 0.9(0.1\theta_{98} + 0.9(0.1\theta_{97} + 0.9(0.1\theta_{96} + \dots)))$
- Equation:  $\frac{0.1}{0.9} \approx 0.11$
- Equation:  $\frac{(1-\epsilon)^{1/\epsilon}}{\epsilon} \approx \frac{1}{e}$
- Equation:  $\epsilon = 0.02 \rightarrow 0.98^{50} \approx \frac{1}{e}$

Andrew Ng

## Implementing exponentially weighted averages

$$v_0 = 0 \text{ 초기값}$$

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

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

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

변수 개수  
변수 개수  
...  
3개

$$V_\theta := 0 \text{ 초기값}$$

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

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

...

$$\rightarrow V_0 = 0$$

Repeat {

Get next  $\theta_t$  → 업데이트

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

이걸 3번 반복해서  
(아래의 100번 반복)

Andrew Ng



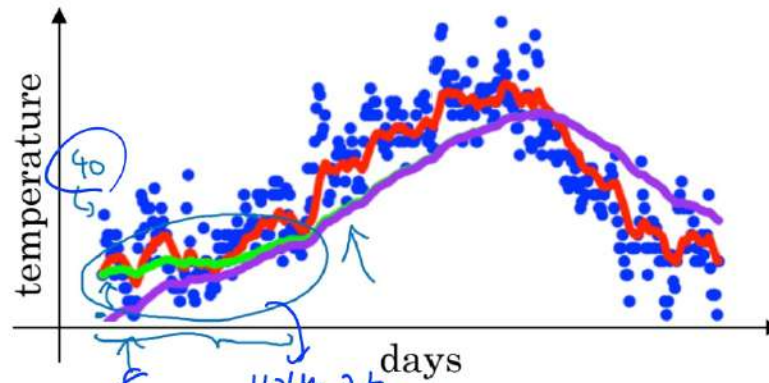
deeplearning.ai

# Optimization Algorithms

## Bias correction in exponentially weighted average

평균 보정.  
→ 평균이 정확히 계산

### Bias correction



$\beta = 0.98$

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

$$v_0 = 0$$

$$v_1 = 0.98 v_0 + 0.02 \theta_1$$

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

$\theta_1$   $\theta_2$  보정  
아직 보정 안함

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

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

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

보정 완료

initial 단계 중요하다.  
t 제곱은 거의 0 ... 항상 보정 완료됨.

Andrew Ng





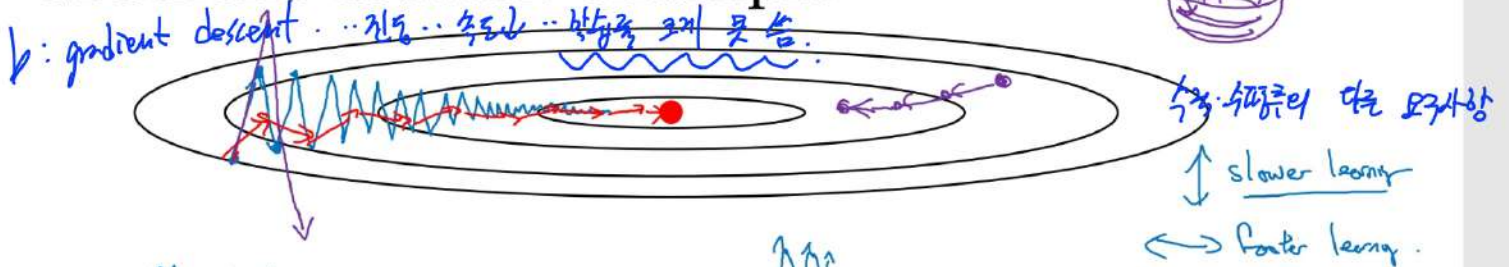
# Optimization Algorithms

모멘텀이 있는 경사 하강법  
→ 경사에 대한 휘저음(흔들림) 개선  
... 가끔씩 업데이트

## Gradient descent with momentum

가이 항상 더 잘 작동한다.

### Gradient descent example



Momentum:

On iteration  $t$ :

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

$$v_{dw} = \beta v_{dw} + (1-\beta) \frac{dW}{dt}$$

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

Friction → Velocity

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



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

가속도

빠르게  
진동 줄이기!

Andrew Ng

### Implementation details

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

$\beta$ : 휘저음(흔들림) 제어 h.p.

On iteration  $t$ :

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

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

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

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

각기 휘저음이 있다

$1/(1-\beta)$ 에 대한 개선을  
시행된다.

$$v_{dw} = \beta v_{dw} + dW \leftarrow ??$$

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

1000 후 휘저음이 점점 감소됨.

Hyperparameters:  $\alpha, \beta$

$$\beta = 0.9$$

average over loss  $\approx 10$  gradient

Andrew Ng

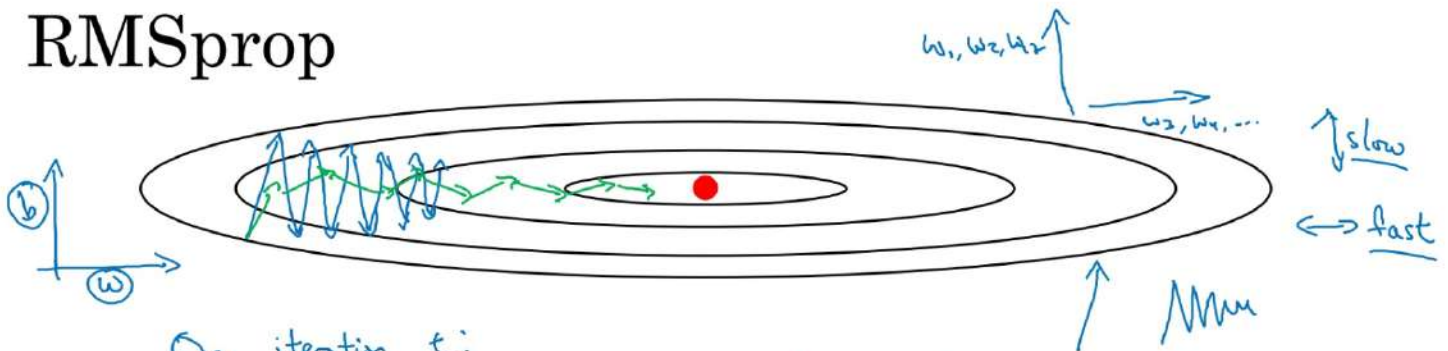


deeplearning.ai

# Optimization Algorithms

## RMSprop

### RMSprop



On iteration  $t$ :

Compute  $dW, db$  on current mini-batch

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

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

$$w := w - \alpha \frac{dw}{\sqrt{S_{dw} + \epsilon}} \leftarrow$$

$$b := b - \alpha \frac{db}{\sqrt{S_{db} + \epsilon}} \leftarrow$$

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

Andrew Ng

장점:

미분값이 큰 곳 : 점동↓

작은 곳 : 더 빠르게 수렴



deeplearning.ai

## Optimization Algorithms

### Adam optimization algorithm

(모멘텀 + 그래디언트)보다  
잘 작동하는 최적화 알고리즘 찾기 어렵다  
→ RMSprop, Adam 잘 안되면...  
RMSprop + 모멘텀

### Adam optimization algorithm

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

On iteration  $t$ :  $\beta$  값

Compute  $dW, db$  using current mini-batch

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

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

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

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

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

### Hyperparameters choice:

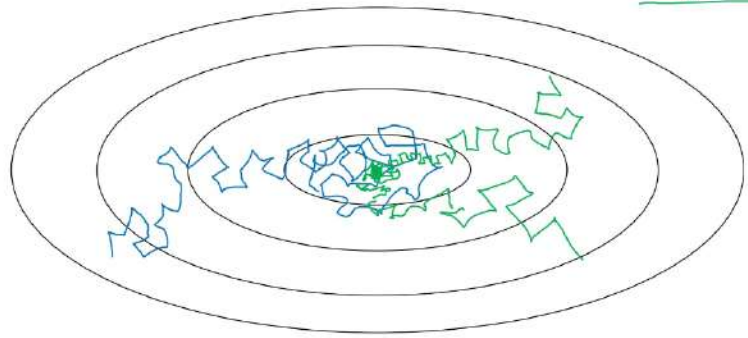
- $\alpha$ : needs to be tune
- $\beta_1$ : 0.9 → ( $dW$ )
- $\beta_2$ : 0.999 → ( $dW^2$ )
- $\epsilon$ :  $10^{-8}$

Adam: Adaptive moment estimation



Adam Coates



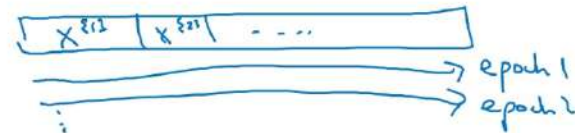


## Learning rate decay

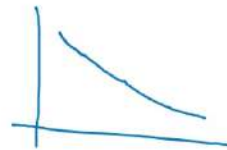
1 epoch = 1 pass through data.

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

Epoch	$\alpha$
1	0.1
2	0.67
3	0.5
4	0.4
$\vdots$	$\vdots$



$\alpha_0 = 0.2$   
decay-rate = 1



## Other learning rate decay methods

Formula

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

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

discrete staircase

Manual decay. step 132  $\alpha$  4321 1322

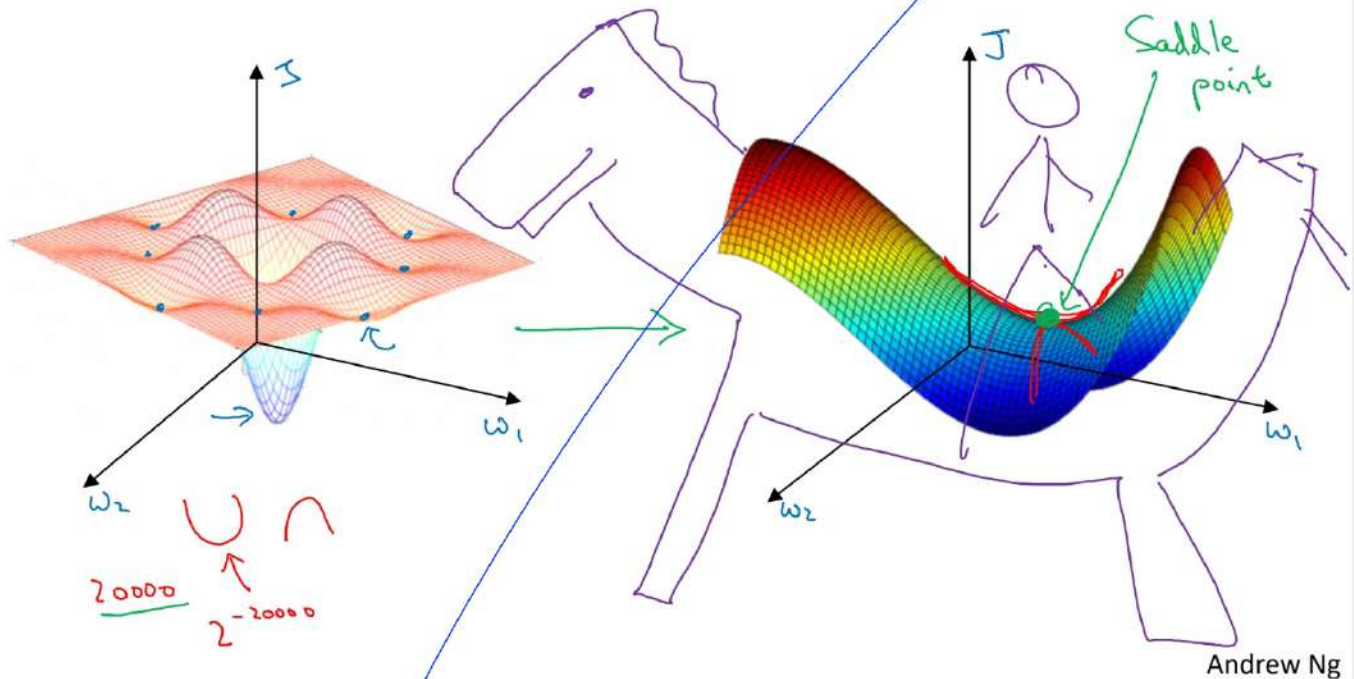


deeplearning.ai

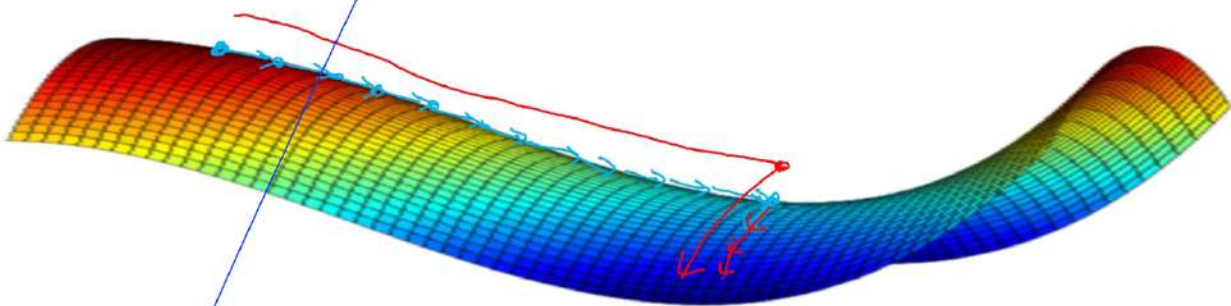
## Optimization Algorithms

### 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

Andrew Ng