

Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



deeplearning.ai

Optimization Algorithms

Mini-batch gradient descent

: 학습률이 너무 크게 설정될 때
일정 속도로 학습을 멈출 때

Batch vs. mini-batch gradient descent

X, Y

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

Vectorization allows you to efficiently compute on m examples.

$$X = \begin{bmatrix} X^{(1)} & X^{(2)} & X^{(3)} & \dots & X^{(500)} & | & X^{(1001)} & \dots & X^{(2000)} & | & \dots & | & \dots & X^{(m)} \end{bmatrix}$$

(n_x, m)

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

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

$(1, m)$

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

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 $t = 1, \dots, 5000$: mini iterations for entire for loop + final.

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

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

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



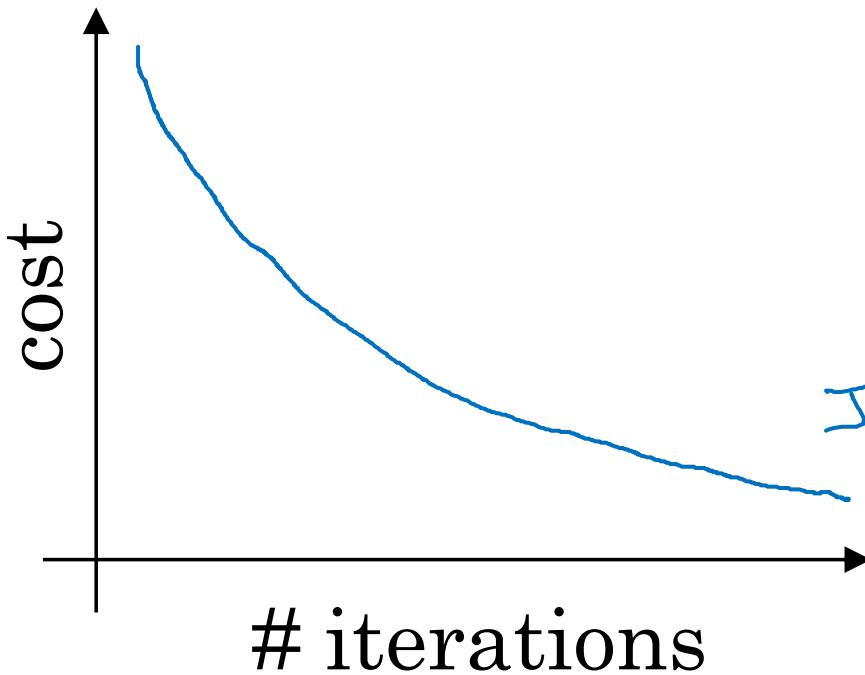
deeplearning.ai

Optimization Algorithms

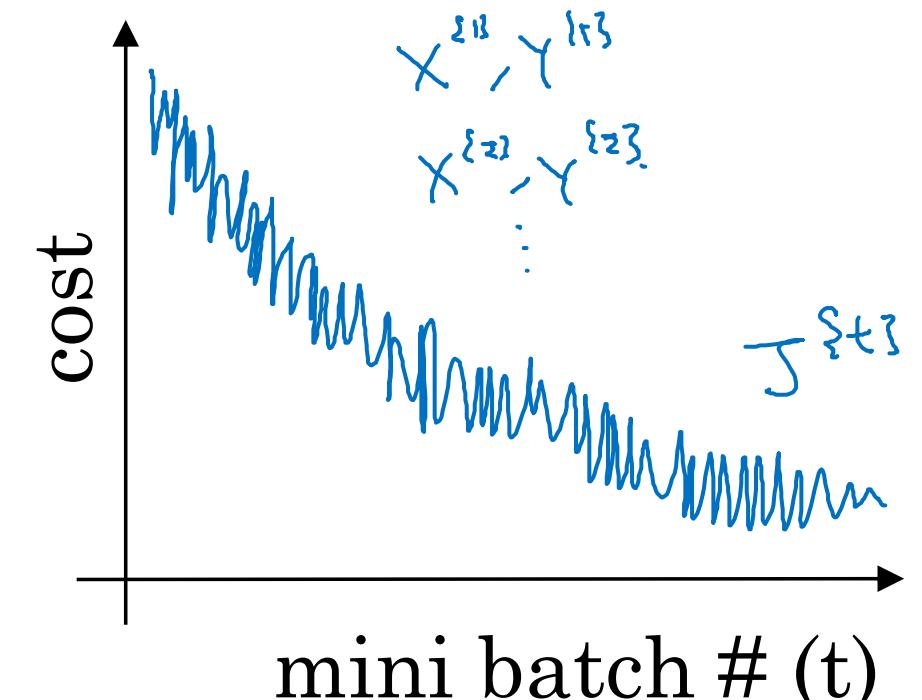
Understanding mini-batch gradient descent

Training with mini batch gradient descent

Batch gradient descent



Mini-batch gradient descent



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

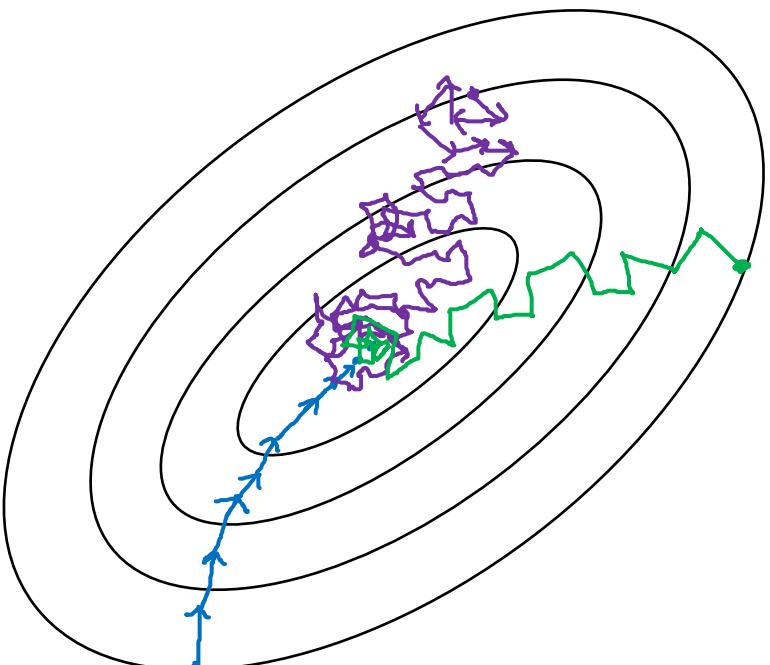
Choosing your mini-batch size

→ 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 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 sparse
feature vectors
many samples mean op. eff.
small batch, larger size vectors

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

Faster learning.

- Vectorization. ($n \times n$)
- Mini batches 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

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

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

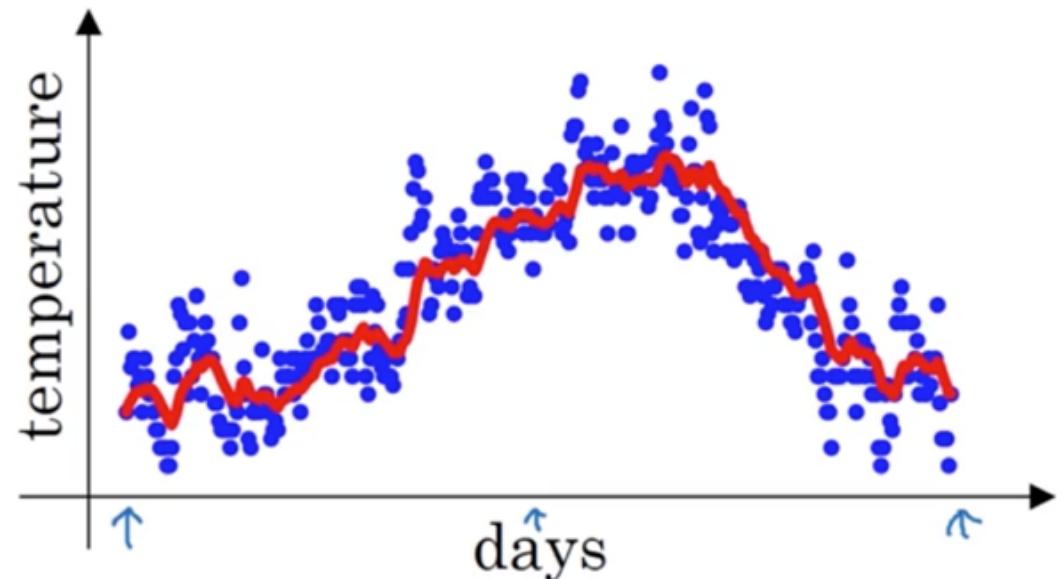
$$\theta_3 = 45^{\circ}\text{F} \quad \vdots$$

\vdots

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

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

\vdots



$$V_0 = 0$$

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

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

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

\vdots

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

Exponentially weighted averages

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

$$\beta = 0.9$$

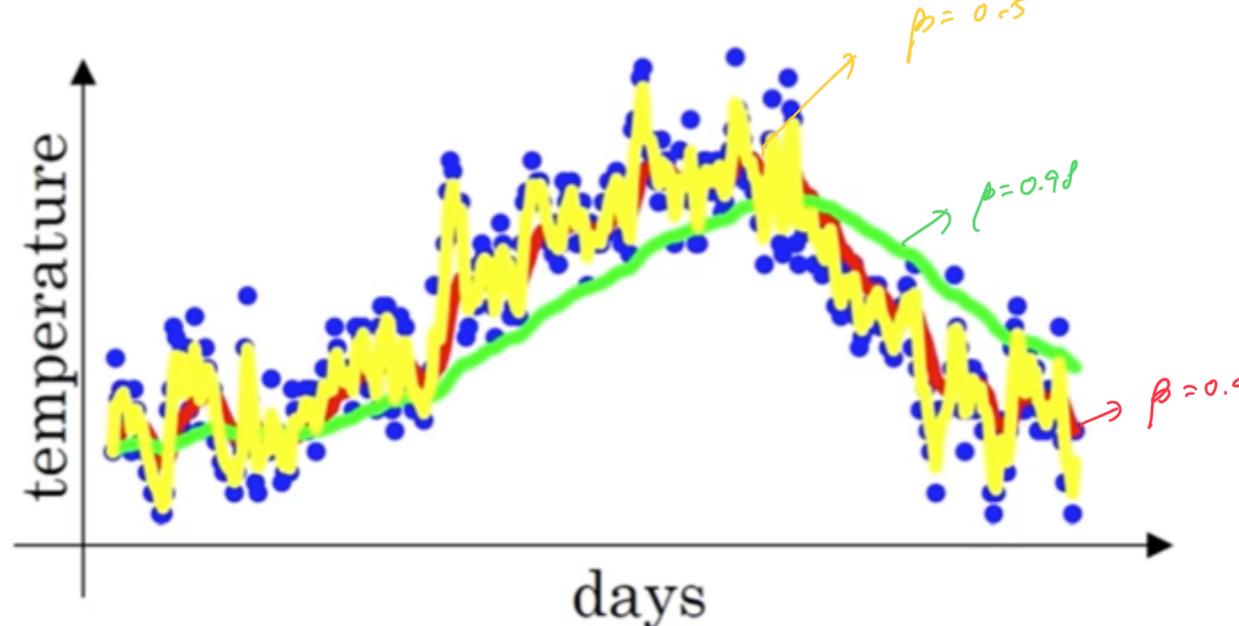
: ≈ 10 days' temper. (= 0.98¹⁰ ≈ 0.37)

$$\beta = 0.98$$

: ≈ 50 days

$$\beta = 0.5$$

: ≈ 2 days



V_t is approximately

average over

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

temperature.

($= V_{t-2}$ until $\frac{1}{1-\beta}$ days ≈ 100% of the data.)

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



deeplearning.ai

Optimization Algorithms

Understanding exponentially weighted averages

Exponentially weighted averages

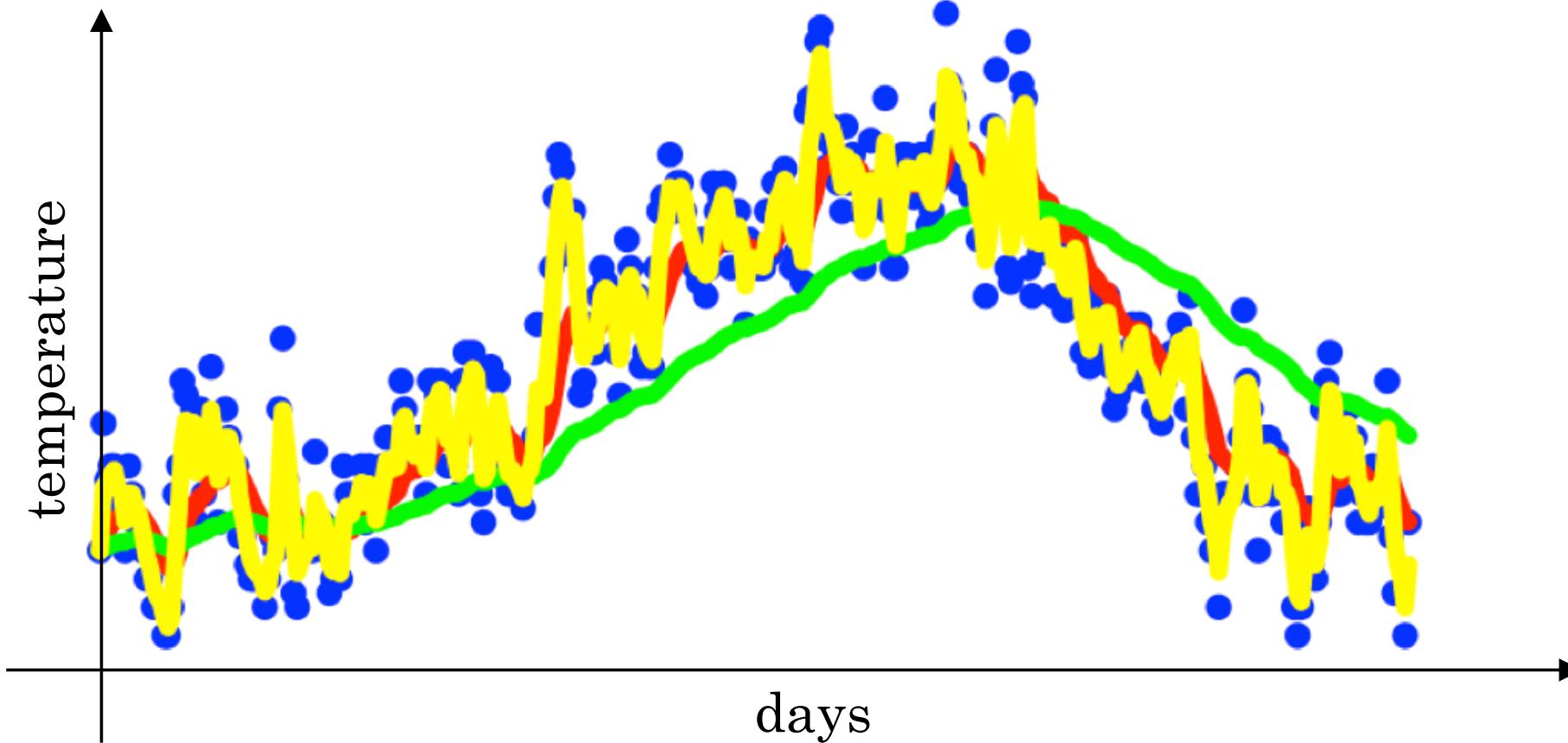
$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

$$\beta = 0.9$$

$$0.98$$

$$0.5$$

기수 가중 평균: 시그널의 흐름을 따라
slowly, 아이디어를
여기서 더 넓힐 때
강화학습
적용된다.



Exponentially weighted averages

$$v_t = \underbrace{\beta v_{t-1}}_{\text{과거의 평균}} + \underbrace{(1 - \beta) \theta_t}_{\text{새로운 정보}}$$

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

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

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

...

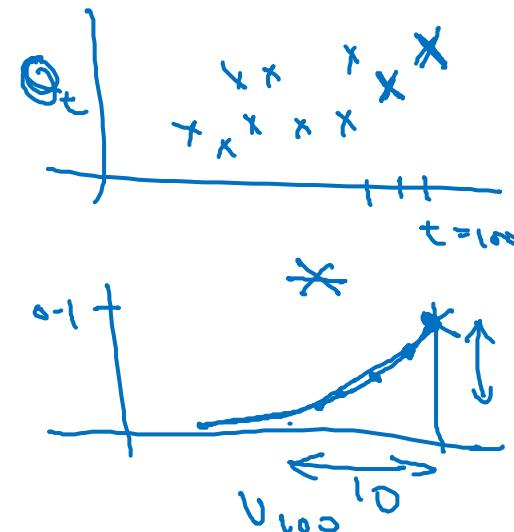
$$\begin{aligned} \overline{v_{100}} &= 0.1\theta_{100} + 0.9 \cancel{(0.1\theta_{99})} + 0.9 \cancel{(0.1\theta_{98})} + 0.9 \cancel{(0.1\theta_{97})} + 0.9 \cancel{(0.1\theta_{96})} \\ &= 0.1\theta_{100} + \dots + 0.1(0.9)^2\theta_{98} + 0.1(0.9)^3\theta_{97} + 0.1(0.9)^4\theta_{96} \end{aligned}$$

각각마다 3분의 2를 빼는 거
 약 100번은 0.02이고, 0.9는 0.9인
 경계에 접근하면서 차이가 줄어들고
 차이가 0.02를 구하는 것 같음.

$\frac{0.9}{0.9} \approx 0.35 \approx \frac{1}{e}$

$\frac{(1-\varepsilon)}{0.9}^{\frac{1}{\varepsilon}} = \frac{1}{e}$
 $\varepsilon = 0.02 \rightarrow \frac{0.98}{0.98-0.02} \approx \frac{1}{e}$

$\frac{1}{1-\beta}$
 $\Sigma = 1 - \beta$
 $\beta = 0.98 \Rightarrow 0.98/50 \approx 1/e$
 이를 50번 더하면 평균에 접근함.



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

$$\Sigma = 1 - \beta$$

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

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

$\varepsilon = 0.02 \rightarrow \frac{0.98}{0.98-0.02} \approx \frac{1}{e}$

Implementing exponentially weighted averages

$$v_0 = 0$$

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

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

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

...

장점: 기울기로 단한줄의 코드 단계에 같은 미분에 적용되는 모든
계수를 omega에 대처 미분을 사용하고, 훈련이 높다.
But, 웨이트는 지속적으로 변화하는 경우, 계산은 힘들다.
하지만 100일짜리 50일의 기울기를 더해 100와 50으로 나누면
더 좋은 계산치를 얻을 수 있지만, 이 방법은 그동안
기울기 변화하는 동안에 더 많은 미분을 필요로 한다.
따라서 그 흐름과 연관된 계수를 적용하는 것이 좋다.

$$V_{\theta} := 0$$

$$V_{\theta} := \beta V + (1-\beta) \theta,$$

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

:

$$\rightarrow V_{\theta} = 0$$

Repeat {

Get next θ_t

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

}

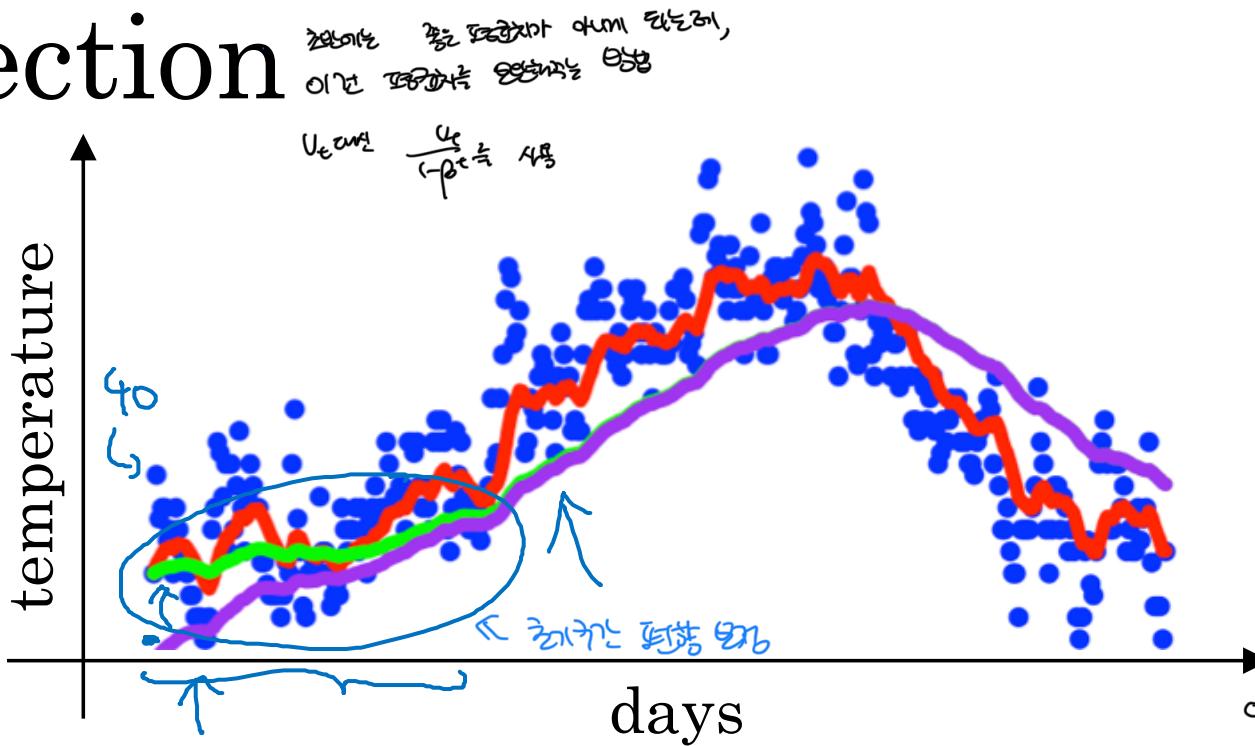


deeplearning.ai

Optimization Algorithms

Bias correction
in exponentially
weighted average

Bias correction



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

$$v_0 = 0$$

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

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

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

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

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

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

Actual v_t vs EMA v_t $1 - \beta^{t-1}$ ≈ 0.02 ≈ 0.0396
Original v_t vs EMA v_t $\beta = 0.98$ ≈ 0.02 ≈ 0.0396
bias correction of EMA v_t ≈ 0.02



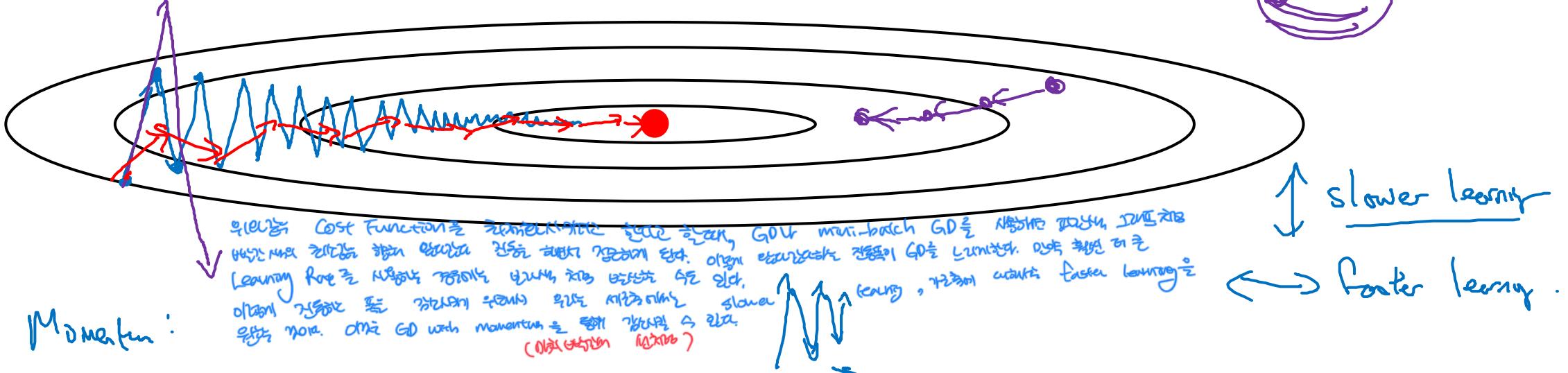
deeplearning.ai

Optimization Algorithms

Gradient descent with momentum

: 온라인 알고리즘 Gradient descent 가 네 터닝 파라미터가 등장한다.
온라인 알고리즘 Gradient descent는 Exponentially Moving Average 를 쓴다.
그것은 온라인 편평한 weight 를 업데이트하는 것이다.

Gradient descent example



Momentum:

On iteration t :

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

$$V_{dw} = \beta V_{dw} + (1-\beta) \frac{\Delta w}{\text{batch size}}$$

$$V_{db} = \beta V_{db} + (1-\beta) \frac{\Delta b}{\text{batch size}}$$

↑ velocity ↑ acceleration

$$w := w - \alpha V_{dw}, \quad b := b - \alpha V_{db}$$

($\Delta w, \Delta b$ 는 가속도 예측을 하고
 V_{dw}, V_{db} 는 속도 예측을 한다.)

" $V_{dw} = \beta V_{dw} + (1-\beta) \theta_t$ "
(θ_t 는 미니 배치)
→ 한 번의 미니 배치에서 $\Delta w, \Delta b$ 의 평균을 계산하고, 이전에 계산한 평균을 업데이트하는
과정에서 w 와 b 를 업데이트하는 과정이다. 이 과정은 Gradient Descent를 더욱 매끄럽게 만든다.
경우에 따라 미니 배치 평균이 0이 되거나, 미니 배치 평균이
반드시 0이 되거나, 아니면 그보다 높거나 낮거나 그 이상으로 추정되는 경우

Andrew Ng

Implementation details

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

On iteration t :

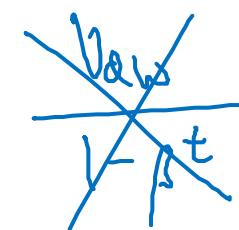
Compute dW, db on the current mini-batch

$$\begin{aligned}\rightarrow v_{dw} &= \beta v_{dw} + (1 - \beta) \underline{dW} \\ \rightarrow v_{db} &= \beta v_{db} + (1 - \beta) \underline{db}\end{aligned}$$

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

여기 $\beta \approx 0.9$ default이다. 초기 값은 대체 Gradient Descent의 학습률을 의미한다.
그리고 bias correction은 학습률을 줄이는 것과 같은 역할을 한다. 단지 100번의
iteration 후에는 학습률이 감소해 bias correction을 적용해 학습률을 초기화하는
방법이다.
마지막 학습률은 $(1 - \beta)$ 가 적용된 값이 됨으로써 $v_{dw} = \beta v_{dw} + dW$ 는
유니언이다.

$$\overbrace{v_{dw}}^{\text{Bias}} = \beta \overbrace{v_{dw}}^{\text{Initial}} + \overbrace{dW}^{\leftarrow}$$



Hyperparameters: α, β

$$\beta = 0.9$$

average over last ≈ 10 gradients



deeplearning.ai

Optimization Algorithms

RMSprop

momentum 이전에 GD를 진행하는 경우에 경사계단을 피할 수 있다.
이전에 진행한 경사계단은, 더 높은 learning Rate α 를 사용해도,
동일한 경사계단을 피하기 어렵다.

RMSprop

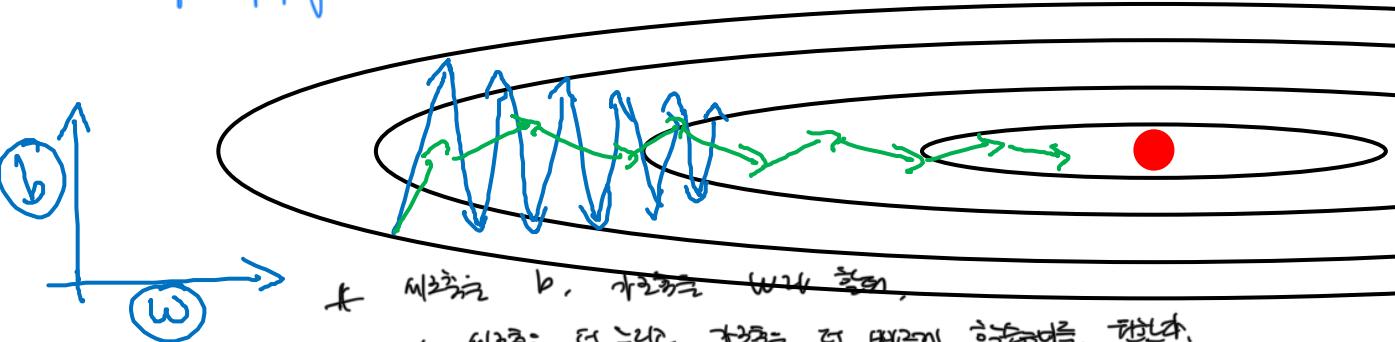
Root Mean Square propagation

이 알고리즘
GD 보다 경사계단 S 안정.

$$w_1, w_2, w_3 \uparrow$$

$$w_3, w_4, \dots \rightarrow$$

slow
↔ fast



On iteration t :

Compute dW, db on current mini-batch

$$\begin{aligned} S_{dw} &= \beta_2 S_{dw} + (1-\beta_2) \underline{dW^2} && \text{element-wise} \\ \rightarrow S_{db} &= \beta_2 S_{db} + (1-\beta_2) \underline{db^2} && \text{small} \end{aligned}$$

EMAX 조정 필요로 함.

large

학습률이 더 빨라 학습률이 적어져서 경지되는 문제를 해결하기 위해 학습률 크기를 비율로 조절할 수 있도록 제안된 방법이다.

학습률은 학습률 크기를 조정하여 학습률을 조정하는 방법을 사용한다.

기본적인 방법은 세션별 학습률을 고정하는 방법이다. 다른 S_{dw} 는 변화하는 것은 아니지만, S_{db} 는 세션별로 값이 달라야 한다. 실제 미분값을 보면, 기본값보다 세션별로 더 큰 경우 불필요하고, 기울기는 0으로 더 크게 된다. 이로 인해 기본방향의 변화폭은 무려하고, 세션별의 변화폭은 차이로 인해 된다. 그래서 그대로 한 번 모아두고 학습률을 조정된다.

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

ϵ 는 학습률 조정을 위한 편차

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

ϵ 는 학습률 조정을 위한 편차

$\sum = 10^{-8}$



deeplearning.ai

Optimization Algorithms

Adam optimization algorithm

Adam optimization algorithm

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

: Momentum or RMSprop 이 초기화된 형태
+ 초기화 Bias Correction은 초기 상태
가정하던 속도에 영향도 주고, 초기 가중치의 곡선적 변화를 고려해 학습률을 알리게
이를 통해 학습률 초기값을 막는 시도

On iteration t :

Compute dw, db using current mini-batch

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

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

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

bias correction

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

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

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

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

Hyperparameters choice:

- α : needs to be tune : 매우 중요한 베타일 필요가 있어서
다양한 값을 시도하여 잘 맞는 값을
찾아내야 한다.
- β_1 : 0.9 → (\underline{dw}) : 모멘트에 만족하는
- β_2 : 0.999 → ($\underline{dw^2}$)
- Σ : 10^{-8} : 불안정성이 있는 경우 초기 편차를 줄이는 데.
반면에, 이 값을 너무 작게 두면 학습 속도가 떨어진다.

Adam : Adaptive moment estimation

- * Adam 알고리즘이 학습률을 자동으로 조절하는
 β_1, β_2 및 초기 편차를 갖는다.
- * 가능한 한 많은 데이터를 사용하여 어떤 값이 잘 학습되는지 확인해야 한다.



Adam Coates



deeplearning.ai

Optimization Algorithms

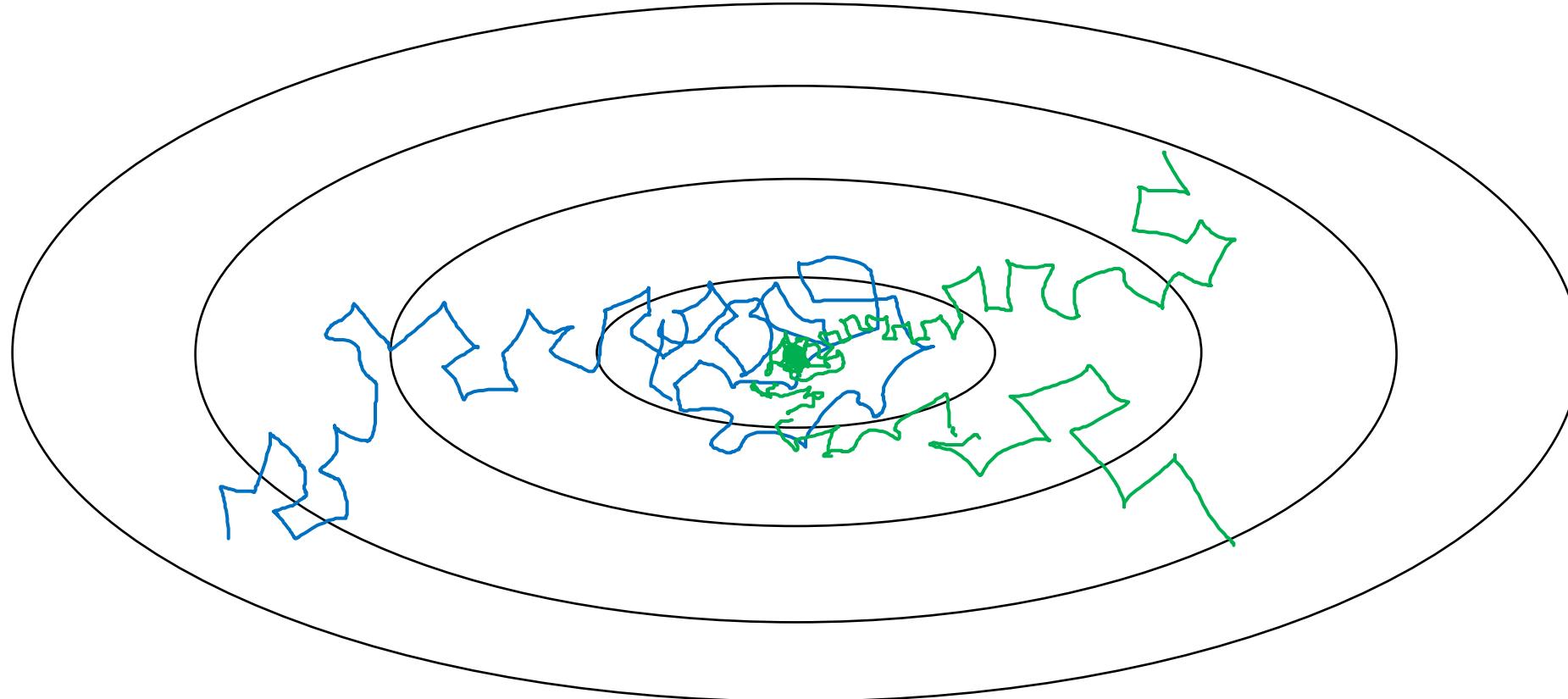
Learning rate decay

Learning rate decay

: 가중치 Learning Rate가 높을 때 loss값은 빠르게 줄을 수는 있지만, 학습의 학습률이 높아지기 때문에, 빠르게 학습하는 학습률을 찾는 데에는 그 단점이 있다.

따라서, 적당한 Learning Rate를 초기 값 후 각 epoch마다 절반
가장 마지막 학습률을 더해나가면 학습률을 찾을 수 있다. 이를 Slowly Reduce 라는 방식.

Slowly Reduce L



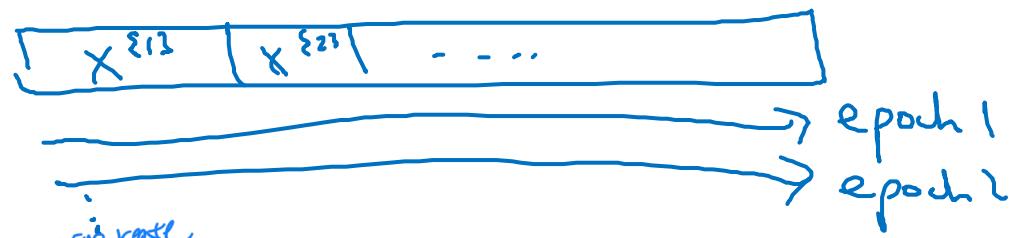
Learning rate decay

1 epoch = 1 pass through data.

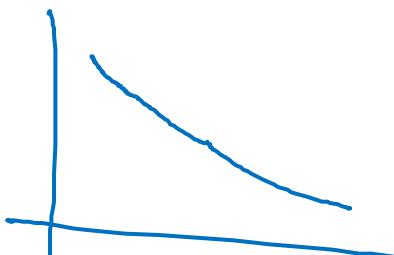
$$\alpha = \frac{\alpha_0}{1 + \text{decay-rate} * \text{epoch-number}}$$

↑ learning rate
↓ mini-batch set size
↑ batch size

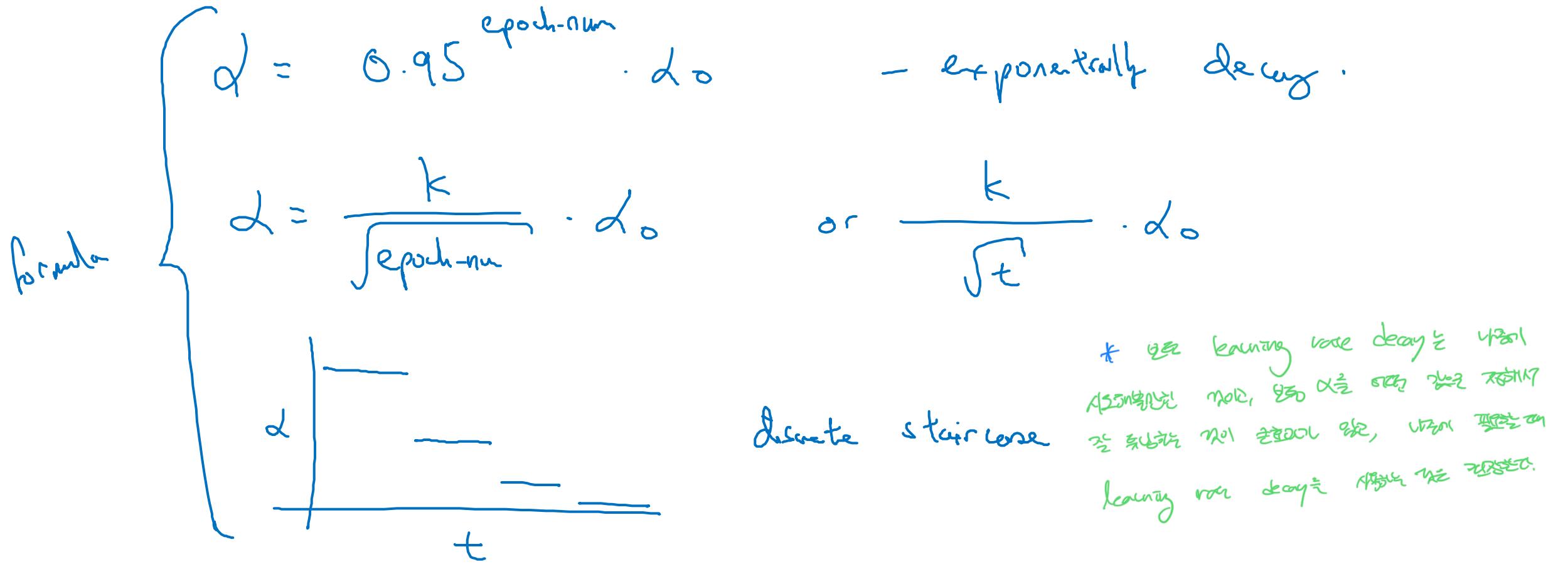
Epoch	α
1	0.1
2	0.067
3	0.05
4	0.04
:	:



$$\alpha_0 = 0.2$$
$$\text{decay-rate} = 1$$



Other learning rate decay methods



Manual decay (수동적 learning rate decay
- 학습률 수동적 조절 사용 가능)

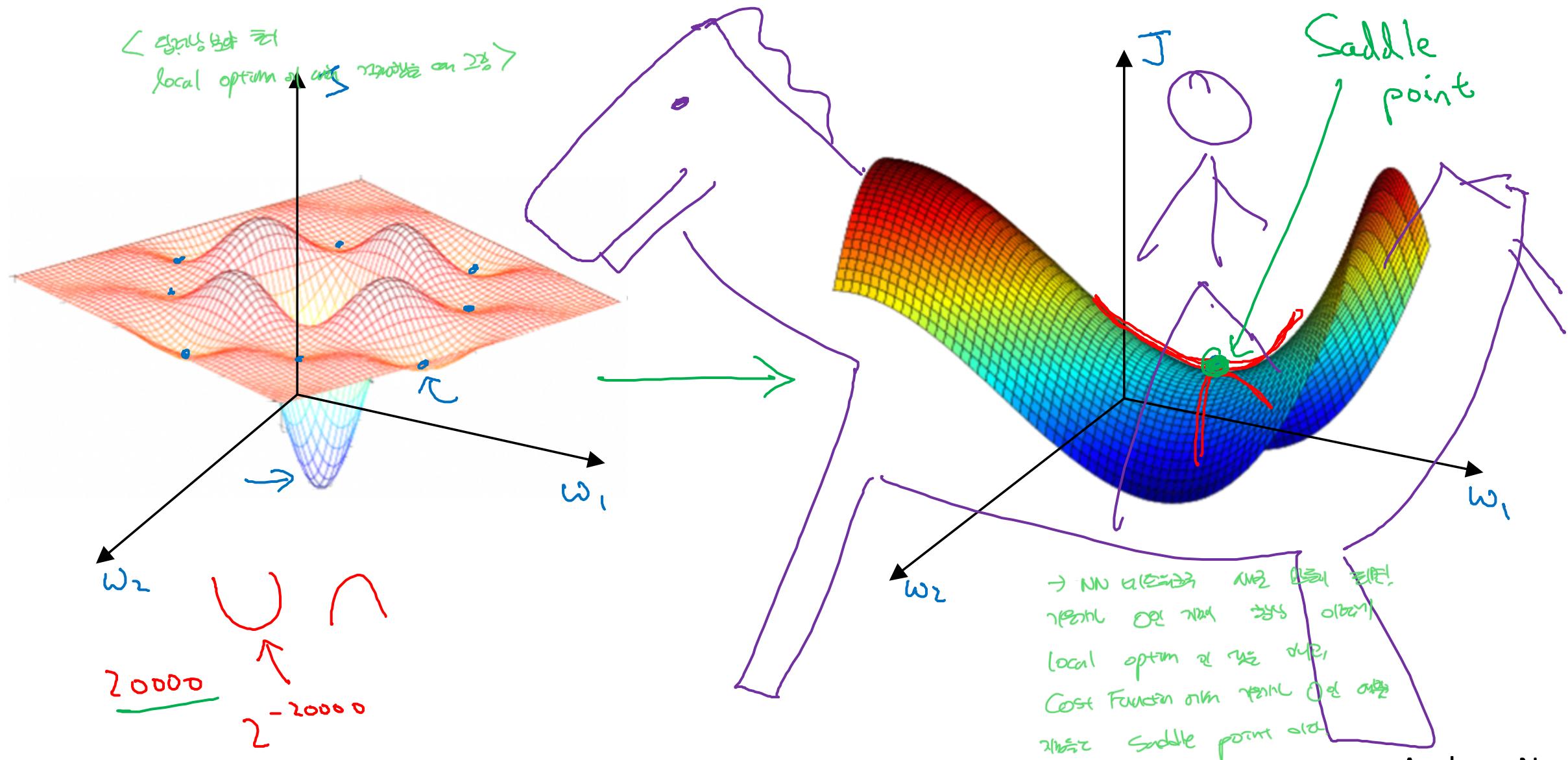


deeplearning.ai

Optimization Algorithms

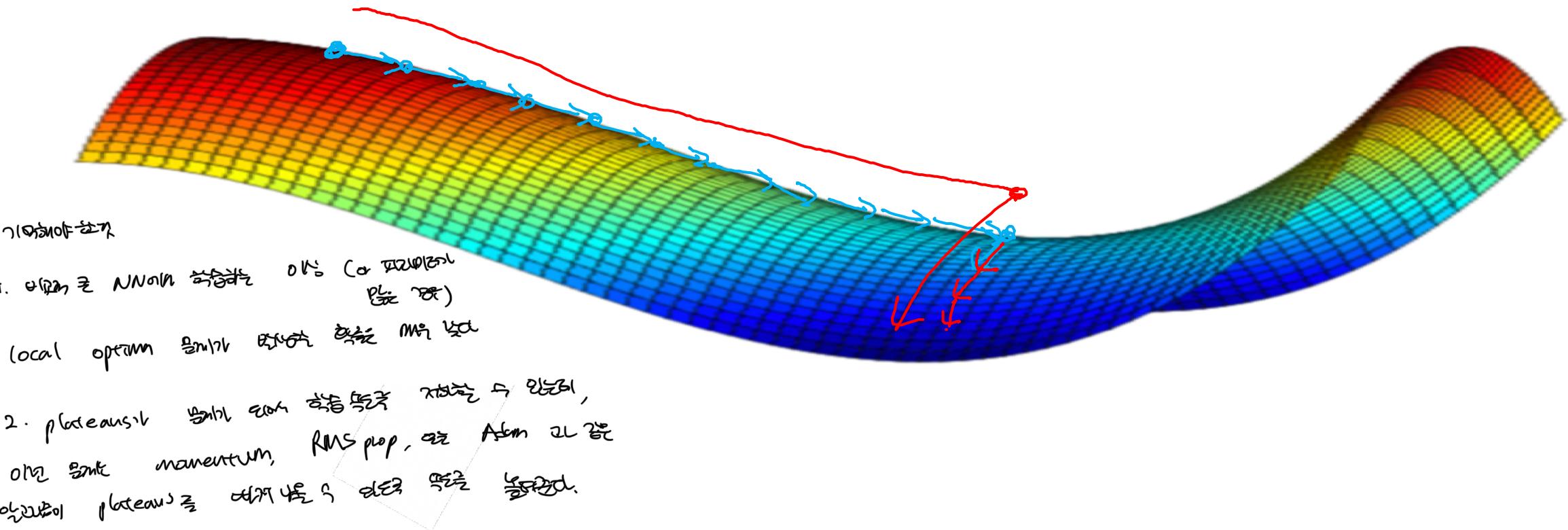
The problem of local optima

Local optima in neural networks



Andrew Ng

Problem of plateaus



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