

WEEK 2. Optimization Algorithms

- Mini-batch gradient descent

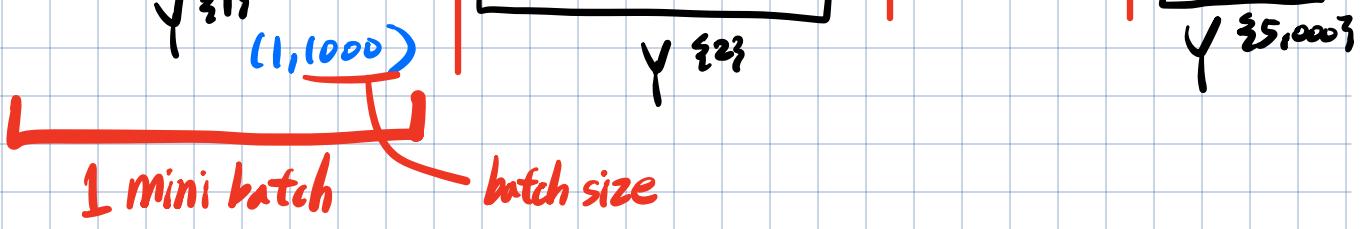
- Deep learning은 highly empirical process, 빅데이터에서 잘 작동
→ 좋은 초기화 알고리즘으로 Speed-up이 필요
 - Batch vs mini-batch gradient descent
→ Vectorization은 M example에 대한 효율적인 계산을 가능하게 해주지만,
매우 큰 M의 대해서는 느려질 수 있음

Batch gradient descent	Mini-batch gradient descent
<p>훈련 샘플의 모든 batch로 훈련</p>	<p>전체 훈련 세트 (X, Y)를 한 번에 진행하기 않고 하나의 mini-batch ($X^{(t)}, Y^{(t)}$)를 동시에 진행</p>
<p>모든 훈련 샘플을 gradient descent 를 할 processing 해야 함</p>	<p>모든 training sample을 processing 하기 전에 gradient descent를 진행시킬 수 있음</p>
<p>1 epoch \rightarrow one gradient descent step</p>	<p>1 epoch \rightarrow 5,000 gradient descent (반복적인 학습)</p>

- What if $M = 5,000,000$?

$$X = \begin{bmatrix} X^{(1)} & X^{(2)} & \dots & X^{(1000)} \end{bmatrix}_{(n_x, m)}$$

$$Y = \begin{bmatrix} Y^{(1)} & Y^{(2)} & \dots & Y^{(1000)} \end{bmatrix}_{(1, m)}$$



* Notation

$a^{[l]} \xi_{t,i}^{[l]}$ \Rightarrow i^{th} sample's activation at l^{th} layer of t^{th} mini batch

- How does it Work?

\Rightarrow 1 step of gradient descent using $X^{[t,i]}, Y^{[t,i]}$ (as if $M=1,000$)

repeat ξ for $t=1, \dots, 5000$ ξ

Forward prop on $X^{[t,i]}$

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

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

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

Vectorized
implementation
(1,000 examples)

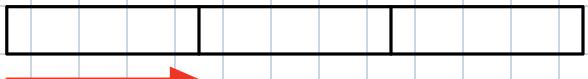
Compute Cost $J^{[t,i]} = \frac{1}{1000} \sum_{i=1}^L \ell(Y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_j \|W^{[j]}\|_F^2$

Backprop to compute gradients w.r.t $J^{[t,i]}$ (using $(X^{[t,i]}, Y^{[t,i]})$)

$$W^{[1]} := W^{[1]} - \alpha dW^{[1]}$$

$$b^{[1]} := b^{[1]} - \alpha db^{[1]}$$

\rightarrow One-epoch: a single pass through training set (full batch / min-batch)



\Rightarrow 1 epoch



\Rightarrow 1 epoch

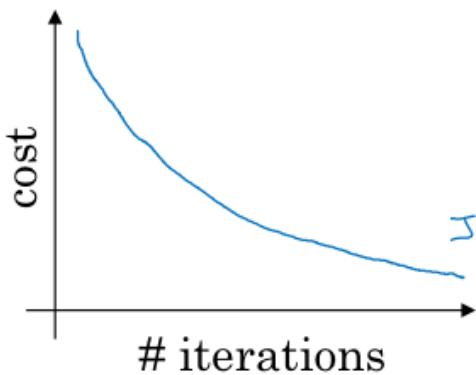


\Rightarrow 1 epoch

• Understanding min-batch gradient descent

— Training with mini batch gradient descent

Batch gradient descent

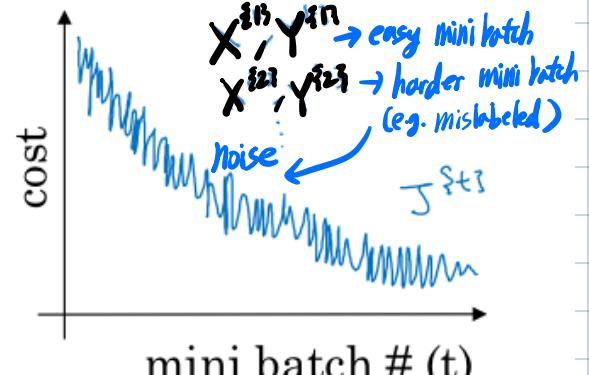


모든 반복에서 전체 샘플 사용 &
각각의 반복마다 비용 감소

⇒ 모든 반복마다 감소해야 함

한번이라도 올라가면 learning rate가
너무 크다는 암울지 같길 필요

Mini-batch gradient descent



모든 반복에서 서로 다른 mini-batch
이용한 흐름

⇒ 모든 반복마다 반드시 J 감소 + noise
(but, 낮아지는 경향)

— Choosing your mini-batch size

⇒ If mini-batch size is...

size = m

Batch Gradient Descent

각 반복마다 huge training
set을 process 해야 함
(X, Y)

$1 < \text{size} < m$

한 번에 batch-size 만큼의
training sample로 진행
($X^{S(t)}, Y^{S(t)}$)

size = 1

Stochastic Gradient descent

한 번에 하나의 training
sample로 진행하는 것
(X^{i1}, Y^{i1})

too long per iter

Fastest learning

- ① Vectorization (~ 1000)
- ② Make progress w.o.
processing entire training set

① Lose benefits of
Vectorization

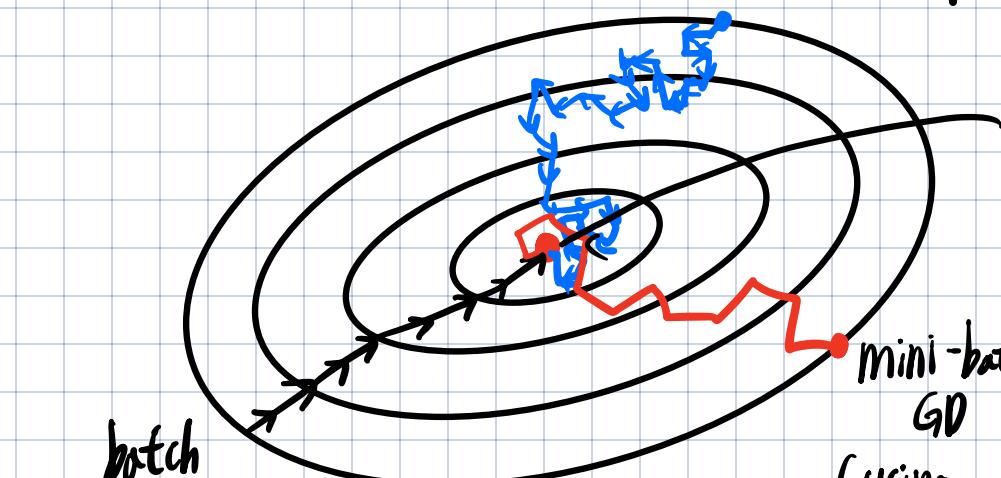
② Will never converge

③ 방향은 맞지만 최소화로

③ always converge / oscillate
within very small region

straight forward 가기
가지는 것 같아.

Stochastic GD (Using 1 training sample)



batch
Gradient descent
(using entire training samples
per iteration)

작은 차례에서
oscillate하는 것은
smaller learning rate
noise는 확률적 %로.

mini-batch
GD

(using $1 < n < m$ training samples
per iteration)

• Exponentially weighted (moving) averages

- Basic concept before learning optimization algorithms that are faster than gradient descent.
- Temperature in London

$$\theta_1 = 40^\circ\text{F}$$

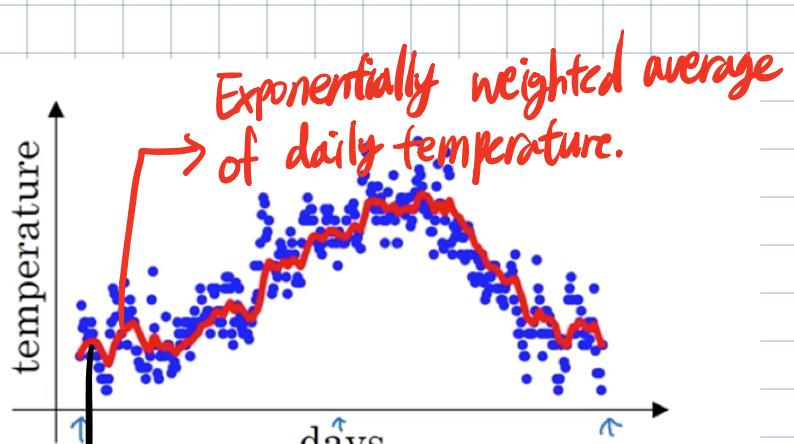
$$\theta_2 = 49^\circ\text{F}$$

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

\vdots

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

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



$$V_0 = 0$$

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

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

\vdots

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

- Formula of exponentially weighted averages

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

β 를 조정하면 약간씩 다른 효과를 얻는다.
parameter

$\beta = 0.9$: ≈ 10 days temp

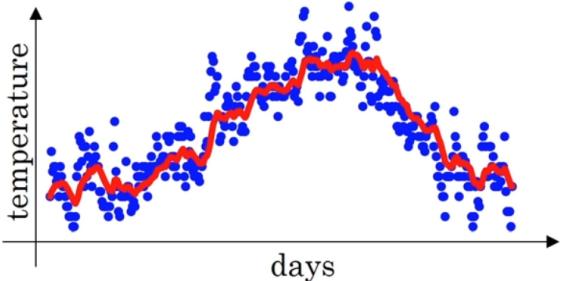
$\beta = 0.98$: ≈ 50 days temp

$\beta = 0.5$: ≈ 2 days temp

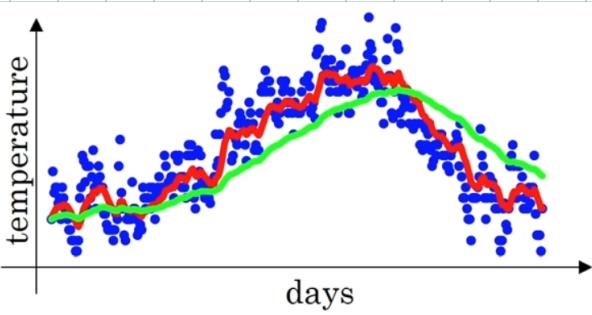
V_t 는 약 $\approx \frac{1}{1-\beta}$ days의

temp를 가지고 평균을 냅

β 가 1에 가까울수록 예전값에 가중치를 한다. : 현재에 변화가 생기면 trend를 따라감

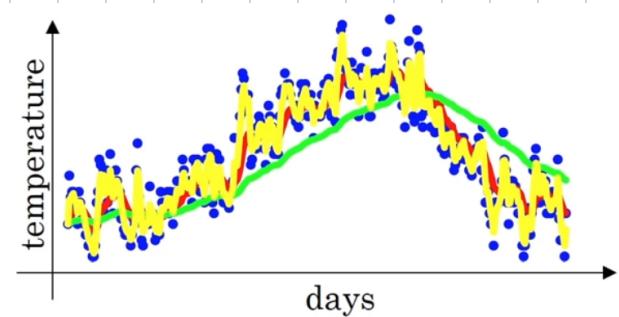


When $\beta = 0.9$ (10 days)



When $\beta = 0.98$ (50 days)

\Rightarrow smoother line but the curve is shifted to right (기온이 바뀌면 느리게 반응)



When $\beta = 0.5$ (2 days)

\Rightarrow Much more noisy & susceptible to the outliers (fast adaptation)

• Understanding exponentially weighted averages

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

$$\Rightarrow V_{100} = 0.9 V_{99} + 0.1 \theta_{100}$$

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

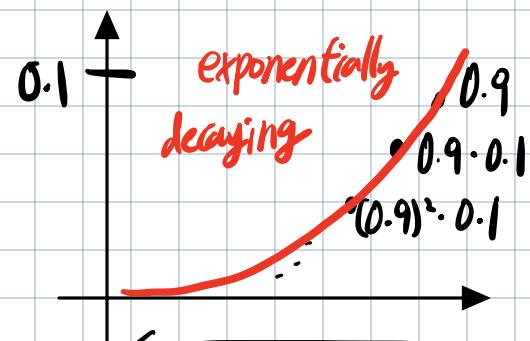
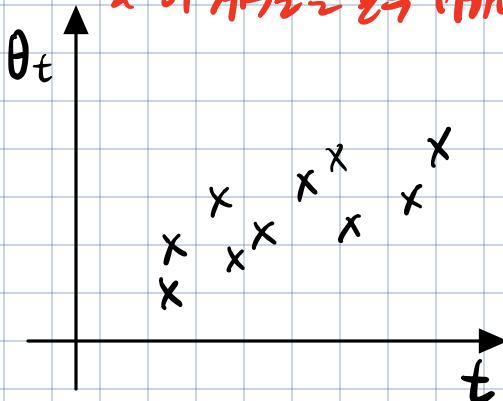
$$V_{98} = 0.9 V_{97} + 0.1 \theta_{98}$$

⋮
⋮

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

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

→ 일의 온도에 기하급수적으로 감소하는 합수를 공한걸 보듯 같아고 있다.
& 이 계수들을 모두 더하면 ≈ 1 (bias correction)



$$10 \text{ days } (\propto \frac{1}{1-\beta} = \frac{1}{\varepsilon})$$

→ 10일 뒤에는 기증자가
현재의 기증자의 $\frac{1}{3}$ 로 줄어든다.

— Implementing exponentially weighted averages

$$V_0 = 0$$

$$V_1 = \beta V_0 + (1-\beta) \theta_1$$

$$V_2 = \beta V_1 + (1-\beta) \theta_2$$

$$V_3 = \beta V_2 + (1-\beta) \theta_3$$

⋮

$$V_0 = 0$$

Repeat ⌂

Get next θ_t

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

(Overwrite the recent value)

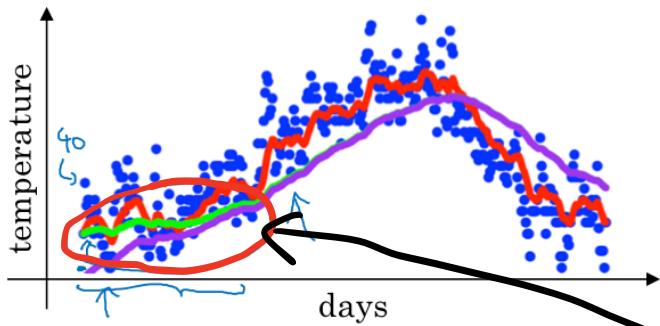
— Advantages of exponentially weighted average

① It takes very little memory & cheap computation

② One line code

• Bias Correction in exponentially weighted average

- Estimate of initial phase



green: theoretical
purple: actual implementation

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

$$V_0 = 0 \rightarrow 0$$

$$V_1 = 0.98 V_0 + 0.02 \theta_1$$

$$V_2 = 0.98 V_1 + 0.02 \theta_2 < \theta_1 \text{ or } \theta_2$$

\Rightarrow 초기값이 0이기 때문에 처음에는
온도를 잘 추정하지 못하기 때문에
초기에 값을 보정해주어야 한다.

- Bias correction

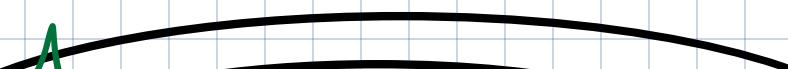
$$\frac{V_t}{1-\beta^t} \Rightarrow \frac{V_2}{1-(0.98)^2} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396}$$

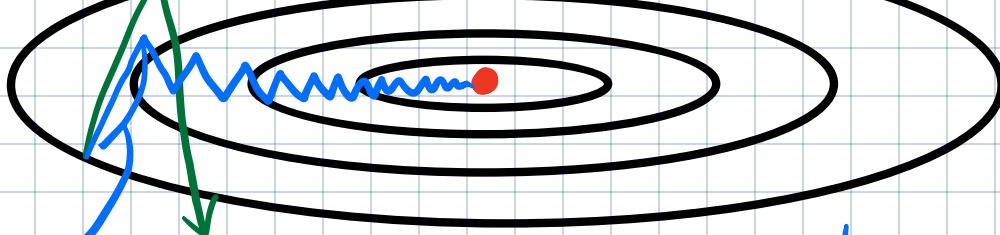
t가 커질수록 $1-\beta^t$ 은 0으로 가까워지므로 시간이 지남수록
보정 효과 \rightarrow (but, 시간이 어느정도 지나면 보정 없어도 추정을 잘한다.)

• Gradient descent with momentum

- min-batch을 이용하는 것은 training set의 subset인 뷔퍼 대로
update direction of variance & oscillation toward convergence
 \Rightarrow Compute an exponentially weighted average of gradients
for updating weights to smooth out direction.

- Gradient descent example





1 iter of
mini batch
of gradient descent

bigger learning rate
makes 'overshooting'

→ slow gradient descent
→ larger learning rate is
not allowed

⇒ We want vertically slow learning ↑
horizontally fast learning ←

- Gradient descent with momentum

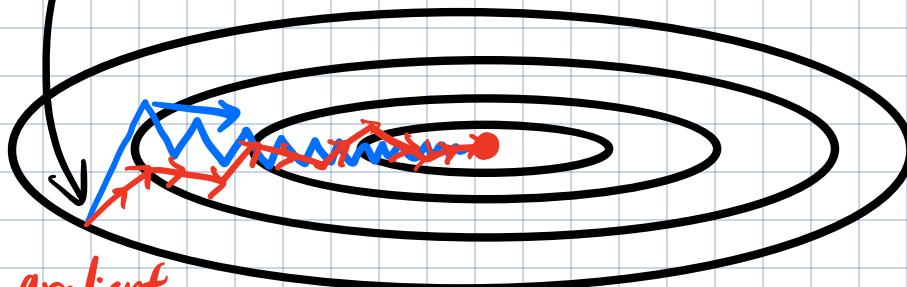
On iter t :

Compute dw, db on current mini-batch

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

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

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



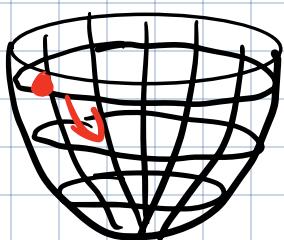
gradient
descent with momentum

average out gradients &
smooth out the steps of
gradient descent

⇒ 속직방향의 진동이 없어 가깝고
수평방향으로는 큰 핵으로
평균이 만족됨

- Bowl-shaped function

⇒ bowl-shaped function을 최적화한다고 했을 때



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

$$V_{db} = \boxed{②} \beta V_{db} + (1-\beta) db$$

- ① friction ($\beta < 1$)
 - ② Velocity
 - ③ Acceleration
- momentum

- Implementation details

$$V_{dw} = 0, V_{db} = 0 \quad (\text{same dimension with } dw)$$

On iteration t :

Compute dw, db on the current mini-batch.

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

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

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

⇒ Hyperparameters : α, β

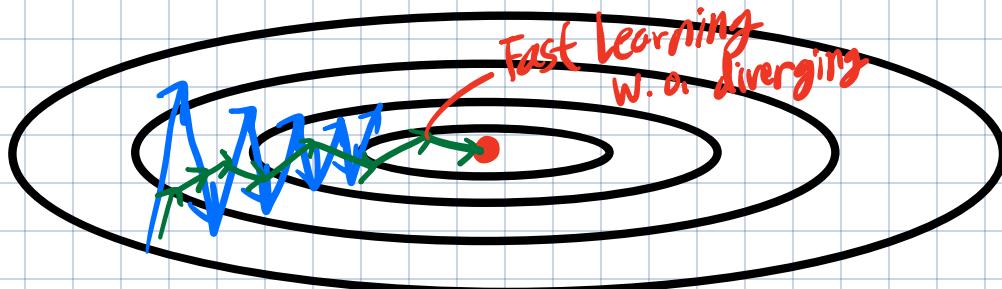
β 0.9 is common

⇒ 비율 10% 정도로 초기 이동 평균 계산(bias) 초기화 진행도 of bias X

∴ $\frac{V_{dw}}{1-\beta^t}$ (bias correction) 은 반드시 필요.

• RMSprop

- to speed up gradient descent



⇒ On iteration t :

Compute dw, db on current mini-batch

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dw^2 \Rightarrow \text{small}$$

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dw^2 \Rightarrow \text{large}$$

speed up $\leftarrow w := w - \alpha \frac{dw}{\sqrt{S_{dw}} + \epsilon}$ $b := b - \alpha \frac{db}{\sqrt{S_{db}} + \epsilon} \Rightarrow$ slow down
 (horizontal) small (vertical) large

$\because \epsilon \Rightarrow$ make sure to not divide by 0, commonly 10^{-8}

- Advantages

① Similar to momentum (진동 \rightarrow)

② Allow us to use larger learning rate \Rightarrow Speed up gradient descent

③ Work well in wide Neural Network.

• Adam optimization algorithm (Adapter moment estimation)

- Wide NN optimization 잘 적용하기 편리

- Momentum + RMS prop

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

On iter t:

momentum Compute dw, db using current mini-batch
 $\beta_1 \Rightarrow V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dw, V_{db} = \beta_1 V_{db} + (1-\beta_1) db$

RMSprop $\Rightarrow S_{dw} = \beta_2 V_{dw} + (1-\beta_2) dw^2, V_{db} = \beta_2 V_{db} + (1-\beta_2) db^2$

bias correction $\left[\begin{array}{l} V_{dw}^{\text{corrected}} = \frac{V_{dw}}{(1-\beta_1^t)}, V_{db}^{\text{corrected}} = \frac{V_{db}}{(1-\beta_1^t)} \\ S_{dw}^{\text{corrected}} = \frac{S_{dw}}{(1-\beta_2^t)}, S_{db}^{\text{corrected}} = \frac{S_{db}}{(1-\beta_2^t)} \end{array} \right]$

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

\Rightarrow Commonly used for Neural Network

- Hyperparameter choice

- α : needs to be tuned
- $\beta_1: 0.9 \rightarrow (dw) \Rightarrow$ momentum-like term
1st moment
- $\beta_2: 0.999 \rightarrow (dw^2)$ 2nd moment
- $\epsilon: 10^{-8}$

- Learning rate decay

- to slowly reduce learning rate

\Rightarrow If you're implementing mini-batch gradient descent with very small mini-batch (e.g. 64, 128...)



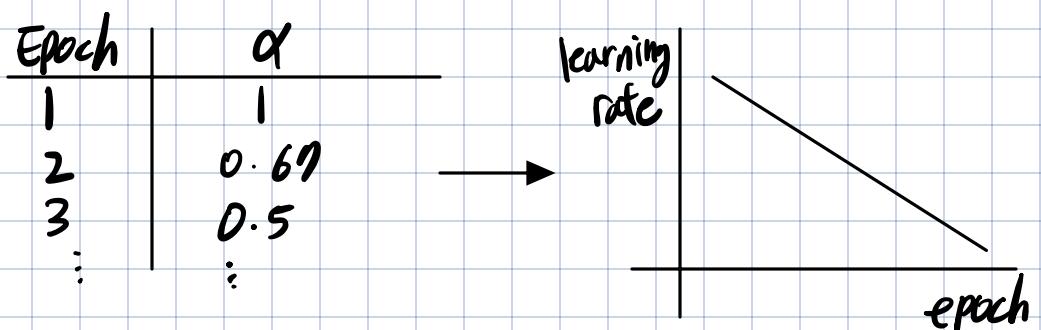
- Decide learning rate

\Rightarrow reduce learning rate as epoch increases

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

Need to be tuned

\Rightarrow if $\alpha_0 = 0.2$, decay-rate = 1,



- Other learning rate methods

① exponential decay

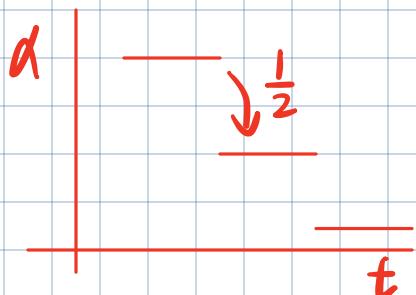
$$\alpha = 0.95^{\text{epoch-num.}} \cdot \alpha_0$$

↳ this should be < 1

②

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

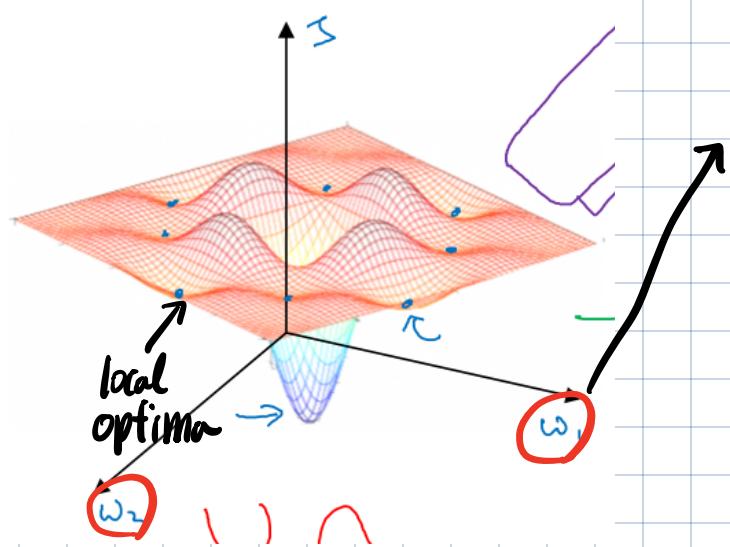
③ discrete stair case



④ Manual decay (when training # is small)

• The problem of local optima

- local optima in neural networks



설계되는 청진 고자원 (n 차원)

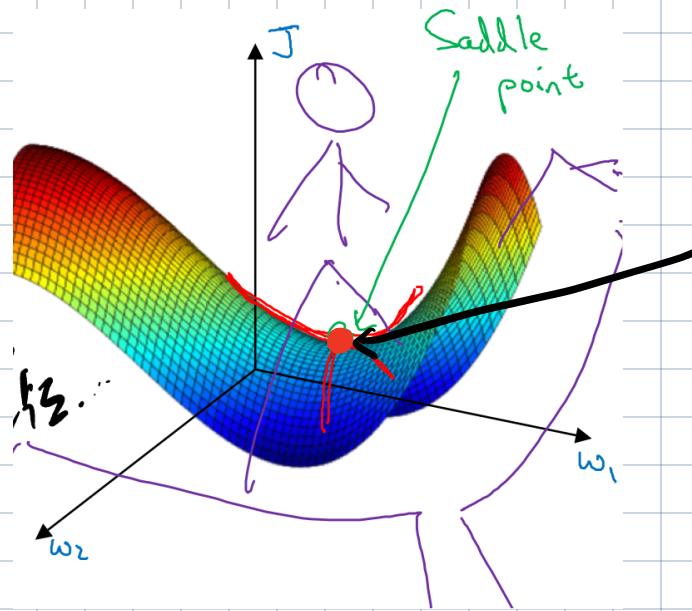
⇒ 모든 n 개의 방향에서 복죽/뾰족
여야하기 때문에 가능성이 적다.

e.g. 20000 차원에서 local optima
가 만들어지는 확률

$$= 2^{-20000}$$

- Saddle point

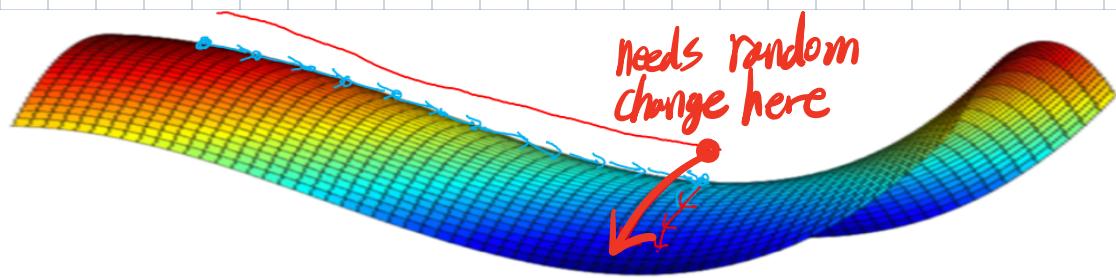
⇒ 심지어는 말 안장과 같은 꼴이야.



0 gradient

⇒ 고장점에서는 local optimality
saddle point가 될 수 있.

- Problem of plateaus: region where derivatives are close to zero for a long time



① Unlikely to get stuck in a bad local optima in high dimension

② Plateaus can make learning slow.

⇒ Momentum, RMSprop or Adam can speed up exiting plateaus