

Courses - 2  
Week - 2

# OPTIMIZATION ALGORITHMS

From prev. discussion

↓

Vectorization allows us to efficiently compute on 'm' examples.

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

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix} \quad 1 \times m$$

If  $m = 50,000,000$  examples

usually forward

✗ to be computed (if considering all m examples)  
and the compute loss  $\rightarrow$  gradient descent  
lot of time consuming

↓ sol  
mini-Batch Gradient  
descent

# mini-batch gradient descent

Let  $m = 50,000$

$$X = \left[ \begin{array}{c|c|c} x^{(1)} & x^{(2)} & \dots & x^{(1000)} \\ \hline & & & \end{array} \right], \left[ \begin{array}{c|c|c} x^{(1001)} & \dots & x^{(2000)} \\ \hline & & \end{array} \right], \dots, \left[ \begin{array}{c|c} x^{(50001)} & x^{(m)} \\ \hline & \end{array} \right]$$

$n \times m$   $\{1\}$   $\{2\}$   $\{5000\}$

$$Y = \left[ \begin{array}{c|c|c} y^{(1)} & y^{(2)} & \dots & y^{(1000)} \\ \hline & & & \end{array} \right], \left[ \begin{array}{c|c|c} y^{(1001)} & \dots & y^{(2000)} \\ \hline & & \end{array} \right], \dots, \left[ \begin{array}{c|c} y^{(50001)} & y^{(m)} \\ \hline & \end{array} \right]$$

$1 \times m$   $\{1\}$   $\{2\}$   $\{5000\}$

If  $m = 50,000$

→ we are not doing all at once as a single batch

→ divide the set of  $m$  examples into mini-batches (each batch of size 1000)

∴ # mini-batches = 5000

Notation:  $X^{\{t\}}$  → mini batch ( $t^{\text{th}}$  mini)

mini batch  $t$ :  $(X^{\{t\}}, Y^{\{t\}})$

dimension of  $X^{\{t\}} = n \times 1000$   
 $Y^{\{t\}} = 1 \times 1000$

# mini-batch GD Algorithm

1 step of GD  
using  $x^{(t)}$ ,  $y^{(t)}$

for  $t = 1, \dots, 5000$

forward pass on  $x^{(t)}$

$$z^{(1)} = w^{(1)} x^{(t)} + b^{(1)}$$

$$A^{(1)} = g^{(1)}(z^{(1)})$$

$$A^{(L)} = g^{(L)}(z^{(L)})$$

Vectorized  
Implementation  
1000 examples

Compute  $\text{cost} = J$  from minibatch  $x^{(t)}, y^{(t)}$

$$= \frac{1}{1000} \sum_{i=1}^{1000} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \times 1000} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

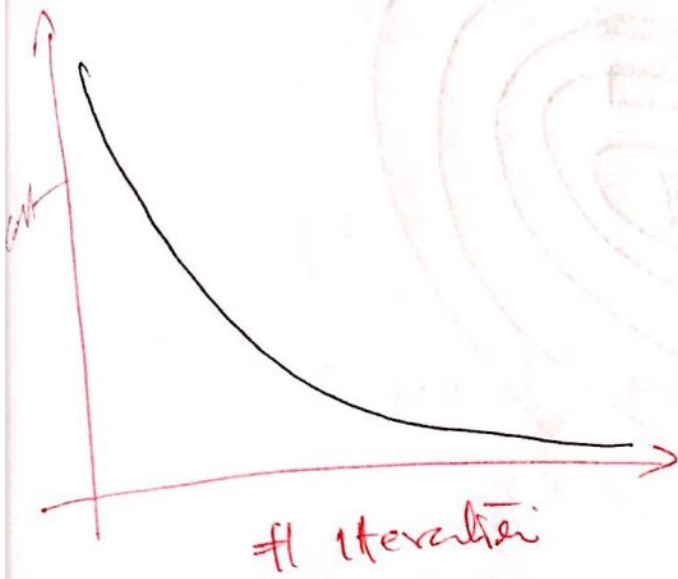
Backprop to compute gradients of cost  
 $J^{(t)}$  using  $(x^{(t)}, y^{(t)})$

$$w^{(l)} := w^{(l)} - \alpha dw^{(l)}; b^{(l)} := b^{(l)} - \alpha db^{(l)}$$

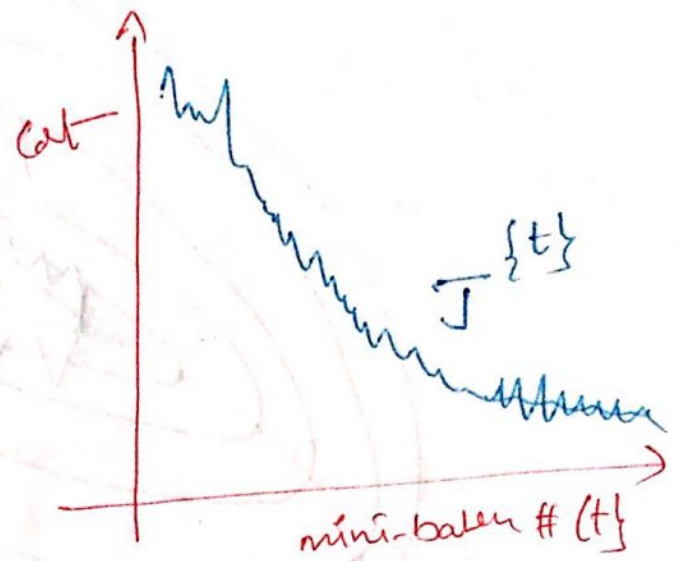
end of 1 epoch  
↓  
are pass through  
training set



Batch Gradient



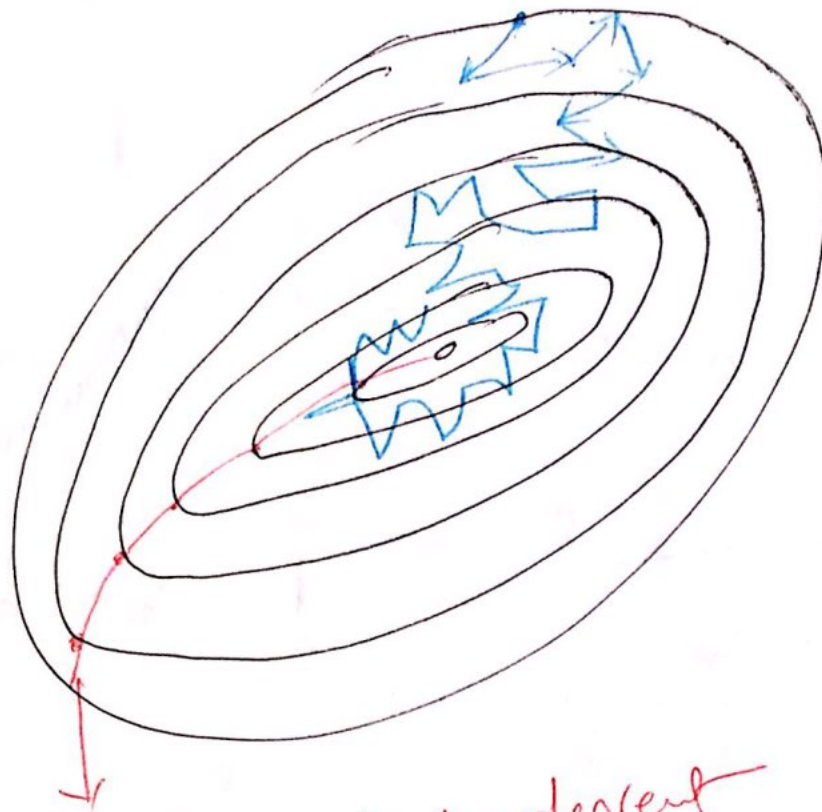
mini-batch



Q. ? minibatch size

① If (each) minibatch size =  $m$  (equal to # examples)  
 $\Rightarrow$  It is same as Batch Gradient

② If minibatch size = 1  $\Rightarrow$  Stochastic Gradient Descent  
 $\Rightarrow$  every single example become mini-batch  
 $\Rightarrow$  # mini batches =  $m$   
 $(x^{(i)}, y^{(i)}) = (x^{(1)}, y^{(1)})$



batch gradient descent

stochastic gradient descent

batch GD  $\rightarrow$  too long per iteration

stochastic GD  $\rightarrow$  very slow convergence

Summary

mini-batch size should not be

TOO LARGE & TOO SMALL

( $\neq m$ )

( $\neq 1$ )  
 $\downarrow$  faster learning

note

Set training set if small

① if  $m < 2000 \rightarrow$  Go for batch GD

② typical mini-batch size

$2^x$  (e.g. 64, 128, 256, 512, ...)   
 ( $2^8$ )

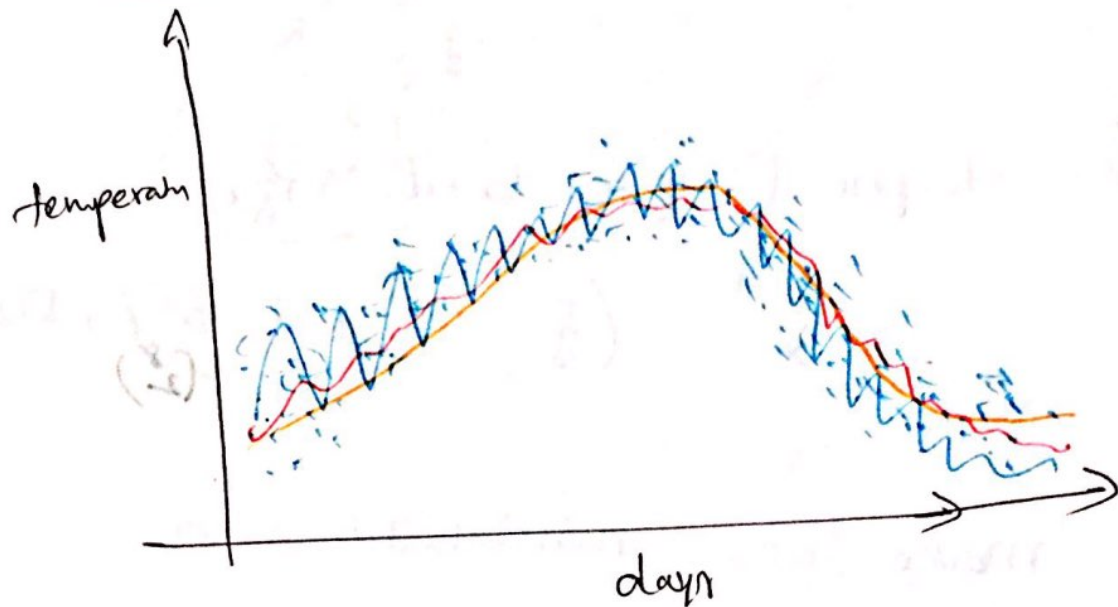
③ make sure mini-batch

$\begin{Bmatrix} x \\ y \end{Bmatrix}$  fits in CPU/RPO  
= (memory)



Algorithm faster than AD

↓  
Exponentially weighted moving averager



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

for  $\beta = 0.9 \Rightarrow V_t$  is approximately averaging  
over  $\frac{1}{1-\beta}$  days temperature

i.e. for  $\beta = 0.9 \Rightarrow$  averaging  $\approx 10$  days temp

for  $\beta = 0.95 \Rightarrow$  averaging  $\approx 30$  days temperature

$\beta = 0.5 \Rightarrow$  averaging  $\approx 2$  days  
temperature

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

exponentially  
weighted  
averaged  
formula

let us compute

$$\beta = 0.9$$

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

$$\theta_t = \text{today's temperature}$$

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

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

Notation

$$V_0 = 0$$

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

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

⋮

$$V_0 = 0$$

$$V_t = \beta V + (1-\beta) \theta_t$$

$$V_t =$$

$$\rightarrow V_0 = 0$$

Repeat {

get next  $\theta_t$ .

$$V_t = \beta V + (1-\beta) \theta_t$$

}



# Bias Correction in exponentially weighted average

$$\beta = 0.98$$

$$t = 2$$

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

$$= 0.0396$$

prob. of  
underestimation

$$\frac{w_t}{1 - \beta^t} \rightarrow \text{bias correction}$$

as  $t \rightarrow \text{large}$

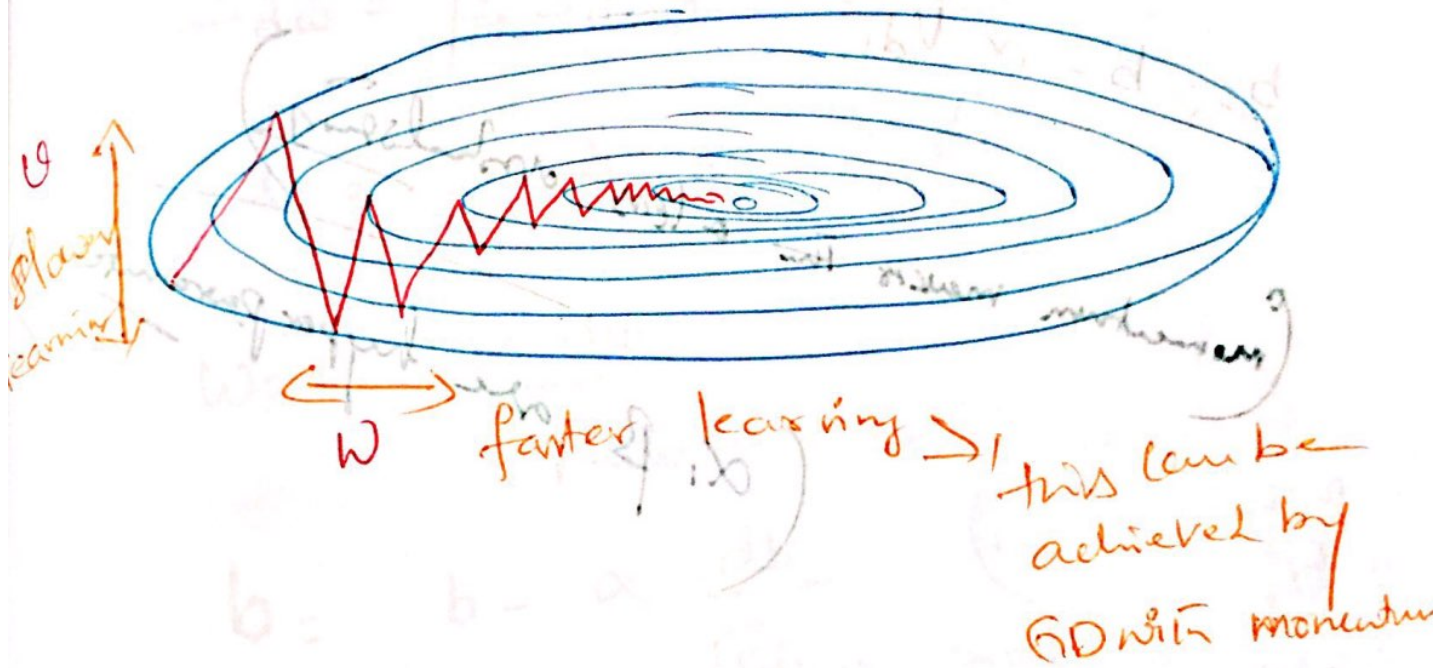
only  $w_t = \beta w_{t-1} + (1 - \beta) \theta_t$

need to be followed by

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

# gradient descent with momentum

key:- Compute exp. weighted average of your gradient and use then to update the new gradients



## Momentum

on each iteration  $t$ :

compute  $dw$ ,  $db$  on current mini-batch

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

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

$$w := w - \alpha \cdot V_{dw}$$

$$b := b - \alpha \cdot V_{db}$$

Momentum term (velocity)

Acceleration

$\beta = 0.9$  is high preference

earlier (usually)  
 $w \leftarrow w - \alpha dw$   
 $b \leftarrow b - \alpha db$

#  $\beta = 0$   
it is standard  
(grad. desc)

often in the literature

$$\begin{cases} v_{dw} = \beta v_{dw} + (1-\beta) d_w \\ v_{db} = \beta v_{db} + (1-\beta) d_b \end{cases} \rightarrow \begin{cases} v_{dw} = \beta^0 d_w + (1-\beta)^0 d_w \\ v_{db} = \beta^0 d_b + (1-\beta)^0 d_b \end{cases}$$

$$w_t = w - \alpha v_{dw}$$

$$b_t = b - \alpha v_{db}$$

(momentum makes the loss oscillations)  
( $\alpha, \beta$  are hyperparameters)

first moment (average)

$$w_b(\eta-1) + w_b \eta = w_b V$$

$$b_b(\eta-1) + b_b \eta = b_b V$$

$$w - \alpha v_{dw} = w$$

$$b - \alpha v_{db} = b$$



RMS prop

on iteration  $t$ :

Compute  $dw$ ,  $db$  on current mini-batch

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

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

$$w = w - \alpha \frac{dw}{\sqrt{S_{dw}}}$$

$\sqrt{S_{dw}} \leftarrow \text{small} + \epsilon$

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

$\sqrt{S_{db}} \leftarrow \text{large} + \epsilon$

$\epsilon$  is added to avoid '0' in denominator

With this algorithm



slow in  $(b)$  direction  
faster in  $(w)$  direction

# Adam optimization algorithm

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

on each iteration  $t$ :

Compute  $dw$ ,  $db$  using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw.$$

$$V_{db} = \beta_1 V_{db} + (1 - \beta_1) db.$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2.$$

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

$$V_{dw}^{\text{corrected}} = \frac{V_{dw}}{(1 - \beta_1^t)}$$

$$V_{db}^{\text{corrected}} = \frac{V_{db}}{(1 - \beta_1^t)}$$

$$S_{dw}^{\text{correct}} = \frac{S_{dw}}{(1 - \beta_2^t)}$$

$$S_{db}^{\text{correct}} = \frac{S_{db}}{(1 - \beta_2^t)}$$

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

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

## hyperparameter choice

$\alpha$ : need to be tuned

$$\beta_1 = 0.9$$

$$\beta_2 = 0.999$$

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

Adam: Adaptive moment estimation

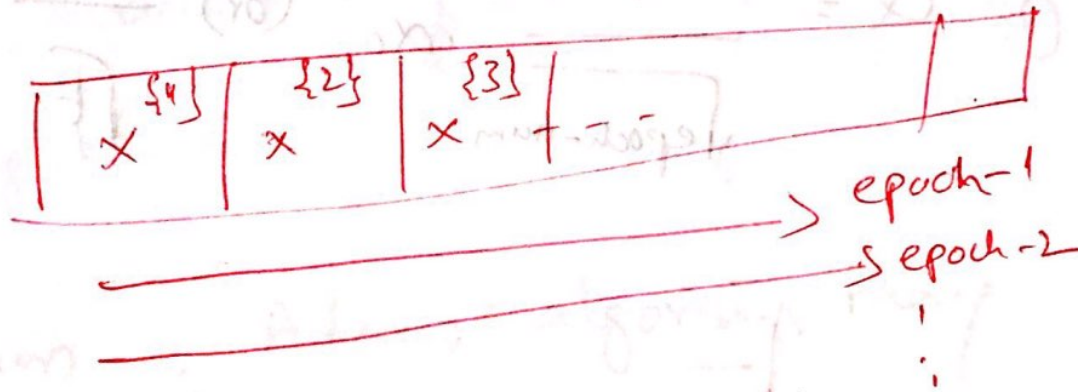


# Learning Rate decay

Idea:- Reduce the learning rate ~~forward~~  
(Convergence) after initial phase

larger learning — initial  
↓  
reduce the learning rate as proceed

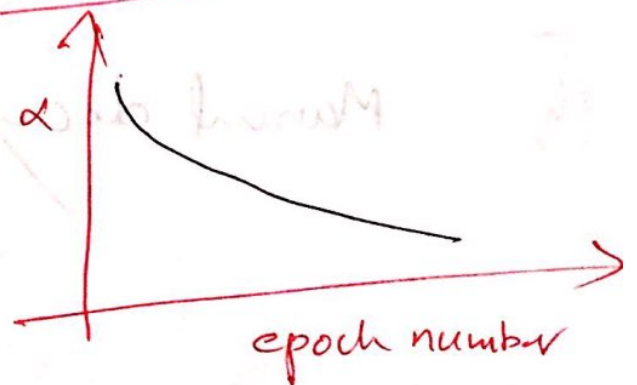
Intuition



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

Ex  $\alpha_0 = 0.2$ , decay-rate = 1

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

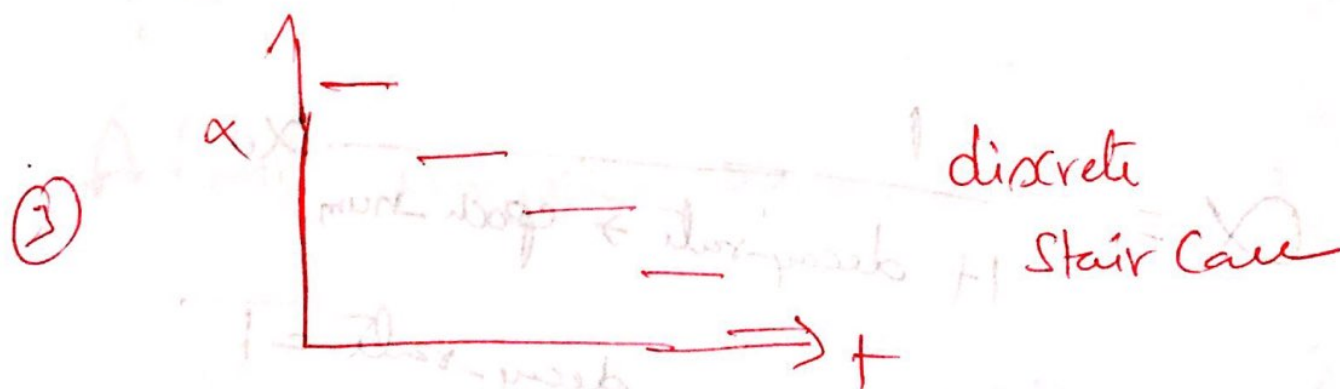


## Other formulas for decay

### ① exponential decay

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

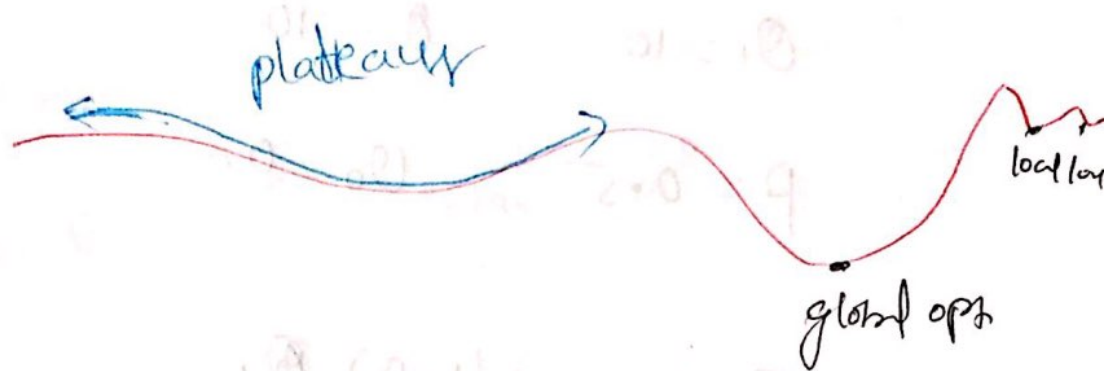
$$\textcircled{2} \alpha = \frac{k}{\sqrt{\text{epoch-num}}} \cdot \alpha_0 \quad (\text{or}) \quad \frac{k}{\sqrt{t}} \cdot \alpha_0$$



### ④ Manual decay

# local optima problem

key points



→ plateaus are the regions where the gradient is close to zero for long duration

\* momentum, Adam algorithm may overcome "plateau" problems

$$\frac{+10}{+9-1} = \frac{10}{8}$$

$$q = \frac{10}{8}$$

$$\frac{7.6}{4.0} = \frac{7.6}{2.0-1} = \frac{7.6}{1} = 7.6$$