



deeplearning.ai

# Optimization Algorithms

---

Mini-batch  
gradient descent

# Batch vs. mini-batch gradient descent

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

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

$(n_x, m)$        $\underbrace{\hspace{10em}}_{X^{\{1\}} (n_x, 1000)}$        $\underbrace{\hspace{10em}}_{X^{\{2\}} (n_x, 1000)}$        $\dots$        $\underbrace{\hspace{10em}}_{X^{\{5,000\}} (n_x, 1000)}$

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

$(1, m)$        $\underbrace{\hspace{10em}}_{Y^{\{1\}} (1, 1000)}$        $\underbrace{\hspace{10em}}_{Y^{\{2\}} (1, 1000)}$        $\dots$        $\underbrace{\hspace{10em}}_{Y^{\{5,000\}} (1, 1000)}$

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

5,000 mini-batches of 1,000 each

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

$$\left| \begin{array}{l} x^{(i)} \\ z^{[l]} \\ \underline{X^{\{t\}}}, \underline{Y^{\{t\}}} \end{array} \right.$$

# Mini-batch gradient descent

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

Forward prop on  $X^{\text{test}}$ .

$$Z^{(1)} = W^{(1)} X^{\text{test}} + b^{(1)}$$

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

$$\vdots$$

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

Vectorized implementation  
(1000 examples)

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

for  $X^{\text{test}}, Y^{\text{test}}$

Backprop to compute gradients w.r.t  $J^{\text{test}}$  (using  $X^{\text{test}}, Y^{\text{test}}$ )

$$W^{(1)} := W^{(1)} - \alpha dW^{(1)}, \quad b^{(1)} := b^{(1)} - \alpha db^{(1)}$$

"1 epoch"

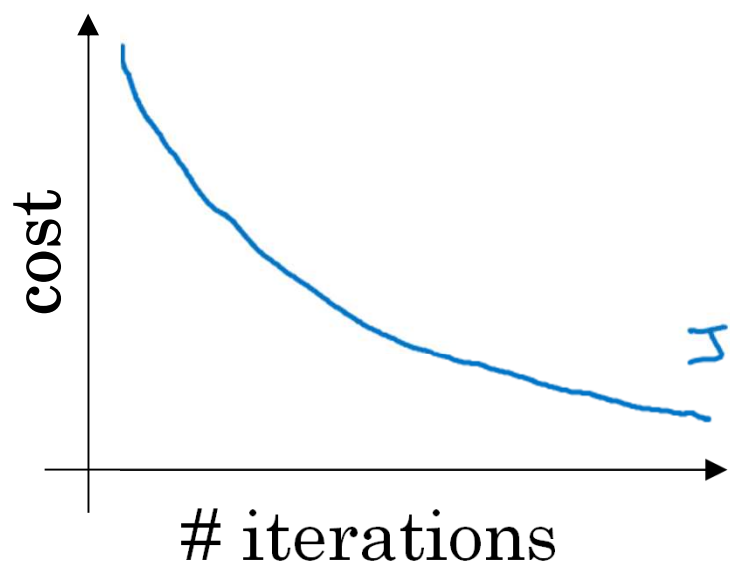
pass through training set.

1 step of grad desc  
using  $X^{\text{test}}, Y^{\text{test}}$ .  
(as if  $m=1000$ )

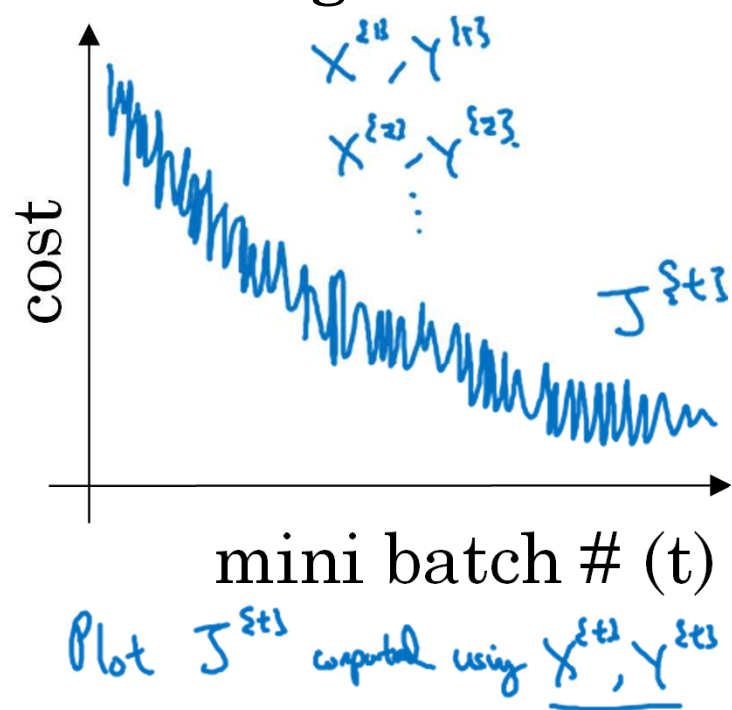
$X, Y$

# Training with mini batch gradient descent

## Batch gradient descent



## Mini-batch gradient descent



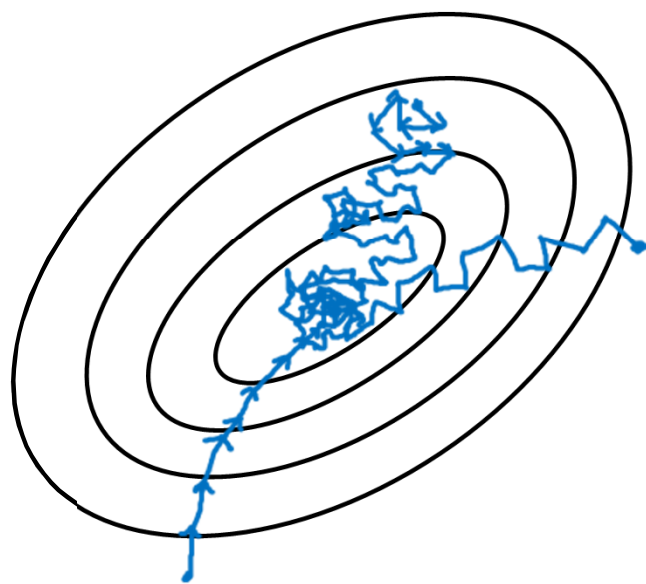
# Choosing your mini-batch size

→ If mini-batch size =  $m$  : Batch gradient descent.

$$(X^{(13)}, Y^{(13)}) = (X, Y).$$

→ If mini-batch size = 1 : Stochastic gradient descent. Every example is its own mini-batch.  
 $(X^{(13)}, Y^{(13)}) = (x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})$  mini-batch.

In practice: Somewhere in-between 1 and  $m$



Stochastic  
gradient  
descent



Loss spreading  
from vectorization

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



Fastest learning.

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

Batch  
gradient descent  
(mini-batch size =  $m$ )



Too long  
per iteration


# Choosing your mini-batch size

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

Typical mini-batch sizes:

→ 64, 128, 256, 512

$2^6$        $2^7$        $2^8$        $2^9$



$\frac{1024}{2^{10}}$

Make sure mini-batch fit in CPU/GPU memory.  
 $X^{(t)}, Y^{(t)}$

# Temperature in London

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

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

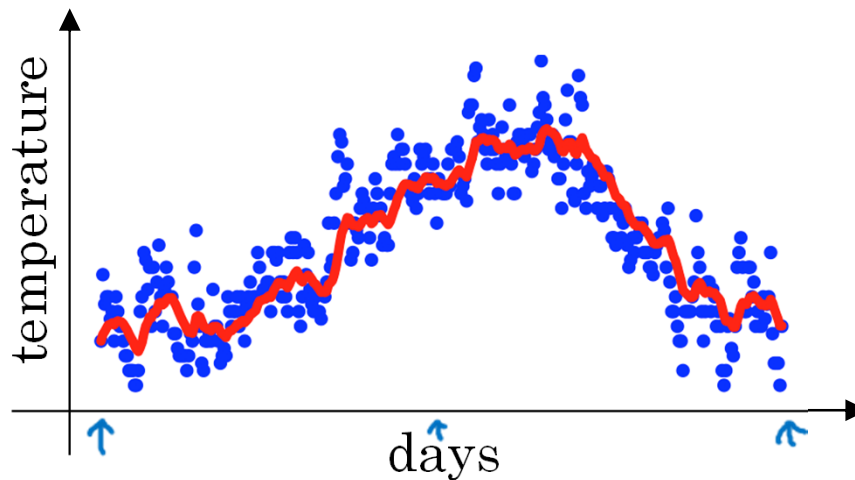
$$\theta_3 = 45^\circ\text{F} \quad \vdots$$

$\vdots$

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

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

$\vdots$



$$V_0 = 0$$

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

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

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

$\vdots$

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

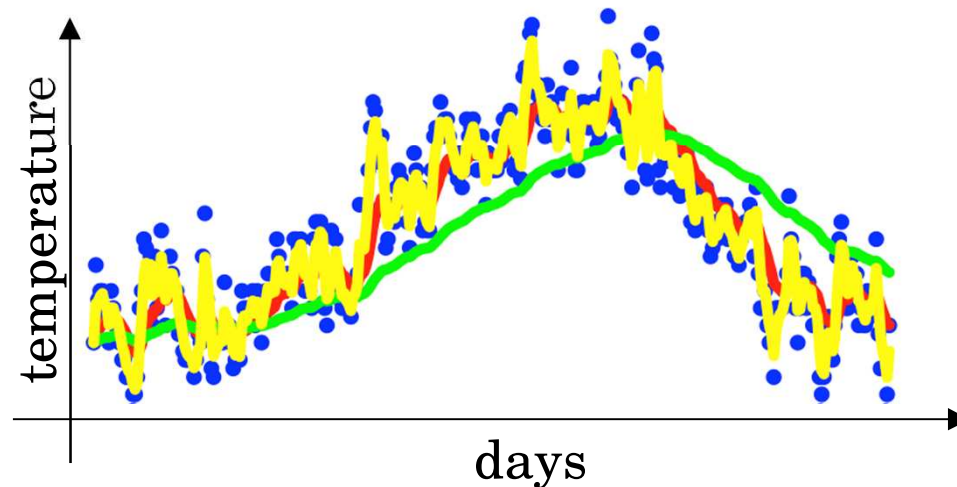
# Exponentially weighted averages <sup>moving</sup>

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

$\beta = 0.9$  :  $\approx 10$  days' temperature.  
 $\beta = 0.98$  :  $\approx 50$  days  
 $\beta = 0.5$  :  $\approx 2$  days

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

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





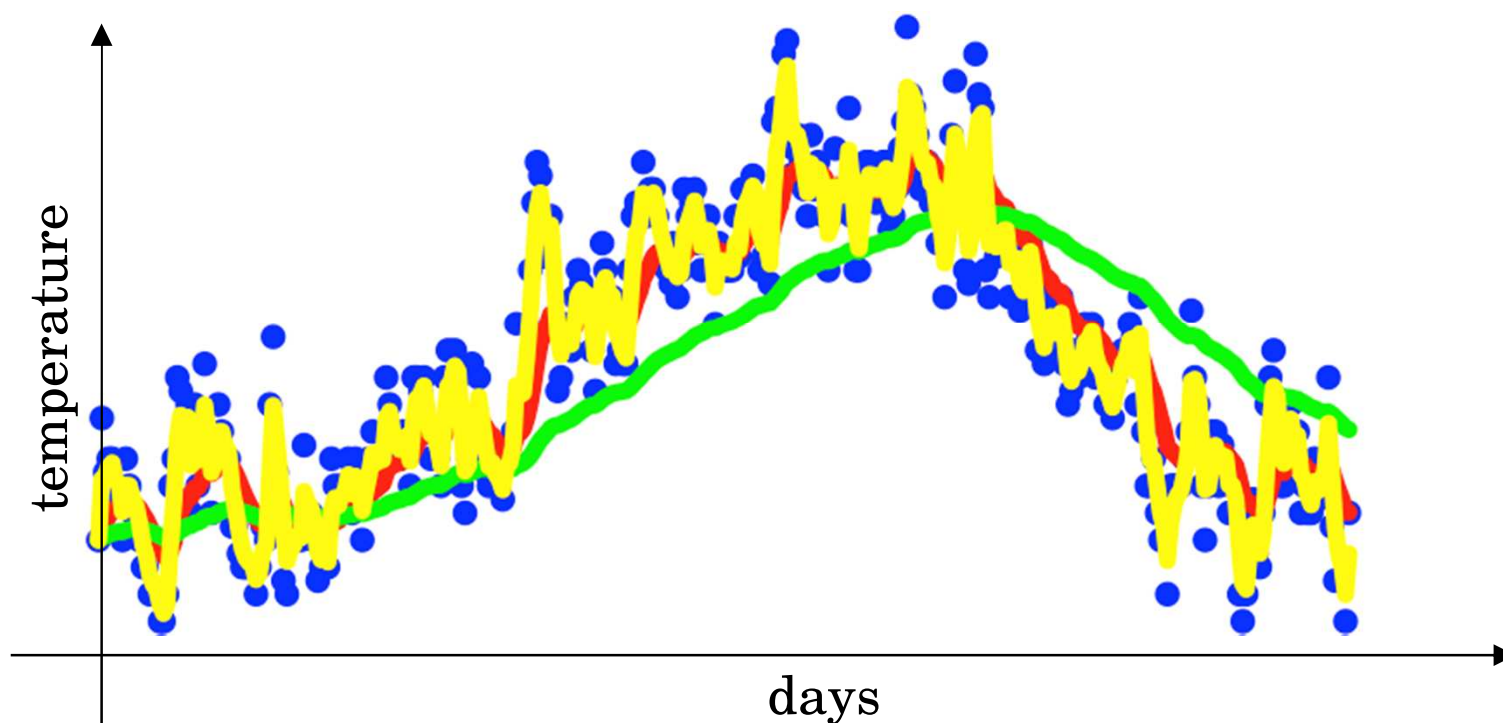
# Exponentially weighted averages

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

$$\beta = 0.9$$

$$0.98$$

$$0.5$$



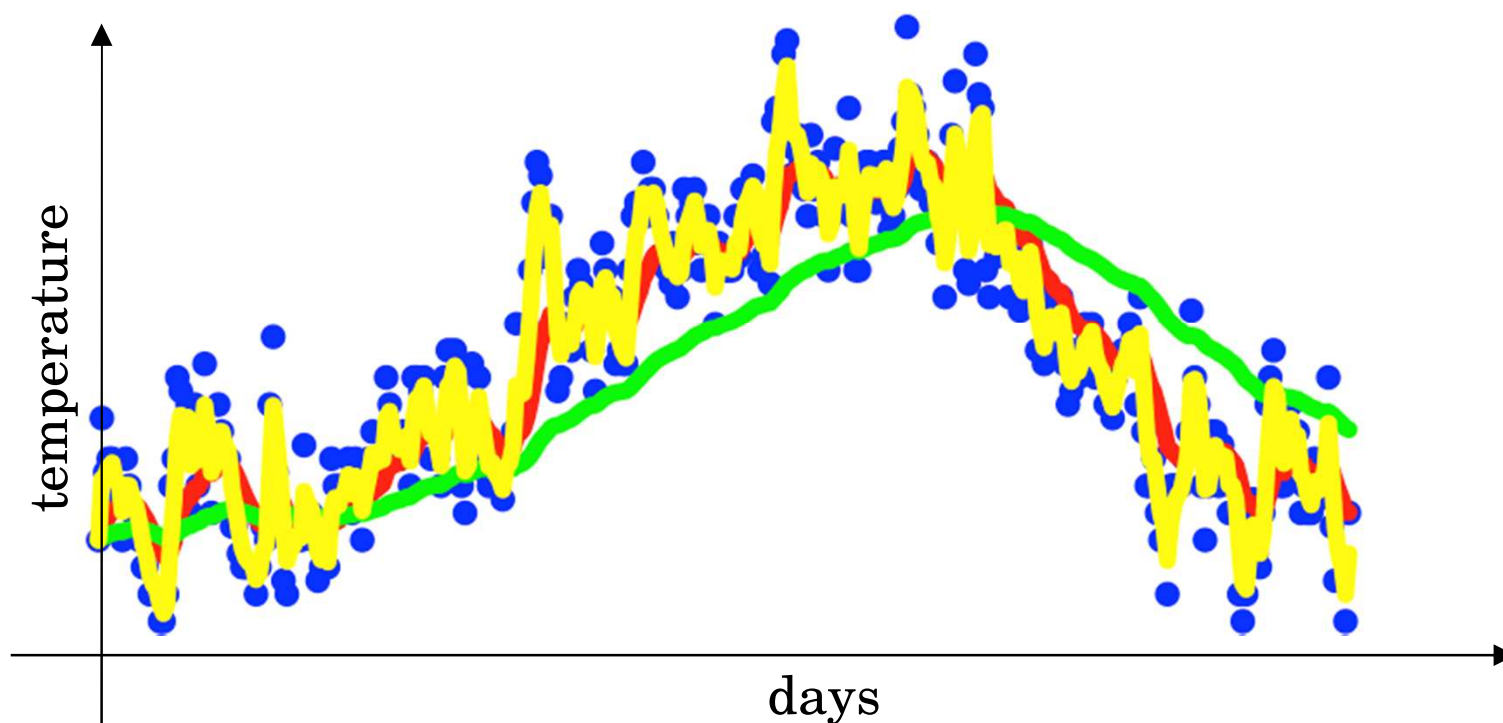
# Exponentially weighted averages

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

$$\beta = 0.9$$

$$0.98$$

$$0.5$$



# Implementing exponentially weighted averages

$$v_0 = 0$$

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

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

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

...

$$V_0 := 0$$

$$V_1 := \beta v + (1 - \beta) \theta_1$$

$$V_2 := \beta v + (1 - \beta) \theta_2$$

⋮

$$\rightarrow V_0 = 0$$

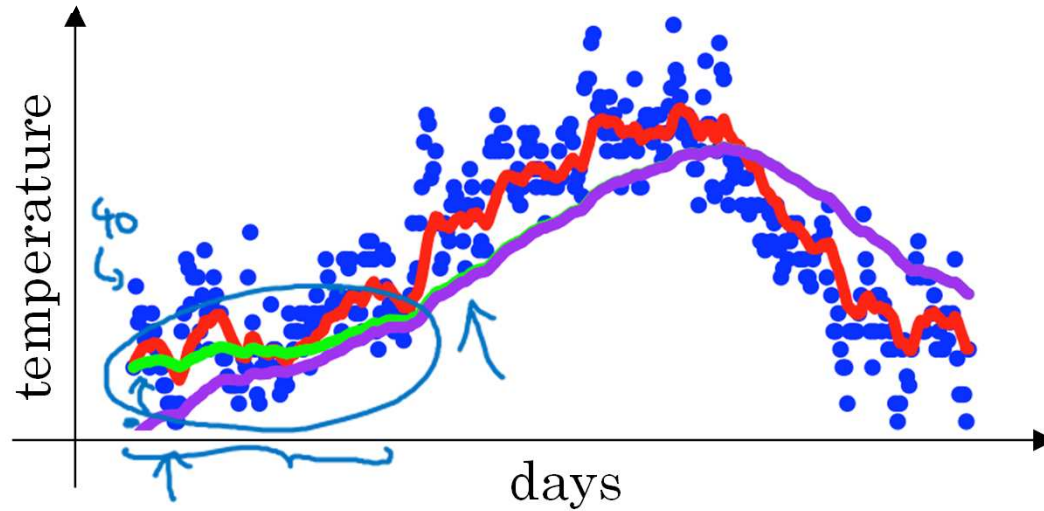
Repeat {

Get next  $\theta_t$

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

}

# Bias correction



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

$$v_0 = 0$$

$$v_1 = \cancel{0.98 v_0} + \underline{0.02 \theta_1}$$

$$v_2 = 0.98 v_1 + 0.02 \theta_2$$

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

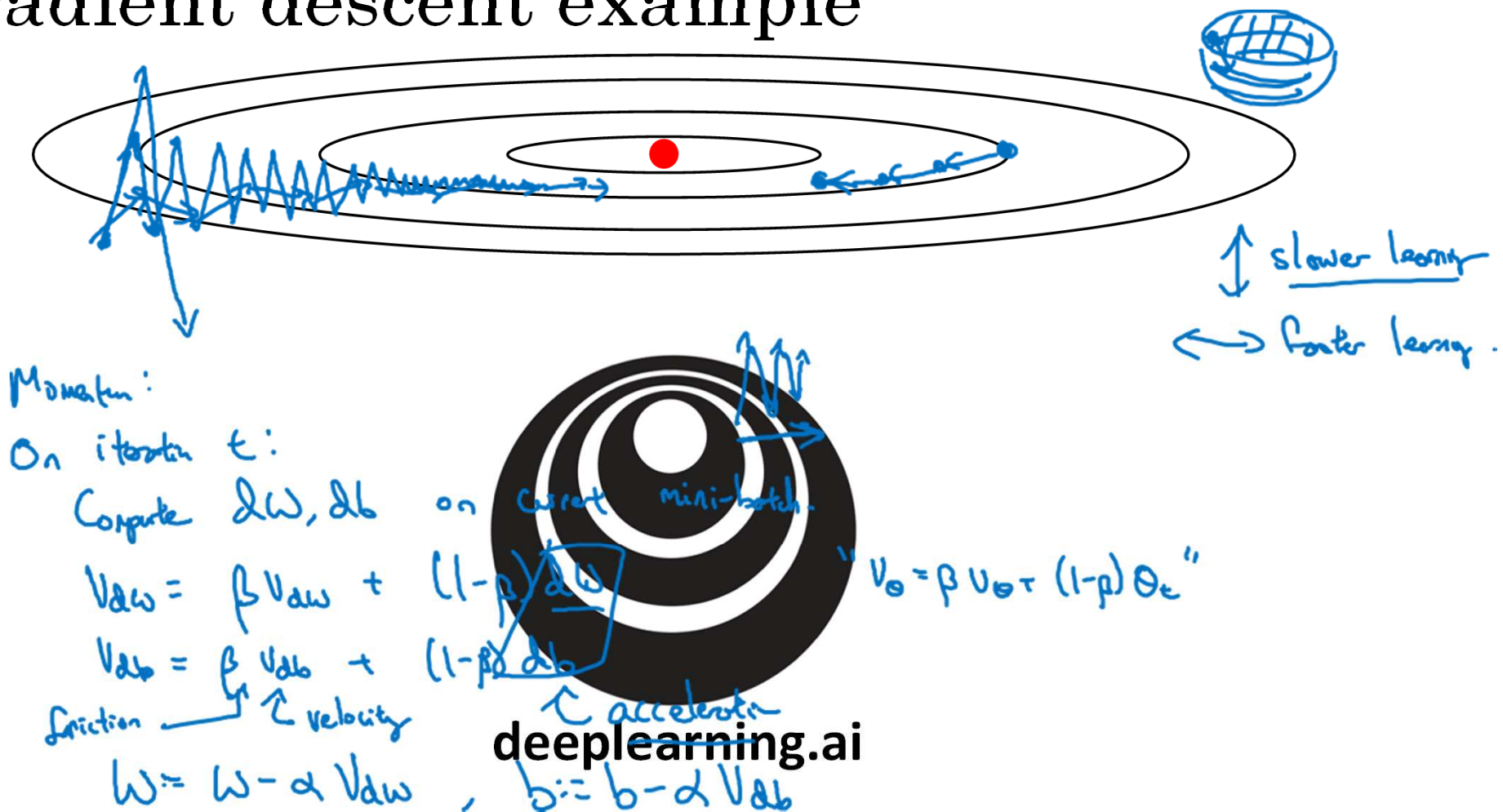
$$= \underline{0.0196 \theta_1} + \underline{0.02 \theta_2}$$

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

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

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

# Gradient descent example



# Implementation details

On iteration  $t$ :   
  $v_{dw} = 0, v_{db} = 0$

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

$$\begin{aligned} v_{dw} &= \beta v_{dw} + (1 - \beta) dW \\ v_{db} &= \beta v_{db} + (1 - \beta) db \\ W &= W - \alpha v_{dw}, \quad b = b - \alpha v_{db} \end{aligned}$$

Handwritten notes:   
 - Blue arrows point to the update equations for  $v_{dw}$  and  $v_{db}$ .   
 - A blue bracket groups the first two equations.   
 - A blue equation  $v_{dw} = \beta v_{dw} + dW \leftarrow$  is written to the right.   
 - A blue equation  $\frac{v_{dw}}{1 - \beta^t}$  is crossed out with a blue 'X'.

Hyperparameters  $\beta$

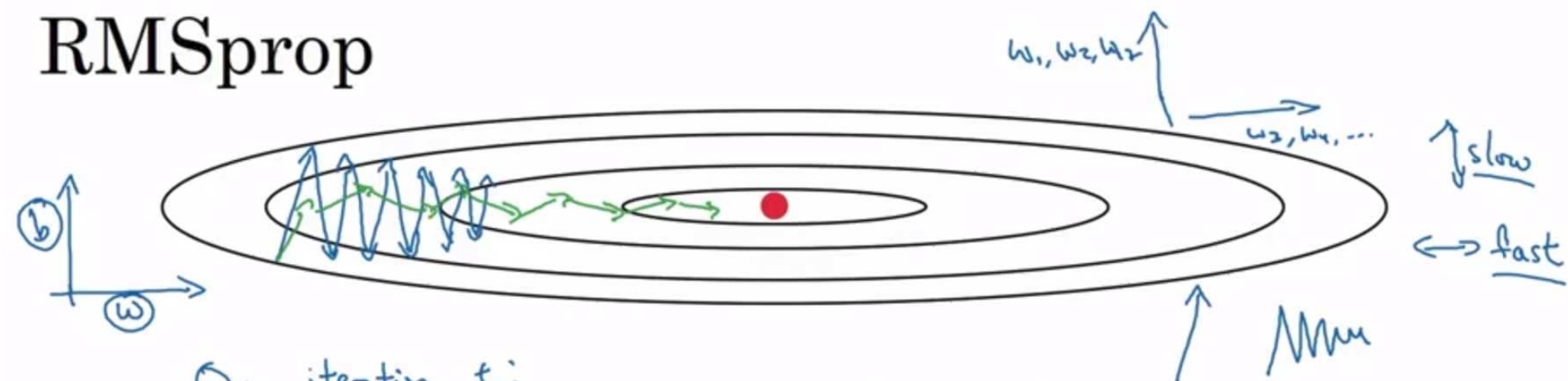
$$\beta = 0.9$$

$\uparrow \uparrow$

over many loss & ln gradients



# RMSprop



On iteration  $t$ :

Compute  $dw, db$  on current mini-batch

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

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

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

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

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

# 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

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

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

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

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

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

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



# Hyperparameters choice:

→  $\alpha$  : needs to be tune

→  $\beta_1$  : 0.9 → ( $dw$ )

→  $\beta_2$  : 0.999 → ( $dw^2$ )

→  $\epsilon$  :  $10^{-8}$

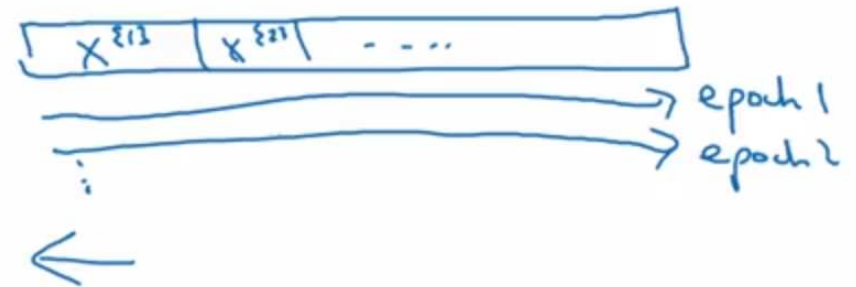
Adam: Adaptive moment estimation

# 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
$\vdots$	$\vdots$




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



# Other learning rate decay methods

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

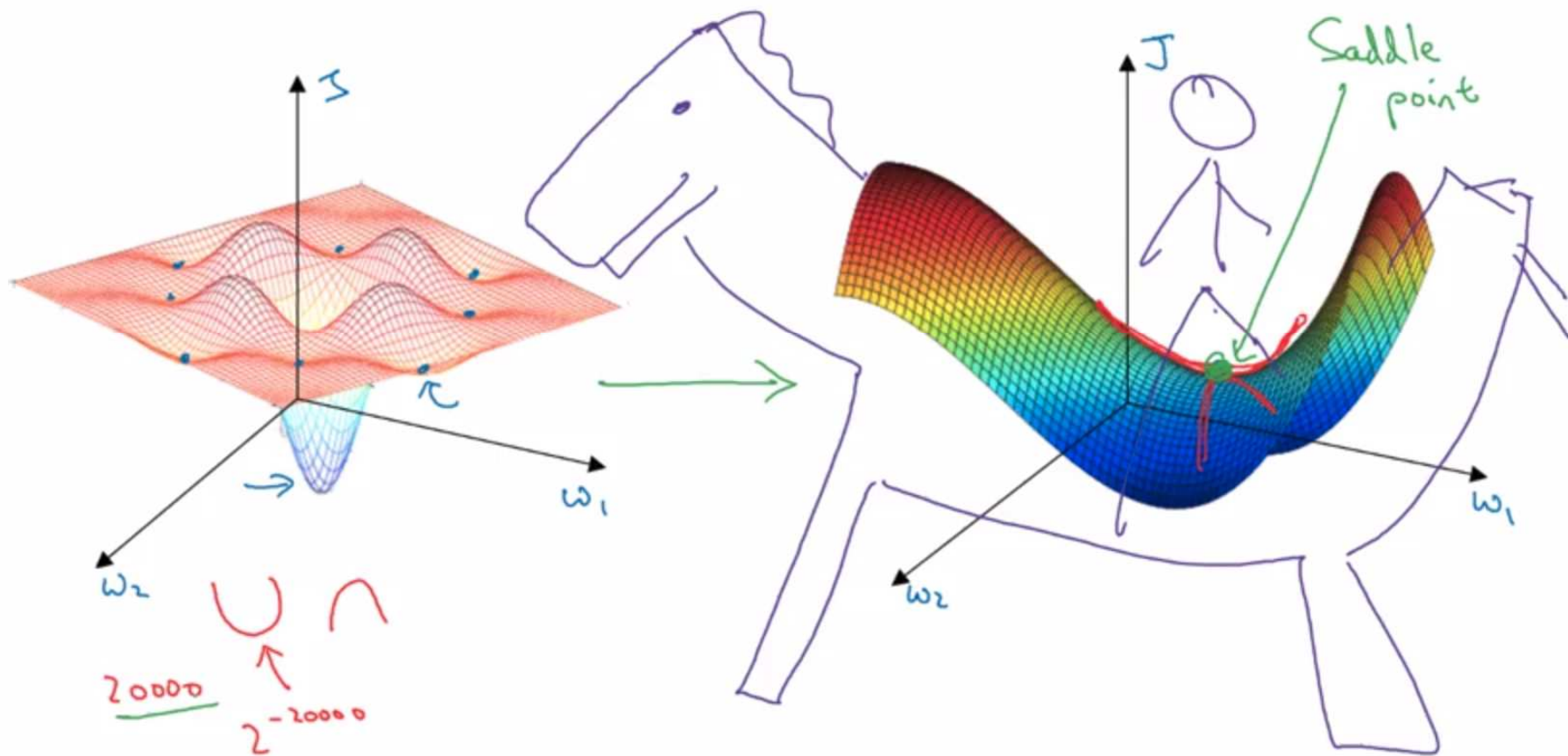
$\alpha = \frac{k}{\sqrt{\text{epoch-num}}} \cdot \alpha_0$  or  $\frac{k}{\sqrt{t}} \cdot \alpha_0$



discrete staircase

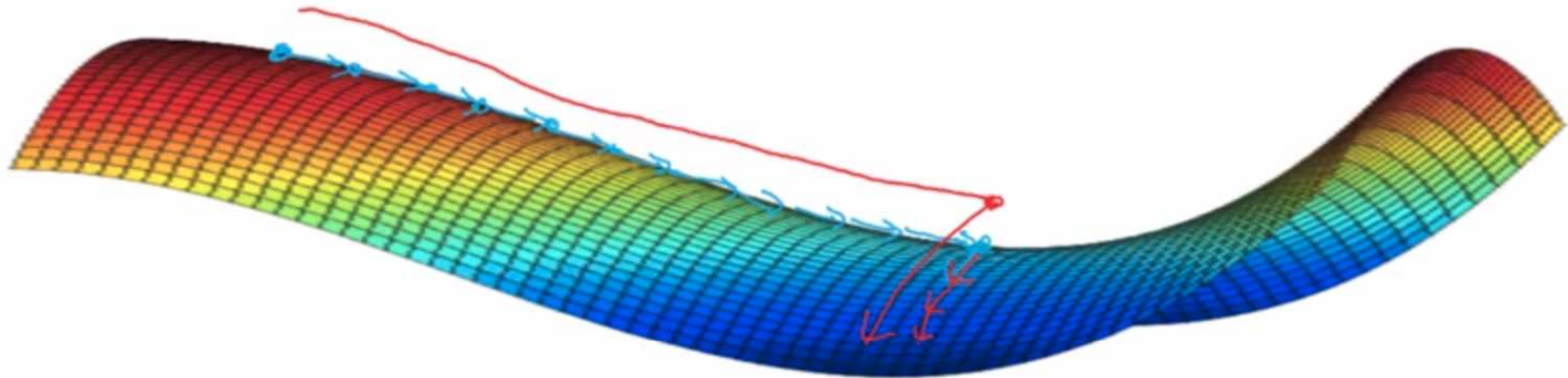
Manual decay.

# Local optima in neural networks



Andrew Ng

# Problem of plateaus



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