

## Other learning rate decay methods

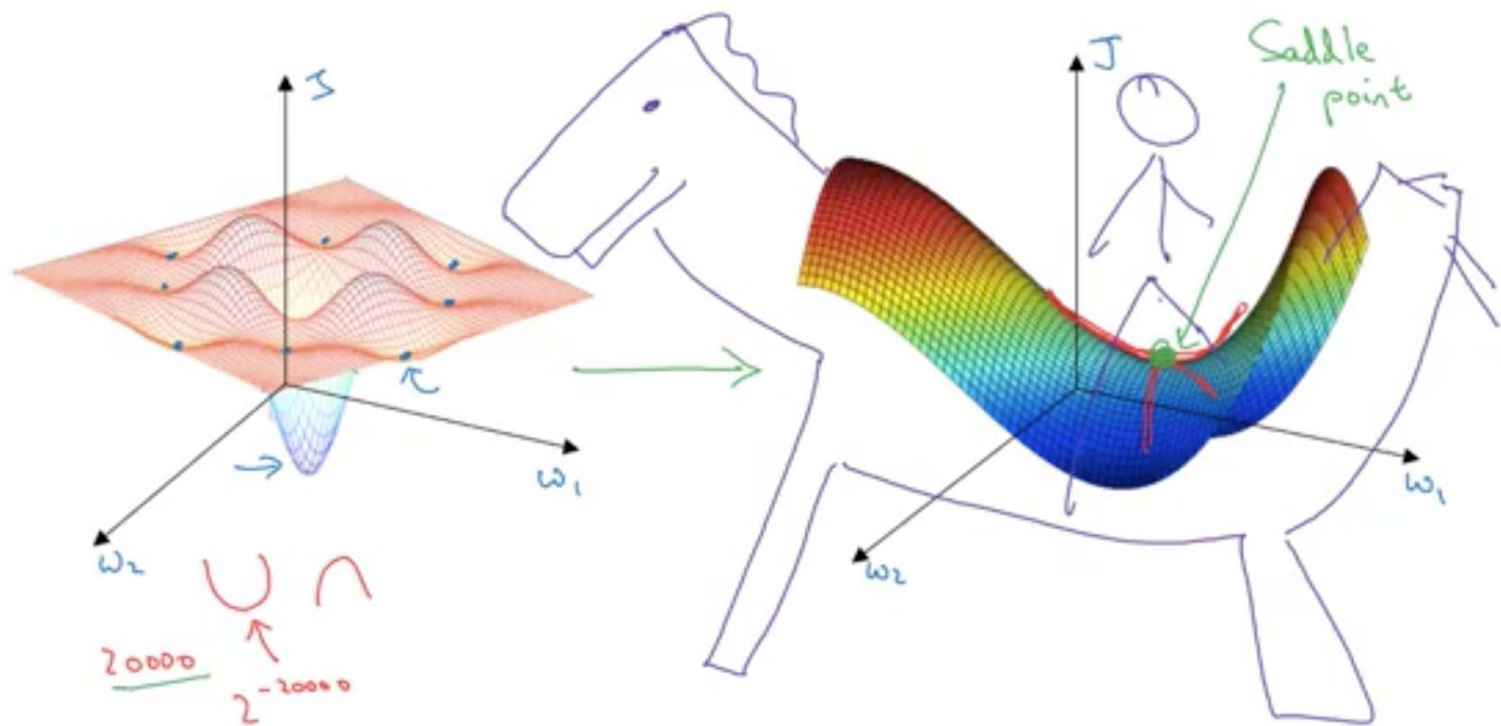
formula

$$\left\{ \begin{array}{l} \alpha = 0.95^{\text{epoch-num}} \cdot \alpha_0 \quad - \text{exponentially decay.} \\ \alpha = \frac{k}{\sqrt{\text{epoch-num}}} \cdot \alpha_0 \quad \text{or} \quad \frac{k}{\sqrt{t}} \cdot \alpha_0 \\ \begin{array}{c} \alpha \\ | \\ \hline | \\ \hline t \end{array} \end{array} \right.$$

discrete staircase

Manual decay.

# Local optima in neural networks

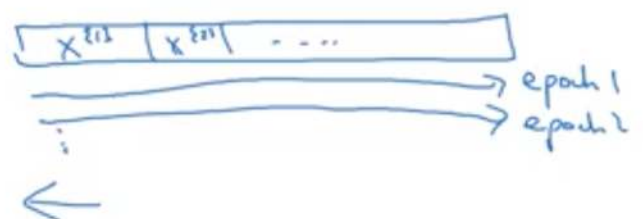


# Learning rate decay

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
$\vdots$	$\vdots$

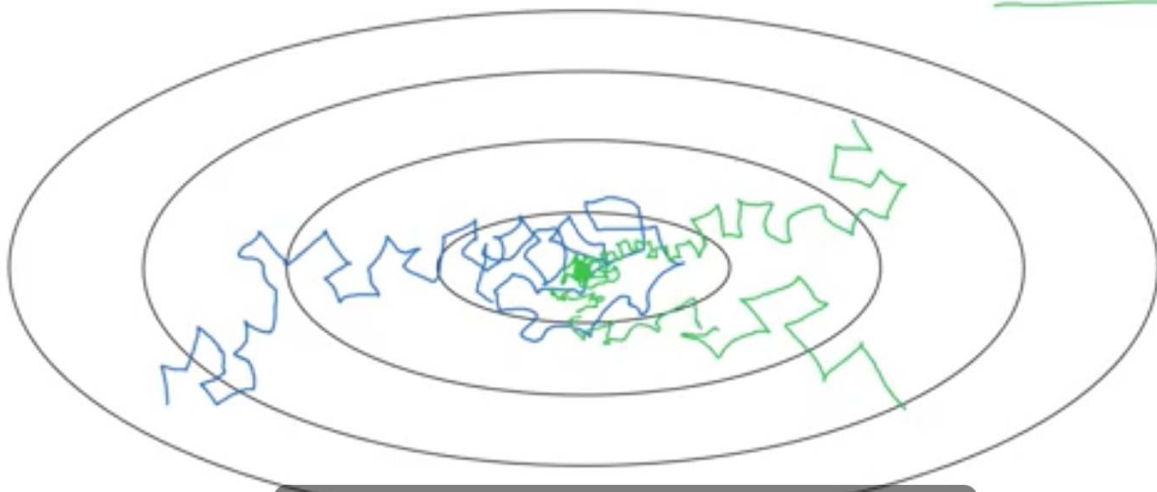


$\alpha_0 = 0.2$   
decay-rate = 1



# Learning rate decay

Slowly reduce  $\alpha$



then having a slower learning rate allows you to take smaller steps.

## 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

# 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

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

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

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

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

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

## Implementation details

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

On iteration  $t$ :

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

$$\left. \begin{aligned} \rightarrow v_{dw} &= \beta v_{dw} + (1 - \beta) dW \\ \rightarrow v_{db} &= \beta v_{db} + (1 - \beta) db \end{aligned} \right\} \quad \left| \quad v_{dw} = \beta v_{dw} + dW \leftarrow$$

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

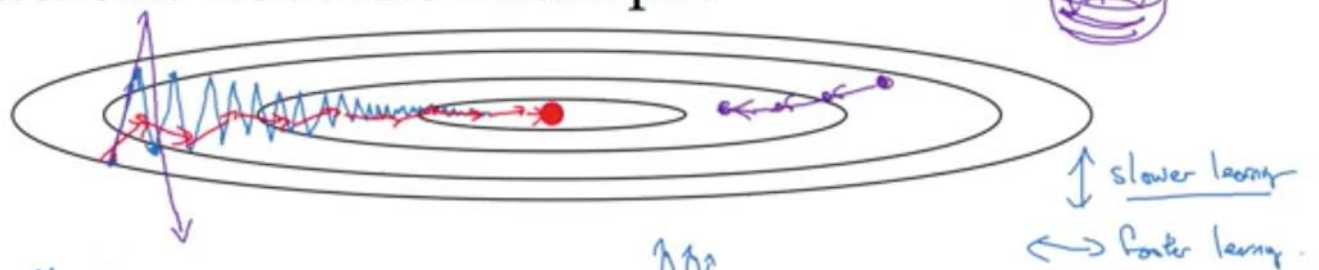
$$\frac{v_{dw}}{1 - \beta^t}$$

Hyperparameters:  $\alpha, \beta$

$$\beta = 0.9$$

average over last  $\approx 10$  gradients

# Gradient descent example



Momentum:

On iteration  $t$ :

Compute  $\partial W, \partial b$  on current mini-batch.

$$V_{dw} = \beta V_{dw} + (1-\beta) \frac{\partial W}{\partial w}$$

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

Friction  $\rightarrow$   $\uparrow$  velocity

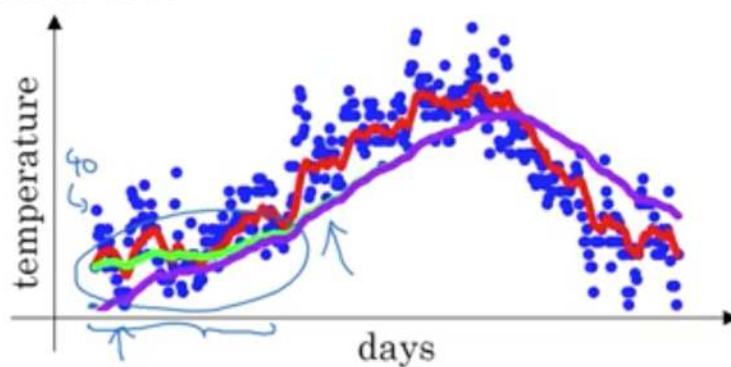
$\uparrow$  acceleration

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

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



## Bias correction



$$\beta = 0.98$$

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

$$v_0 = 0$$

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

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

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

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

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

## Implementing exponentially weighted averages

$$v_0 = 0$$

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

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

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

...

$$V_0 := 0$$

$$V_0 := \beta v + (1 - \beta) \theta_1$$

$$V_0 := \beta v + (1 - \beta) \theta_2$$

⋮

→  $V_0 = 0$

Repeat {

Get next  $\theta_t$

$$V_0 := \beta V_0 + (1 - \beta) \theta_t \leftarrow$$

}

## Exponentially weighted averages

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

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

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

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

...

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

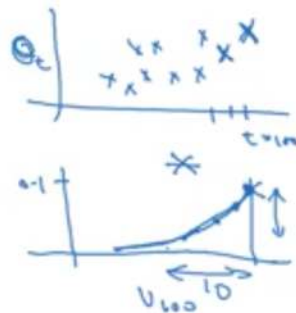
$$= 0.1 \theta_{100} + 0.1 \times 0.9 \cdot \theta_{99} + 0.1 (0.9)^2 \theta_{98} + 0.1 (0.9)^3 \theta_{97} + 0.1 (0.9)^4 \theta_{96}$$

$$0.9^{10} \approx 0.35 \approx \frac{1}{e}$$

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

$$0.98^{50} \approx \frac{1}{e}$$

Andrew Ng

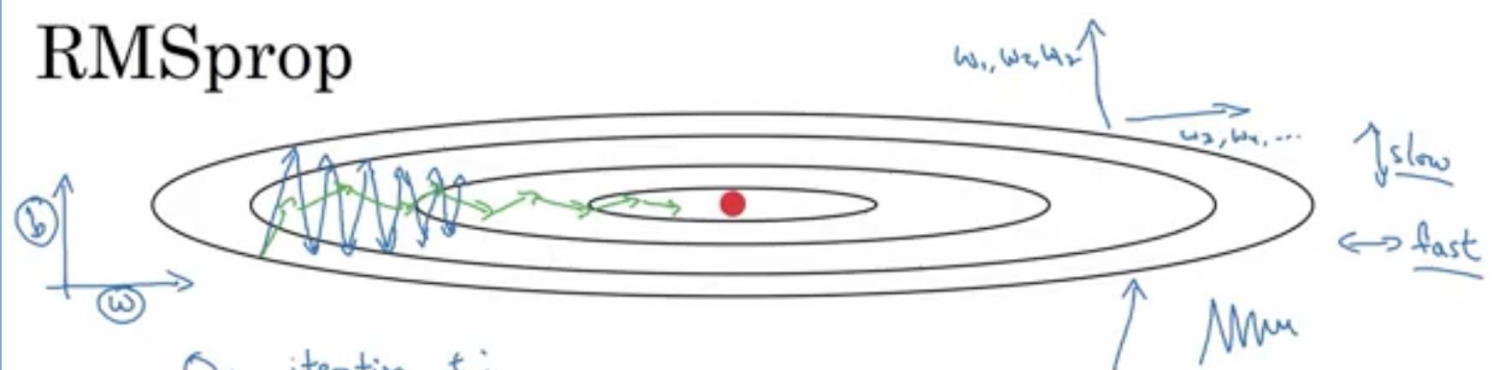


$$0.1 \theta_{97} + 0.9 \nu_{97}$$

6.48?

$$\rightarrow \underline{0.98^{50}} \approx \frac{1}{e}$$

# RMSprop



On iteration  $t$ :

Compute  $dw, db$  on current mini-batch

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

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

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

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

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

# Exponentially weighted averages

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

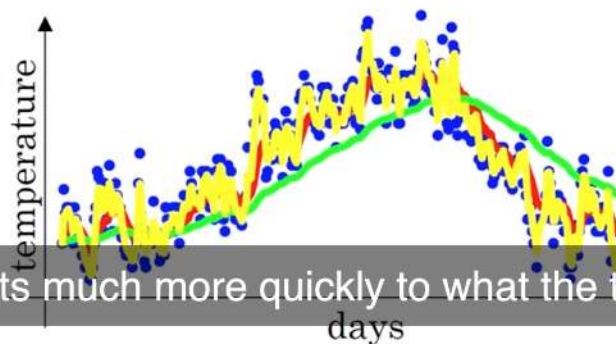
$\beta = 0.9$  :  $\approx 10$  days' temperature

$\underline{\beta = 0.98}$  :  $\approx 50$  days

$\beta = 0.5$  :  $\approx 2$  days

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

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



But this adapts much more quickly to what the temperature changes.

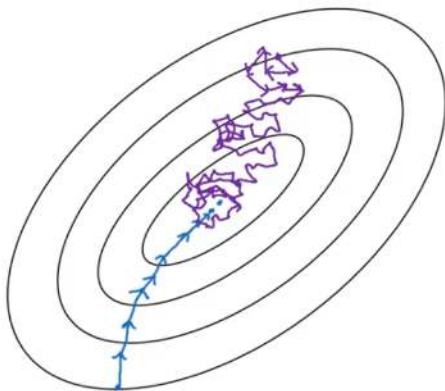
# Choosing your mini-batch size

→ 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 mini-batch.  
 $(X^{(1)}, Y^{(1)}) = (x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})$  mini-batch.

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



Stochastic  
gradient  
descent



Less sensitive  
to vectorization

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



Fastest learning.

- Vectorization.  
(~1000)
- Make passes without  
processing entire training set:

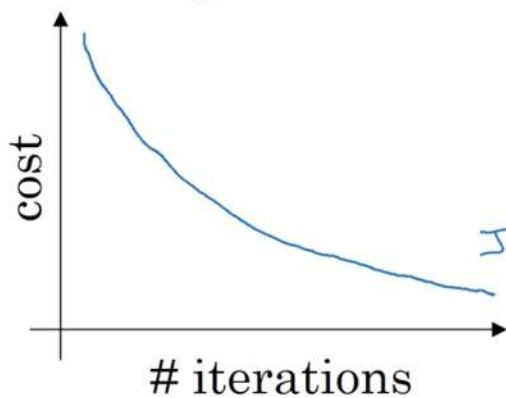
Batch  
gradient descent  
(mini-batch size =  $m$ )



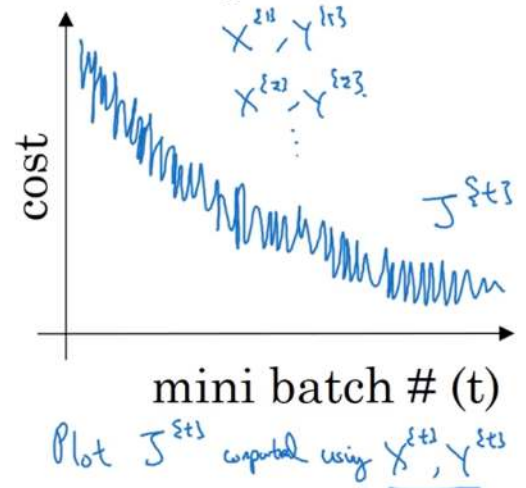
Too long  
per iteration

# Training with mini batch gradient descent

Batch gradient descent



Mini-batch gradient descent





# Mini-batch gradient descent

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

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

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

$$A^{(t)} = \sigma^{(t)}(Z^{(t)})$$

$\vdots$

$$A^{(t)} = \sigma^{(t)}(Z^{(t)})$$

Vectorized implementation  
(1000 examples)

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

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

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

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

"1 epoch"

pass through training set.

1 step of gradient descent  
using  $X^{\text{test}}, Y^{\text{test}}$ .  
(as if  $n=1000$ )

$X, Y$

