

# Optimization Algorithms

Most of this material is from Prof. Andrew Ng's and Chang's slides.

## Introduction

- We will learn the optimization algorithms enable you to train your neural network much faster -> **training efficiency!**
- Applying machine learning is a highly empirical process (i.e., highly iterative process)
  - You should train a lot of models to find one that works really well
- Deep learning works best in a regime of big data
  - We should train our model on a huge dataset
- So, having fast optimization algorithms can really help the above problems

## Batch vs. mini-batch gradient descent

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

$$\underset{(n,m)}{X} = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & \dots & x^{(m)} \end{bmatrix}$$

$$\underset{(1,m)}{Y} = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & \dots & y^{(m)} \end{bmatrix}$$

## Batch vs. mini-batch gradient descent

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

$$\underset{(n,m)}{X} = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & \dots & x^{(m)} \end{bmatrix}$$

$$\underset{(1,m)}{Y} = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & \dots & y^{(m)} \end{bmatrix}$$

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

## Batch vs. mini-batch gradient descent

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

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

$(n, m)$        $\underbrace{\hspace{10em}}_{X^{1\{1\}}}$        $\underbrace{\hspace{10em}}_{X^{2\{2\}}}$        $\underbrace{\hspace{10em}}_{X^{5000\{5000\}}}$

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

$(1, m)$

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

5,000 mini-batches of 1,000 samples

## Batch vs. mini-batch gradient descent

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

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

$(n, m)$        $\underbrace{\hspace{10em}}_{X^{1\{1\}}_{(n, 1000)}}$        $\underbrace{\hspace{10em}}_{X^{2\{2\}}_{(n, 1000)}}$        $\underbrace{\hspace{10em}}_{X^{5000\{5000\}}}$

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

$(1, m)$        $\underbrace{\hspace{10em}}_{Y^{1\{1\}}_{(1, 1000)}}$        $\underbrace{\hspace{10em}}_{Y^{2\{2\}}_{(1, 1000)}}$        $\underbrace{\hspace{10em}}_{Y^{5000\{5000\}}}$

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

5,000 mini-batches of 1,000 samples

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

## Batch vs. mini-batch gradient descent

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

$$\begin{array}{l}
 X = \underbrace{\begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(1000)} \end{bmatrix}}_{X^{\{1\}}_{(n,1000)}} \quad \underbrace{\begin{bmatrix} x^{(1001)} & \dots & x^{(2000)} \end{bmatrix}}_{X^{\{2\}}_{(n,1000)}} \quad \dots \quad \underbrace{\begin{bmatrix} \dots & x^{(m)} \end{bmatrix}}_{X^{\{5000\}}} \\
 (n,m) \\
 Y = \underbrace{\begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(1000)} \end{bmatrix}}_{Y^{\{1\}}_{(1,1000)}} \quad \underbrace{\begin{bmatrix} y^{(1001)} & \dots & y^{(2000)} \end{bmatrix}}_{Y^{\{2\}}_{(1,1000)}} \quad \dots \quad \underbrace{\begin{bmatrix} \dots & y^{(m)} \end{bmatrix}}_{Y^{\{5000\}}} \\
 (1,m)
 \end{array}$$

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

5,000 mini-batches of 1,000 samples

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

$x^{(i)} ?$

$z^{[l]} ?$

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

## Batch vs. mini-batch gradient descent

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

$$\begin{array}{l}
 X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & \dots & x^{(m)} \end{bmatrix} \\
 (n,m)
 \end{array}$$

$$\begin{array}{l}
 Y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & \dots & y^{(m)} \end{bmatrix} \\
 (1,m)
 \end{array}$$

$$\begin{aligned}
 dW^{[l]} &= \frac{\partial J}{\partial W^{[l]}} = \frac{\partial \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\mathcal{Y}, y)}{\partial W^{[l]}} = \frac{1}{m} \sum_{i=1}^m \frac{\partial \mathcal{L}(\mathcal{Y}, y)}{\partial W^{[l]}} \\
 db^{[l]} &= \frac{\partial J}{\partial b^{[l]}} = \frac{\partial \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\mathcal{Y}, y)}{\partial b^{[l]}} = \frac{1}{m} \sum_{i=1}^m \frac{\partial \mathcal{L}(\mathcal{Y}, y)}{\partial b^{[l]}}
 \end{aligned}$$

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

5,000 mini-batches of 1,000 samples

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

$x^{(i)} ?$

$z^{[l]} ?$

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

## Batch vs. mini-batch gradient descent

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

$$\begin{array}{l}
 X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(1000)} & | & x^{(1001)} & \dots & x^{(2000)} & | & \dots & | & \dots & x^{(m)} \end{bmatrix} \\
 (n, m) \quad \underbrace{\hspace{10em}}_{X^{\{1\}}_{(n, 1000)}} \quad \underbrace{\hspace{10em}}_{X^{\{2\}}_{(n, 1000)}} \quad \dots \quad \underbrace{\hspace{10em}}_{X^{\{5000\}}_{(n, 1000)}} \\
 \\
 Y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(1000)} & | & y^{(1001)} & \dots & y^{(2000)} & | & \dots & | & \dots & y^{(m)} \end{bmatrix} \\
 (1, m) \quad \underbrace{\hspace{10em}}_{Y^{\{1\}}_{(1, 1000)}} \quad \underbrace{\hspace{10em}}_{Y^{\{2\}}_{(1, 1000)}} \quad \dots \quad \underbrace{\hspace{10em}}_{Y^{\{5000\}}_{(1, 1000)}}
 \end{array}$$

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

5,000 mini-batches of 1,000 samples

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

$x^{(i)}$  ?

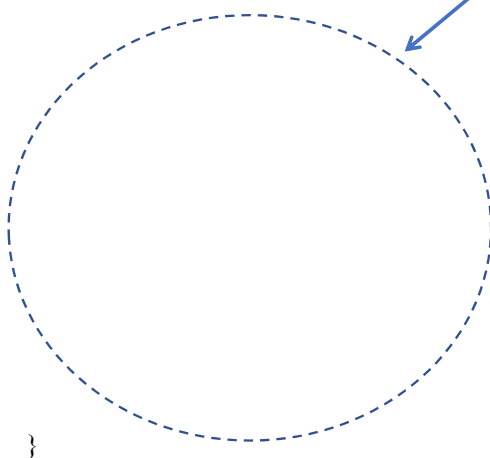
$z^{[l]}$  ?

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

## Mini-batch gradient descent

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

one step of gradient descent  
using  $X^{\{t\}}, Y^{\{t\}}$  (as if  $m = 1000$ )



## Mini-batch gradient descent

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

one step of gradient descent  
using  $X^{\{t\}}, Y^{\{t\}}$  (as if  $m = 1000$ )

Forward prop on  $X^{\{t\}}$

$$\begin{aligned} Z^{[1]} &= W^{[1]}X^{\{1\}} + b^{[1]} \\ A^{[1]} &= g^{[1]}(Z^{[1]}) \\ &\vdots \\ A^{[L]} &= g^{[L]}(Z^{[L]}) \end{aligned} \quad \left. \vphantom{\begin{aligned} Z^{[1]} &= W^{[1]}X^{\{1\}} + b^{[1]} \\ A^{[1]} &= g^{[1]}(Z^{[1]}) \\ &\vdots \\ A^{[L]} &= g^{[L]}(Z^{[L]}) \end{aligned}} \right\} \begin{array}{l} \text{vectorized implementation} \\ (1000 \text{ examples}) \end{array}$$

Compute cost :  $J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^{1000} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_l \|W^{[l]}\|_F^2$

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

}

## Mini-batch gradient descent

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

one step of gradient descent  
using  $X^{\{t\}}, Y^{\{t\}}$  (as if  $m = 1000$ )

Forward prop on  $X^{\{t\}}$

$$\begin{aligned} Z^{[1]} &= W^{[1]}X^{\{1\}} + b^{[1]} \\ A^{[1]} &= g^{[1]}(Z^{[1]}) \\ &\vdots \\ A^{[L]} &= g^{[L]}(Z^{[L]}) \end{aligned} \quad \left. \vphantom{\begin{aligned} Z^{[1]} &= W^{[1]}X^{\{1\}} + b^{[1]} \\ A^{[1]} &= g^{[1]}(Z^{[1]}) \\ &\vdots \\ A^{[L]} &= g^{[L]}(Z^{[L]}) \end{aligned}} \right\} \begin{array}{l} \text{vectorized implementation} \\ (1000 \text{ examples}) \end{array}$$

Compute cost :  $J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^{1000} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_l \|W^{[l]}\|_F^2$

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

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

$$\begin{aligned} W^{[l]} &:= W^{[l]} - \alpha \cdot dW^{[l]} \\ b^{[l]} &:= b^{[l]} - \alpha \cdot db^{[l]} \end{aligned}$$

}

## Mini-batch gradient descent

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

one step of gradient descent  
using  $X^{(t)}, Y^{(t)}$  (as if  $m = 1000$ )

Forward prop on  $X^{(t)}$

$$\begin{aligned} Z^{[1]} &= W^{[1]}X^{(1)} + b^{[1]} \\ A^{[1]} &= g^{[1]}(Z^{[1]}) \\ &\vdots \\ A^{[L]} &= g^{[L]}(Z^{[L]}) \end{aligned}$$

vectorized implementation  
(1000 examples)

Compute cost :  $J^{(t)} = \frac{1}{1000} \sum_{i=1}^{1000} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_l \|W^{[l]}\|_F^2$

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

Backprop to compute gradients w.r.t.  $J^{(t)}$  (using  $X^{(t)}$  and  $Y^{(t)}$ )

$$\begin{aligned} W^{[l]} &:= W^{[l]} - \alpha \cdot dW^{[l]} \\ b^{[l]} &:= b^{[l]} - \alpha \cdot db^{[l]} \end{aligned}$$

One Epoch

}

## Mini-batch gradient descent

repeat {

one step of gradient descent  
using  $X^{(t)}, Y^{(t)}$  (as if  $m = 1000$ )

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

Forward prop on  $X^{(t)}$

$$\begin{aligned} Z^{[1]} &= W^{[1]}X^{(1)} + b^{[1]} \\ A^{[1]} &= g^{[1]}(Z^{[1]}) \\ &\vdots \\ A^{[L]} &= g^{[L]}(Z^{[L]}) \end{aligned}$$

vectorized implementation  
(1000 examples)

Compute cost :  $J^{(t)} = \frac{1}{1000} \sum_{i=1}^{1000} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_l \|W^{[l]}\|_F^2$

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

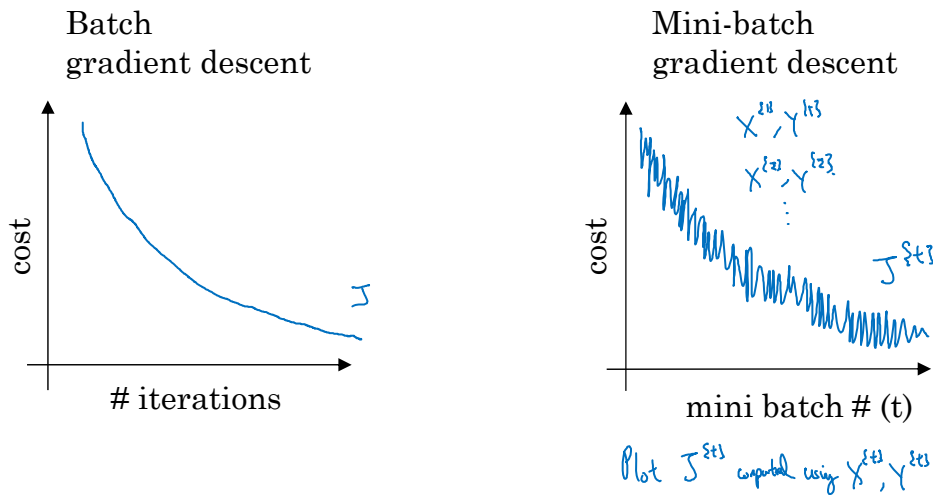
Backprop to compute gradients w.r.t.  $J^{(t)}$  (using  $X^{(t)}$  and  $Y^{(t)}$ )

$$\begin{aligned} W^{[l]} &:= W^{[l]} - \alpha \cdot dW^{[l]} \\ b^{[l]} &:= b^{[l]} - \alpha \cdot db^{[l]} \end{aligned}$$

One Epoch

}

## Training with mini-batch gradient descent



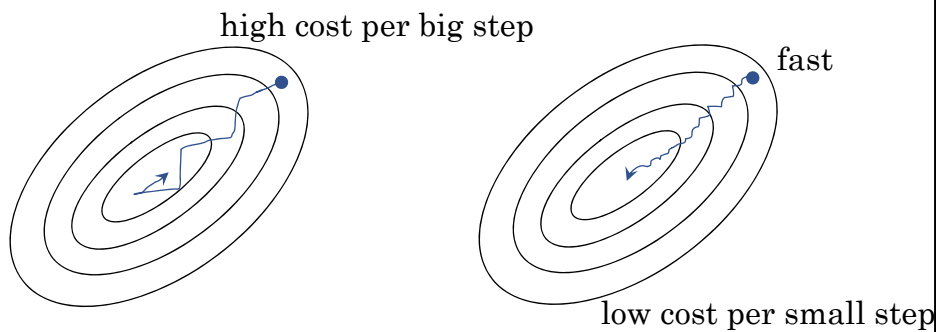
## Choosing your mini-batch size

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



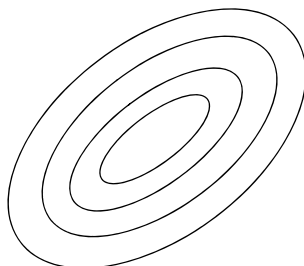
# Choosing your mini-batch size

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




# 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:  
 $(X^{(1)}, Y^{(1)}) = (x^{(1)}, y^{(1)}), \dots, (X^{(m)}, Y^{(m)}) = (x^{(m)}, y^{(m)})$   
→ Every example is its own mini-batch
- In practice, somewhere in-between 1 and m



stochastic gradient descent (mini-batch size = 1)	In-between (mini-batch size not too big or small)	batch gradient descent (mini-batch size = m)
Lose speedup from vectorization (=parallelization)	Fastest learning - Vectorization - Make progress without processing entire training set	Too long per iteration(=epoch)

## Choosing your mini-batch size

- if small training set : use batch gradient descent  
(e.g.,  $m \leq 2000$ )
- else, typical mini-batch size :  
→ 64, 128, 256, 512, 1024, ...  
 $2^6 \quad 2^7 \quad 2^8 \quad 2^9 \quad 2^{10}$   
  
power of 2
- Make sure mini-batch fits in CPU/GPU memory

## Exponentially weighted averages for understanding momentum

- We need to use exponentially weighted averages to understand more sophisticated optimization algorithms than gradient descent

## Temperature in London

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

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

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

$\vdots$

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

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

$\vdots$

## Temperature in London

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

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

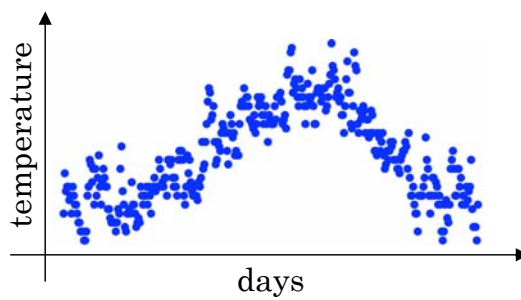
$$\theta_3 = 45^\circ\text{F}$$

$\vdots$

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

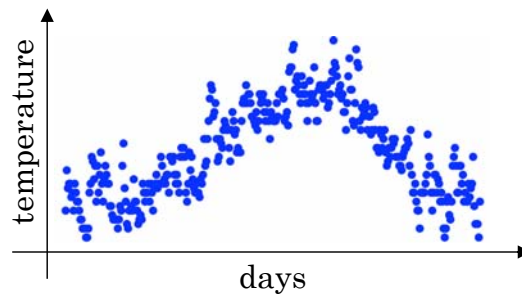
$$\theta_{181} = 56^\circ\text{F}$$

$\vdots$



## Temperature in London

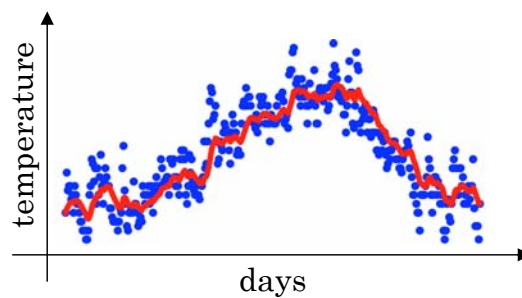
$\theta_1 = 40^\circ\text{F}$   
 $\theta_2 = 49^\circ\text{F}$   
 $\theta_3 = 45^\circ\text{F}$   
 $\vdots$   
 $\theta_{180} = 60^\circ\text{F}$   
 $\theta_{181} = 56^\circ\text{F}$   
 $\vdots$



$$\begin{aligned}v_0 &= 0 \\v_1 &= 0.9v_0 + 0.1\theta_1 \\v_2 &= 0.9v_1 + 0.1\theta_2 \\v_3 &= 0.9v_2 + 0.1\theta_3 \\&\vdots \\v_t &= 0.9v_{t-1} + 0.1\theta_t\end{aligned}$$

## Temperature in London

$\theta_1 = 40^\circ\text{F}$   
 $\theta_2 = 49^\circ\text{F}$   
 $\theta_3 = 45^\circ\text{F}$   
 $\vdots$   
 $\theta_{180} = 60^\circ\text{F}$   
 $\theta_{181} = 56^\circ\text{F}$   
 $\vdots$

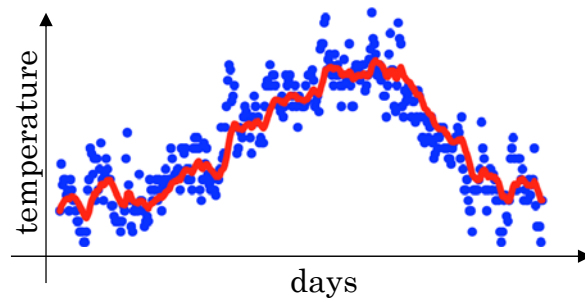


$$\begin{aligned}v_0 &= 0 \\v_1 &= 0.9v_0 + 0.1\theta_1 \\v_2 &= 0.9v_1 + 0.1\theta_2 \\v_3 &= 0.9v_2 + 0.1\theta_3 \\&\vdots \\v_t &= 0.9v_{t-1} + 0.1\theta_t\end{aligned}$$

## Exponentially weighted averages

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

$$\beta = 0.9$$

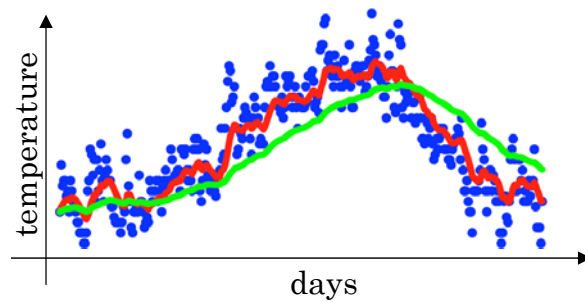


## Exponentially weighted averages

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

$$\beta = 0.9$$

$$\beta = 0.98$$



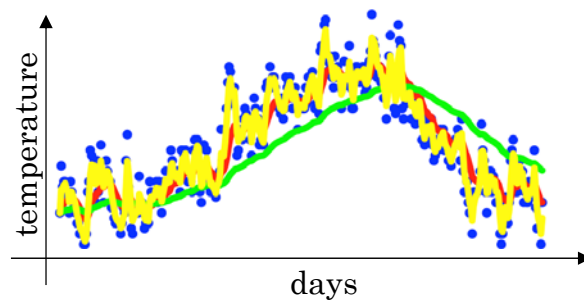
## Exponentially weighted averages

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

$$\beta = 0.9$$

$$\beta = 0.98$$

$$\beta = 0.5$$



## Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1 - \beta)\theta_t \quad \text{When } \beta = 0.9, \text{ what is } v_{100} ?$$

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

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

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

...

## Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1 - \beta)\theta_t \quad \text{When } \beta = 0.9, \text{ what is } v_{100} ?$$

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

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

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

...

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

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



## Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1 - \beta)\theta_t \quad \text{When } \beta = 0.9, \text{ what is } v_{100} ?$$

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

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

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

...

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

## Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1 - \beta)\theta_t \quad \text{When } \beta = 0.9, \text{ what is } v_{100} ?$$

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

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

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

...

$$v_{100} = 0.1\theta_{100} + 0.9(0.1\theta_{99} + 0.9(0.1\theta_{98} + 0.9v_{97}))$$

## Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1 - \beta)\theta_t \quad \text{When } \beta = 0.9, \text{ what is } v_{100} ?$$

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

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

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

...

$$v_{100} = 0.1\theta_{100} + 0.9(0.1\theta_{99} + 0.9(0.1\theta_{98} + 0.9v_{97}))$$



exponentially decaying weight!

$$v_{100} = 0.1\theta_{100} + 0.1 \times 0.9\theta_{99} + 0.1 \times (0.9)^2\theta_{98} + 0.1 \times (0.9)^3\theta_{97} + \dots + 0.1 \times (0.9)^{100}\theta_0$$



## Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t \quad \text{When } \beta = 0.9, \text{ what is } v_n ?$$

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

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

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

...

$$v_{100} = 0.1\theta_{100} + 0.9(0.1\theta_{99} + 0.9(0.1\theta_{98} + 0.9v_{97}))$$

$$v_{100} = 0.1\theta_{100} + 0.1 \times 0.9\theta_{99} + 0.1 \times (0.9)^2\theta_{98} + 0.1 \times (0.9)^3\theta_{97} + \dots + 0.1 \times (0.9)^{100}\theta_0$$



$$v_n = 0.1\theta_n + 0.1 \times 0.9\theta_{n-1} + 0.1 \times (0.9)^2\theta_{n-2} + 0.1 \times (0.9)^3\theta_{n-3} + \dots + 0.1 \times (0.9)^n\theta_0$$

$\therefore$  sum of all weights  $\rightarrow 1$  if  $n \rightarrow \infty$

## Implementing exponentially weighted averages

$$v_0 = 0$$

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

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

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

...

## Implementing exponentially weighted averages

$$v_0 = 0$$

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

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

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

...

$$v_\theta := 0$$

$$v_\theta := \beta v_\theta + (1 - \beta) \theta_1$$

$$v_\theta := \beta v_\theta + (1 - \beta) \theta_2$$

⋮

## Implementing exponentially weighted averages

$$v_0 = 0$$

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

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

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

...

$$v_\theta := 0$$

$$v_\theta := \beta v_\theta + (1 - \beta) \theta_1$$

$$v_\theta := \beta v_\theta + (1 - \beta) \theta_2$$

⋮

---

$$v_\theta := 0$$

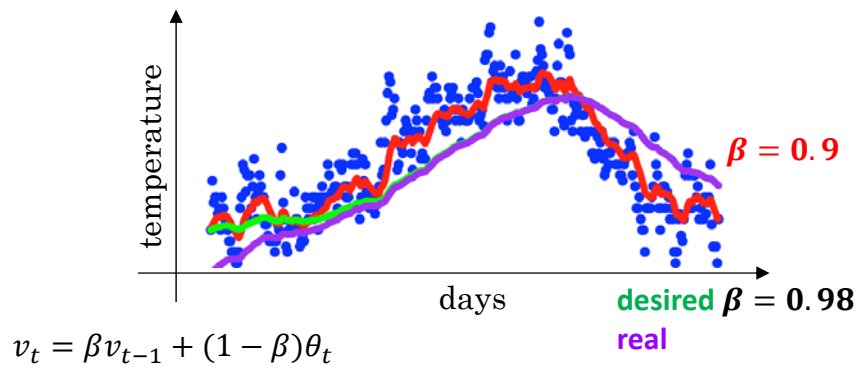
Repeat {

    Get next  $\theta_t$

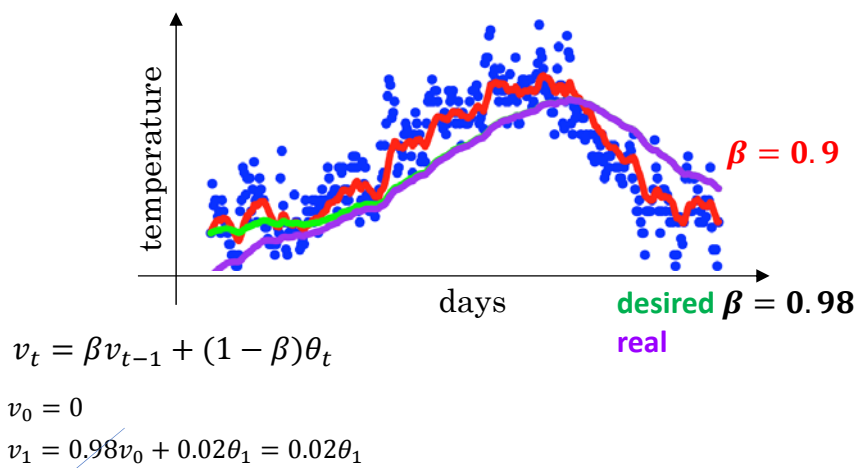
$$v_\theta := \beta v_\theta + (1 - \beta) \theta_t$$

}

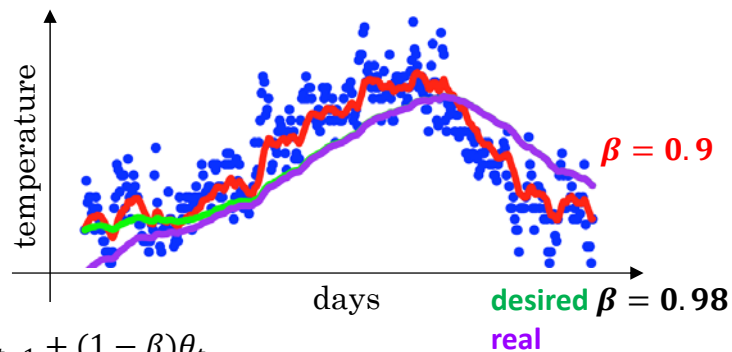
## Bias correction



## Bias correction



## Bias correction



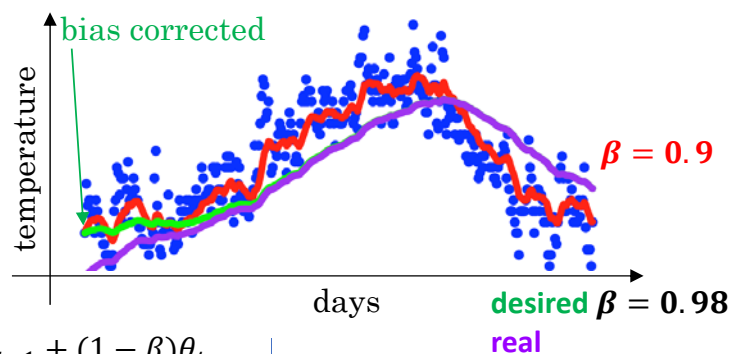
$$v_t = \beta v_{t-1} + (1 - \beta)\theta_t$$

$$v_0 = 0$$

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

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

## Bias correction



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

$$v_0 = 0$$

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

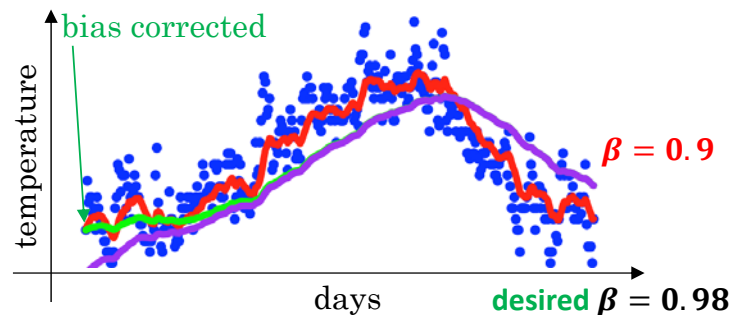
$$\begin{aligned} v_2 &= 0.98v_1 + 0.02\theta_2 \\ &= 0.98 \times 0.02\theta_1 + 0.02\theta_2 \\ &= 0.0196\theta_1 + 0.02\theta_2 \end{aligned}$$

$$v_t \leftarrow \frac{v_t}{1 - \beta^t} \quad \leftarrow \text{bias correction!}$$

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

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

## Bias correction exponentially weighted average!



$$v_t = \frac{\beta v_{t-1} + (1 - \beta) \theta_t}{1 - \beta^t} \leftarrow \text{bias correction!}$$

$$v_n = \frac{(1 - \beta) \theta_n + (1 - \beta) \times \beta \theta_{n-1} + (1 - \beta) \times \beta^2 \theta_{n-2} + \dots + (1 - \beta) \times \beta^n \theta_0}{1 - \beta^{n+1}}$$

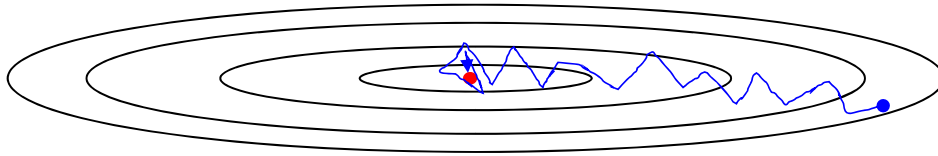
$$v_n = \frac{(1 - \beta) \theta_n + (1 - \beta) \times \beta \theta_{n-1} + (1 - \beta) \times \beta^2 \theta_{n-2} + \dots + (1 - \beta) \times \beta^n \theta_0}{(1 - \beta)(1 + \beta + \beta^2 + \dots + \beta^n)}$$

$$v_n = \frac{\theta_n + \beta \theta_{n-1} + \beta^2 \theta_{n-2} + \dots + \beta^n \theta_0}{1 + \beta + \beta^2 + \dots + \beta^n} = \frac{1 \theta_n + \beta \theta_{n-1} + \beta^2 \theta_{n-2} + \dots + \beta^n \theta_0}{1 + \beta + \beta^2 + \dots + \beta^n}$$

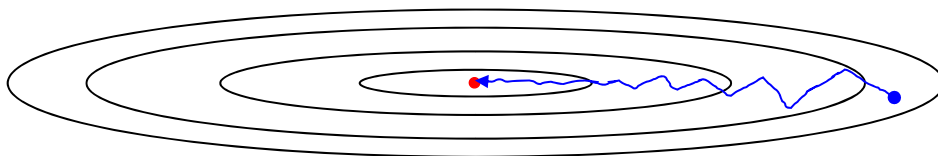
## Optimization algorithms

- Gradient descent with momentum
  - almost always works faster than the standard gradient descent
  - use exponentially weighted average of gradients
- RMSProp
- ADAM

## Gradient descent example



## Gradient descent example



On iteration  $t$  :

Compute  $dW, db$  on current minibatch

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

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

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

## Implementation details

On iteration  $t$ :

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

$$v_{dW} := \beta v_{dW} + (1 - \beta) dW \quad \text{initial } v_{dW} = 0$$

$$v_{db} := \beta v_{db} + (1 - \beta) db \quad \text{initial } v_{db} = 0$$

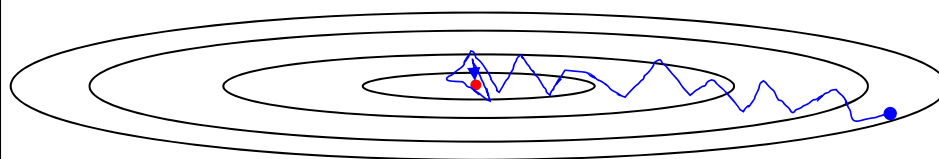
$$W := W - \alpha v_{dW}, \quad b := b - \alpha v_{db}$$

In this case, usually **not use bias correction**,  
because initial gradients are not important.

Hyperparameters:  $\alpha, \beta$        $\beta = 0.9$

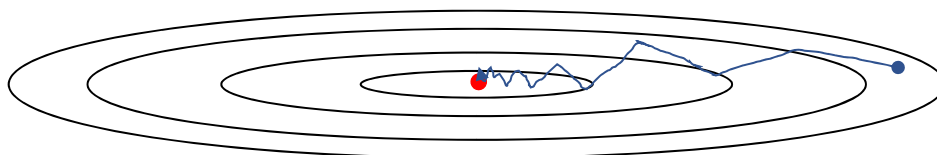
### RMSProp

: root mean square propagation



oscillation of gradient descent

## RMSPProp



On iteration  $t$  :

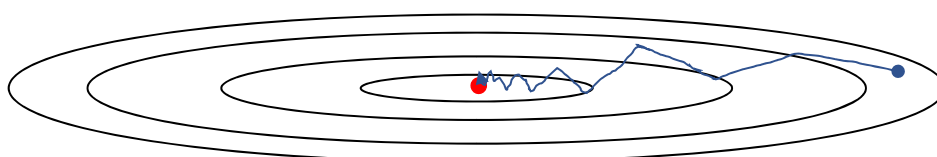
Compute  $dW, db$  on current minibatch

$$S_{dW} = \beta S_{dW} + (1 - \beta) dW^2 \quad \begin{array}{l} \text{element-wise} \\ \text{squaring} \end{array} \quad \text{initial } S_{dW}=0, \text{ initial } S_{db}=0$$

$$S_{db} = \beta S_{db} + (1 - \beta) db^2 \quad \begin{array}{l} \text{element-wise} \\ \text{division by squared root} \end{array}$$

$$W := W - \alpha \frac{dW}{\sqrt{S_{dW}}} \quad b := b - \alpha \frac{db}{\sqrt{S_{db}}}$$

## RMSPProp



On iteration  $t$  :

Compute  $dW, db$  on current minibatch

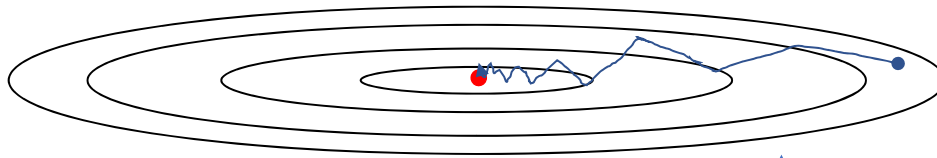
$$S_{dW} = \beta S_{dW} + (1 - \beta) dW^2 \quad \text{Suppose it is 'small',}$$

$$S_{db} = \beta S_{db} + (1 - \beta) db^2 \quad \text{Suppose it is 'large',}$$

$$W := W - \alpha \frac{dW}{\sqrt{S_{dW}}} \quad b := b - \alpha \frac{db}{\sqrt{S_{db}}}$$



# RMSProp



On iteration  $t$  :

Compute  $dW, db$  on current minibatch

$$S_{dW} = \beta S_{dW} + (1 - \beta) dW^2 \quad \text{Suppose it is 'small',}$$

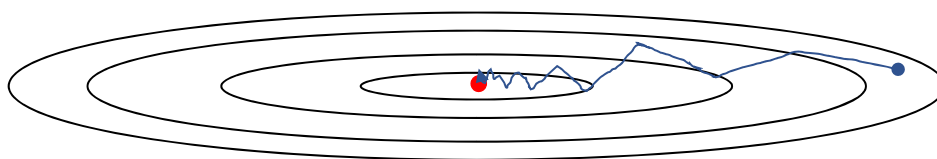
$$S_{db} = \beta S_{db} + (1 - \beta) db^2 \quad \text{Suppose it is 'large',}$$

$$W := W - \alpha \frac{dW}{\sqrt{S_{dW}}}$$

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

then 'large' (speed-up)   then 'small' (slow-down)

# RMSProp



On iteration  $t$  :

Compute  $dW, db$  on current minibatch

$$S_{dW} = \beta S_{dW} + (1 - \beta) dW^2 \quad \text{Suppose it is 'small',}$$

$$S_{db} = \beta S_{db} + (1 - \beta) db^2 \quad \text{Suppose it is 'large',} \quad \epsilon = 10^{-8} \text{ for numerical stability}$$

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

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

then 'large' (speed-up)   then 'small' (slow-down)

## Adam optimization algorithm

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

On iteration  $t$  :

Compute  $dW, db$  on current minibatch

$$v_{dW} = \beta_1 v_{dW} + (1 - \beta_1) dW, \quad v_{db} = \beta_1 v_{db} + (1 - \beta_1) db$$

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

## Adam optimization algorithm

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

On iteration  $t$  :

Compute  $dW, db$  on current minibatch

"momentum" with  $\beta_1$

$$v_{dW} = \beta_1 v_{dW} + (1 - \beta_1) dW, \quad v_{db} = \beta_1 v_{db} + (1 - \beta_1) db$$

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

"RMSProp" with  $\beta_2$

# Adam optimization algorithm

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

On iteration t :

Compute  $dW, db$  on current minibatch

"momentum" with  $\beta_1$

$$v_{dW} = \beta_1 v_{dW} + (1 - \beta_1) dW, \quad v_{db} = \beta_1 v_{db} + (1 - \beta_1) db$$

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

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

"RMSProp" with  $\beta_2$

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

bias correction

# Adam optimization algorithm

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

On iteration t :

Compute  $dW, db$  on current minibatch

"momentum" with  $\beta_1$

$$v_{dW} = \beta_1 v_{dW} + (1 - \beta_1) dW, \quad v_{db} = \beta_1 v_{db} + (1 - \beta_1) db$$

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

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

"RMSProp" with  $\beta_2$

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

$$W := W - \alpha \frac{v_{dW}^{\text{corrected}}}{\sqrt{S_{dW}^{\text{corrected}} + \epsilon}} \quad b := b - \alpha \frac{v_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}} + \epsilon}}$$

## Hyperparameters choice

$\alpha$ : needs to be tuned

$\beta_1$ : 0.9 (weighted average for  $dW$ )

$\beta_2$ : 0.999 (weighted average for  $dW^2$ )

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

Adam? Adaptive momentum estimation

## Hyperparameters choice



Adam Coates

$\alpha$ : needs to be tuned

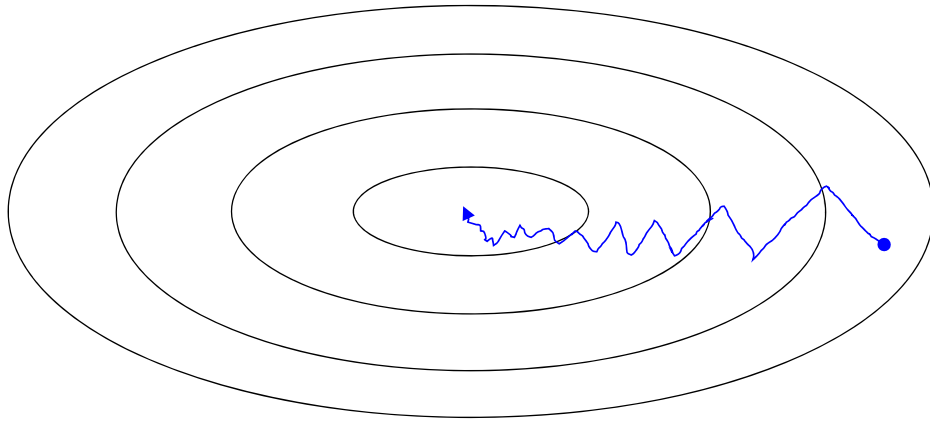
$\beta_1$ : 0.9 (weighted average for  $dW$ )

$\beta_2$ : 0.999 (weighted average for  $dW^2$ )

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

Adam? Adaptive momentum estimation

## Learning rate decay



slowly reduce the learning rate over time

## Implementing learning rate decay

1 epoch = 1 pass through data



## Implementing learning rate decay

1 epoch = 1 pass through data



$$\alpha = \frac{1}{1 + \text{decay\_rate} \times \text{epoch\_num}} \alpha_0$$

## Implementing learning rate decay

1 epoch = 1 pass through data

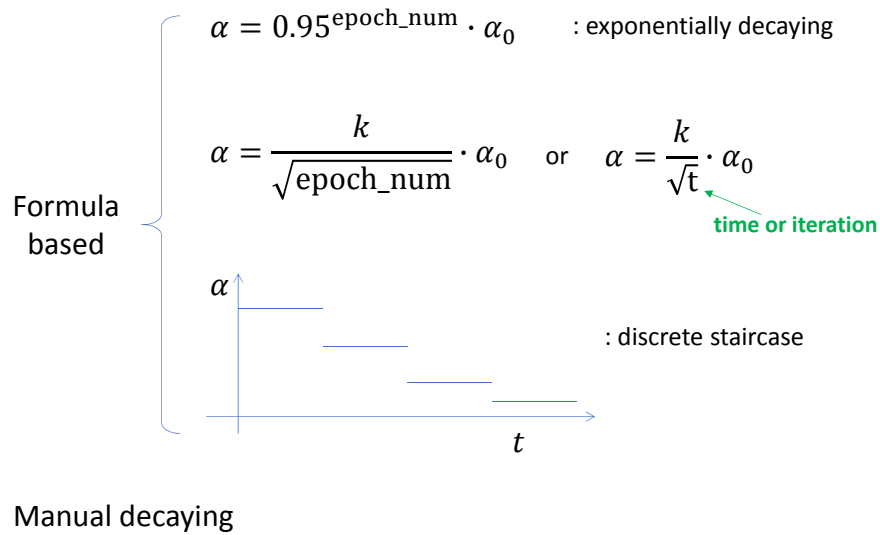


$$\alpha = \frac{1}{1 + \text{decay\_rate} \times \text{epoch\_num}} \alpha_0$$

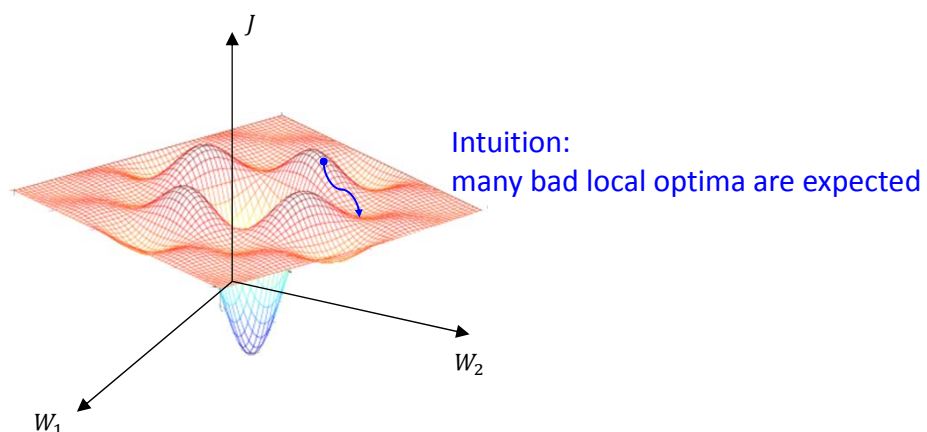
$\alpha_0 = 0.2$   
 $\text{decay\_rate} = 1$

Epoch	$\alpha$
1	0.1
2	0.067
3	0.05
4	0.04
$\vdots$	$\vdots$

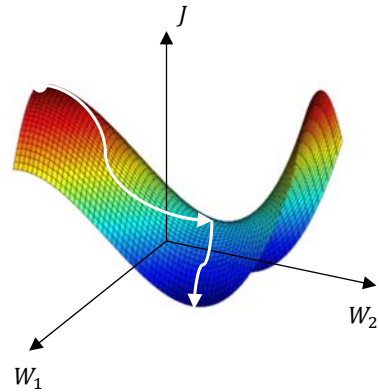
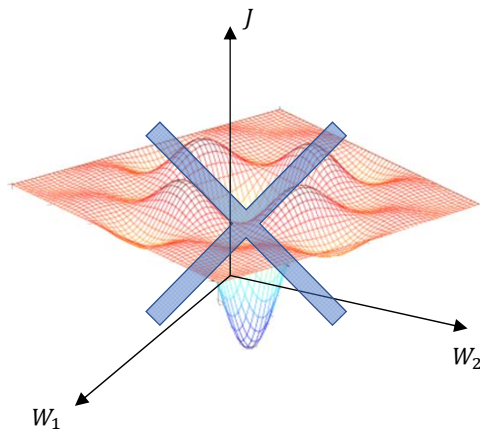
## Other learning rate decay methods



## Local optima in neural networks

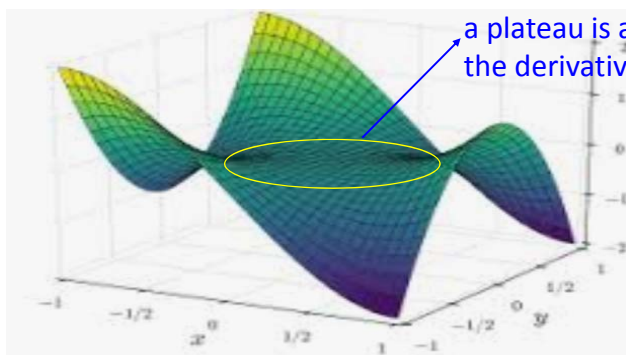


## Local optima in neural networks



Real: many saddle points exists

## Problem of plateaus



- Unlikely to get stuck in a bad local optima
- Plateaus can make learning slow
- momentum, rmsprop, adam optimization algorithm may be very helpful



- End -