

$x^{(i)}$   $\rightarrow$   $i$ th training example  
 $z^{[l]}$   $\rightarrow$  value of  $z$  for layer  $l$   
 $x^{(t)}, y^{(t)}$   $\rightarrow$  mini-batch  $t$

## Week 2 : Optimization Algorithms

Mini-batch Gradient descent

$$X = [x^{(1)} \ x^{(2)} \ \dots \ x^{(m)}] \rightarrow (n, m)$$

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}] \rightarrow (1, m)$$

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

Split training sets into smaller parts  $\rightarrow$  (mini-batches)

minibatches each of 1000 examples

$$X = [ \underbrace{x^{(1)} \ x^{(2)} \ \dots \ x^{(1000)}}_{X^{(1)} \rightarrow (n, 1000)} \mid \underbrace{x^{(1001)} \ \dots \ x^{(2000)}}_{X^{(2)} \rightarrow (n, 1000)} \mid \dots \mid \underbrace{\dots \ x^{(m)}}_{X^{(3000)} \rightarrow (n, 1000)} ]$$

$$Y = [ \underbrace{y^{(1)} \ y^{(2)} \ \dots \ y^{(1000)}}_{Y^{(1)} \rightarrow (1, 1000)} \mid \underbrace{y^{(1001)} \ \dots \ y^{(2000)}}_{Y^{(2)} \rightarrow (1, 1000)} \mid \dots \mid \underbrace{\dots \ y^{(m)}}_{Y^{(3000)} \rightarrow (1, 1000)} ]$$

Minibatch  $t$  :  $X^{(t)}, Y^{(t)}$

$\downarrow$  check assignment code for creating mini-batch

Batch vs minibatch gradient descent



all  $(X, Y)$   
at once are  
processed



$X^{(t)}, Y^{(t)}$  are  
processed at once

#  $S$  = no. of samples in the mini-batch

Page No.

Date

Mini-batch gradient descent (algorithm)

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

forward prop. on  $x^{\{t\}}$  ✓

$$z^{[1]} = w^{[1]} x^{\{t\}} + b^{[1]}$$

$$a^{[1]} = g^{[1]}(z^{[1]})$$

⋮

$$a^{[L]} = g^{[L]}(z^{[L]})$$

} this vectorized implementation processes 1000 examples rather than whole training set

Compute cost  $J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^S L(\hat{y}^{(i)}, y^{(i)}) +$

for  $x^{\{t\}}, y^{\{t\}}$  ↑ ↑

$$\frac{\lambda}{2 \cdot 1000} \sum_i \|w^{[1]}\|_F^2$$

Back prop to compute gradients w.r.t  $J^{\{t\}}$  ↘

(using  $x^{\{t\}}, y^{\{t\}}$ )

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

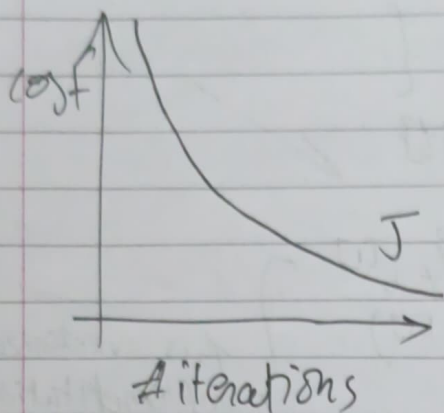
}

"1 epoch"

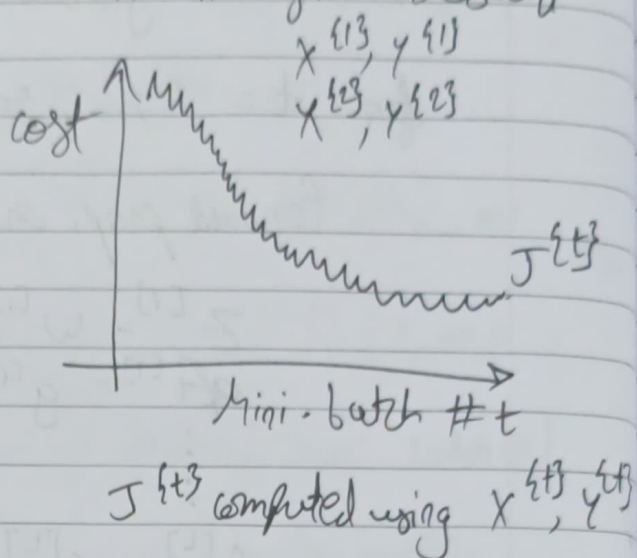
↪ single pass through training set

When you have a large training set, a mini-batch algorithm is preferred as it would work faster than the regular algo.

Batch grad. descent



Mini-batch grad. descent



The graph in mini-batch gets noisy because we are computing  $J^{(t)}$  for different different mini-batches i.e. for different  $x^{(t)}, y^{(t)}$  in each iteration so that is why some  $J^{(t)}$ 's might go up and some might go down, but overall trend should be seen as decreasing only.

(choosing mini-batch size (hyperparameter))

If mini-batch size =  $m$ : Batch grad. descent  $(x^{(1)}, y^{(1)})$   
 $= (x, y)$

If mini-batch size = 1: Stochastic grad. descent  
 $\rightarrow$  every example is its own mini batch.

$(x^{(1)}, y^{(1)}) = (x^{(1)}, y^{(1)}) \Leftarrow, (x^{(2)}, y^{(2)}) = (x^{(2)}, y^{(2)})$

In practise: mini-batch size is between 1 and  $m$ .



Stochastic grad  
descentIn-between  
(mini-batch size)Batch grad  
descent

mini-batch size = 1

lose speed up from  
Vectorizationbecause processing one em. w/o training  
at a time is inefficient. entire set

Fastest learning:

- vectorization
- make progress

mini-batch size = m

Too long to  
process per iterationIf small training set: Use batch grad. descent  
→ ( $m \leq 2000$ )

Typical mini-batch sizes:

$$64, 128, 256, 512$$

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

Also make sure mini-batch size fits in CPU/GPU memory

Other optimization Algos: Exponentially weighted  
averages

→ Faster than grad. descent

em:  $\theta_1 = 4^\circ\text{C}$

$\theta_2 = 9^\circ\text{C}$

$\theta_3 = \vdots$

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

 $\vdots$ 

$V_0 = a$

$V_1 = 0.9V_0 + 0.1\theta_1$

$V_2 = 0.9V_1 + 0.1\theta_2$

$V_3 = 0.9V_2 + 0.1\theta_3$

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

$$V_t = \beta V_{t-1} + (1-\beta) \Theta_t, \quad \beta = 0.9$$

$V_t$  can be inferred as approximately averaging temperature over  $\approx \frac{1}{1-\beta}$  days.

$\beta = 0.9 : \approx 10$  days' average temp.

$\beta = 0.98 : \approx 50$  days' average  
 → when plotted, the curve gets more smoother and it shifts more to the right.

$\beta = 0.5 : \approx 2$  days → more noisy curve

$$V_t = \beta V_{t-1} + (1-\beta) \Theta_t, \quad \beta = 0.9$$

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

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

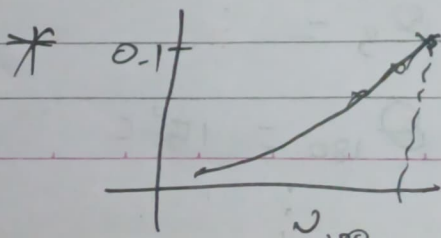
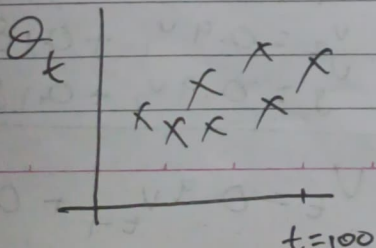
$$V_{98} = 0.9 V_{97} + 0.1 \Theta_{98}$$

⋮

$$\rightarrow V_{100} = 0.1 \Theta_{100} + 0.9 \cancel{V_{99}} (0.1 \Theta_{99} + 0.9 \cancel{V_{98}})$$

$$0.1 \Theta_{99} + 0.9 V_{97}$$

$$= 0.1 \Theta_{100} + 0.1 \times 0.9 \times \Theta_{99} + 0.1 (0.9)^2 \Theta_{98} + \dots$$



$$0.9^{10} \approx 0.35 \approx \frac{1}{e} \quad (1-\epsilon)^{1/\epsilon} = \frac{1}{e}$$

Implementing exponentially weighted averages

$$\begin{aligned} v_0 &= 0 \\ v_1 &= \beta v_0 + (1-\beta)\theta_1 \\ v_2 &= \beta v_1 + (1-\beta)\theta_2 \\ &\vdots \end{aligned}$$

$$\begin{aligned} v_0 &= 0 \\ v_{10} &= \beta v + (1-\beta)\theta_1 \\ v_0 &:= \beta v + (1-\beta)\theta_2 \\ &\vdots \end{aligned}$$

$$v_0 = 0$$

Repeat {

get next  $\theta_t$

$$v_0 := \beta v_0 + (1-\beta)\theta_t$$

\* Very efficient as it takes only one line of code and less memory (as we are overwriting the same variable), still not the best way to find ~~average~~ to compute average.

Bias correction in exponentially weighted averages

$$v_t = \beta v_{t-1} + (1-\beta)\theta_t, \quad \beta = 0.98$$

$$v_0 = 0$$

$$v_1 = 0.98 \overset{0}{v_0} + 0.02\theta_1, \quad \beta = 0.98$$

If  $\theta_1 = 40^\circ\text{C}$ ,  $v_1 = 0.02 \times 40 = 0.8$  which is not a good estimate on first day's temperature.



$$V_2 = 0.98 V_1 + 0.02 \Theta_2$$

$$= 0.98 (0.02) \times \Theta_1 + 0.02 \Theta_2$$

$$= 0.0196 \Theta_1 + 0.02 \Theta_2$$

→ Even  $V_2$  will not be a very good estimate of first 2 days' temperature.

To solve this:

→ insted of ~~but~~ considering  $V_t$ , take  $\frac{V_t}{1-\beta^t}$

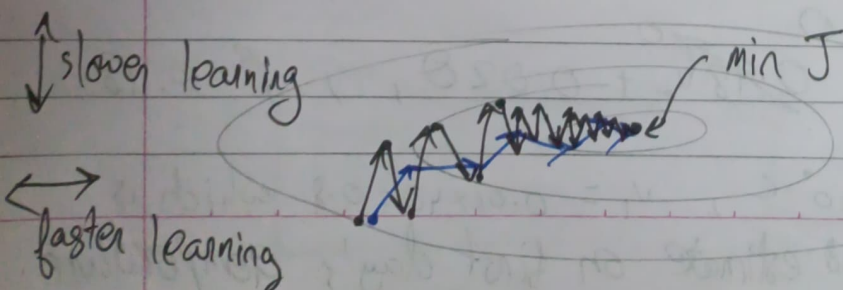
eg:  $t=2$  :  $1-\beta^t = 1-(0.98)^2 = 0.0396$

$$\frac{V_2}{0.0396} = \frac{0.0196 \Theta_1 + 0.02 \Theta_2}{0.0396}$$

so as seen above, this becomes a weighted average of  $\Theta_1$  &  $\Theta_2$  ( $\because 0.0196 + 0.02 = 0.0396$ ) and this removes the bias.

so as  $t$  becomes large,  $\frac{1}{1-\beta^t}$  becomes  $\beta^t$  will approach 0, so bias correction will make no difference.

- Gradient descent with Momentum:



- mini-batch /  
- ~~stochastic~~ stochastic  
grad desc.

- With momentum

#  $v['dw' + \text{str}(i)] = \text{np.zeros}((\text{parameter}['w' + \text{str}(i)]).shape)$

Page No.

Date

The up and down oscillations slows down gradient descent and prevents us from using larger learning rates.

If learning rate is large:

it ends up diverging

Momentum:  $\beta$

On iteration  $t$ :

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

friction ( $< 1$ )

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

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

velocity → acceleration

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

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

Implementation details:

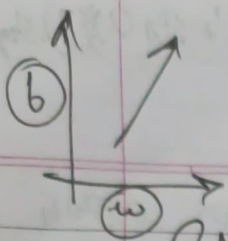
$$\# v_{db} = 0, v_{dw} = 0 \rightarrow \text{shape}(dw) = \text{shape}(w)$$

Hyperparameters:  $\alpha, \beta$

$\beta = 0.9$  (robust value) / common value  
i.e. average of ~~temp~~ days  $\approx 10$

⚡ No need to use bias correction for  $v_{dw}$  &  $v_{db}$





$\Rightarrow db$  is large  
 $dw$  is small

$\updownarrow$  slow  
 $\longleftrightarrow$  fast

Page No.

Date

RMS prop: Root Mean Square prop

On iteration  $t$ :

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

small  $\rightarrow S_{dw} = \beta S_{dw} + (1 - \beta) dw^2$   
 $\because dw$  is small  
 $\#$  "dw"  
 $\#$  keeps the exponentially weighted average of the squares of the derivatives.  
 $\uparrow$  element-wise squaring

large  $\rightarrow S_{db} = \beta S_{db} + (1 - \beta) db^2$   
 $\because db$  is large  
 $\#$

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

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

Similar effects as seen in momentum as, vertical oscillations damp out and it moves faster horizontally, so we can also increase the learning rate to further make the process and computations faster without any issue.

# Adam optimization algorithm:

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

On iteration  $t$ :

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

$$\left. \begin{aligned} V_{dw} &= \beta_1 V_{dw} + (1 - \beta_1) dw \\ V_{db} &= \beta_1 V_{db} + (1 - \beta_1) db \end{aligned} \right\} \begin{array}{l} \text{Momentum} \\ \beta_1 \end{array}$$

$$\left. \begin{aligned} S_{dw} &= \beta_2 S_{dw} + (1 - \beta_2) dw^2 \\ S_{db} &= \beta_2 S_{db} + (1 - \beta_2) db^2 \end{aligned} \right\} \begin{array}{l} \text{RMSprop} \\ \beta_2 \end{array}$$

In Adam's optimization, we include correction factors

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

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

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

→ Very effective for wide variety of architectures

Hyperparameters choice:

$\alpha$  : needs to be tuned

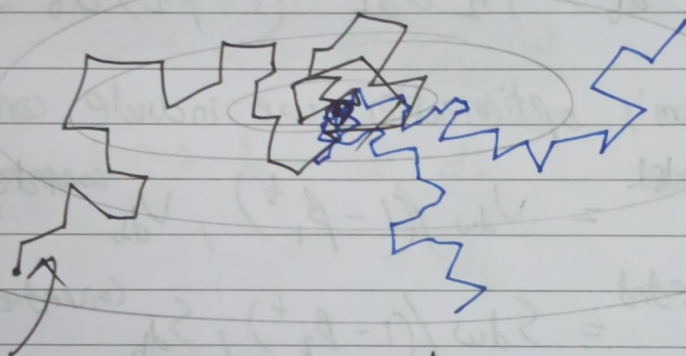
$\beta_1 : 0.9$  ( $dw$ )  $\rightarrow$  mean of the derivatives

$\beta_2 : 0.999$  ( $dw^2$ )  $\rightarrow$  compute exponentially weighted average of squares

$\epsilon : 10^{-8}$

Adam : Adaptive moment estimation

Learning rate decay :

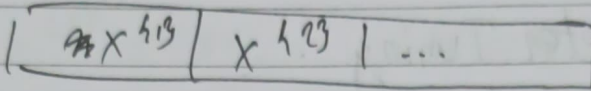


This does not exactly converge at the minimum hence it wanders around the minimum  $J$  but never really reaches there.

By slowly reducing  $\alpha$ , we will take smaller and smaller steps and as the min  $J$  comes closer, it will finally converge with smaller steps and not wander around.



1 epoch = 1 pass through the data



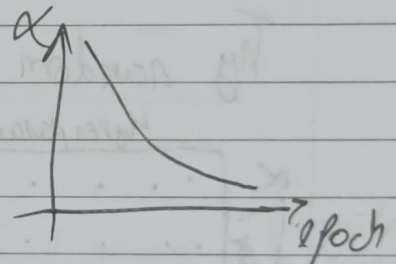
epoch 1  
epoch 2 (2<sup>nd</sup> pass)

$$\alpha = \frac{\alpha_0}{1 + \text{decayrate} * \text{epoch-number}}$$

example:

$$\alpha_0 = 0.1, \text{ decay Rate} = 0.1$$

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



Other learning rate decay methods:

1)  $\alpha = 0.95^{\text{epoch-num}} \cdot \alpha_0 \rightarrow$  exponentially decay the value of  $\alpha$

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

