

## Week-2 Mini-batch gradient descent

for  $t = 1, \dots, 5000$

Forward prop on  $x^{[t]}$

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

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

:

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

} Vectorized implementation

for  $x^{[t]}, y^{[t]}$

$$\text{Compute cost } J^{[t]} = \frac{1}{1000} \sum_{i=1}^n (g^{[i]}, y^{[i]}) + \frac{\lambda}{2 \times 1000} \sum_{l} \|w^{[l]}\|_F^2$$

Backprop to compute gradients

$$w^{[l]} = w^{[l]} - \alpha \delta w^{[l]}, \text{ by for } b.$$

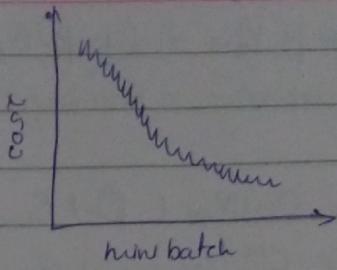
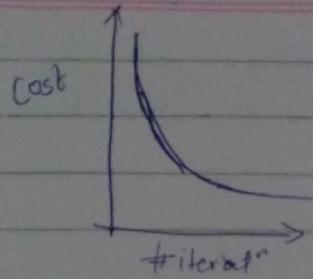
}

Here, 1 epoch will indicate 1 gradient descent step.

So, instead of waiting for 8 million training examples, we

so, 1 epoch of training set = 1 gradient descent step

$$1 \text{ " } \text{ " mini" } \text{ " } s = 5000 \text{ " } \text{ " } \text{ " } s$$



minibatch size

$m$

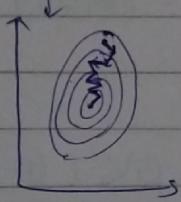
Name

Batch gradient descent  
stochastic "

$$(x^{(1)}, y^{(1)}) = (x, y)$$

$$(x^{(1)}, y^{(1)}) = (x^{(1)}, y^{(1)}),$$

every e.g. is  
a minibatch



In practice, stochastic gd will take more time & as each entry will need to be calculated. Thus, "vectorizat" meaning.

minibatch size  $\neq 1 \neq m \Rightarrow$  fastest learning.

- How to choose mini-batch size
- small training set : Use batch gradient descent  
m.s 2020
- Typical minibatch size:  $2^6, 2^7, 2^8$  ...
- Make sure mini batch fits the CPU/GPU memory.

In practice, minibatch size is one, <sup>hyper</sup>parameter to adjust for better performance

Exponentially weighted Average

- For cases when the data is moving. temp. v/s days. Average (trendy) sh needs to be calculated.
- to build more sophisticated optimizat" algorithms

We take weighted avg.

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

approx. average over  $\frac{1}{1-\beta}$  days

$$V_{100} = 0.9 V_{99} + 0.1 O_{100}$$

$$V_{99} = 0.9 V_{98} + 0.1 O_{99}$$

$$V_{98} = 0.9 V_{97} + 0.1 O_{98}$$

$$\hookrightarrow V_{100} = 0.9 (0.9 V_{97} + 0.1 O_{98}) + 0.1 O_{100}$$

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

$$V_0 = 0$$

$$V_1 = \beta V_0 + (1-\beta) O_1$$

$$V_2 = \beta V_1 + (1-\beta) O_2$$

$$V_3 = \beta V_2 + (1-\beta) O_3$$

$$V_O = 0$$

$$V_O := \beta V_0 + (1-\beta) O_0$$

$$V_O := \beta V_O + (1-\beta) O_2$$

?

vectorized

## Bias Correction

To compute weighted average efficiently in  
 Instead of  $V_t$ , if we In the previous formula, the initial stage  
 the graph will start with a bit smaller value.  
 So, to correct it, we can use  $\frac{V_t}{1-\beta^t}$  instead of  $V_t$   
 for  $t=2$ ,

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

$$V_2 = \beta \left( \frac{V_1}{1-\beta^t} \right) + (1-\beta) O_2$$

$$= 0.0196 O_1 + 0.02 O_2$$

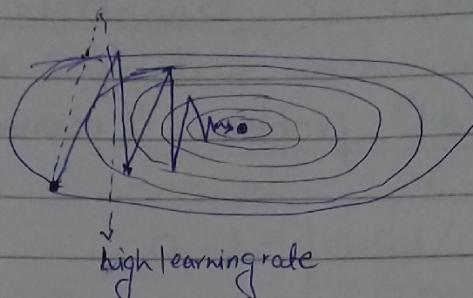
$$= 0.0396$$

$$V_1 = 0.98 \overset{\circ}{V}_0 + 0.02 O_0$$

## Gradient Descent with Momentum

As shown in figure, the gradient descent oscillates on the contour in order to find minima.

Also, this prevents from using larger  $\alpha$ , as the the oscillation gradient descent will go out of the contour



What we want:

$\downarrow$  slow learning vertically  
 $\leftarrow$  fast " horizontally

$$V_{dw} = \beta V_{dw} + (1-\beta) dw, \quad \text{from } (V_0 = \beta V_0 + (1-\beta) V_{0-1})$$

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

$$W := W - \alpha V_{dw}$$

instead of  $dw$ , we take weighted average.

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

Thus, the oscillations will become more smooth

Now, we have two hyperparameters  $\alpha, \beta$

Generally  $\beta=0.9$ , nobody changes it, it is a robust value.

In practice, many people do not use  $(1-\beta)$  term, but Andrew prefers to use it. It may happen that, formula with and w/o  $1-\beta$  terms can need different  $\alpha$

### RMS prop

Like momentum, it is also used to speed up gradient descent for iteration  $t$ ,

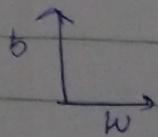
Compute  $dw, db$  on current mini-batch

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

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

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

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



Assume that  $S_{dw}$  is very small, that is change in  $w$  direction is small, thus,  $w$  won't update much. If  $S_{dw}$  is very large, in that case  $\sqrt{S_{dw}}$  will be large, and hence more change can be expected. Instead of  $w$  &  $b$ , vectors there might be high dimensional vectors, that can change be modified, based on avg change. Thus, root mean square prop!

Hence, now we can have large learning rate to speed up gradient descent

DisAdvantage, what if  $S_{dw} = 0$  ??

### Adam Optimizat<sup>n</sup> problem (momentum + RMS prop)

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

On iteration t:

Compute  $dw, db$  using correct <sup>mini-batch</sup> ~~implemented~~

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

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

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

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

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

choice of hyperparameters' values:

$\alpha$  : needs to be tuned

$$\beta_1 : 0.9$$

$$\beta_2 : 0.999$$

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

Adam = Adaptive Moment Estimation

Learning rate Decay

For a fixed value of alpha, the convergence to minima won't take place.

Incase, if the value of alpha is gradually decreasing, the convergen won't be on minima still, but it will oscillate to a more lighter space

So, by reducing alpha, we can take big steps initially and slowly take small steps

1 epoch = 1 pass through decay rate

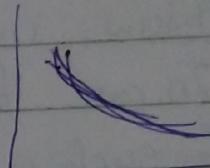
$$\alpha = \frac{1}{1 + \text{decay rate} * \text{epoch\_num}} \alpha_0$$

$$\alpha = 0.2$$

$$\text{decay rate} = 1$$

here, decay  
rate is another  
hyper-parameter

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



Exponential decay

$$\alpha = 0.95^{\text{epoch\_num}} \times \alpha_0$$

discrete staircase decay

$$\alpha = \frac{k}{\text{epoch\_num}} \times \alpha_0 \quad \text{or} \quad \frac{k \cdot \alpha_0}{\sqrt{t}} \quad \left[ \begin{array}{l} t \leftarrow \text{iteratn no.} \\ z \end{array} \right] - \frac{z}{z}$$

manually decay after observation