

dec1: Mini batch Gradient descent:

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}_{(n_x, m)} \quad Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix}_{(1, m)}$$

If $m = 5,000,000$, the G.D is very slow.

sd. mini batches of 1000 each

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(1000)} & \dots & x^{(1000)} & \dots & x^{(m)} \end{bmatrix}_{(n_x, 1000)} \quad | \quad \dots \quad | \quad \dots \quad |$$

$\underbrace{x^{(1)} \dots x^{(1000)}}_{(n_x, 1000)}$ $\underbrace{x^{(1001)} \dots x^{(1000)}}_{(n_x, 1000)}$ $\underbrace{x^{(1000)} \dots x^{(m)}}_{(n_x, 1000)}$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(1000)} & \dots & y^{(1000)} & \dots & y^{(m)} \end{bmatrix}_{(1, 1000)} \quad | \quad \dots \quad | \quad \dots \quad |$$

$\underbrace{y^{(1)} \dots y^{(1000)}}_{(1000, 1)}$ $\underbrace{y^{(1001)} \dots y^{(1000)}}_{(1000, 1)}$ $\underbrace{y^{(1000)} \dots y^{(m)}}_{(1000, 1)}$

minibatch t : $x^{t\frac{1}{2}}, y^{t\frac{1}{2}}$

$$\begin{bmatrix} x^{(1)} & y^{(1)} \end{bmatrix}_{(n_x, 1000)} \quad \begin{bmatrix} y^{(1)} \end{bmatrix}_{(1000, 1)}$$

Algorithm:- (Mini Batch GD)

implement 1 step of GD using x^{st}, y^{st}

for t = 1, ..., 5000:

Forward Prop on x^{st}

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

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

Vectorized implementation (1000 examples)

compute cost function

$$J = \frac{1}{1000} \sum_{i=1}^m d(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2(m=1000)} \sum \|w^{[l]}\|_F^2$$

Backprop to compute gradient w.r.t J^{st} (using x^{st}, y^{st})

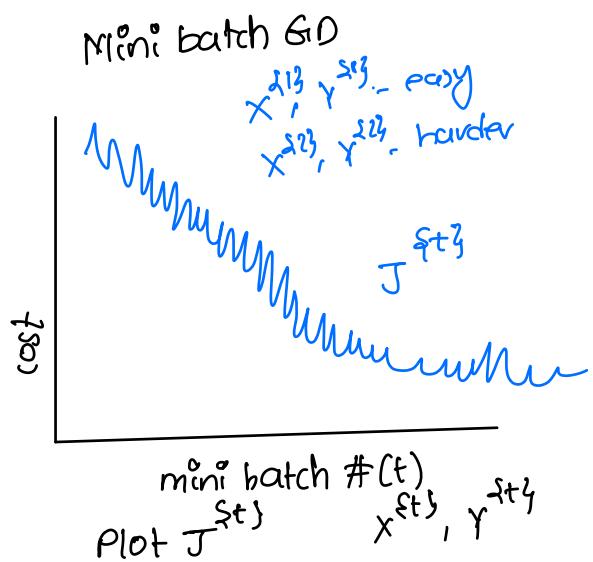
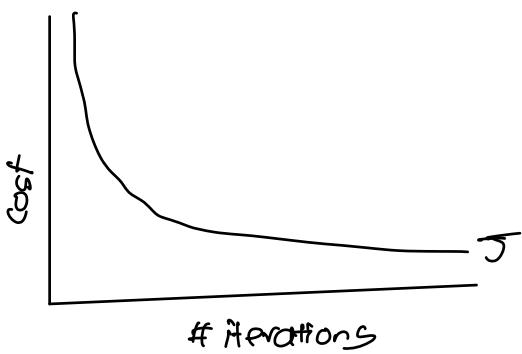
$$w^{[l]} = w^{[l]} - \alpha d w^{[l]}$$

$$b^{[l]} = b^{[l]} - \alpha d b^{[l]}$$

y

dec: understanding mini batch GD:-

Batch GD



choose mini batch size

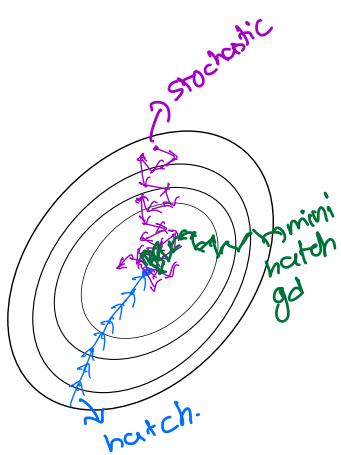
if mini batch size = $m \Rightarrow$ batch gradient descent.

$$(x^{s1}, y^{s1}) = (x, y)$$

if mini batch size = 1 \Rightarrow stochastic gradient descent

$$(x^{s1}, y^{s1}) = (x^{(1)}, y^{(1)})$$

$$(x^{s2}, y^{s2}) = (x^{(2)}, y^{(2)})$$



In practice mini batch size in between 1 and m

batch gd \rightarrow too long per iteration

stochastic gd \rightarrow worse speed from vectorization.

minibatch gd \rightarrow fast learning, vectorisation (1000)

\rightarrow Make progress without needing to
wait to iterate entire training set

each epoch will update 5000 times w

if you have small training set \rightarrow use batch gradient descent
($m \leq 2000$)

typical minibatch sizes:

64, 128, 256, 512, 1024

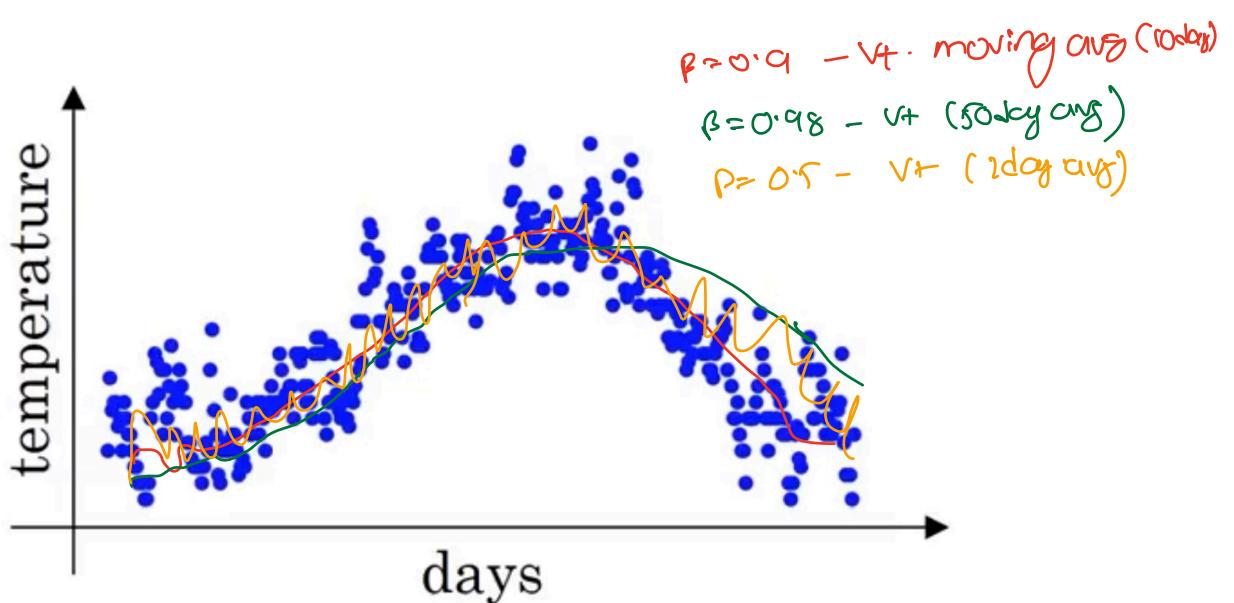
$2^6, 2^7, 2^8, 2^9, 2^{10}$

Make sure minibatch fits in CPU/GPU memory (x^{s1}, y^{s1})

dec3. Exponentially Weighted Averages:-

temperature in London

$$\theta_1 = 40^\circ\text{F} \quad \theta_2 = 40^\circ\text{F} \quad \theta_3 = 45^\circ\text{F} \dots \theta_{180} = 60^\circ\text{F}$$



Moving average (exponentially moving average)

$$v_0 = 0$$

$$v_1 = 0.9 v_0 + 0.1 \theta_1$$

$$v_2 = 0.9 v_1 + 0.1 \theta_2$$

$$v_3 = 0.9 v_2 + 0.1 \theta_3$$

⋮

$$v_t = 0.9 v_{t-1} + 0.1 \theta_t$$

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

$\beta : 0.9 \approx 10 \text{ day avg temp}$

v_t = approximately averaging over $\approx \frac{1}{1-\beta}$ days temperature

dec9: understanding moving avg

$$v_t = \beta v_{t-1} + (1-\beta) \theta_t \quad \beta=0.9$$

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

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

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

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



$$0.9^{10} \approx 0.35 \approx \frac{1}{e} \quad (10 \text{ days avg})$$

$$(0.9)^{10} = (1-\varepsilon)^{\frac{10}{\varepsilon}} = \frac{1}{e}$$

if $\beta=0.98$ $(0.98)^{50} = \frac{1}{e}$ (50 days avg).

Implementation

$$v_0 = 0$$

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

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

:

:

:

$$v_{\theta} = 0$$

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

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

$$v_{\theta} = 0$$

Repeat $\Leftarrow \theta_t$

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

}

decs: Bias correction in exponential Avg:

$$\Theta_1 = 40$$

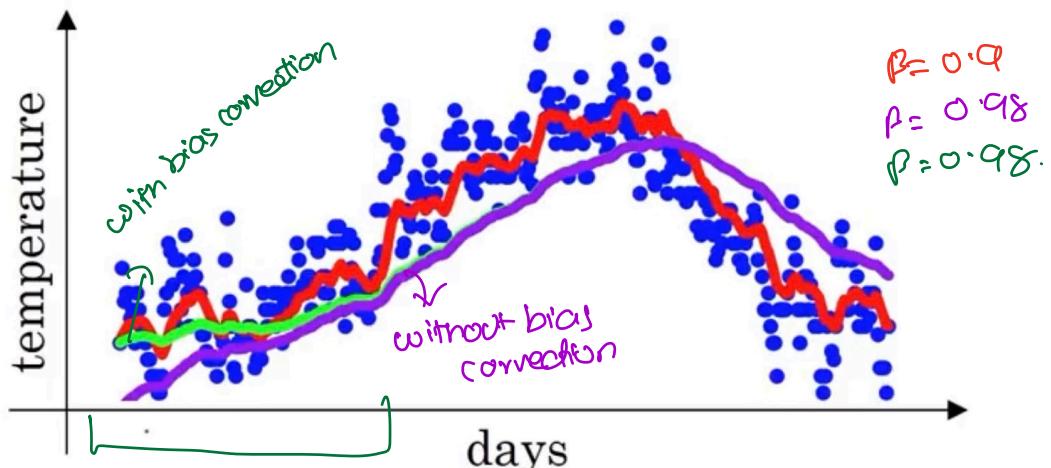
$$V_t = \beta V_{t-1} + (1-\beta) \Theta_t$$

$$V_0 = 0$$

$$V_1 = 0.98 V_0 + 0.02 \Theta_1 = 0.02 \Theta_1 = 0.8$$

$$V_2 = 0.98 V_1 + 0.02 \Theta_2 = (0.98)^2 \Theta_1 + 0.02 \Theta_2$$

very less

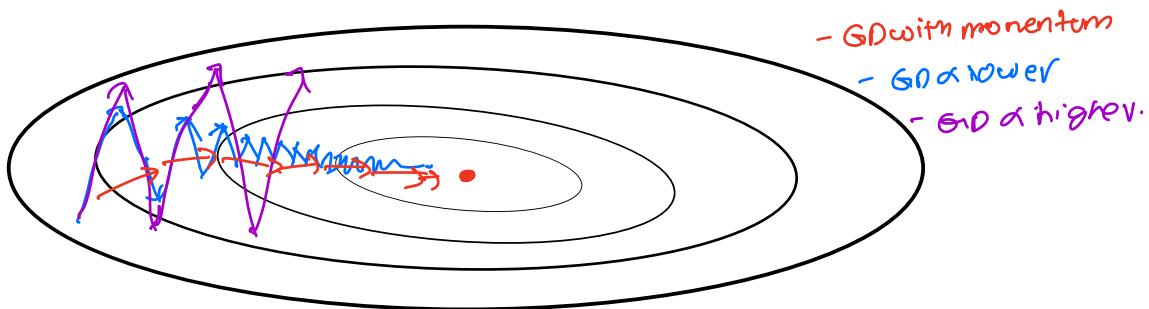


$$\frac{V_t}{1-\beta^t} \quad t \text{ large} \approx V_t$$

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

$$\frac{V_2}{0.0396} = \frac{0.0196 \Theta_1 + 0.02 \Theta_2}{0.0396} = \text{removes the bias.}$$

lec 6: Gradient Descent with Momentum :-



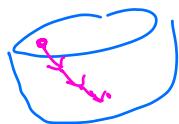
vertical axis - slower learning

horizontal axis - faster learning

Momentum:

on iteration t

compute $d\omega, db$ on current mini batch.



$$vd\omega = \beta vd\omega + (1-\beta) d\omega$$

$$v\theta = \beta v\theta + (1-\beta) \theta t$$

$$vd\theta = \beta vd\theta + (1-\beta) d\theta$$

$$\begin{aligned} \omega &= \omega - \alpha vd\omega \\ b &= b - \alpha vd\theta \end{aligned}$$

smooth steps of gd

Implementation:

on iteration t

compute $d\omega, db$ on current mini batch

$$vd\omega = \beta vd\omega + (1-\beta) d\omega$$

Hyper parameter

$$\alpha, \beta$$

$$vd\theta = \beta vd\theta + (1-\beta) d\theta$$

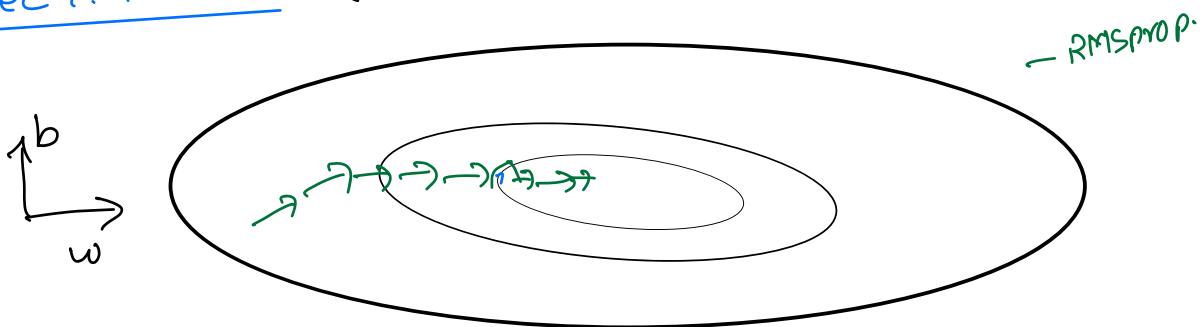
$$\omega = \omega - \alpha vd\omega$$

$$\beta = 0.9 \text{ (avg over 10 iter)}$$

$$b = b - \alpha vd\theta$$

bias correction
= $\frac{vd\omega}{(1-\beta^t)}$ can be included.

dec 7: RMSprop :- (root mean square prop)



on iteration t:

compute d_w, d_b on current mini batch.

elementwise

$$(small) \quad Sd_w = \beta_1 Sd_w + (1-\beta_1) d_w^2$$

$$(large) \quad Sd_b = \beta_1 Sd_b + (1-\beta_1) d_b^2$$

$$w = w - \alpha \frac{d_w}{\sqrt{Sd_w + \epsilon}} \quad (\text{large})$$

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

numerical stability :

$$b = b - \alpha \frac{d_b}{\sqrt{Sd_b + \epsilon}} \quad (\text{reduced})$$

lec8: Adam optimization Algorithm:-

Adam = momentum + RMSprop

$$v_{dw} = 0 \quad s_{dw} = 0 \quad v_{db} = 0 \quad s_{db} = 0$$

on iteration t :

compute dw, db using current mini batch

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

momentum β_1

$$s_{dw} = \beta_2 s_{dw} + (1 - \beta_2) dw^2 \quad s_{db} = \beta_2 s_{db} + (1 - \beta_2) db^2$$

RMSprop β_2

$$\hat{v}_{dw} = \frac{v_{dw}}{(1 - \beta_1^t)} \quad \hat{v}_{db} = \frac{v_{db}}{(1 - \beta_1^t)}$$

$$\hat{s}_{dw} = \frac{s_{dw}}{(1 - \beta_2^t)} \quad \hat{s}_{db} = \frac{s_{db}}{(1 - \beta_2^t)}$$

$$w = w - \alpha \frac{\hat{v}_{dw}}{\sqrt{\hat{s}_{dw}} + \epsilon}$$

$$b = b - \alpha \frac{\hat{v}_{db}}{\sqrt{\hat{s}_{db}} + \epsilon}$$

Hyperparameters:

α : needs to be true

β_1 : 0.9 (momentum) (dw) (db) (1^{st} momentum)

β_2 : 0.99 (RMSprop) (dw^2) (db^2) (2^{nd} momentum)

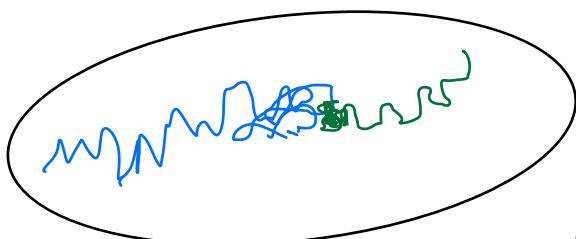
ϵ : 10^{-8}

Adam - adaptive moment estimation.

decay: learning rate decay :-

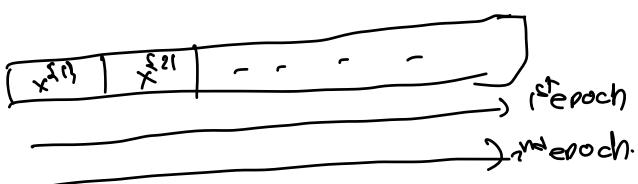
- slowly reduce learning rate overtime.

64, 128 minibatch



epoch = 1 pass through data.

- a - without decay
- a - with decay.

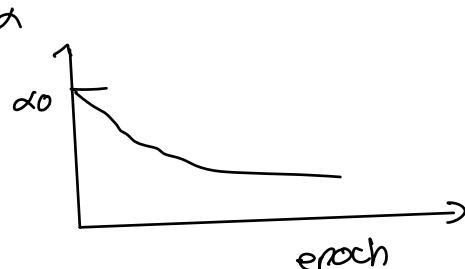


$$\alpha = \frac{1}{1 + (\text{decay rate}) \times \text{epoch number}} \times \alpha_0$$

$$\text{epoch } \alpha_0 = 0.2 \quad \text{decay rate} = 1$$

$$\frac{1}{1 + (1 \times 1)} \cdot 0.2$$

epoch	α
1	0.1
2	0.067
3	0.05
4	0.04



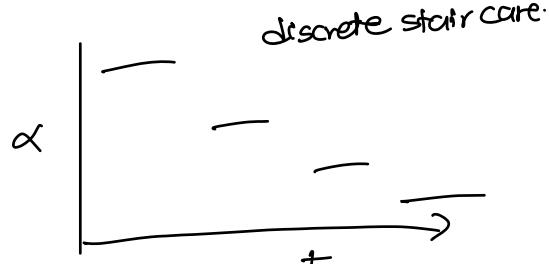
exponential decay

$$\alpha = 0.95^{\text{epochnum}} \cdot \alpha_0$$

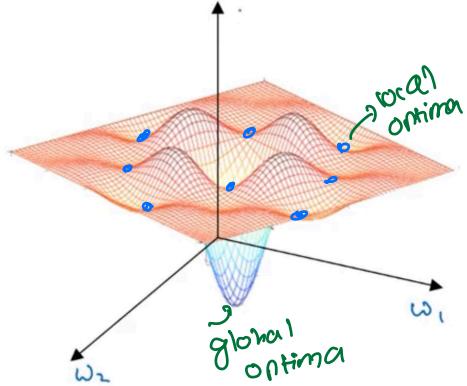
$$\alpha = \frac{k}{\sqrt{\text{epochnum}}} \cdot \alpha_0$$

$$\alpha = \frac{k}{\sqrt{t}} \cdot \alpha_0 \quad (t - \text{minibatch iter})$$

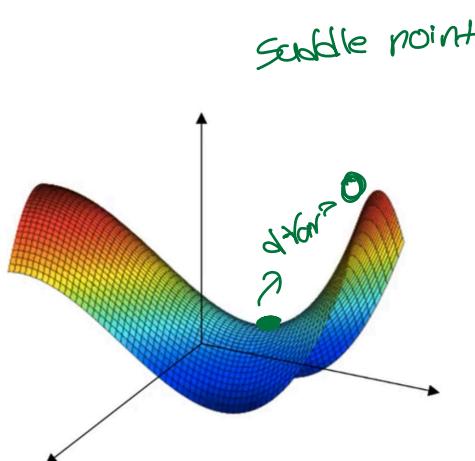
Manual decay is also possible.



dec10: Problem of local optima:-



- most point where gradient zero
one not like this.

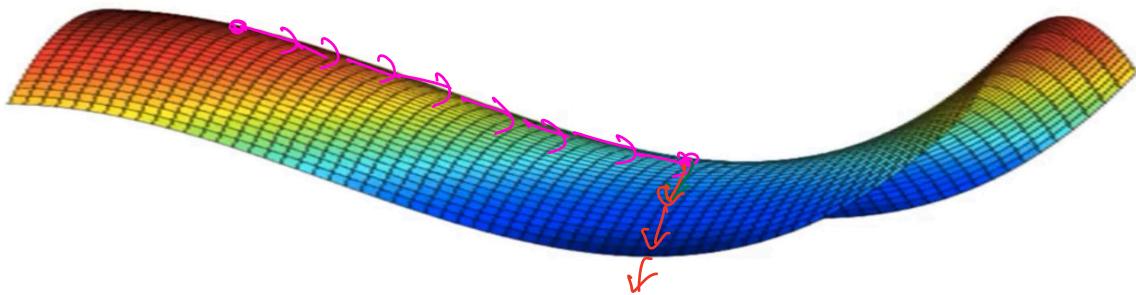


- This is ~~more~~ important problem-

if $n = 20000$ dim

all should be like this \rightarrow

Problem of plateaus:-



- unlikely to get stuck in had local optima
- plateaus can make learning slow
- adam, RMSprop, momentum can speed up GD