

WEEK 2

→ you need fast algo's coz its small & big problems

★ OPTIMIZATION ALGORITHMS

Mini Batch Gradient

vectorization allow to efficiently calculate m examples

what if m is big?

→ calculate $\frac{1}{m}$ calculate gradient of m examples

→ do gradient descent on m examples

→ calculate again $\frac{1}{m}$

do gradient descent for m examples

if m is large, this is very slow

Batch vs mini-batch Gradient descent.

↓

you divide your training set (m) into multiple examples batches.

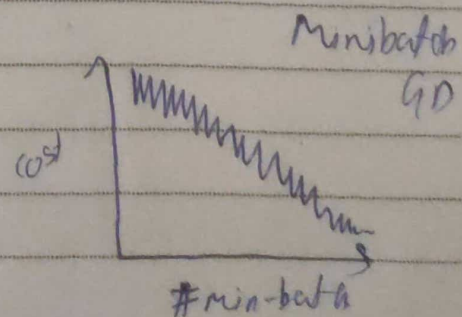
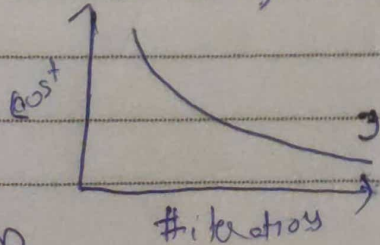
for $m = 5000$, you divide it into 5000 batches of 1000 each.

epoch $\left\{ \begin{array}{l} \text{Now you run gradient descent on batches of 1000} \end{array} \right.$

so you'll need a total of 5000 epochs.

★ minibatch works faster than batch gradient descent for big m .

understanding Mini batch



hyperparameter.

→ Choosing your mini-batch size.

mini-batch size = m : Batch grad Descent.

mini-batch size = 1 : Stochastic GD

big steps

large steps

straight up goes to minima.

but slow

small steps

very noisy

takes zig-zag turn to minima

So usually you'll have something in between 1 and m .

if small training set : use batch gradient descent.
otherwise,

Typical mini-batch size : $\rightarrow 64, 128, 256, 512$
Apparently for powers of 2
it is faster

Make sure mini-batch size fits in CPU/GPU memory

Exponentially weighted average.

is used in "optimization" algo which works faster than GD

eg:- Say you want to approximate a very noisy data

predict value

$$v_0 = 0$$

$$v_1 = 0.9v_0 + 0.1d_1$$

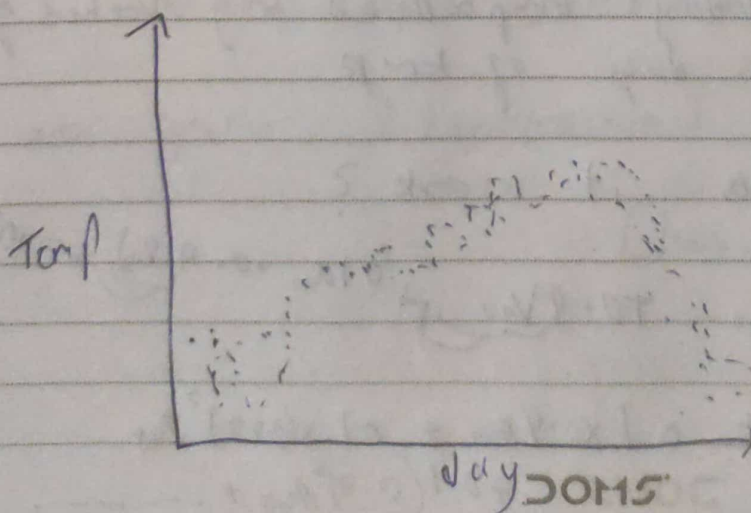
$$v_t = 0.9v_{t-1} + 0.1d_t$$

temp on day t

$v \approx$ our average approximation

$$\beta = 0.9$$

we can choose
diff values of β



$V \rightarrow$ exponentially weighted avg.

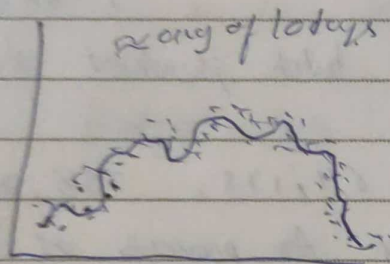
11

β is actually a hyperparameter
 β approximately represents, the no. of days one
 combined with V which we're averaging
 values to predict today's value

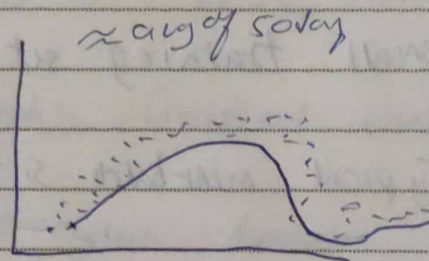
$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

\hookrightarrow \approx avg over $\frac{1}{1-\beta}$ days of temp

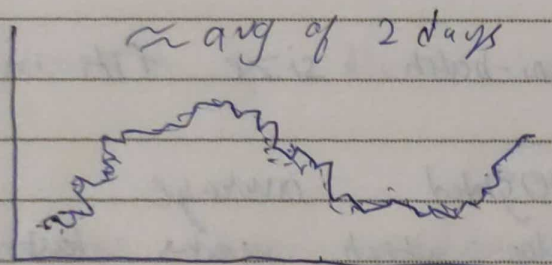
if $\beta = 0.9$ \approx 10 days of temp taken into
 consideration to predict today's temp



$\beta = 0.9$



$\beta = 0.95$



$\beta = 0.5$

But what are we actually doing here?

\hookrightarrow approximating today's temp based on avg values of
 some other days of temp

how does it work?

eg:-

$$V_{100} = 0.1 \theta_{100} + 0.9 V_{99} \rightarrow (0.1 \theta_{99} + 0.9 V_{98}) \rightarrow \text{and so on.}$$

exponentially
 decaying
 fn.

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

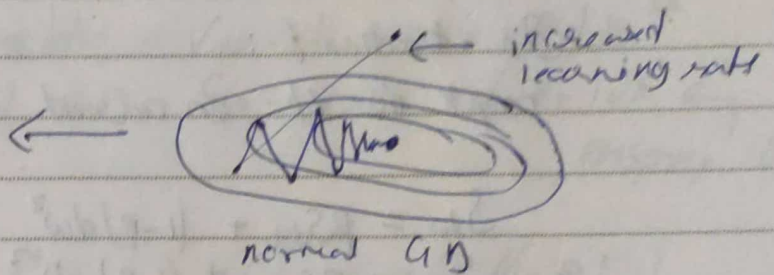
Bias correction in exponentially weighted avg

↳ makes it more accurate

→ divide V_t by $1 - \beta^t \rightarrow V_t = \frac{V_t}{(1 - \beta^t)}$

Gradient Descent with momentum

lots of unnecessary computation, bcz it doesn't go in the direction of minima



★ → we want to give it the right direction but if we increase learning rate, it might overshoot

what we want :-

↕ slow learning in vertical direction
↔ fast learning in horizontal direction

→ so we introduce momentum.

Algorithm: On iteration t :

compute dW, db on current minibatch

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

$$V_{db} = \beta V_{db} + (1 - \beta) db \rightarrow \text{provides accel}^n \text{ in } G, N$$

provides direction in G, N

provides velocity in G, N

$$W = W - \alpha V_{dw} ; b = b - \alpha V_{db}$$

now you've 2 hyperparameters = α, β

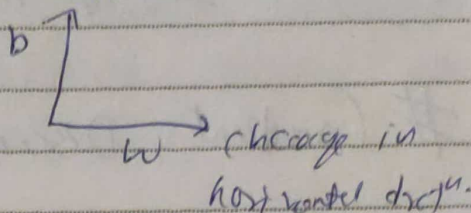
usually 0.9

[bias correction isn't usually considered a hyperparameter bcz it doesn't matter that much]

RMS prop

↳ speeds up GD

change in vertical direction



On iteration f :

compute dw , db on current mini-batch

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

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

$$W := W - \frac{dw}{\sqrt{S_{dw}}}$$

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

we hope $\rightarrow S_{dw}$ is relatively small
so change in w (horizontal) is big.

$\rightarrow S_{db}$ is v large so change in b (vertical direction) is small.

Adam Optimization

\rightarrow there have been alot optimization algo's in DL history

\rightarrow but most of them don't generalize in all cases

\rightarrow RMSprop and GD with momentum are among v rare ones which generalize.

Adam optimization \rightarrow combines GD with momentum and RMSprop

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

On iteration t :

compute dw, db using current mini-batch.

momentum β_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, β_2

bias correction

$$\begin{cases} V_{dw}^{corrected} = V_{dw} / (1-\beta_1^t) \\ S_{dw}^{corrected} = S_{dw} / (1-\beta_2^t) \end{cases}, \begin{cases} V_{db}^{corrected} = V_{db} / (1-\beta_1^t) \\ S_{db}^{corrected} = S_{db} / (1-\beta_2^t) \end{cases}$$

$$w := w - \alpha \frac{V_{dw}^{corrected}}{\sqrt{S_{dw}^{corrected}} + \epsilon}$$

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

hyperparameters choice:

α : needs to be tuned

$\beta_1 = 0.9$

(dw)

Suggested by

$\beta_2 = 0.999$

(dw^2)

authors of
Adam optimization
paper

$\epsilon = 10^{-8}$

why is it called adam?

→ Adaptive Moment Estimation

Learning Rate decay

↳ slowly reduces learning-rate over time

-yK why we need this, you actually proposed it

1 epoch = 1 pass through the data.

$$\alpha = \frac{1}{1 + \text{decay-rate} \times \text{epoch-num}} \alpha_0$$

decay-rate \rightarrow another hyperparameter

for $\alpha_0 = 0.2$, decay rate = 7

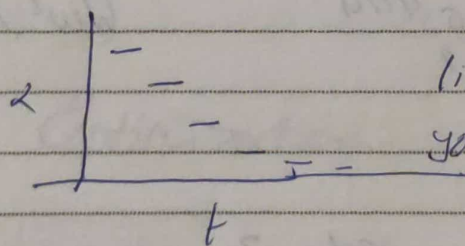
epoch	α
1	0.1
2	0.067
3	0.05
4	0.04

other ways of calculating α

$$\alpha = (0.95)^{\text{epoch-num}} \alpha_0$$

- exponential decay

$$\alpha = \frac{K}{\sqrt{\text{epoch-num}}} \alpha_0 \quad \text{or} \quad \frac{K}{\sqrt{t}} \alpha_0$$



like a stair case

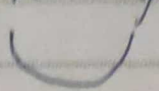
you keep dividing α by 2

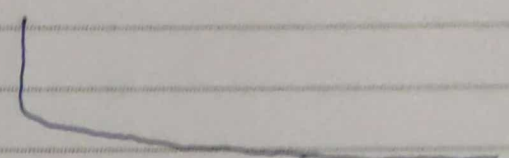
in each iteration

OR you can control α , manually !!

the problem of local optima

Our intuitions about local and principal minima in 2, 3 dimensions don't really hold up in higher dimension functions

we prefer to choose a convex fⁿ
but what if we're in 20000 dimensional fⁿ?
it's crazy hard to get a function
which looks like  this in 20000 dimensions
lots of saddle points
the point of ^{what you get in high dimensions} principal ~~minima~~ = saddle point

another thing which slows GD is "plateau" regions
the slope remains close to zero for a long
period of time
something like 

- we're
- unlikely to get stuck in local optima in high dimensions
 - Plateaus are a problem