

Momentum based GD $\rightarrow w_{t+1} = w_t - \eta \nabla W_t$

but takes a lot of time to navigate a simple slope
if I'm moving in one dirⁿ w/ many small steps
I can take a big step.

dude is capturing my history!
update₀ = 0

$$\left. \begin{aligned} \text{update}_t &= r \cdot \text{update}_{t-1} + \eta \nabla W_t \\ w_{t+1} &= w_t - \text{update}_t \end{aligned} \right\} \text{recursive tool!}$$

$$\text{update}_1 = r \text{update}_0 + \eta \nabla W_1 = \eta \nabla W_1$$

$$\text{update}_2 = r \text{update}_1 + \eta \nabla W_2 = r \eta \nabla W_1 + \eta \nabla W_2$$

$$\text{update}_t = r \cdot \text{update}_{t-1} + \eta \nabla W_t$$

exponentially weighted avg.
as $t \uparrow$, 1st term \downarrow weightage
 $r < 1$ too.

- In regions w/ gentle slope — momentum based GD can take longer steps due to momentum.
- Sometimes moving fast = overshooting \rightarrow momentum based GD may oscillate in + out of minimum valley.
 \rightarrow better than vanilla GD. \downarrow lots of u-turns before convergence.

Nesterov Accelerated Gradient: (NAG)

"look before you leap."

Before: $\text{update}_t = \underbrace{r \text{update}_{t-1}}_{\text{move by this much}} + \underbrace{\eta \nabla W_t}_{\text{and then some more.}}$

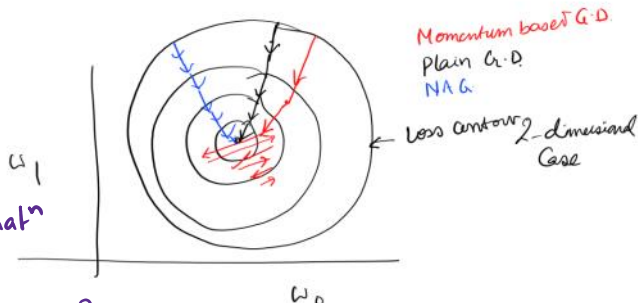
$$\left. \begin{aligned} w_{\text{look-ahead}} &= w_t - r \cdot \text{update}_{t-1} \\ \text{update}_t &= r \text{update}_{t-1} + \eta \nabla W_{\text{look}} \\ w_{t+1} &= w_t - \text{update}_t \end{aligned} \right\}$$

But
What if we calculated gradient @
partially updated value of w ?

$$\rightarrow w_t - r \text{update}_{t-1} \text{] look ahead}$$

- dude corrects momentum GD's oscillations
- here, oscillations smaller
- chances of overshooting minima smaller.

- update @ current step = avg. of grad + prev step's estimatⁿ



Here's where it kinda sucks:

- treat w_0, w_1 equally but loss surface never symmetrical
- slope high in one dimⁿ, low in another.

$\left. \begin{aligned} &\text{Dim}^n w \text{ large grad} = \text{small update} \\ &\text{vice versa.} \end{aligned} \right\}$

Adjusting Learning Rate: Gradient Descent

\rightarrow goes over entire data before
doing one update.

Stochastic
G.D.
update after
data point

Minibatch
G.D.
update after
a batch

Adjusting Learning Rate

Gradient Descent	Stochastic G.D.	Minibatch G.D.
<ul style="list-style-type: none"> → goes over entire data before doing one update. → True Gradient of Loss ↓ No approximation ↓ Theoretical guarantee holds ↓ very very slow → <exreme> 	<ul style="list-style-type: none"> update after data point → <exreme> 	<ul style="list-style-type: none"> update after a batch → <in-between>

1 epoch = one pass over entire dataset
 1 step = one update of the parameters (\vec{w}, \vec{b})
 N = no of data points
 B = minibatch size

Algorithm	# steps in 1 epoch
Vanilla G.D	L
Stochastic G.D	N
Batch G.D	N/B

Heuristics based learning rate optimizations:

① Tune learning rate:

- try diff values on log scale for η
 - run a few epochs and find which best
 - do finer search.
- } linear search.

② Step decay:

- Halve η after every 5 epochs
- Halve η after epoch if validation error > prev epoch.

③ Exponential decay:

- $\eta = \eta_0 e^{-kt}$ $t = \text{step no.}, k, \eta_0 = \text{hyperparams.}$

④ Y_+ decay: $\eta = \frac{\eta_0}{1+kt}$

Adagrad — gradient update integrating adaptive learning rate.

accumulated squares of gradients of each weight (cumulative)

$$V_t = \nabla w_1^2 + \nabla w_2^2 + \dots + \nabla w_t^2$$

→ used for normalization.

$$V_t = V_{t-1} + (\nabla w_t)^2$$

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{V_{t+1}}} \nabla w_t$$

Issues!

- decays η aggressively (denom $\uparrow \uparrow$)
- freq. param = slow update



fix by decaying denom!

$$V_t = \beta V_{t-1} + (1-\beta)(\nabla w_t)^2$$

$$v_t = \beta v_{t-1} + (1-\beta)(\nabla w_t)^2$$

$w_{t+1} = \text{same}$

ADAM — momentum gradient descent + RMS prop moving avg. of sq. grad. } 2 decay hyperparams β_1 β_2

Now, the final update Rule

$$w_{t+1} = w_t - \eta \frac{\bar{v}_t}{\sqrt{s_t + \epsilon}} \nabla_{\bar{w}} \mathcal{L}(w_t + \bar{w})$$

Benefit of momentum approach

Benefit of RMSprop approach

used in this eqn.

Bias Correction in ADAM

Sum of weights of past values of state vectors \bar{v}_t, s_t will tend to ∞ if $t \rightarrow \infty$. This will happen only at very very large value of t .

So, for smaller values of t , we need to correct \bar{v}_t & s_t .

$$\hat{\bar{v}}_t = \frac{\bar{v}_t}{(1-\beta_1^t)} \quad \hat{s}_t = \frac{s_t}{(1-\beta_2^t)}$$

looks out what is small !!

