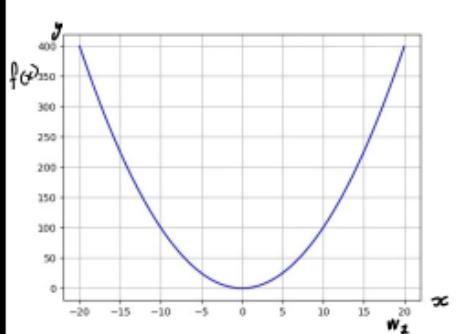


SGD with Momentum



→ How loss function is changing as a single weight is getting changed.

→ X axis = w_2 , Y axis = $f(x)$

$$\hookrightarrow \text{Loss} = f(x) = f(w_2)$$

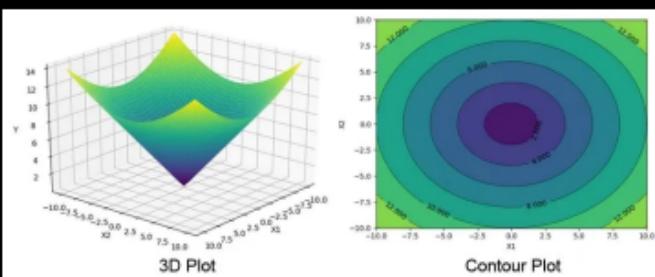
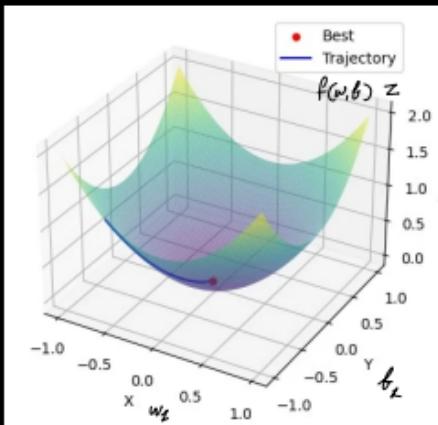
→ Our aim is to minimize loss as much as possible, by finding the optimal value of w_2 .

→ This is 3D representation of loss function w.r.t weight (w) & bias (b).

→ X axis = w , Y axis = b , Z axis = Loss function = $f(w, b)$

→ How loss is changing with respect to ' w ' & ' b '.

→ Our aim is to find such optimal value of w & b which brings down loss function as much as possible.



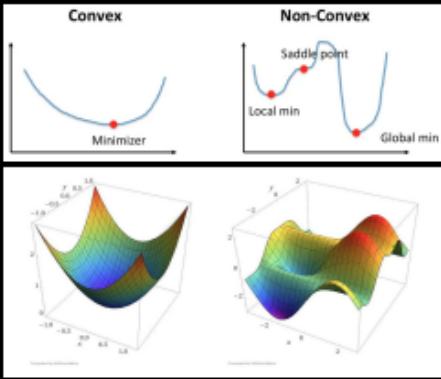
→ Contour plot is a simple 2D visualization of 3D plot from the top angle.

→ Contour plot Shows lost dimension with the help of colors.

→ As given image shows, blue region shows down side.

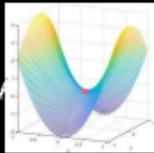


Convex vs Non-convex optimization



Problems in Non-convex opt.

- local minimum trap
- Saddle point →
↳ in one direction, slope is upwards and another direction, slope is downwards for long distance.
- High curvature



- ↳ Small circle = Small radius = High curvature
- ↳ Big circle = Large radius = Low curvature

Momentum Optimization

Target Problems

- High curvature
- Consistent gradient.
- Noisy gradients (local minimum trap)

Basic idea

→ It accelerates the convergence of gradient descent in that direction which previous gradient strongly suggests.

→ If our gradients are sending us in one direction then we will accelerate the convergence speed.

Single most benefit of momentum optimization

Speed

Mathematics

$$W_{t+1} = W_t - \alpha \Delta W_t \quad \text{if } V = \text{Velocity} \text{ then , now}$$

$$W_{t+1} = W_t - V_t \quad \text{where}$$

→ We will use all our past velocity as momentum and that will provide acceleration.

$$V_t = \beta V_{t-1} + \alpha \Delta W_t$$

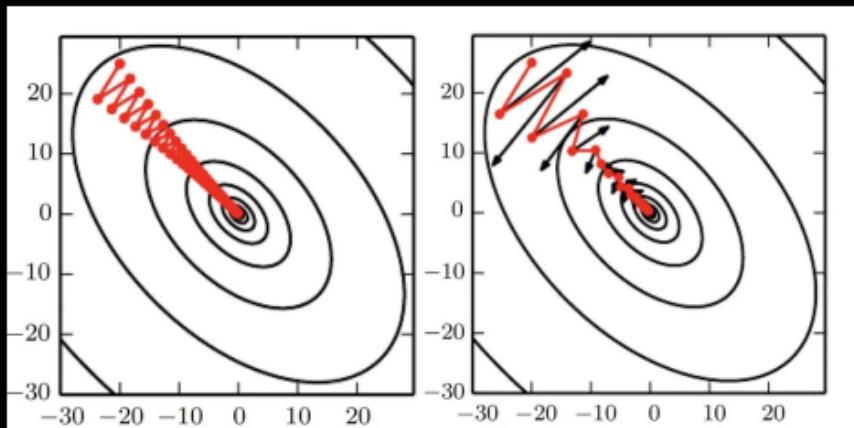
(past velocity) \rightarrow $\underbrace{\beta}_{0 \leq \beta \leq 1}$ ΔW_t \rightarrow $\underbrace{\alpha}_{0 < \alpha < 1}$



Created with Notewise

SGD

SGD with momentum



Pole of β in momentum

$$W_{t+1} = W_t - V_t \quad \text{if } \beta = 0 \text{ then}$$

$$V_t = \beta V_{t-1} + \alpha \nabla W_t \quad V_t = \alpha \nabla W_t$$

→ Hence if $\beta = 0$, momentum will act like SGD.

$$\begin{aligned} W_{t+1} &= W_t - V_t \\ &= [W_t - \alpha \nabla W_t] = SGD \end{aligned}$$

→ Hence, β is a kind of decaying factor

→ Newer velocity will have more impact than older.

→ Even $\underbrace{\beta = 1}_{\text{No decay}}$ won't perform as expected, hence β should be 0.9 or 0.99...



Nesterov Accelerated Gradient (NAG)

- NAG is a different technique to generate momentum.
- NAG is a technique in which we reduces oscillations.

Muth intution

Momentum

$$W_{t+2} = W_t - V_t$$

where

$$V_t = \beta V_{t-2} + \alpha \nabla W_t$$

In momentum, next position depends on 2 factors

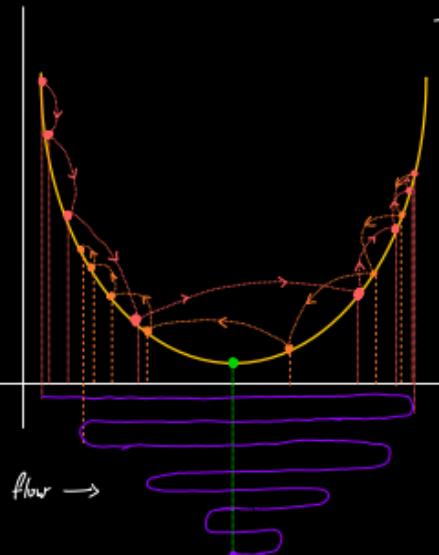
similation process

1. History velocity (βV_{t-2})

2. Gradient at that point ($\alpha \nabla W_t$)

In **NAG**, first we will update our position based on momentum (look ahead term)

→ Based on new position, we will calculate gradient from that location.



$$W_{t+2} = W_t - \beta V_{t-2} \quad \text{where}$$

$$V_t = \beta V_{t-2} + \alpha \nabla W_{t+2} \quad W_{t+2} = \text{look-ahead}$$

$$W_t = W_t - V_t$$

In simple words...

Step : 1 follow the momentum & go to that direction.

Step : 2 after reaching at the location calculate the gradient at new location.

Disadvantages

- local minima trap.

```
sgd_optimizer = tf.keras.optimizers.SGD(
    learning_rate=lr_schedule,
    momentum=0.9,
    nesterov=True,
    clipnorm=1.0 # Gradient clipping to avoid exploding gradients
)
```

SGD	SGD + momentum	NAG
momentum = 0.0	momentum = 0.9 (decay factor)	momentum = 0.9
nesterov = False	nesterov = False	nesterov = True



AdaGrad (adaptive gradient)

→ learning-rate will keep changing based on situation.

Usefull scenarios

- Scale of i/p features are different. ex. (age, salary)
- Sparse data (Most values are 0)

→ Problem with Sparse data.

• Elongated bowl Problem.

→ Slope changes only on one axis but keeps constant on another axis.

Why it happens..?

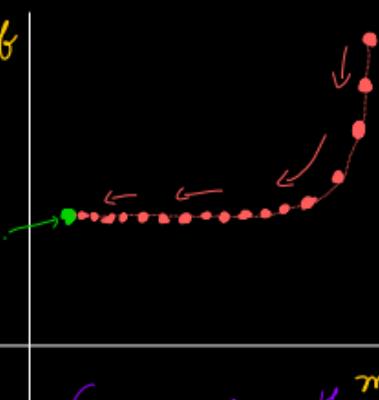
→ Suppose we have a dataset like α, b where α is sparse column by making prediction $m\alpha + b$, the terms $m\alpha$ will become zero in majority cases & b will dominate in some data-points. hence weight (w) will converge slowly.

$$w = w - \alpha \frac{\partial L}{\partial w} \text{ (slow convergence)}$$

$$b = b - \alpha \frac{\partial L}{\partial b} \text{ (fast convergence)}$$

→ As graph shows, initially b will converge very fast, whereas due to sparse values, convergence in m (horizontal) direction becomes too slow until b reaches its optimal location.

→ hence majority of values in α are zero, the convergence in m (horizontal direction) will be very slow.



Convergence path. m



- This is the problem, we don't want such behaviour, we want that convergence in both direction must be equal.
- In this situation we have one last option learning-rate, so in Adagrad there will be different learning rate for individual parameters.
- In this situation, to balance convergence of each parameter (w, b) we will use different learning rate for individual one.

Exe $w = w - \alpha \frac{\partial L}{\partial w}$ (slow convergence) → Set high learning rate

$b = b - \alpha \frac{\partial L}{\partial b}$ (fast convergence) → Set low learning rate

Math formulae

$$w_{t+1} = w_t - \frac{\alpha \nabla w_t}{V_t + \epsilon} \quad \xrightarrow{\frac{\partial L}{\partial w}}$$

$\underbrace{V_t}_{\text{Small number}} + \epsilon$

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

(squared sum of past gradients)

Same for b also..

Disadvantages

- Not useful in NN. (Can be useful in Linear Regression but not for complex problems)
- It can reach near to solution but never to exact solution. due to decreasing learning rate.

Solution
 \hookrightarrow RMSProp
 \hookrightarrow Adam



RMSProp (Root Mean Squared Propagation)

- Improved version of AdaGrad
- AdaGrad reduces learning-rate continuously but at some point learning-rate becomes too small & due to this reason, conversion can't happen, due to this problem RMSProp is here to solve the problem.

Adagrad

$$W_{t+1} = W_t - \frac{\alpha}{\sqrt{V_t + \epsilon}} \nabla W_t$$

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

for b

$$b_{t+1} = b_t - \frac{\alpha}{\sqrt{V_t + \epsilon}} \nabla b_t$$

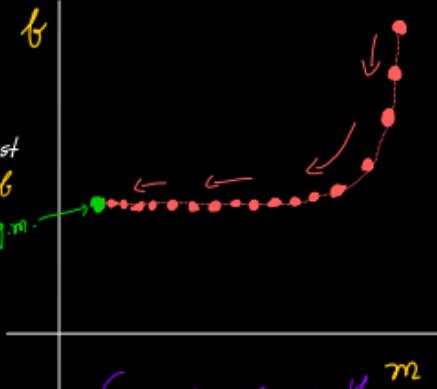
$$V_t = V_{t-1} + (\nabla b_t)^2$$

- In beginning within few epochs

convergence of W will be accelerated & b there will be almost no conversion in b .

→ But after some time due to too many past convergences of W , eventually learning rate of b will become too smaller. hence no convergence will occur in b .

→ So we need to control the explosion of V_t . that's what RMSProp does.



RMSProp

$$W_{t+1} = W_t - \frac{\alpha}{\sqrt{V_t + \epsilon}} \nabla W_t$$

$$\beta = 0.95 \text{ (Most cases)}$$

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

→ Exponentially decaying average

→ Giving more importance to recent gradient instead of giving more importance to older gradients also.

→ This mechanism won't reduce learning-rate gradually after so



Created with

Notewise

→ hence V_t never shoot.

$$V_0 = 0$$

$$V_1 = 0.95 \times 0 + 0.05 (\nabla w_1)^2$$

$$V_2 = 0.95 \times 0.05 (\nabla w_1)^2 + 0.05 (\nabla w_2)^2$$

$$V_3 = \underbrace{0.95 \times 0.95 \times 0.05 (\nabla w_1)^2}_{\text{epoch-1}} + \underbrace{0.95 \times 0.05 (\nabla w_2)^2}_{\text{epoch-2}} + \underbrace{0.05 (\nabla w_3)^2}_{\text{epoch-3}}$$

→ This is one of the best optimization technique with no major disadvantages.

→ Widely used optimizer before Adam.



Adam (Adaptive Moment Estimation) Optimizer

→ Top performance
ANN
CNN
RNN

→ Till now whichever optimizers we have seen, they have used 2 main techniques 1. Momentum, 2. learning rate decay but Adam uses both of them.

Math intuition

$$W_{t+1} = W_t - \frac{\alpha}{V_t + \epsilon} \times M_t$$

where

$$M_t = \beta_2 M_{t-1} + (1-\beta_2) \nabla W_t \quad \} \text{ momentum}$$

$$V_t = \beta_2 V_{t-1} + (1-\beta_2) (\nabla W_t)^2 \quad \} \text{ adagrad}$$

Bias correction

$$\begin{cases} \beta_2 = 0.9 \\ \beta_2 = 0.99 \end{cases} \quad \} \text{ By default in libraries}$$

$$\hat{M}_t = \frac{M_t}{1-\beta_2^t}$$

$$\hat{V}_t = \frac{V_t}{1-\beta_2^t}$$

Where

$t = \text{number of epochs}$

