

Optimization algo used to find \min^m of a $\frac{f(x)}{f(x)}$.

$$x_{\text{new}} = x_{\text{old}} - \alpha \vec{\nabla} f(x_{\text{old}})$$

↓
learning rate (step size)

{ Epoch = One full pass over the entire dataset.
Batch = A small subset of dataset used in one step of training
Iteration = One update of the model using a batch.

ex: 1000 samples ; batch size = 100

$$1 \text{ epoch} = 10 \text{ iterations } (1000 / 100 = 10)$$

Batch Size =
Updates per epoch =
Speed =
Stability =

Batch GD	Stochastic GD	Mini-batch GD
All data	1 sample	Small batch
1 update per epoch	Many updates	Medium updates
Slow	Fast	Balanced
High	low (noisy)	Balanced

$$1 \text{ epoch} = \frac{\text{dataset size}}{\text{batch size}}$$

<https://colab.research.google.com/drive/1PKOoRT76HTA-5WrbDY6jAiH87c5KVw> t?authuser=2#scrollTo=q3zASsyttNsN (For code and visualization purpose)

→ Code of all 5 optimizers, we studied.

SGD with Momentum

Wednesday, 19 February 2025 11:30

Why we need momentum?

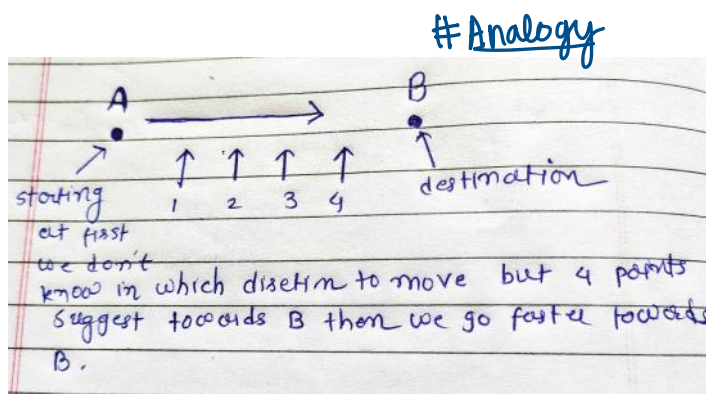
① non-convex cost funcⁿ.

② saddle pts

A saddle point is a point where in one direction the surface goes in the upward direction and in another direction it goes downwards

From <<https://paperswithcode.com/method/sgd-with-momentum>>

③ low performance
||
Speed + Stability



EWMA (Expo Wt Moving Avg)

↳ Tech used to find trends in time-series data.

"velocity"

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

wt assigned to past value

$$0 < \beta < 1$$

ex: if $\beta = 0.5$ then

$$\frac{1}{1-\beta} = 2$$

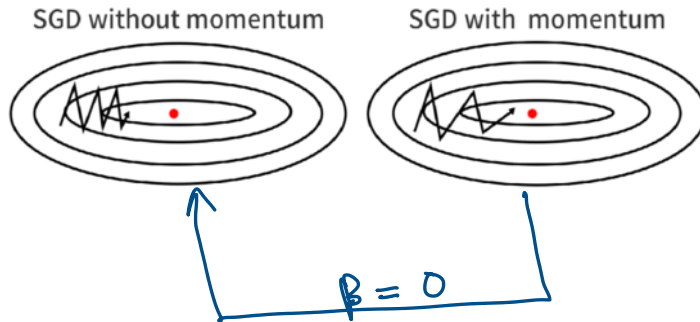
calc. avg was from previous 2 readings.

high $\beta \Rightarrow$ Avg of more past data.

where,

$$W_{t+1} = W_t - v_t$$
$$v_t = \beta v_{t-1} + \eta \Delta W_t$$

Decaying factor



Advantages :

- ①. Faster convergence.
- ②. Escaping local minima (goal \rightarrow global minima)

Adagrad adapts the learning rate individually for each parameter based on past gradients.

Wednesday, 19 February 2025 11:57

Let, $\epsilon = 10^{-8}$ (to avoid div. by zero)

let $h_0 = 0$

Sum of squared gradients :

$$h_t = h_{t-1} + g_t^2$$

element-wise sq. of grad vec

Initial learning rate

$$\theta_t = \theta_{t-1} - \frac{\eta}{\sqrt{h_t + \epsilon}} g_t$$

Adaptive term that scales the learning rate.

* Large gradients \Rightarrow small l.r.

* Small " \Rightarrow large l.r.

* Good for sparse data.

Root Mean Square Propagation

↳ Improvement over
Adagrad

Instead of accumulating all past squared gradients (which causes the learning rate to decay too much), RMSProp uses a **moving average of squared gradients**. This makes it more suitable for **non-convex optimization problems**, such as training deep neural networks.

↓
expo dec. avg \Rightarrow More wt. to recent gradients.
Initialize θ_0 ; $E[g^2]_0 = 0$ "Moving avg of sq gradients"
 $\gamma \rightarrow$ expo decay factor (typically, $\gamma = 0.9$)

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (1 - \gamma) g_t^2$$

$$\theta_t = \theta_{t-1} - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

ADAM = Mom + RMSprop

Wednesday, 19 February 2025 12:01

(Adaptive Moment Estimation)

* ADAM uses 2 moving averages \Rightarrow

① mom term (m_t)

② Adaptive scaling (v_t)

$g_t \rightarrow$ gradient

since m_0 & $v_0 = 0$
 \Rightarrow biased towards smaller values in early steps.

Bias correction \Rightarrow

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

\Downarrow
Too small updates in beginning.

$$\begin{aligned} m_t &= \beta_1 m_{t-1} + (1 - \beta_1) g_t \\ v_t &= \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \\ \theta_t &= \theta_{t-1} - \frac{\alpha}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t \end{aligned}$$

$$\eta \sim 0.001 ; \quad \beta_1 \sim 0.9 ; \quad \beta_2 \sim 0.999$$
$$\epsilon \sim 1e-8$$