

Session-8

NN - OPTIMIZERS

Feb 23, 2024



Will you please listen? I am not the messiah.

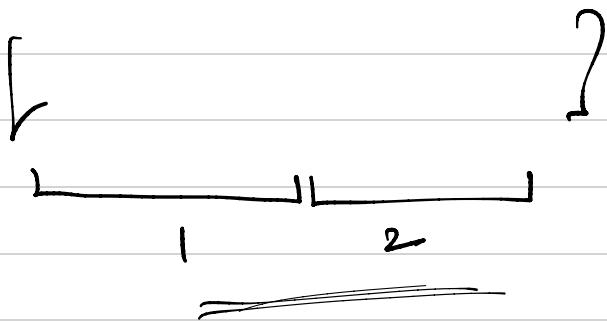
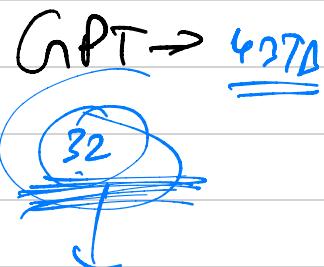
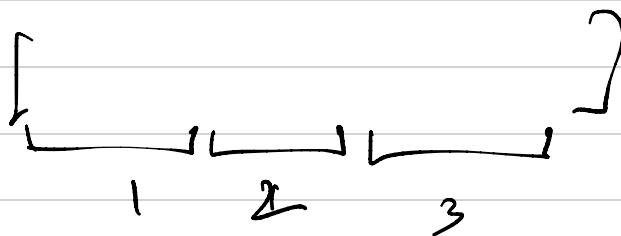


imgflip.com

AGENDA

- ① Some problems with SGD
- ② EMA
- ③ RMS Prop
- ④ ADAM

How SGD Works



2.217

$$\omega_{ij}^k = \omega_{ij}^k - \eta \sum_{i=1}^n \frac{\partial L}{\partial \omega_{ij}^k}$$

no. of data points

Diagram illustrating the update rule for weights ω_{ij}^k using Stochastic Gradient Descent (SGD):

$$\omega_{ij}^k = \omega_{ij}^k - \eta \sum_{i=1}^n \frac{\partial L}{\partial \omega_{ij}^k}$$

ω_{ij}^k → Layer
From Layer → to Layer

ω_{34}^2 → 2nd layer,
3rd neuron → 4th neuron

→ Stochastic Gradient descent

↑ single samp. (update weight)

→ Batch / Mini-batch G. descent (You take avg of grad → batch)

a batch of dataset (3-2 MB - GPU)

$$32 \text{ - } 64 \text{ - } 128$$

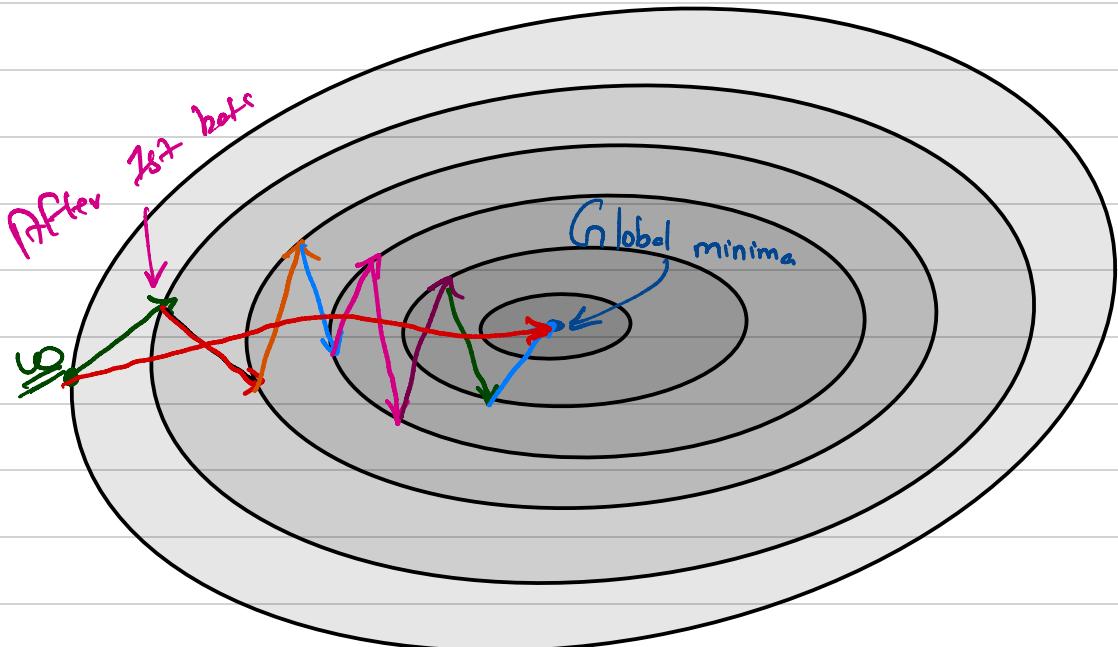
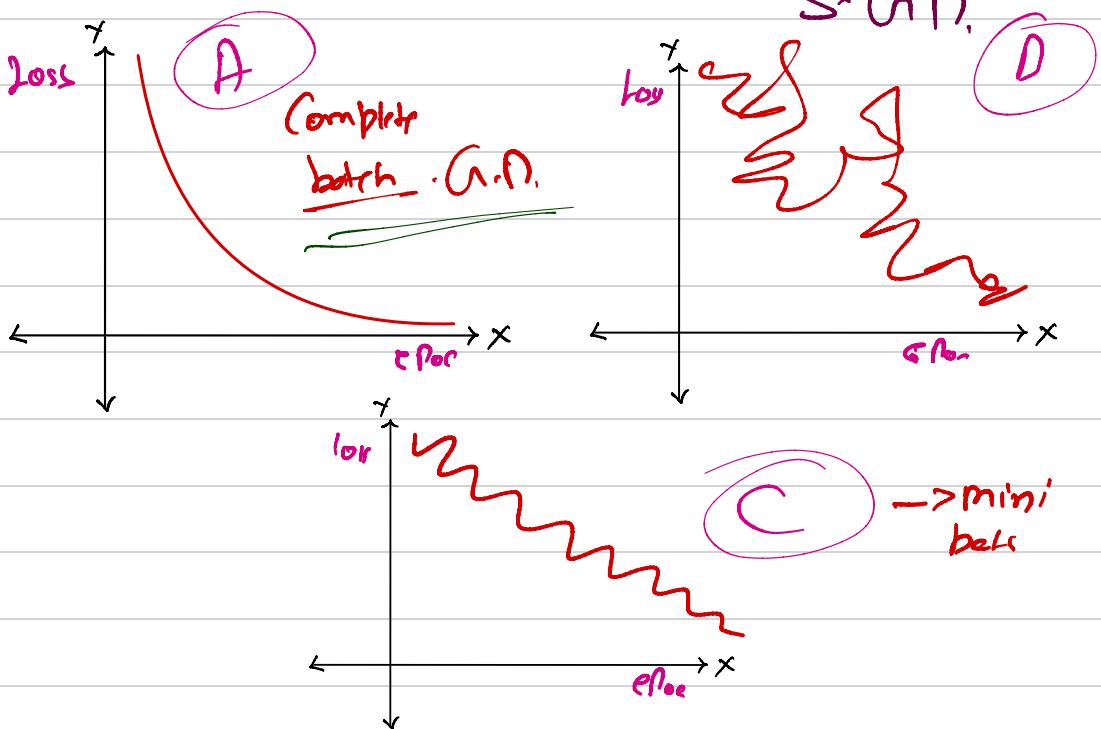
Replacement - not recommended

$$\left. \begin{array}{l} 1MB = 1024 KB \\ 1GB = 1024 MB \end{array} \right\}$$

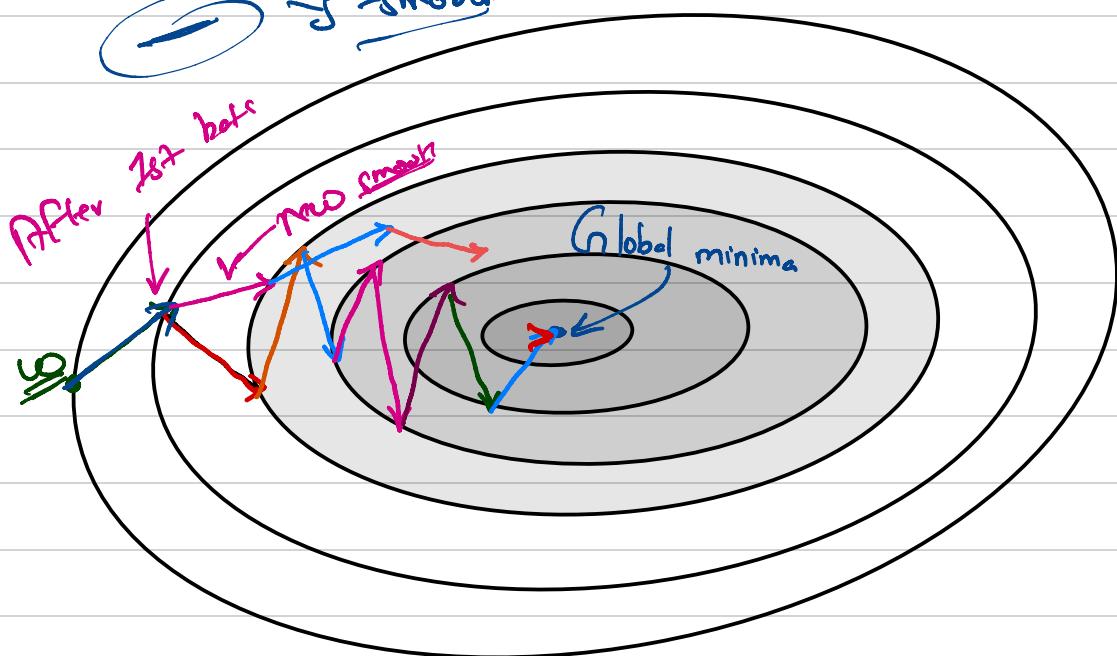
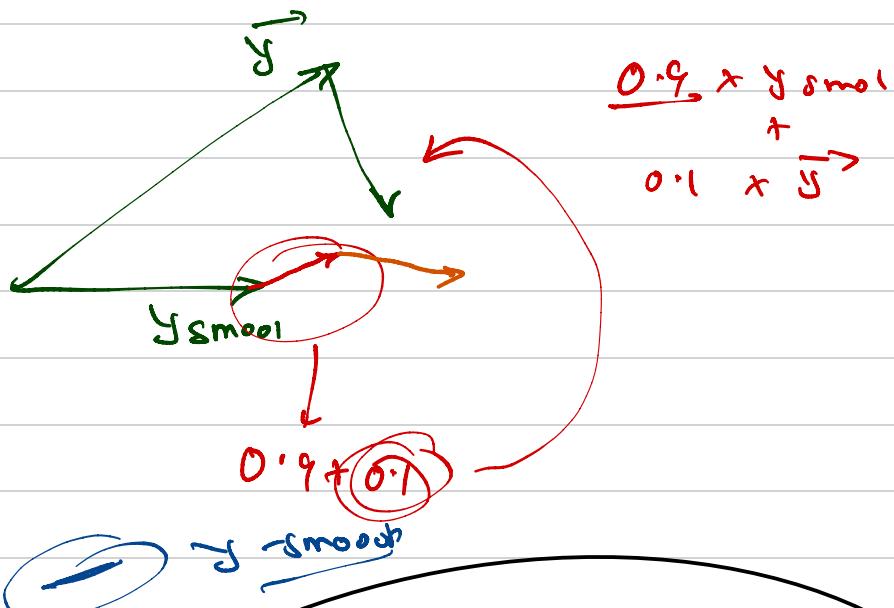
only handle an array $\rightarrow 2^n$

(10) (GPU → 16 Samples)

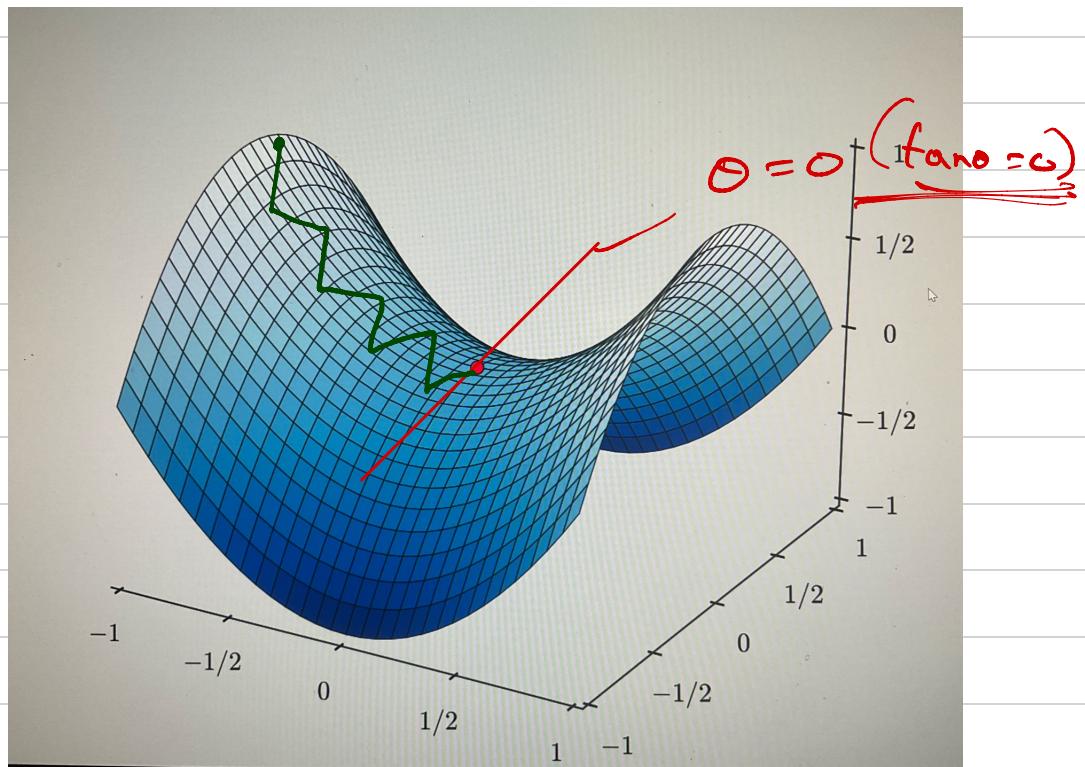
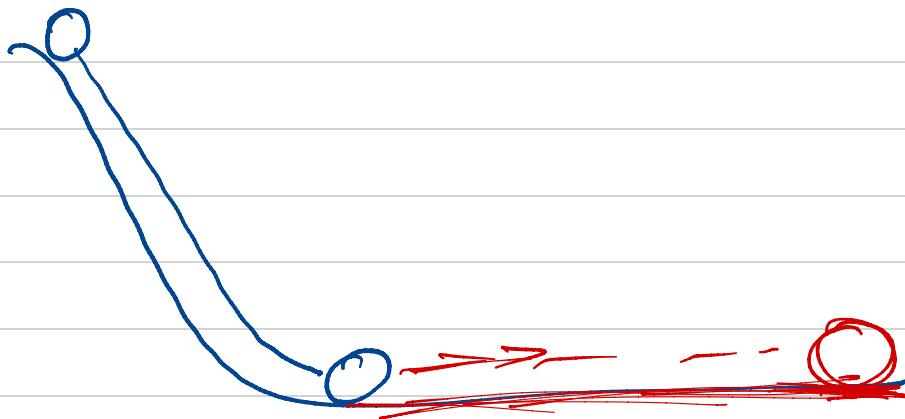
Using → 10 samp.



Exponential moving average



Momentum



When using EM

$$\frac{\partial \mathcal{L}}{\partial w_{ij}^k} \rightarrow \text{at time step } t \rightarrow g^t$$

$$w^{t+1} = w^t - \eta \times g^t$$

In EM

$$g^t = \beta \times g^{t-1} + (1-\beta) \times \frac{\partial \mathcal{L}}{\partial w_{ij}^k}$$

batch gradient of
Prev step
t-1

batch current step

How momentum helps

$$v_0 = 0, \beta = 0.9, w_0 = 0.1$$
$$v_1 = \beta \times v_0 + (1-\beta) \times w_0$$
$$= 0.9 \times 0 + (1-0.9) \times 0.1 = \underline{0.01}$$

$$\begin{aligned} g^{t-1} &= 0.04 \\ \beta &= 0.9 \\ g^t &= 0.06 \end{aligned}$$

$$0.06 = 0.9 \times 0.04$$

$$+ (1-0.9) \times \underline{\Delta \omega} \quad g^{t+1} = \frac{\partial L}{\partial \omega} \quad (t=1)$$

$\approx 0.24\omega$

$$90\% \times g^t + 10\% \rightarrow g^{t+1}$$

0.06

$$\omega = \omega - \eta \times \frac{\partial L}{\partial \omega} \quad \begin{matrix} \nearrow \text{gradient} \\ t=0 \end{matrix}$$

gradient at $t=1$

$$\beta \times g^{t-1} + (1-\beta) \frac{\partial L}{\partial \omega}$$

gradient step

Vanilla gr

$$\theta = \omega + \text{at}$$



Rms Prop



$$\Rightarrow \sqrt{\sum_{t=1}^n g_t^2 + G} \rightarrow 10^{-8}$$

as time inc

$$\omega_{\text{new}} = \omega_{\text{old}} - \frac{n}{\sqrt{\sum_{t=1}^n g_t^2 + G}} \times \frac{\partial L}{\partial \omega}$$

$$\rightarrow t \uparrow \Rightarrow \left(\sum_{t=1}^n g_t^2 \right) \rightarrow \downarrow$$

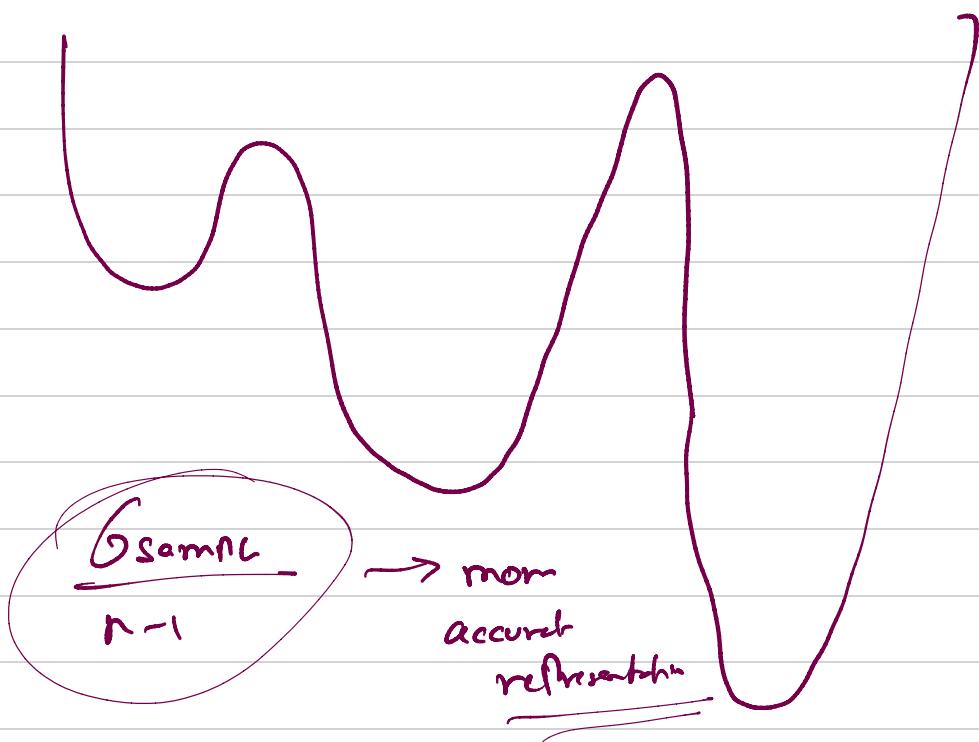
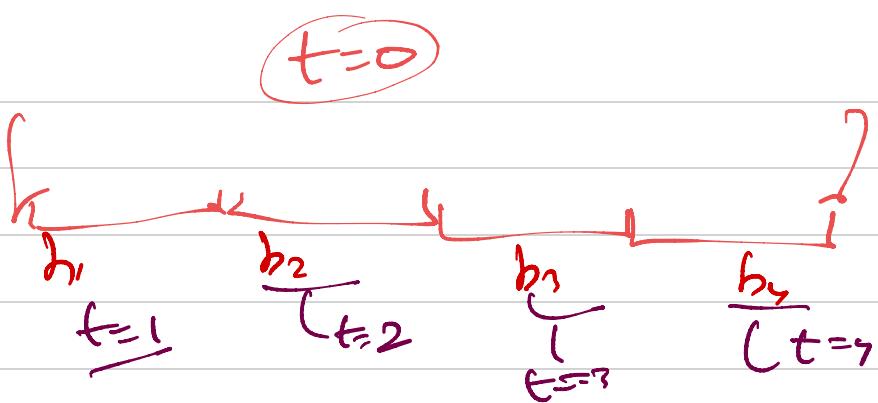
Adam



Algorithm 1: Adam, our proposed algorithm for stochastic optimization. See section 2 for details, and for a slightly more efficient (but less clear) order of computation. g_t^2 indicates the elementwise square $g_t \odot g_t$. Good default settings for the tested machine learning problems are $\alpha = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$. All operations on vectors are element-wise. With β_1^t and β_2^t we denote β_1 and β_2 to the power t .

Require: α : Stepsize → Learning Rate
Require: $\beta_1, \beta_2 \in [0, 1]$: Exponential decay rates for the moment estimates
Require: $f(\theta)$: Stochastic objective function with parameters θ → loss function
Require: θ_0 : Initial parameter vector → weight param
 $m_0 \leftarrow 0$ (Initialize 1st moment vector) → momentum vector
 $v_0 \leftarrow 0$ (Initialize 2nd moment vector) → learning rate
 $t \leftarrow 0$ (Initialize timestep)
while θ_t not converged **do**
 $t \leftarrow t + 1$
 $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep t) → $\frac{\partial L}{\partial w_t}$
 $m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate)
 $v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (Update biased second raw moment estimate)
 $\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$ (Compute bias-corrected first moment estimate)
 $\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$ (Compute bias-corrected second raw moment estimate)
 $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$ (Update parameters)
end while
return θ_t (Resulting parameters)

$$g_t^t = \beta \times g^{t-1} + (1-\beta) \frac{\partial L}{\partial w_t}$$



Train \rightarrow Standardize

$$\frac{x - \bar{M}}{\sigma} \rightarrow Z\text{-score}$$

Fit \rightarrow trainer \rightarrow ~~$M \& G$~~
transforming test