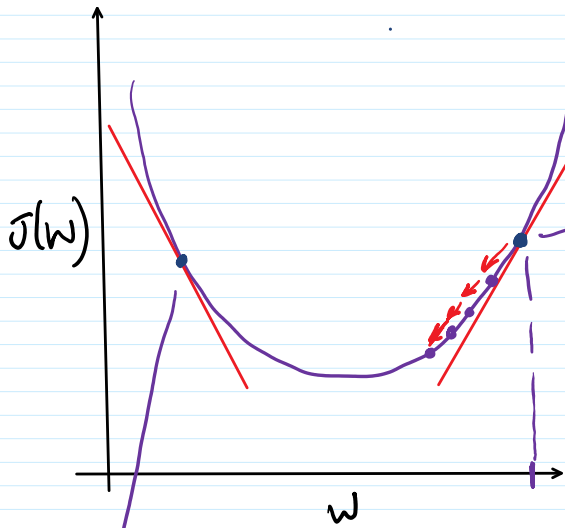


for training faster and getting max/min faster.

Gradient Descent :-



$$w = w - \alpha \underbrace{\frac{dJ(w)}{dw}}_{> 0}$$

$$w = w - \alpha \cdot (+ve)$$

$$w = w - (+ve)$$

w ↓ decreasing

$$w = w - \alpha \underbrace{\frac{dJ(w)}{dw}}_{< 0}$$

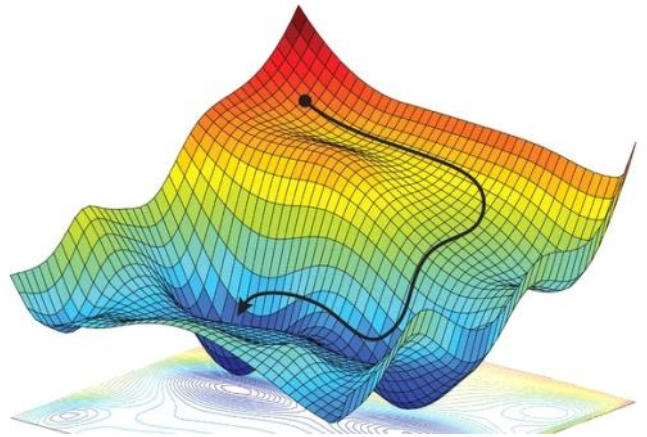
$$w = w - \alpha \cdot (-ve)$$

$$= w - (-ve)$$

w ↑ increasing

||y for "b".

Choosing "α".



In general,

$$\theta_{t+1} = \theta_t - \alpha \cdot \left. \frac{\partial L(\theta)}{\partial \theta} \right|_{\theta = \theta_t}$$

$$\theta_{t+1} = \theta_t - \alpha \cdot \nabla_{\theta} L(\theta)$$

Let MSE loss for images.

$$L(\theta) = \sum_{j=1}^L \sum_{i=1}^M (\hat{y}_{ij} - y_{ij})^2$$

① Batch Gradient Descent :-

The entire training set is used in each iteration.

(slow, too heavy)

② Stochastic Gradient Descent :-

The weights are updated after computing gradients for every training sample.

(random directions)

③ Mini-batch Gradient Descent :-

Update weights using batches from the whole dataset.

$\left(\frac{c}{m} \right)$
of samples # of batches.

① Momentum-based Gradient Descent :-
 Nestron Momentum / (1983)
 [Nestron Accelerated Gradient]

$$\underline{v_t = \gamma v_{t-1} + \alpha \nabla_t L(\theta)}$$

$$\theta_{t+1} = \theta_t - v_t$$

①

$$\Rightarrow \theta_{t+1} = \theta_t - [\gamma v_{t-1} + \alpha \nabla_t L(\theta)]$$

$$\theta_{t+1} = \theta_t - \underbrace{\gamma v_{t-1}} - \underbrace{\alpha \nabla_t L(\theta)}$$

Momentum factor (to get momentum to move faster)

\Rightarrow Exponentially decaying average

\Rightarrow Exhibits oscillations near minimum/maximum.

Modify learning rate instead?
 (α)

\rightarrow reduce ' α ' after every 'n' iterations

\rightarrow Exponential decay
 $\alpha = \alpha_0 e^{-kt}$

$$\alpha = \alpha_0 e^{-kt}$$

② Ada Grad Optimization (Adaptive Gradient)

$$\underline{x_t = x_{t-1} + (\nabla \theta_t)^2}, \quad x_0 \geq 0$$

$$\theta_{t+1} = \theta_t - \frac{\alpha}{\delta + \sqrt{x_t}} \nabla \theta_t$$

the number
to avoid divide by 0.



⇒ As time progresses, x_t will get to a large value, hence, the moment does not happen.

③ Ada M Optimization (Adaptive Moments)

$$b_t = \rho_1 b_{t-1} + (1 - \rho_1) \nabla \theta_t$$

$$0 \leq \rho_1 \leq 1$$

$$x_t = \rho_2 x_{t-1} + (1 - \rho_2) (\nabla \theta_t)^2$$

$$0 \leq \rho_2 \leq 1$$

$$\theta_{t+1} = \theta_t - \frac{\alpha}{\delta + \sqrt{\tilde{x}_t}} \cdot \tilde{b}_t$$

$$\tilde{x}_t = \frac{x_t}{1 - \rho_2^t}, \quad \tilde{b}_t = \frac{b_t}{1 - \rho_1^t}$$