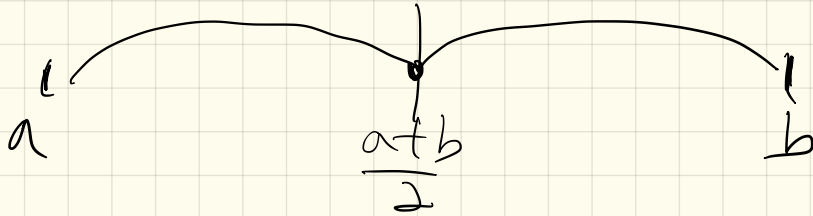


$$\Delta \beta(t+1) = (1+r) \left(\delta \Delta \beta(t) + (1-\delta) \nabla f(t) \right)$$



$$\hat{M}_t \quad \text{vs} \quad M_t$$

$$M_0 = 0$$

$$M_1 = (1 - \eta_1) \nabla L(B(1))$$

$$M_2 = \eta_1 (1 - \eta_1) \nabla L(B(1)) + (1 - \eta_1) \nabla L(B(2))$$

$$M_3 = \underbrace{\eta_1^2 (1 - \eta_1)}_{\uparrow} \nabla L(B(1)) + \underbrace{\eta_1 (1 - \eta_1)}_{\uparrow} \nabla L(B(2)) + \underbrace{(1 - \eta_1)}_{\uparrow} \nabla L(B(3))$$

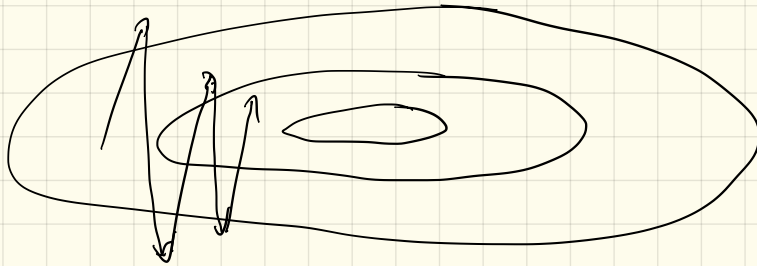
a, b, c

$$\frac{1}{3}a + \frac{1}{3}b + \frac{1}{3}c$$

$$\Rightarrow \frac{(1 - \eta_1)(1 - \eta_1^2)}{1 - \eta_1} = \underbrace{1 - \eta_1^2}$$

Optimization

⇒ GD \Leftarrow 1st order optimization



zig-zag

$$\boxed{\gamma}$$

↑ step size
learning rate.

⇒ 2nd order Newton's method.

$$M_t = \frac{1}{t} \sum_{i=1}^t \nabla L(\beta(i))$$

$$G_t = \frac{1}{t} \sum_{i=1}^t (\nabla L(\beta(i))) \quad \leftarrow \text{element wise}$$



Momentum : M_t

AdaGrad : G_t

RMSProp : G_t

Adam : M_t, G_t

$$\beta(t+1) = \beta(t) - \frac{\gamma M(t)}{\sqrt{G_t}}$$

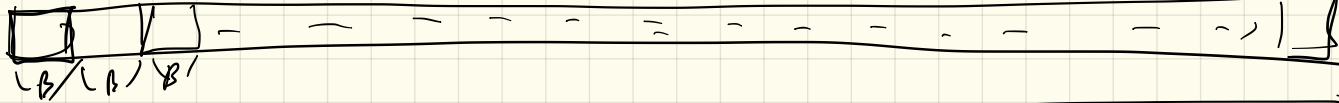
$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, X^2 = \begin{bmatrix} x_1^2 \\ \vdots \\ x_n^2 \end{bmatrix}, \sqrt{X^2} = \begin{bmatrix} \sqrt{x_1} \\ \vdots \\ \sqrt{x_n} \end{bmatrix} \quad \sqrt{G_t} \leftarrow \text{element wise}$$

① random shuffle.

125, 769, 3, 11, - -

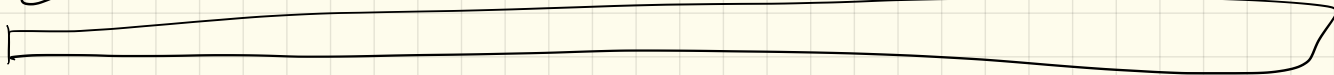
② pick the next B data points

1st 2nd



epoch.

↩



$B \uparrow$

\approx

$r \downarrow$