

# SGDWM

15 September 2023

07:07

## Stochastic gradient descent

with momentum

$$\beta = 0.9$$

Batch

$$w_{new} = w_{old} - \alpha \frac{\partial L}{\partial w_{old}} \quad \left. \begin{array}{l} \text{Deriv of loss} \\ \text{w.r.t. deriv} \\ \text{of weight} \end{array} \right\}$$

$$b_{new} = b_{old} - \alpha \frac{\partial L}{\partial b_{old}} \quad \left. \begin{array}{l} \text{Deriv of loss} \\ \text{w.r.t. deriv} \\ \text{of bias} \end{array} \right\}$$

$$w_{new} = w_{old} - \alpha v_{old} \quad \rightarrow \text{velocity component.}$$

Random

$$\text{Data} \rightarrow \text{NN} \rightarrow \hat{y} \rightarrow \text{Loss} \leftarrow wB$$

$$P = B \times \text{previous Days} + (1-B) \times \text{Today}$$

$$= B \times PD + (1-B) \times T.$$

$$v_{dw} = \underbrace{B \times v_{dw}}_{\text{previous weights}} + \underbrace{(1-B) \times \text{dw}}_{\text{current weight.}}$$

$$\text{①} \quad t_2 \quad t_3 \quad t_4 \quad \dots \quad t_n$$

$$\text{②} \quad \underbrace{q_1}_{\text{rand}} \quad q_2 \quad q_3 \quad q_4 \quad \dots \quad q_n.$$

$$B = 0.8$$

$$v_{t1} = B \times PD + (1-B) \times T.$$

$$= 0.8 \times 0 + (1-0.8) \times q_1$$

$$= 0.2 \times q_1 \quad \checkmark \quad \text{①}$$

$$v_{t2} = B \times PD + (1-B) \times T.$$

$$= B \times (0.2 \times q_1) + (1-B) \times q_2$$

$$= \frac{0.2 \times 0.8 + 0.2 \times q_2}{\checkmark \quad \rightarrow 0.8} \quad \checkmark \quad \text{②}$$

$$v_{t3} = B \times PD + (1-B) \times T.$$

$$= B \times (0.2 \times 0.8 + 0.2 \times q_2) + (1-B) \times q_3$$

$$=$$

$$\beta = 0.8 \rightarrow \text{③}$$

$$\left. \begin{array}{l} 49 \\ 48 \\ 47 \\ 46 \\ 45 \end{array} \right\} \text{50th iter}$$

By using Exponential moving

avg we tend out to avg oscillation.

in vertical direction close to zero.

$$\beta \uparrow \quad \text{smoothing} \uparrow$$

$$B \downarrow \quad \text{smoothing} \downarrow$$

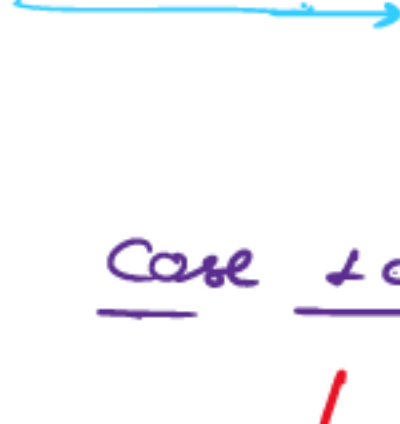
It enables the algorithm to

take straight path in forward

direction to global minima

and dampen out vertical oscillation.

## Adagrad.



$$\begin{array}{ll} \eta & \propto \text{constant} \\ \text{sgd} & \propto \text{constant} \\ \eta \text{sgd} & \propto \text{constant} \\ \text{gdwm} & \propto \text{constant} \end{array}$$

Case 1 step.

$\rightarrow$  dynamic learning rate

$$\alpha = 0 \text{ to } 1$$

$$\alpha = 0 \quad \text{model is not learning anything}$$

$$\alpha = 1 \quad \text{overshooting}$$

$$\alpha = 0.01 \quad \text{preferred}$$

$$w_{new} = w_{old} - \alpha \frac{\partial L}{\partial w_{old}}.$$

$$w_{new} = w_{old} - \alpha$$

$$w_{new} = w_{old}.$$



$$\begin{array}{lll} \text{age} & \text{salary} & \text{gender} \\ 0-120 & 0-\infty & M/F \\ \text{Dense} & \text{Dense} & \text{Sparse} \end{array}$$

$$\checkmark \quad \checkmark \quad \checkmark$$

In pool world of datasets

model learn 1st dense feature

and then sparse feature

$$w_{new} = w_{old} - \alpha \frac{\partial L}{\partial w_{old}}.$$

$$\rightarrow \text{dynamic lr}$$

$$w_{new} = w_{old} - \alpha' \frac{\partial L}{\partial w_{old}}. \quad \checkmark$$

$$\text{New LR: } \alpha' = \frac{\alpha}{\sqrt{nL + \epsilon}} \quad \text{old LR} \quad \text{Epsilon } 10^{-7}, 10^{-9}$$

$$\epsilon \quad \text{small positive no.}$$

used to avoid zero

Division Error.

$$\frac{1}{\sqrt{L}}$$

$$L = \sum_{i=1}^n \left( \frac{\partial L}{\partial w} \right)^2 \quad \checkmark$$

Random

$$0.01 \quad 0.009$$

$$\text{Data} \rightarrow \text{NN} \rightarrow \hat{y} \rightarrow \text{Loss} \leftarrow wB\alpha$$

$$\text{w.p.a. Data} \rightarrow \text{NN} \rightarrow \hat{y} \rightarrow \text{Loss} \leftarrow wB\alpha.$$

$$\alpha' = \frac{\alpha}{\sqrt{nL + \epsilon}}$$

$$\text{①} \quad \sum_{i=1}^1 \left( \frac{\partial L}{\partial w} \right)^2 = 5$$

$$\text{②} \quad \sum_{i=1}^2 \left( \frac{\partial L}{\partial w} \right)^2 + \left( \frac{\partial L}{\partial w} \right)^2 = 10$$

$$\text{③} \quad \sum_{i=1}^3 \left( \frac{\partial L}{\partial w} \right)^2 + \left( \frac{\partial L}{\partial w} \right)^2 + \left( \frac{\partial L}{\partial w} \right)^2 = 15$$

Batch

at

$\alpha'$

$$\begin{array}{llll} \beta_1 & \uparrow & \uparrow \uparrow \uparrow \uparrow & 0.01 \\ \beta_2 & \uparrow \uparrow & \uparrow \uparrow \uparrow & 0.008 \\ \beta_3 & \uparrow \uparrow \uparrow & \uparrow \uparrow & 0.004 \\ \beta_4 & \uparrow \uparrow \uparrow \uparrow & \uparrow & 0.001 \end{array}$$

$$\left. \begin{array}{ll} \beta_1 \rightarrow 5 \\ \beta_2 \rightarrow 10 \\ \beta_3 \rightarrow 15 \end{array} \right\} \quad \text{①}$$

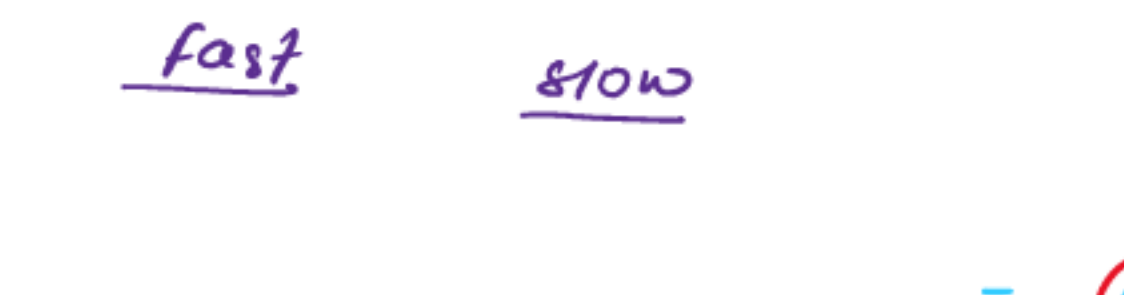
Increasing decreasing

$$\text{Data} \rightarrow \text{NN} \rightarrow \hat{y} \rightarrow \text{Loss} \leftarrow wB\alpha. \quad 0.01$$

$$\text{Data} \rightarrow \text{NN} \rightarrow \hat{y} \rightarrow \text{Loss} \leftarrow wB\alpha. \quad 0.008$$

$$\beta_3$$

$$\beta_4$$



Fast slow

$$= 10 - 0.01 \times 10$$

$$= 10 - 1$$

$$= 9$$

$$= 10 - 0.0001 \times 10$$

$$= 10 - 0.001$$

$$= 9.999$$

① Global minima is not reached

in case of Adagrad.

② Learning rate reduces monotonically

③  $w_{new} = w_{old}.$

Learning rate decay phenomenon.