

Adagrad optimizers

GD uses the same step size (learning rate) for each parameter or I/p feature. This leads to have different amounts of curvature in different ~~step size~~ dimensions. so, algorithm needs a different step size to each parameter.

Need of Adagrad

- Adagrad is an extension of the GD optimization.
- It allows different learning rate (η_i) for each weight.
- In real-world dataset, some I/p features are sparse (most of the features are zero) and some I/p features are dense (most features are non-zero).
- In this scenario, same learning rate (η) for all the weights is not good for optimization.
- So, Adagrad algorithm adapts the learning rate to the parameters.
 - * It is using low learning rates for parameters associated with frequently occurring features.
 - * High learning rates for parameters associated with rare features.
- This is well suitable for the sparse data.

e.g.

During training, some parameters might reach near to converging stage where only fine adjustment is required, but some parameters need to be adjusted significantly due to small no. of matching samples i.e., far away from convergence.

→ The weight update rule for adagrad is

$$w_t = w_{t-1} - \eta_t \frac{\partial E}{\partial w_{t-1}}$$

$$\eta_t = \frac{v}{\sqrt{\alpha_t + \epsilon}}$$

where α_t = different learning rates for each weight at each iteration.

ϵ - small constant to avoid division by 0 error

where $v_0 \rightarrow 0$, $\epsilon = 10^{-8}$

* If α_t becomes 0 then η_t will become 0 which lead to slow convergence

i.e., no change in w .

$$\alpha_t = \sum_{t=1}^T \left(\frac{\partial E}{\partial w_{t-1}} \right)^2$$

$\frac{\partial E}{\partial w_{t-1}}$ is derivative of loss w.r.t. weight and

$\left(\frac{\partial E}{\partial w_{t-1}} \right)^2$ will always be the since it is a square term.

so, α_t is always the due to squaring the gradients, hence $\alpha_t \geq \alpha_{t-1}$

→ From this, α_t and η_t is inversely proportional to one another.

$$\eta_t = \frac{v}{\sqrt{\alpha_t + \epsilon}}$$

→ This implies that when α_t will increase, η_t will decrease.

* So, when the no. of iterations will increase the learning rate will reduce adaptively. Hence, no need to ~~redu~~ select the learning rate manually.

Advantages

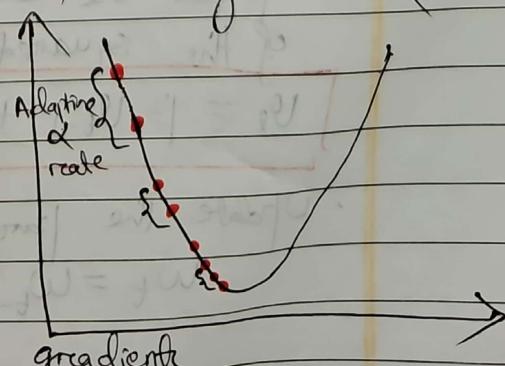
- No manual tuning of the learning rate is required.
- Faster convergence and more reliable.

Disadvantages

when α_t increases, it will decrease when α_t becomes large, increases the no. of iterations. Due to this, α_t will decrease at the larger rate. This will make the old weight almost equal to the new weight which may lead to slow convergence.

② RMSprop optimizer:

- It is an adaptive learning rate optimization algorithm.
- It is a variant of SGD and Adagrad algorithms.
- It improves the convergence speed and stability of the model training process.
- It uses moving average of the squared gradients i.e., $(\nabla_w E(w_t))^2$ to scale (exponentially decay) the learning rate for each parameter.
- It makes the learning rate less than a certain threshold.
- RMSprop is designed to dramatically reduce the computations in training neural networks.
- This helps in reducing the learning rate by forgetting each gradients, and focus on the most recently observed partial gradients.



During the progress of the search, overcoming the limitation of Adam.

→ It leads to smoothly adjust the learning rate for each parameter in the network providing a better performance than batch GD.

→ RMSProp algorithm utilizes exponentially weighted moving averages of squared gradients to update the parameters
ie, $(\nabla_w E(w_t))^2$

RMSProp Algorithm

→ Initialize parameters :

- Learning rate : α , exponential decay rate for averaging B .
- Small constant for numerical stability ϵ , initial parameter value w_0 .

→ Initialize accumulated gradients (Exponentially weighted avg.)

Accumulated squared gradient for each parameter : $v_t = 0$

→ Repeat until convergence or maximum iterations :

• Compute the gradient of the objective function w.r.t. the parameters

$$g_t = \nabla_w E(w_t)$$

• Update the exponentially weighted average of the squared gradients

$$v_t = \beta v_{t-1} + (1-\beta) g_t^2$$

• Update the parameters

$$w_t = w_{t-1} - \alpha \cdot \frac{g_t}{\sqrt{v_t + \epsilon}}$$

where $g_t \rightarrow$ gradient of the loss function w.r.t.
the parameters at time t .

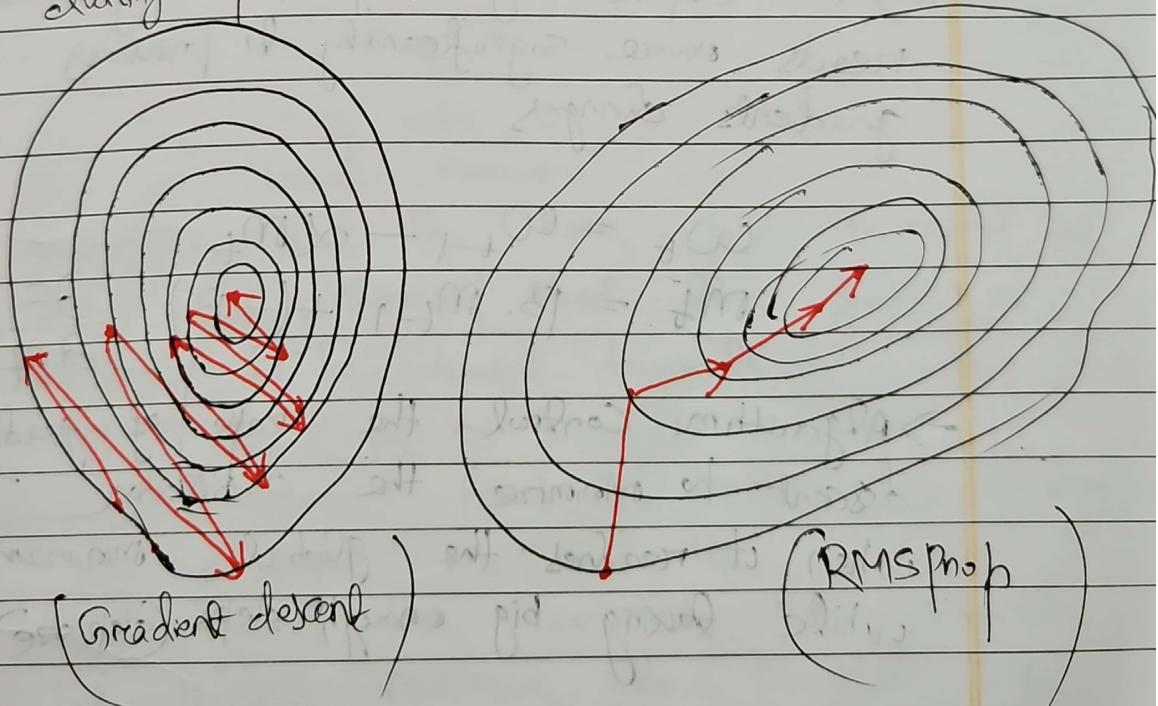
$\beta \rightarrow$ decay factor

$v_t \rightarrow$ exponentially weighted average of
the squared gradients.

$\epsilon \rightarrow$ small constant to prevent division
by 0

→ this process is repeated for each parameter,
and it helps adjust the learning rate for
each parameter based on the historical gradients.

→ The exponential moving average allows $(t-1)$.
the algorithm to give more importance to
recent gradients, and clampen the effect
of older gradients, providing stability
during optimization.



③ Adam (Adaptive Momentum estimation) Optimizer

- Adam is an extension to SGD algorithm
- \hat{g}_t is a 1st order derivative optimization algorithm.
- \hat{g}_t uses the SGD extensions of Adam and RMSprop, to deal with ML and DL problems involving large datasets and high dimensional parameter spaces.

Momentum :

- Momentum is used to accelerate the GD algorithm by considering the exponentially weighted average of the gradients.
- These averages makes the convergence towards the minima faster.
- An exponentially weighted moving average reacts more significantly to previous gradients changes.

$$w_t = w_{t-1} - \alpha m_t$$

$$m_t = \beta_1 \cdot m_{t-1} + (-\beta_1) \cdot \frac{\partial E}{\partial w_t}$$

- Algorithm control the rate of gradient descent to minimize the oscillation when it reaches the global minimum while taking big enough steps (step size)

$$v_t = \beta_2 \cdot v_{t-1} + (1-\beta_2) \left(\frac{\partial E}{\partial w_t} \right)^2$$

- \hat{g}_t makes the algorithm to pass the local minima hurdles in a search space.
- Hence, combine the features of the above methods to reach the global minimum.

Mathematical Aspect of Adam optimizer

→ First momentum update

$$m_t = \beta_1 \cdot m_{t-1} + (1 - \beta_1) \left(\frac{\partial E}{\partial w_t} \right)$$

→ Second momentum update

$$v_t = \beta_2 \cdot v_{t-1} + (1 - \beta_2) \left(\frac{\partial E}{\partial w_t} \right)^2$$

where,

β_1 & β_2 : exponential decay rates average of gradients ($\beta_1 = 0.9$, $\beta_2 = 0.999$)

These are estimating the moments.

→ Exponential Decay → Learning rate schedule where decay (decrease) the α with more iterations using the exponential function:

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

k : decay rate

t : iteration number

since m_t and v_t have both initialized as 0, it is observed that they gain a tendency to be biased towards 0 by both β_1 & $\beta_2 \approx 1$

→ This optimizer fixes this problem by computing bias-corrected m_t and v_t :

this controls the weights while reaching the global minimum to prevent high oscillations when near it.

$$\boxed{\begin{aligned} \hat{m}_t &= m_t \\ &\quad \cdot \frac{1 - \beta_1^t}{1 - \beta_1} \\ \hat{v}_t &= v_t \\ &\quad \cdot \frac{1 - \beta_2^t}{1 - \beta_2} \end{aligned}}$$

→ Intuitively, learning is adapting to the error after every iteration so that it remains controlled and unbiased throughout the training.

→ Now, instead of our normal weight parameters m_t and v_t , we take the bias-corrected weight parameters \hat{m}_t and \hat{v}_t .

→ Use these into general equation of

$$w_t = w_{t-1} - \hat{m}_t \frac{\alpha}{\sqrt{\hat{v}_t + \epsilon}}$$

ϵ → small the constant to avoid division by zero error when $v_t \rightarrow 0$, $\epsilon = 10^{-8}$

Algorithm :- Adam

β_1 and β_2 = decay rates of average of momentum and RMSprop methods ($\beta_1=0.9$, $\beta_2=0.999$).

w_0 : Initial parameter vector

$f(w)$: stochastic objective function

→ Initialise $m_0 = 0$ // First momentum term.

- $v_0 = 0$ // 2nd moment term.
- $t = 0$ // timestep

Repeat {

$$t = t + 1$$

$$g_t = \nabla f(w) \cdot f_t(w_{t-1})$$

$$m_t = \beta_1 \cdot m_{t-1} + (1-\beta_1) \cdot g_t \quad // \text{update 1st moment estimate}$$

$$v_t = \beta_2 \cdot v_{t-1} + (1-\beta_2) \cdot g_t^2 \quad // \text{update 2nd raw moment estimate}$$

$$\hat{m}_t = \frac{m_t}{1 - p_1^t} \quad // \text{bias-corrected 1st moment estimate}$$

$$\hat{v}_t = \frac{v_t}{1 - p_2^t} \quad // \text{bias-corrected 2nd moment estimate}$$

$$w_t = w_{t-1} - \hat{m}_t \left(\frac{\alpha}{\hat{v}_{t+e}} \right) \quad // \text{update the parameters}$$

→ until w_t reach convergence

①

ML

Data:- Requires less amount of data.

H/w requirements:- work on low-end machines

Training:- Takes less time.

Data:- Fails to interpret the interpretation of the result.

Feature engineering:- Identifies by expert Hand-coded by experts as per domain and database.

③ Ideal value for $\eta \rightarrow 0.40$

(ii) Application of DL:

- Image recognition
- Natural language processing
- Medical diagnosis
- Self driving cars
- Information retrieval

DL

Requires large amount of data.

Heavily depend on high-end machines.

Takes more time for training.

Easy to interpret the result.

Learn automatically high-level features from data.

- ④ Can handle structured & unstructured data
- High degree of accuracy
- Automatic feature extraction
- or models can generalize well to new unseen classes.