

Understanding Optimization

ALGORITHMS

Let us quantify the quality of any particular set of weights w . The goal of optimization is to find " w " such that it minimize the loss function.

- ① Finding set of param i.e w corresponding to global minima.

Strategy 01: A first VERY BAD idea try out many different values randomly. keep the track of the best. Procedure look like below

```
# assume X-train (each column) an example (3073 x 50,000)  
# assume Y-train labels (1D array of 50,000)  
# L evaluate my loss function
```

```
bestLoss = float("inf")  
for num in range(1000):  
    w = np.random.randn(...)  
    loss = L(X-train, Y-train, w)  
    if loss < bestLoss:  
        bestLoss = loss  
        bestW = w
```

This approach gives 15% accuracy, whereas completely Random will give ~10% accuracy.

Strategy 02: (Random local search)

Try to take small step in random direction, take further step only if cost is getting lesser.

$w = np.random(\dots)$

best loss = float('inf')

Step size = .0001

for(
 $x=1 \rightarrow 10^{10}$)
 $w_{try} = w + np.random(\dots) * \text{Step_size}$

loss = $L(X_{\text{tra}}, Y_{\text{tra}}, w_{try})$

if loss < best loss:

$w = w_{try}$

best loss = loss

classification accuracy $\approx 21.4\%$.

Strategy 03: - Stop guessing direction use "Gradient"

* The gradient tells us direction in which function

has the steepest "Rate of Increase". So to mini
-mize we move in opposite (so -ve direction)

* Gradient tells only the direction not the step size.

Gradient Descent:- Compute gradient loss

of the loss function. perform the parameter update.

while True :

weights-grad = evaluate-grad(loss-fn, data, weights)

weights += - step-size * weights-grad

⊗ Now it's clear step-size = learning-rate.

Mini-Batch-Gradient-Descent:- Perform the calculation

over small batches rather than complete dataset.

while True :

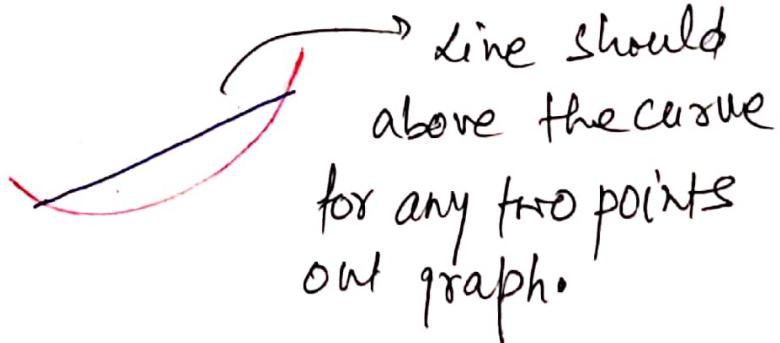
data-batch = sample-training-data(data, 256) $\xrightarrow{\text{L} \rightarrow \text{batch size}}$

weights-grad = evaluate-gradient(loss-fn, data-batch, wts)

weights += - weights-gradient * step-size.

challenges for plain Mini-Batch-Gradient Descent

Convex Function:-



Saddle point:- Points where one dimension

slopes up & another goes down. On these points Gradient is almost ≈ 0 .

④ Refer: [blog.paperspace.com/intro to optimi
zation in deep learning - gradient descent](http://blog.paperspace.com/intro-to-optimization-in-deep-learning-gradient-descent/)

⑤ Do visit:- Complicated landscapes

[cs.umd.edu/~ntomg/projects/landscapes]

→ Surrounded by a plateau.

Plateau := Also called high plain
or tableland, is an area of
highland.



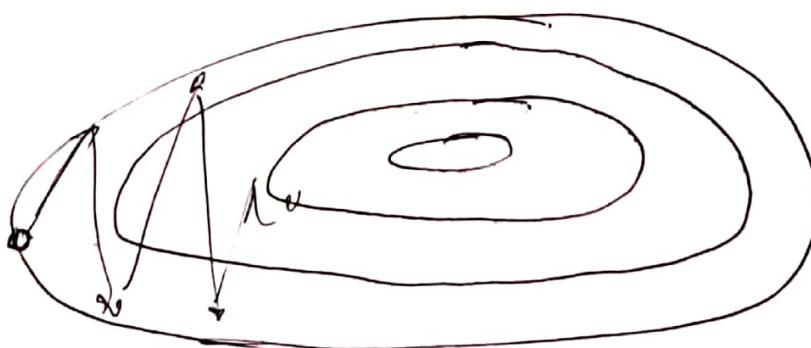
P₀ T₀ '6

Challenges Continue:-

Same learning rate applies to all parameters. If data is sparse & features have updates. If data is sparse & features have different frequencies, we might not want to update all of them with same extent. perform larger update to rarely occurring features.

Gradient Descent + Momentum :- The momentum

term increases for dimensions whose gradients point in the same direction & reduces updates for dimensions whose gradient direction changes. As a result we converge faster & reduced oscillation. Let's take it by an example. Visualize below loss function.



↓ what we want is
↓ slow learning rate VERTICAL

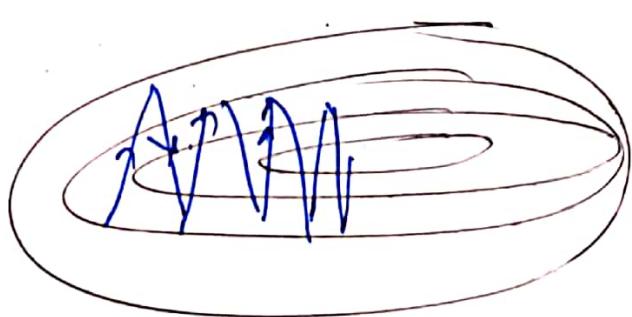
↔ fast learning rate in HORIZONTAL.

(P.T.D)

$$\text{So } V_t = V_{t-1} \gamma + V_{t-1} + \eta \nabla_\theta J(\theta) \quad (\text{param smooth})$$

$$\theta = \theta - V_t \leftarrow (\text{paramUpdate}) \quad \gamma = 0.9$$

So by averaging over last iterations (gradients in the vertical direction will be ≈ 0 for prev example) (bcz in ONE iter it will go up & say +ve gradient in other go down -ve gradient hence averagin it will ≈ 0) for vertical. Where as in horizontal gradient is almost always fve) hence oscillation will reduce in that direction it will go faster horizontally with less oscillation.



Before Momentum



After Momentum.

$\gamma = 0.9$ most of the cases.

Adagrad: — It adapt learning rate to the parameters, performing larger updates for infrequent parameters & smaller updates for frequent parameters. This is well-suited for to deal with SPARSE DATA.

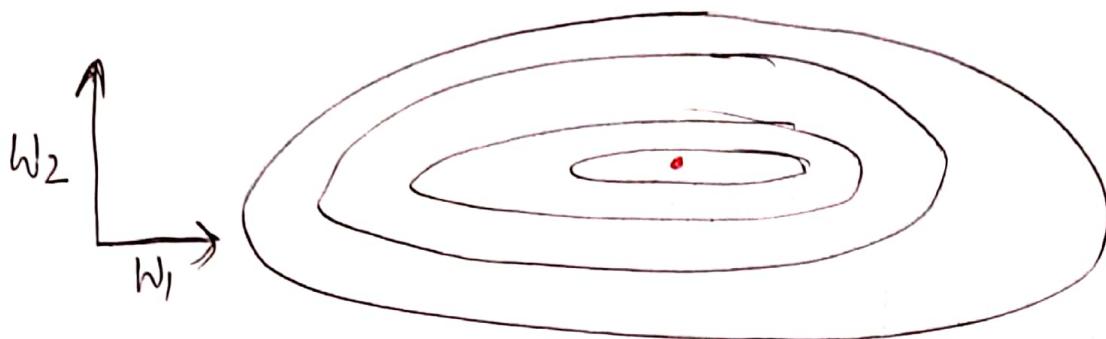
$g_{t,i} = \nabla_{\theta_t} J(\theta_{t,i}) \rightarrow$ gradient of objective function w.r.t to the parameter θ_i at the time step t !

$$\theta_{t+1,i} = \theta_{t,i} - \eta \cdot g_{t,i}$$

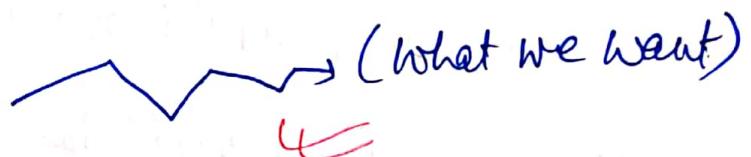
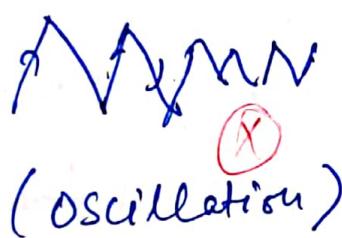
F.T.O
=

RMSProp

α



Assume above is our loss function & we want to make learning faster in (w_1) direction (w_2) slow to avoid lot of oscillation.



On iteration ' t ' will compute $d\mathbf{w}_1$, $d\mathbf{w}_2$ on minibatch

$$S_{d\mathbf{w}_1}^t = \beta S_{d\mathbf{w}_1}^{t-1} + (1-\beta) d\mathbf{w}_1^2 \quad (\text{element wise square})$$

$$S_{d\mathbf{w}_2}^t = \beta S_{d\mathbf{w}_2}^{t-1} + (1-\beta) d\mathbf{w}_2^2$$

$$\mathbf{w}_1 = \mathbf{w}_1 - \alpha \frac{d\mathbf{w}_1}{\sqrt{S_{d\mathbf{w}_1}} + \epsilon}$$

$$\mathbf{w}_2 = \mathbf{w}_2 - \alpha \frac{d\mathbf{w}_2}{\sqrt{S_{d\mathbf{w}_2}} + \epsilon}$$

$\epsilon \approx$ very small number to avoid divide by ZERO $\approx 10^{-8}$

Intuition why it Works :-

In horizontal direction we want fast learning,
vertical we want slow down.

So we dW_1^2 will be relatively small $\sqrt{\sum dW_1^2}$
whereas dW_2^2 will be relatively big $\sqrt{\sum dW_2^2}$
by dividing by big number we \therefore slow down
in vertical direction (W_2). Hinton suggested
 $\beta = 0.9$ good default value for learning rate = .001.

Adam Optimization

It is a combination of (Gradient Descent with Momentum + RMS Prop) Algorithm.

Steps

$$V_{dw1}, S_{dw1}, V_{dw2}, S_{dw2} = 0, 0, 0, 0$$

on iteration 'T' compute dw_1 & dw_2 with current mini-batch.

$$\left. \begin{array}{l} V_{dw1} = \beta_1 \times V_{dw1} + (1 - \beta_1) \times dw_1 \\ V_{dw2} = \beta_1 \times V_{dw2} + (1 - \beta_1) \times dw_2 \end{array} \right\} \begin{array}{l} \text{momentum like} \\ \text{update} \end{array}$$

$$\left. \begin{array}{l} S_{dw1} = \beta_2 \times S_{dw1} + (1 - \beta_2) dw_1^2 \\ S_{dw2} = \beta_2 \times S_{dw2} + (1 - \beta_2) dw_2^2 \end{array} \right\} \begin{array}{l} \text{RMS prop like} \\ \text{update} \end{array}$$

Bias Correction

$$V_{dw1}^{\text{corrected}} = V_{dw1} / (1 - \beta_1^t)$$

$$V_{dw2}^{\text{corrected}} = V_{dw2} / (1 - \beta_1^t)$$

$$S_{dw1}^{\text{corrected}} = S_{dw1} / (1 - \beta_2^t)$$

$$S_{dw2}^{\text{corrected}} = S_{dw2} / (1 - \beta_2^t)$$

P.T.O

$$w_1 = w_1 - \alpha * \frac{V_{dw_1}^{\text{Corrected}}}{\sqrt{S_{dw_1}^{\text{Corrected}} + \epsilon}}$$

α = learning rate

$$\beta_1 = .9$$

$$\beta_2 = .99$$

$$\epsilon \approx 10^{-8}$$

$$w_2 = w_2 - \alpha * \frac{V_{dw_2}^{\text{Corrected}}}{\sqrt{S_{dw_2}^{\text{Corrected}} + \epsilon}}$$

11

RADAM

- ① Recently many new methods have been proposed to accelerate optimization by applying "adaptive learning rate". Example ADAGRAD (2011) - RMSProp
- ② RADAM proposed new variant of ADAM, which ~~eg~~ explicitly rectifies the variance of the "adaptive learning rate" based on derivation.
- Contribution: — Convergence issue is due to undesirable big variance of the adaptive learning rate in the early stage of model training.

③ Proposing RADAM as solution.

$$\phi(q_1, q_2, \dots, q_t) = (1 - \beta_1) \sum_{i=1}^t \beta_1^{t-i} q_i$$

$\psi(q_1, q_2, \dots, q_t) = \frac{1 - \beta_2^t}{1 - \beta_2} \sqrt{\epsilon + \left((1 - \beta_2) \left(\sum_{i=1}^t \beta_2^{t-i} q_i^2 \right) \right)}$

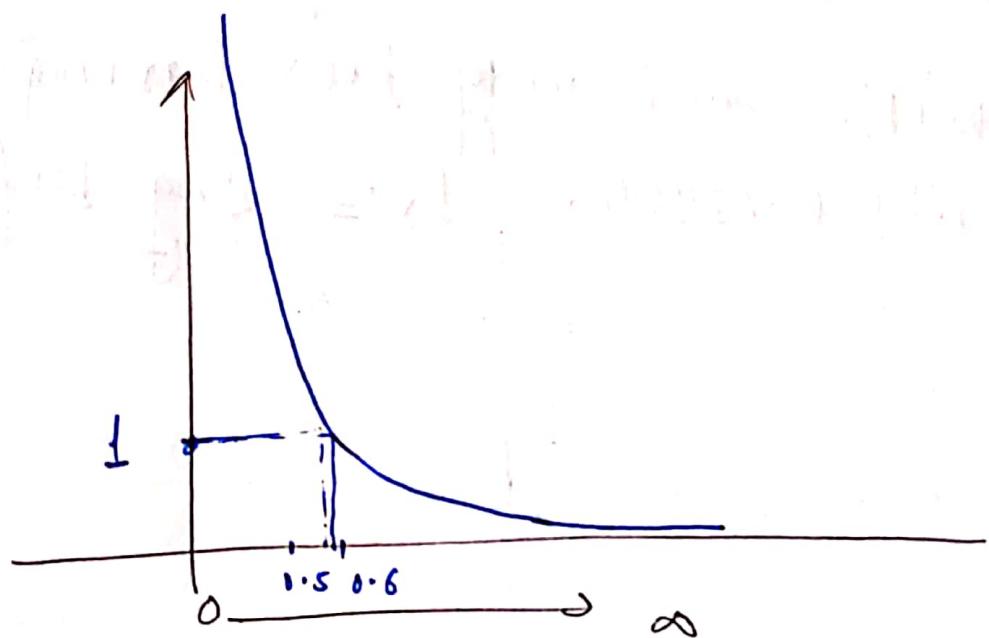
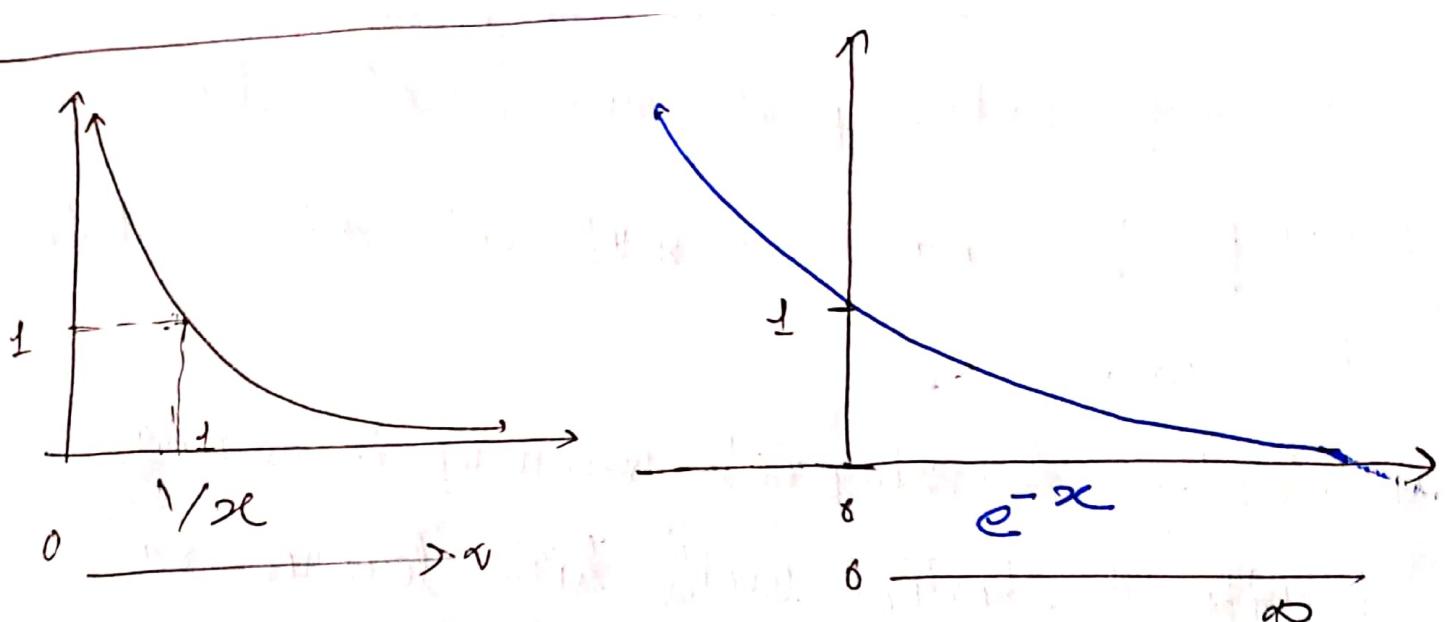
ψ = RMS type update.

ϕ momentum type update.

Learning Rate Warm Up :- Instead of setting learning rate (α_t) as constant or in decreasing order. Learning Rate Warmup strategy sets α_t as some small value.

Example:- Linear warmup set $\alpha_t = t\alpha_0$ $t < T_w$.

= "My Notes" :- Learning rate warm up is saying inspite of starting with high learning rate start with some very low learning rate + gradually increase. Ex:- Jog before Run.



$$\underline{\underline{(\bar{x}' \bar{e}^x)}}$$

Variance of ADAPTIVE RATE

"Due to lack of samples in the early stage, the adaptive learning rate has an undesirably large variance."

Say $t=1$, $\psi(q_1) = \sqrt{q_2}$, we view $\{q_1, q_2, \dots, q_t\}$ as random variables drawn from Normal Distribution $N(0, \sigma^2)$. Therefore $\frac{1}{q_2}$ is subject to scale

Scale-inv- $\chi^2(1, 1/q_2)$. Note that

$$\text{Var}[\sqrt{1/q_2}] \propto \int_0^\infty x^2 e^{-x} dx.$$

adaptive learning rate can have very

high variance ~~at~~ first stage of learning

& setting small learning rate can reduce

Variance \Rightarrow why ↓

$$\text{Var}[dx] = \alpha^2 \text{Var}[x]$$

= ==

Therefore it is unbounded variance that cause problem.

Estimation of P or f_t :-

find the length is SMA(Simple Moving Average)
such it has same "center of mass" with EMA(Exponential Moving Average). In other length $f(t, \beta_2)$ satisfies below equation.

$$\frac{(1-\beta_2) \sum_{t=1}^t \beta_2^{t-1} (t+1-i)}{1 - \beta_2^t} = \sum_{i=1}^t \frac{(t+1-i)}{f(t, \beta_2)}$$

solve above equation :-

$$f(t, \beta_2) = \frac{2}{1-\beta_2} - 1 - \frac{2t\beta_2^t}{1-\beta_2^t} = (P_t) \leftarrow \text{refer}$$

Rectification Term :-

$$\gamma_t = \frac{(p_t - 4)(p_t - 2)p_0}{(p_0 - 4)(p_0 - 2)p_t}$$

< RADAM >

while $t = \{1, \dots, T\}$ do

$$g_t \leftarrow \Delta_\theta f_t(\theta_{t-1}) \text{ (gradient at } t)$$

$$v_t \leftarrow \beta_2 v_{t-1} + (1-\beta_2) g_t^2 \text{ (2nd moment)}$$

$$m_t \leftarrow \beta_1 m_{t-1} + (1-\beta_1) g_t \text{ (1st moment)}$$

$$\hat{m}_t \leftarrow m_t / (1-\beta_1^t) \text{ (Bias correction)}$$

$$\text{Calculate } p_t \\ \leftarrow p_0 - 2t \beta_2^t / 1 - \beta_2^t$$

if $p_t > 4$:

$$\hat{v}_t \leftarrow \sqrt{v_t / (1 - \beta_2^t)} \text{ (Bias correction)}$$

$r_t \leftarrow$ calculate variance rectification

$$\theta_t \leftarrow \theta_{t-1} - \alpha r_t \hat{m}_t / \hat{v}_t \text{ (with adaptive moment)}$$

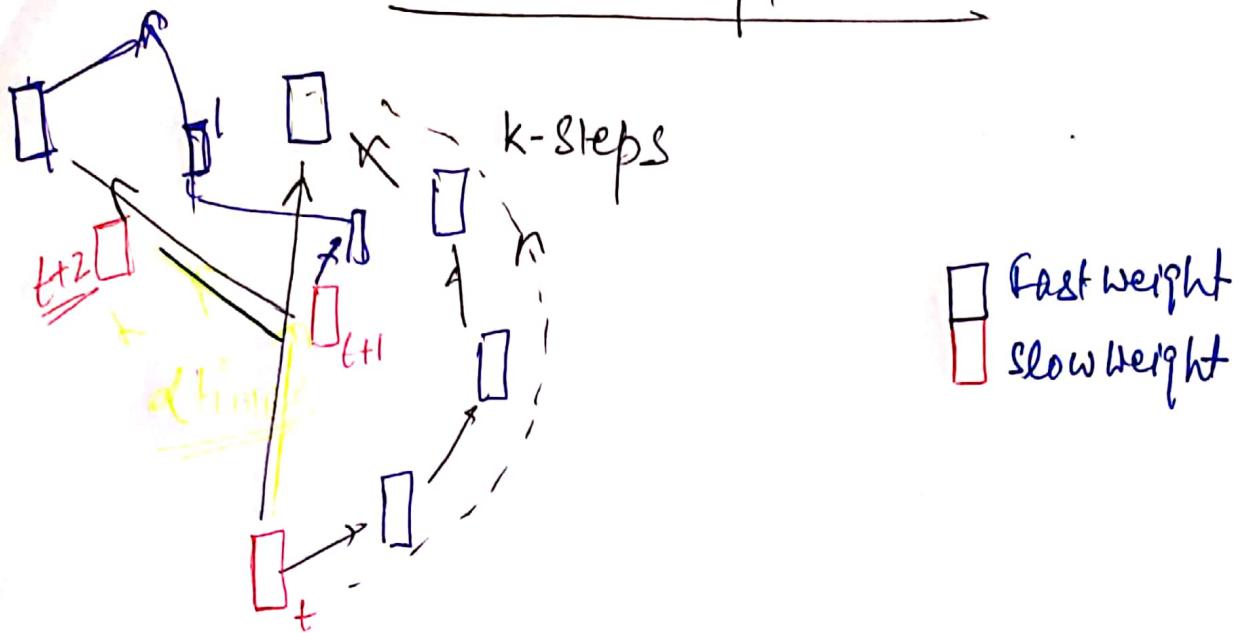
else:

$$\cancel{\theta_t \leftarrow \theta_{t-1} - \alpha r_t \hat{m}_t}$$

$$\theta_t \leftarrow \theta_{t-1} - \alpha r_t \hat{m}_t$$

(unadapted moment).

LookAhead Optimizer



RADAM

- ④ Recently many new methods have been proposed to accelerate optimization by applying "adaptive learning rate". Example ADAGRAD (2011).
- ④ RMSProp
- ④ RADAM proposed new variant of ADAM, which ~~eg~~ explicitly rectifies the variance of the "adaptive learning rate" based on deviation.
- ④ Contribution — Convergence issue is due to undesirable big variance of the adaptive learning rate in the early stage of model training.

④ Proposing RADAM as solution.

$$\phi(q_1, q_2, \dots, q_t) = (1 - \beta_1) \sum_{i=1}^t \beta_1^{t-i} q_i$$

$\psi(q_1, q_2, \dots, q_t) = \frac{1 - \beta_2^t}{1 - \beta_1^t}$

$\psi = \sqrt{\epsilon + \left(1 - \beta_2\right) \left(\sum_{i=1}^t \beta_2^{t-i} q_i^2\right)}$

ψ = RMS type update.

ϕ momentum type update.

Learning Rate Warm Up :- Instead of setting learning rate (α_t) as constant or in decreasing order. Learning Rate Warmup strategy sets α_t as some small value.

Example:- Linear warmup set $\alpha_t = t \alpha_0$
 $t < T_w$.

= My Notes! - Learning rate warm up is saying
inspite of starting with high learning rate
start with some very low learning rate
& gradually increase. Ex! = Jog before Run.

Variance of ADAPTIVE RATE

"Due to lack of samples in the early stage, the adaptive learning rate has an undesirably large variance."

Say $t=1$, $\Psi(g_1) = \sqrt{\frac{1}{g_1^2}}$, we view $\{g_1, g_2, \dots, g_t\}$ as random variables drawn from Normal Distribution

$N(0, \sigma^2)$. Therefore $\frac{1}{g_1^2}$ is subject to scale

Scale-inv- $\chi^2(1, \frac{1}{g_1^2})$. Note that

$$\text{Var}\left[\sqrt{\frac{1}{g_1^2}}\right] \propto \int_0^\infty x^{-1} e^{-x}.$$

adaptive learning rate can have very

high variance ~~at~~ first stage of learning

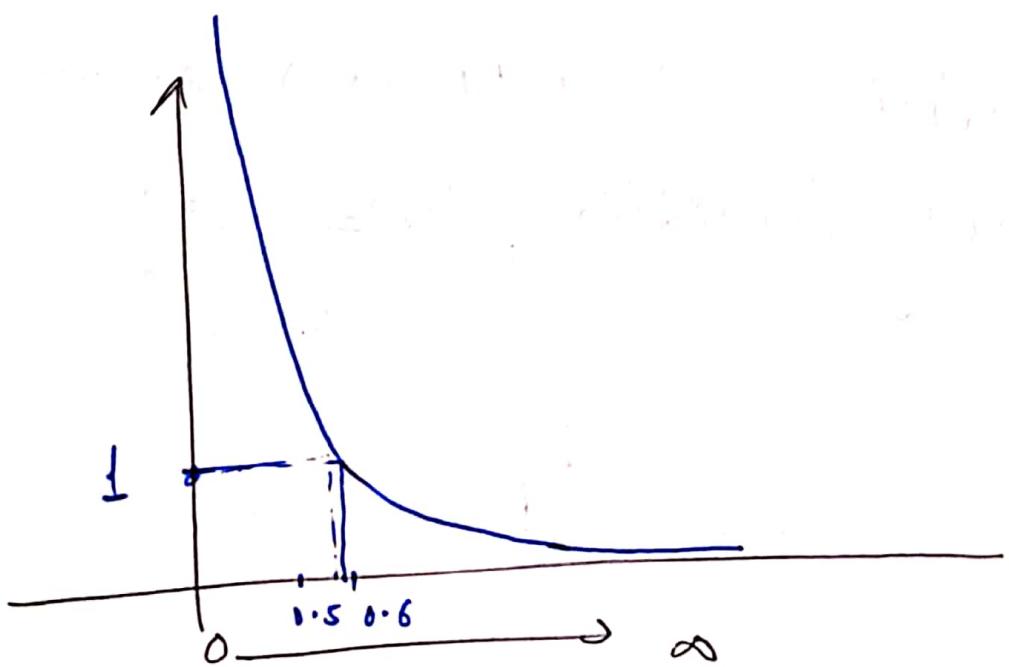
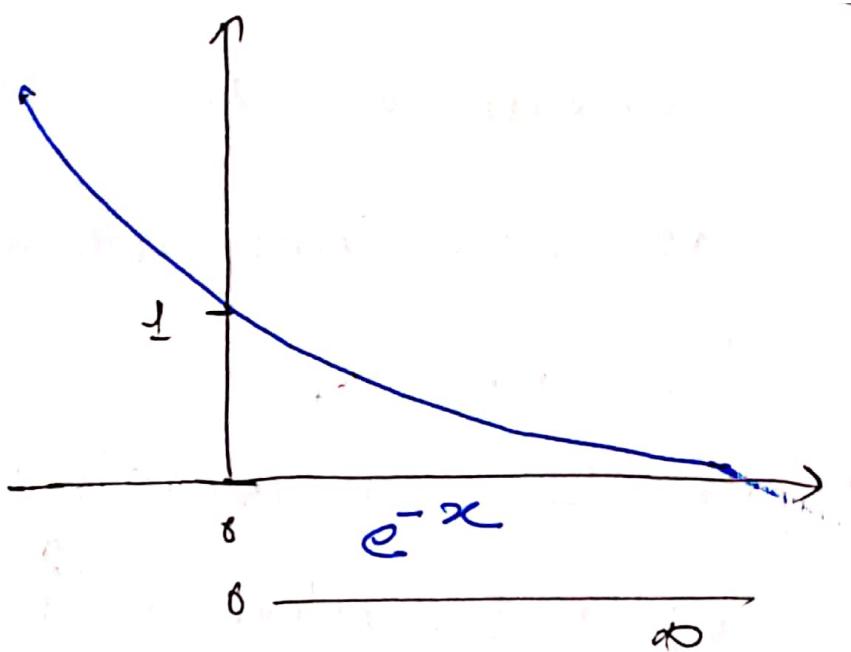
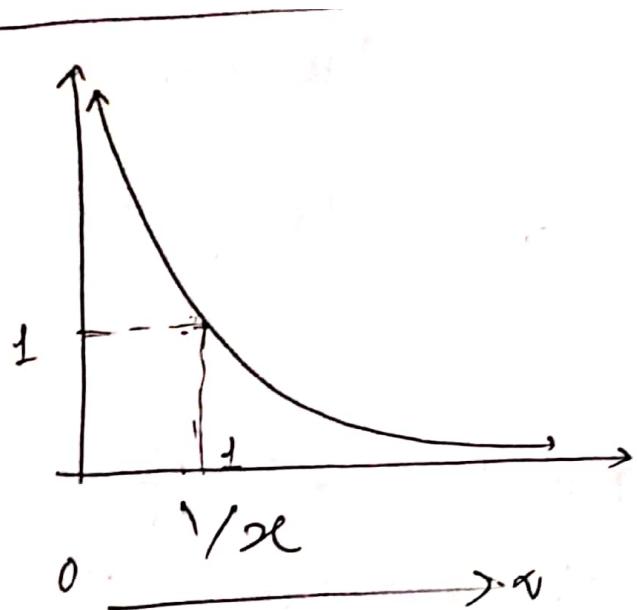
& setting small learning rate can reduce

Variance \Rightarrow why ↓

$$\text{Var}[\alpha x] = \alpha^2 \text{Var}[x]$$

= _____

Therefore it is unbounded variance that cause problem.



$$\underline{\underline{(x^{-1}e^{-x})}}^{(x^{-1}e^{-x})}$$

Estimation of P_t or f_t :-

find the length is SMA(Simple Moving Average)
such it has same "center of mass" with EMA(Exponential Moving Average). In other length $f(t, \beta_2)$ satisfies below equation.

$$\frac{(1-\beta_2) \sum_{t=1}^t \beta_2^{t-1} (t+1-i)}{1 - \beta_2^t} = \sum_{i=1}^t \frac{(t+1-i)}{f(t, \beta_2)}$$

Solve above equation :-

$$f(t, \beta_2) = \frac{2}{1-\beta_2} - 1 - \frac{2t\beta_2^t}{1-\beta_2^t} = (P_t) \leftarrow \text{refer as}$$

Rectification Term:-

$$r_t = \frac{(p_{t-4})(p_{t-2})p_0}{(p_0-4)(p_0-2)p_t}$$

RADAM

while $t = \{1, \dots, T\}$ do

$$g_t \leftarrow \Delta_\theta f_t(\theta_{t-1}) \quad (\text{gradient at } t)$$

$$v_t \leftarrow \beta_2 v_{t-1} + (1-\beta_2) g_t^2 \quad (2^{\text{nd}} \text{ moment})$$

$$m_t \leftarrow \beta_1 m_{t-1} + (1-\beta_1) g_t \quad (1^{\text{st}} \text{ moment})$$

$$\hat{m}_t \leftarrow m_t / (1 - \beta_1^t) \quad (\text{Bias correction})$$

Calculate p_t

$$\leftarrow p_0 - 2t \beta_2^t / (1 - \beta_2^t)$$

if $p_t > 4$:

$$\hat{v}_t \leftarrow \sqrt{v_t / (1 - \beta_2^t)} \quad (\text{bias correction})$$

$r_t \leftarrow$ calculate variance rectification

$$\theta_t \leftarrow \theta_{t-1} - \alpha r_t \hat{m}_t / \hat{v}_t$$

(with adaptive moment)

else:

$$\cancel{\theta_t \leftarrow \theta_{t-1} - \alpha r_t \hat{m}_t}$$

$$\theta_t \leftarrow \theta_{t-1} - \alpha r_t \hat{m}_t$$

(unadapted moment).