# Adagrad Optimization - SPSSe

$$W_n = W_0 - \eta \sqrt{\frac{\partial L}{\partial w_0}}$$

## Adaptive Gradient Algorithm

→ learning rates

→ sparse

$\frac{1}{0} = \infty$

$$W_{t+1} = W_t - \frac{\eta}{\sqrt{\sum_{i=1}^{n} g_i^2 + \epsilon}} * \frac{\partial L}{\partial W_t}$$

gradient @ current iteration

very very small No. $= 10^{-8}$ ⇒ more 0.00 f non-zero = Dense

Text Data ← OHE

BOW/TF-IDF/
N-GRAMS/Word2ve
Glove

| 0 | 0 | 0 | 1 | 0 |
|---|---|---|---|---|
| 0 | 0 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 2 | 0 |
| 0 | 0 | 3 | 0 | 0 |

| 1 | 2 | 3 | 4 |
|---|---|---|---|
| 5 | 6 | 0 | 7 |
| 8 | 9 | 3 | 4 |
| 0 | 1 | 3 | 5 |

⇒ More No. of Zeros = Sparse

Ravinder raised hand

very very smarly

BackProp

$$w_t = w_{t-1} - \eta' \left( \frac{\partial L}{\partial w_0} \right) 2$$

$$\eta' = \frac{\eta}{\sqrt{\sum_{i=1}^{D} \alpha_i + \epsilon}}$$

0.001

$\frac{1}{10}$  $\frac{1}{100}$  $\frac{1}{1000}$  $\frac{1}{10000}$  1 to 22

$\alpha = 1, 1.5, 0.5, 0.6, 2, 2.5 \cdots$

$\frac{1}{10} \downarrow$

Smaller

$$658 \cdot \sum_{i=1}^{n} (1 + 1.5 + 0.5 + 0.6 + 2 + 2.5 \cdots)$$

$$\eta' = 0.00000001 * 2$$

$$= 0.00000002$$

Webinar Chat

Read twice but difficult to digest..Anyway read it again

PRAMOD K. to Hosts and panelists
PK  yes

abhishek to Hosts and panelists
A   what is alpha?
    yes sir
    no
    1/10

Urmil Shah to Hosts and panelists
US  The idea of very small eta value is to address exploding gradient problem but if it gets very small it never reach local minima? is my understanding correct.
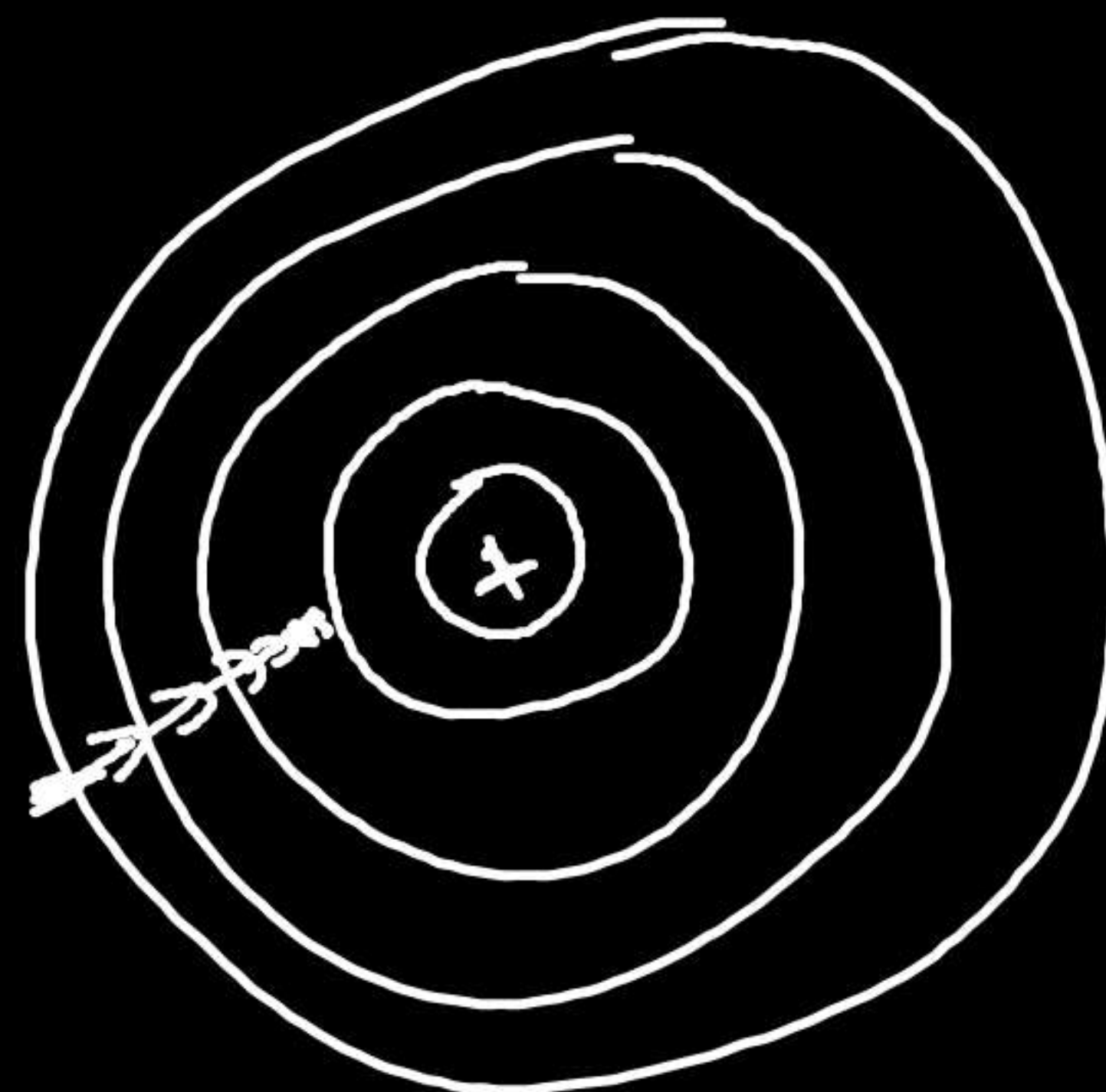
PRAMOD K. to Hosts and panelists
PK  Bigger denominator means very small value will receive

Who can see your messages? Recording on

To:  Hosts and panelists

Type message here...

2 / 2

Ravinder raised hand

Learnvista Priva...

Learnvista Private Limited

$$W_n = W_0 - \eta' \frac{\partial L}{\partial W_0}$$

$$\eta' = \frac{\eta}{\sqrt{\sum_{i=1}^{n} \alpha_i} + \epsilon}$$

$$W = 50$$
$$W_1 = 48$$
$$W_2 = 53$$
$$W_4 = \frac{}{} = \frac{50.000001}{49.999555} \quad 0.00060$$

**Webinar Chat**

A — what is alpha?

yes sir

no

1/10

Urmil Shah to Hosts and panelists
US — The idea of very small eta value is to address exploding gradient problem but if it gets very small it never reach local minima? is my understanding correct.

PRAMOD K. to Hosts and panelists
PK — Bigger denominator means very small value will receive

abhishek to Hosts and panelists
A — but how does it help to solve sparse data?

Ravinder to Hosts and panelists
R — How is the formula of Alpha ?

Who can see your messages? Recording on

To: Hosts and panelists ⌄

Type message here...

3 / 3

# RMSProp → Root mean Squared Propagation

↳ extension of gradient descent and the

sparse
+
Dense

Adagrad version of GD that uses a decaying average of partial gradient in the adaptation of the step size for each parameter.

## RMSprop

$$V_t = \beta V_{t-1} + (1-\beta)(\nabla w_t)^2$$

, $\beta = 0.55$

, exponential
decaying Avg

4 / 4

$$W_1 = 1, 2, 3, 4, 5, 6, 7, 8, 9, 10$$

$$\sum W_i = 55$$

$$(W_i)_{Avg} = \frac{55}{10} = 5.5$$

$$\frac{1}{5.5}$$

$$\frac{\partial L}{\partial W_0}$$

$$\nabla W_t^2$$

$$W_{t+1} = W_t - \frac{\eta}{\sqrt{V_t + \epsilon}} \sqrt{(W_t)^2}$$

Avg current Broelup

exp

---

**Webinar Chat**

**Urmil Shah** to Hosts and panelists
Thanks

**Usha Kumari** to Hosts and panelists
but where is squaring happening in RMS prob?

**abhishek** to Hosts and panelists
yes sir

**Rohan** to Hosts and panelists
Yes Clear

**Usha Kumari** to Hosts and panelists
yes sir clear

**Urmil Shah** to Hosts and panelists
Understood Sir

**Rohan** to Hosts and panelists
Adagrad - Can handle Sparse Matrix Only. RMS Prop - both Sparse and Dense

Who can see your messages? Recording on

To: Hosts and panelists

Type message here...

# RMSprop - DIS-adv → NO

2014 - Adam - most significant

2024 - Almost 10 years completed

Adaptive Momentum Estimation,
→ it requires less memory
→ large data set
→ combination of Momentum + RMSprop $\beta = 0.9$

$$\frac{1}{1-\beta}$$

Adam :-

SGD/MBGD/BGD

Momentum

NAG

Adagrad

↳ RMSprop/Adadelta → $W_n = W_0 - D' \frac{\partial L}{\partial W_0}$,

most of the algorithm

ANN/MLP/CNN/RNN/LSTM/GRU etc.

↳ Adam

→ Adam ⟨ Dense
         Sparse

→ RMSprop ⟩ combine
→ Momentum ⟩ merge

DIS-

$$W_{t+1} = W_t - \frac{\eta}{\sqrt{V_t + c}} \times m_t$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

← Formula

$$\hat{V}_t = \frac{V_t}{1 - \beta_2^t}$$

$$\beta_1 = 0.9$$

$$\beta_2 = 0.99 \quad - \text{ by default}$$

-Keras

where,

$$m_t = \beta_1 m_{t-1} + (1-\beta) \nabla W_t \quad - \text{ momentum}$$

$$V_t = \beta_2 V_{t-1} + (1-\beta_2)(\nabla W_t)^2 \quad - \text{ RMSprop}$$

How to improve a neural Network

-> fine tuning Neural Network hyperparameters.

epochs   # hidden layer   # Neurons   Learning Rate   Activation

Optimization

-> By Solving Problem   -> vashing Gradient Prob

-> Less data -> Slow training

-> overfitting

→ NO. of hidden layer

→ NO. of Neurons per layer

→ learning rate

→ best optimization — Adam/RMSprop

→ Batch Normalization | early stopping | Regularization | Dropout

→ batch size
- BGD
- SGD
- MBGD

→ Activation

→ epochs