

## [NESTED LEARNING]

PAGE No:



Tundarani

DATE:

### \* CORE IDEA :

- In deep learning, there is multiple layers of successive memory compression, i.e depth  $\neq$  deep computation
- In NL  $\rightarrow$  multiple optimisers nested inside each other. Every layer, optimizer, memory unit has its own optimization problem.
  - .. It is a stack of optimizers & not functions

#### ① Deep Optimizers

↳ Adam, SGD + Momentum  $\rightarrow$  memory system that learns to compress gradients usefully

∴ Expressive optimizers with deep memory are designed which has multiple layers of memory inside optimizer itself.

↳ Deep momentum

↳ Better preconditioning

↳ Better gradient compression.

#### ② Self modifying model

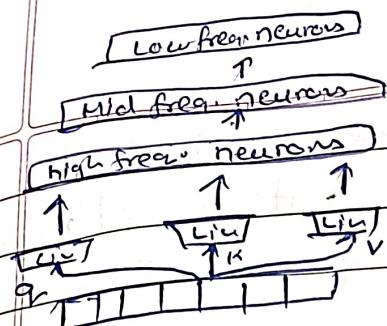
↳ sequence model that learns to modify itself by learning to update its own algo.

#### ③ Continuum memory system (CMS)

↳ diff parts of memory update at different freq.

like long / mid / short term memory

→ each compresses at own time scale



each level update at different frequency & doesn't rely on centralised clock.

\* ASSOCIATIVE MEMORY: Ability to form & retrieve connection

2.1 IMP! Learning = neural update

Memory = neural update via I/P

Learning = acquiring effective & useful memory.

Every neural net layer & every optimize is an associative memory system

↳ It stores associations b/w things it has seen.

Let associative memory  $M$  as a function

$$M : K \rightarrow V$$

↳ value  $\downarrow$  key

$$\therefore M^* = \arg \min_M \tilde{L}(M(K); v)$$

↳ objective that measures quality of mapping

Eg 1: Training 1 Layer MLP = Associative memory

→ model is  $y = w x$

& gradient descent is  $w_{t+1} = w_t - \eta \nabla_w L(w_t; x_{t+1})$

↳ standard backprop

now,

$$\min_w (w x_{t+1}, u_{t+1}) + \frac{1}{2\eta} \|w - w_t\|^2$$

$$\cdot u_{t+1} = \nabla_y L(w_t; x_{t+1}) \quad \begin{matrix} \langle w x, u \rangle \\ \text{dot prod.} \end{matrix}$$

what model produces ( $w x$ )  
& what it should.

Now,  $u_{t+1}$  = error signal

aka local surprise signal

$$+ \|w - w_t\|^2$$

If  $u_{t+1} \uparrow \uparrow \Rightarrow$  bad pred.

$u_{t+1} \downarrow \downarrow \Rightarrow$  good pred.

↳ stabilizer.

like gradient descent.

∴ It is clear from fig 1 that neural net is storing compressed map b/w each data point  $x$  & its error/surprise signal  $u$ . ie training = memory compression

Eg 2: Momentum = 2-level memory

$$m_{t+1} = m_t - \eta \nabla_w L(w_t; x_{t+1})$$

$$w_{t+1} = w_t + m_{t+1}$$

Here  $m_t$  is a second associative memory that stores & compresses past gradients.

$$m_{t+1} = \arg \min_w (-\langle m, \nabla_L(w_t; x_{t+1}) \rangle + \eta \|m - m_t\|^2)$$

Eg 3: Linear attention = Memory

$$M_{t+1} = M_t + v_t k_t^T$$

$$\text{aka } M_{t+1} = M_t - \nabla_L(M_t; k_t, v_t)$$

$M_t$  is associative memory mapping key → value

by gradient descent

∴  $w_k, w_q, w_v \rightarrow$  trained by GD

&  $M \rightarrow$  updated by inner GD, nested optimization

### # SUMMARY of 2.1

	Component	Key	Value	Memory
①	MLP training	$\times$	$\nabla L$	Long-term
②	Momentum	grad index	grad	medium
③	Linear attention	$k$	$v$	short
④	Optimizer	weights	grad stats	meta

## \* NESTED OPTIMIZATION PROBLEMS

Q.2

- Update freq. → for any component A:

$f_A = \text{no. of updates of component A per unit time}$

In transforms:

- ① Attention is fast memory.
- ② Adam momentum is medium memory
- ③ Weights are slow memory.

If A faster than B then  $A > B$

↳ If  $f_A > f_B$

or  $f_A = f_B$  but A must be computed faster

Hierarchy

LVL 1 → fastest → optimization prob. → scatter flow

2 → slower

:

K → slowest

or  $A \equiv B$

when no one is faster & are independent to each other

Eg: Gradient descent

$$w_{t+1} = w_t - \eta \nabla L(w_t)$$

only W updates:

freq. lvl = slow

& only one optimization problem.

Eg: GD + Momentum

$$m_{t+1} = \alpha m_t - \eta \nabla L(w_t)$$

$$w_{t+1} = w_t + m_{t+1}$$

m updates faster & W updates slower

∴ Inner optim. → updates m to compress recent gradients

Outer optim. → updates W based on  $\eta$ .

Eg: Linear Attention

$$M_{t+1} = M_t + V_t K_t^T, \text{ happens every token} \rightarrow \text{very fast}$$

$w_k, w_v, w_q \rightarrow \text{slow}$

∴ attention M is fastest  
optimizer is mid

& weights is long term

## \* OPTIMIZERS as LEARNING MODULES

2.3 ↳ Optimizer is a multi lvl learner, with it's own memory.

Eg 1: GD as a learning Module

$$w_{t+1} = w_t - \eta g_t$$

in NL

$$w_{t+1} = \arg \min_w \left[ \langle w_t, g_t \rangle + \frac{1}{2\eta} \|w - w_t\|^2 \right]$$

→ here optimizer solving tiny problem at ever step

key = gradient direction

value = weight correction

Eg 2: SGD+ Momentum

$$m_{t+1} = \beta m_t + (1-\beta) g_t$$

$$w_{t+1} = w_t - \eta m_{t+1}$$

→ momentum memory  $m_t$  stores a compressed history of past gradients.

→ this is faster memory

Inner optim.

$$\text{Update momentum: } m_{t+1} = \arg \min_m \left[ \langle m_t, g_t \rangle + \frac{1}{2(1-\beta)} \|m - m_t\|^2 \right]$$

$$\begin{aligned} \text{momentum} &= \|m\|_2 \\ \text{weights} &= \|w\|_2 \end{aligned}$$

Eg 3: Adam

$$m_t = \beta_1 m_{t-1} + (1-\beta_1) g_t$$

→  $m_t$  = fast surprise memory

$m_t$  = first moment  
(mean)

$v_t$  = slower variance memory

$$v_t = \beta_2 v_{t-1} + (1-\beta_2) g_t^2$$

$w_t$  = slow long term memory

$v_t$  = second moment  
(variance)

$$w_{t+1} = w_t - \eta \frac{m_t}{\sqrt{v_t} + \epsilon}$$

$w_t$  = weights

Eg 4: Shampoo / Muon (Preconditioners)

$$P_t \approx \sum_i g_i g_i^T$$

→ preconditioner store correlations b/w GRDs.

→ It is a rich associative memory

→ Adam stores scalar memory

Shampoo stores matrix memory

## CONTINUUM MEMORY SYSTEM (CMS)

\* Fast memory  $\rightarrow$  update every token

medium  $\rightarrow$  every few token

slow  $\rightarrow$  per sequence

v. slow  $\rightarrow$  only during training

$T \rightarrow$  time const.

$$\therefore M_i(t+1) = \left(1 - \frac{1}{\tau_i}\right) M_i(t) + \frac{1}{\tau_i} F_i(x_i)$$

$M_i \rightarrow$  memory at  $i$ ;  $F_i \rightarrow$  update signal;  $\tau_i =$  update speed.

Fast memory receives  $I_{IP}$ , previous  $O_P$  & its own memory state.

$$h_i(t) = f(W_i h_{i-1}(t), M_i(t))$$

$$M_i(t+1) = \alpha_i M_i(t) + (1 - \alpha_i) h_i(t)$$

$\alpha_i$ near 0 $\rightarrow$ fast memory
$\alpha_i$ near 1 $\rightarrow$ slow memory

### \* HOPE Arch. [Hierarchical Optimization Process Engine]

3.2 HOPE = CMS + Self modifying sequence model.

1. CMS  $\rightarrow$  memory hierarchy

2. Titas-style self modifying blocks

3. Linear attention  $\rightarrow$  for long context

4. Deep optimizer inside forward pass

5. Model modifies its own internal fns.

①  $I_{IP}$  goes through encoder

$$h_0(t) = \sum E(x_t)$$

④ Linear Attention

$\hookrightarrow$  to handle long sequences

$$A(t+1) = A(t) + v(t) \kappa(t)^T$$

$\kappa(t)$  is fastest

② Feature passes CMS

$$h_i(t) = f_i(h_{i-1}(t), M_i(t))$$

$$M_i(t+1) = (1 - \alpha_i) h_i(t) + \alpha_i M_i(t)$$

⑤ Decoder head produces next token

$$y_t = D(h_{final}(t))$$

③ Self modifying 'Neural learning module'

$\hookrightarrow$  predicts how to change part of model's own weight

$$\Delta w_i(t) = U_i(h_i(t))$$

$$w_i(t+1) = w_i(t) + \Delta w_i(t)$$

$U_i$  is learned