

- How to optimize deep learning algorithms for hardware : focus
- Presentation - Motivation

Problem Statement }
Method } for research paper
Results }

- Coding - part : present some part of the paper (1 or 2 experiments)
- Autograd - What is it? : something differentiation etc

29/07/25

Topic: ML Basics

- a) Focus on: (i) Gradient Descent
(ii) Backprop (BP)

- b) Covering high level motivation in ML

Q1 • MNIST - handwritten digits (e.g. 0 written in various ways)
data set

\Leftarrow Q1 How can we identify it with ML?
Q2 Using DSA? - (hard)

Q2 • Supervised ML : training set of images with known labels

↓

no one ML algorithm

↑ (Starts with random h₀)

Foundational

Models

features: (i) Hypothesis class - why? Processing?

(ii) loss function - for assessing how well model works

(iii) optimization method - helps in determining set of params

which help in reducing the sum of losses

G3 Logistic Regression - Theory

- What's the setting? - k classes to learn

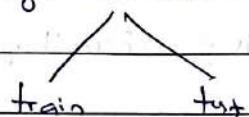
ex. for MNIST digit identifying

$n = 28 \times 28 = 784$ - dimension of input data

$k = 10$ - no. of different classes / 10 digits

$m = 60,000$ - # points in the training set

- Is random splitting of data into always good? (No)



for time series: no

it is fluctuating

→ 60% data for training

Step 1:

- Hypothesis function selected: Linear H.F.

$$\check{h}_{\theta}(n) = \Theta^T n$$

Θ : random

proportion of n

→ don't start with 0

n is input

directly - for

complicated optimisation

$$h_{\theta}(n) = \Theta^T \begin{bmatrix} n \\ 1 \end{bmatrix} \quad \text{--- (1)}$$

n is projected from n

dimensional to k dimensional

→ it takes time to compute

slow

→ we need to do parallelize it

→ matrix notation used

- Matrix based notation: (Parallel computation)

$$x \in \mathbb{R}^{m \times n}, X = \left[\begin{array}{c} x^{(1)} \\ \vdots \\ x^{(m)} \end{array} \right] \quad (i) \quad y \in \{1, \dots, k\}^m$$

$$\text{The dimension of } y \text{ is } \left[\begin{array}{c} y^{(1)} \\ \vdots \\ y^{(m)} \end{array} \right]$$

$$\checkmark h_{\theta}(x) = X \theta \quad \text{--- } ②$$

$m \times k \quad m \times n \quad n \times k$

What's the advantage of ② over ①?

→ faster

→ GPU - parallel hardware (else no benefit in CPU)

Step 2: (This is done to evaluate)

- Loss function selection #1: classification error

$$\checkmark l_{\text{err}}(h(x), y) = \begin{cases} 0 & \text{correct} \\ 1 & \text{incorrect} \rightarrow i.e. \text{ if } \arg \max_i h_i(x) \neq y \end{cases}$$

What's wrong with this?

$$\Rightarrow \text{ex. } n = 1000 \quad \left. \begin{array}{l} \text{(random if we scale)} \\ \text{if correct} = 250 \\ \text{incorrect} = 750 \end{array} \right\} \begin{array}{l} \text{loss function is} \\ \text{not optimiser friendly} \Rightarrow \text{as not differentiable} \end{array}$$

↓

What should we do now?

Loss #2: → to make all positive & sum to 1

• Can we do with probabilities?

$$\checkmark z_i = P(\text{label} = i) = \frac{\exp(h_i(x))}{\sum_{j=1}^k \exp(h_j(x))} \quad \begin{array}{l} h_i(x) : \\ -0.5 \quad -0.2 \quad 0.1 \quad 0.5 \end{array}$$

= \underbrace{\exp(\text{normalise}(h_i(x)))}_{\exp(h_i(x))}

This is → softmax function

Right!!

To handle with this, we

(as minimising) for numerical consistency go for

$$\checkmark L(h(x), y) = -\log P(\text{label} = y)$$

$$\cancel{\text{softmax}} \rightarrow \text{cross entropy} = -h_y(x) + \log \sum_{j=1}^k \exp(h_j(x))$$

(*)

(our focus)

→ exponential function
(w) absⁿ → (pos) out

optimism friendly loss $\rightarrow (\bar{w}) > (\bar{z} + \bar{w})$

Note: softmax & cross entropy same!!.

minimize w.r.t

Step 3: (Optimiser) (This is done to optimise the optional)

What do we do with loss function L ?

\Rightarrow minimize:

$$\frac{1}{m} \sum_{i=1}^m L(h_\theta(x^{(i)}), y^{(i)})$$

for cross entropy: L_{ce}

MISSING

C-4 Gradient Descent Discussed

(We start
some
with loss
first)

- We get optimal $\Theta_{n \times k}$

- for optimization: we use gradient descent

Backprop

Gradient Descent
(By hand)

$L(w)$: When can we apply GD? $\text{cost} = w \approx \Theta$

not loss

• Convex

 $\rightarrow \text{GD} \rightarrow \text{minimise}$

• Continuous

 $\rightarrow \text{GD} = \text{minimise}$

• Differentiable

How to do it? (optimisation part)

Algo: \rightarrow how do we find (Θ) ? Start with random w_0

$$w_0 \downarrow \rightarrow w_1$$

$$w_1 \downarrow \rightarrow w_2$$

$$w_2 \downarrow \rightarrow w_3$$

$$w_3 \downarrow \rightarrow w_4$$

$$w_4 \downarrow \rightarrow w_5$$

$$w_5 \downarrow \rightarrow w_6$$

$$w_6 \downarrow \rightarrow w_7$$

$$w_7 \downarrow \rightarrow w_8$$

$$w_8 \downarrow \rightarrow w_9$$

$$w_9 \downarrow \rightarrow w_{10}$$

$$w_{10} \downarrow \rightarrow w_{11}$$

$$w_{11} \downarrow \rightarrow w_{12}$$

$$w_{12} \downarrow \rightarrow w_{13}$$

$$w_{13} \downarrow \rightarrow w_{14}$$

$$w_{14} \downarrow \rightarrow w_{15}$$

$$w_{15} \downarrow \rightarrow w_{16}$$

$$w_{16} \downarrow \rightarrow w_{17}$$

$$w_{17} \downarrow \rightarrow w_{18}$$

$$w_{18} \downarrow \rightarrow w_{19}$$

$$w_{19} \downarrow \rightarrow w_{20}$$

$$w_{20} \downarrow \rightarrow w_{21}$$

$$w_{21} \downarrow \rightarrow w_{22}$$

$$w_{22} \downarrow \rightarrow w_{23}$$

$$w_{23} \downarrow \rightarrow w_{24}$$

$$w_{24} \downarrow \rightarrow w_{25}$$

$$w_{25} \downarrow \rightarrow w_{26}$$

$$w_{26} \downarrow \rightarrow w_{27}$$

$$w_{27} \downarrow \rightarrow w_{28}$$

$$w_{28} \downarrow \rightarrow w_{29}$$

$$w_{29} \downarrow \rightarrow w_{30}$$

$$w_{30} \downarrow \rightarrow w_{31}$$

$$w_{31} \downarrow \rightarrow w_{32}$$

$$w_{32} \downarrow \rightarrow w_{33}$$

$$w_{33} \downarrow \rightarrow w_{34}$$

$$w_{34} \downarrow \rightarrow w_{35}$$

$$w_{35} \downarrow \rightarrow w_{36}$$

$$w_{36} \downarrow \rightarrow w_{37}$$

$$w_{37} \downarrow \rightarrow w_{38}$$

$$w_{38} \downarrow \rightarrow w_{39}$$

$$w_{39} \downarrow \rightarrow w_{40}$$

how do we know we are improving?

ii) \rightarrow if $\| \vec{w}^{t+1} - \vec{w}^t \| < \epsilon \Rightarrow$ converging

$$(1 = 1.101)^2 = 1.21 = (1.101)^2$$

• If we don't know loss: L

\rightarrow if we want to do: minimize $L =$

\rightarrow How would we do it? \rightarrow (we can't if we don't know, we need to know it)

\rightarrow μ ridge(\vec{w})

$L(\vec{w} + \vec{s}) < L(\vec{w}) \rightarrow$ How to ensure? during optimisation?

Taylor Series Expansion

Taylor series

$$l(\vec{w} + \vec{s}) \approx l(\vec{w}) + \underbrace{\vec{g}(\vec{w})^T \vec{s}}_{\text{gradient}} + \frac{1}{2} \vec{s}^T H(\vec{w}) \vec{s}$$

High order terms: req different methods to
Optimize (not GD)

$$\vec{g}(\vec{w}) = \nabla l(\vec{w})$$

04/08/25

C1 (Continuing Gradient Descent)

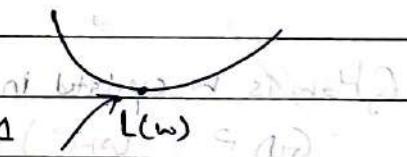
- If I don't know $\log(\cdot)$, how do I ensure the loss always decreases? \rightarrow minimize $L = \text{loss}$ but not much aware about L

\rightarrow Step 1: Construct simpler/tractable functions

\rightarrow Approach: Taylor Approximation

- Take small step(s) in $L(w)$

(S not defined yet) $\|S\|_2$ is small



\rightarrow Then what will be loss? (Break it down)

\Rightarrow Initial Loss = $L(w)$

\rightarrow Now = $L(w + S)$

S is vector (or scalar product)

V1: $L(\vec{w} + \vec{s}) \approx L(\vec{w}) + \vec{g}(\vec{w})^T \vec{s}$ if $\|\vec{s}\|_2$ is small

(GD) Where: $\vec{g}(\vec{w}) = \nabla L(\vec{w})$ notes on $L(w)$

V2: $L(\vec{w} + \vec{s}) \approx L(\vec{w}) + \vec{g}(\vec{w})^T \vec{s} / 2$ in 2D

Newton's Method \rightarrow (unit step) $\frac{1}{2} \vec{s}^T H(\vec{w}) \vec{s}$

Condition 2:

L twice differentiable(b)

Why: $H(\vec{w}) = \nabla^2 L(\vec{w})$

loop for \vec{w}

Hessian

- Which one is more appropriate?

→ (2) is proper but there is a pb — presence of H which requires function to be double differentiable (expensive)

↓

We go with VI (so henceforth gradient descent)

- What is exactly gradient descent in VI?

→ Our goal is $\min(L)$

→ We need to ensure $L(\vec{w} + \vec{s}) < L(\vec{w})$

→ for this we need to do some trick in VI

NOTE: $\vec{g}(\vec{w})$ is fixed we can't do it

\vec{s} needs to be ~~minimized~~ changed

↓ $\Rightarrow \vec{s} = -\alpha \vec{g}(\vec{w})$ (force it)

$$L(\vec{w} + \vec{s}) = L(\vec{w} - \alpha \vec{g}(\vec{w})) \quad (B)$$

$$\approx L(\vec{w}) - \alpha \vec{g}(\vec{w})^T \cdot \vec{g}(\vec{w})$$

Q: How is w updated in

GD? $(\vec{w} + \vec{s})$ || $\vec{s} = -\alpha \vec{g}(\vec{w})$ $\Rightarrow \alpha > 0$ (as eq of gradient)

↓ \vec{s} is, in our hand

$(\vec{w})_1 = \vec{w}$ we can force it to be > 0

$\Rightarrow (\vec{s}_{RHS})_1$ is always greater than LHS

Dashed blue line shows how to choose α (but = 1)

- What's still the problem?

(but it can be many) \rightarrow only 1 local minima: used in classical model

\rightarrow how can we determine α ? $\vec{s}_1 = (\vec{w})_1 - \vec{g}(\vec{w})$ (eq)

- Since (B) is an iterative process

as we start with $\vec{s}(\vec{w})_1 = (\vec{s}_1 + \vec{w})_1 = \vec{w}$

$\Rightarrow \vec{s}_1 = 1 / \vec{g}(\vec{w})_1$ (first time) \leftarrow ensures initially high α &

What if α is always long? \rightarrow not good

bottom first then lower

finally converge

- Note: previous section works for only 1 dimension

What if $n \times n$ dimensions? \leftarrow each dimension will have a
(input) \rightarrow curve associated with it

\downarrow
it should be independent \leftarrow • (α should be associated with
a dimension)

(i) coordinate instead of gradient descent, we go for AdaGrad
learning rate) \leftarrow adaptive gradient \rightarrow
adaptation # if multidimensional, how to
descent

C2 # AdaGrad do we update? \rightarrow multidimension! = multiply

w_0 and \bar{z} $\nparallel d \equiv w_d^0 := D, z_d = 0$ minima
 \uparrow (running avg)

Assuming initially

- Repeat until convergence

$$S1: \vec{g} = \nabla f(\vec{w})$$

$$\sum_{i=1}^t g_i^2 = Z_{i,t} \quad \text{if prev } g_i \text{ was high}$$

$$w_d^{t+1} = Z_{d,t} \leftarrow Z_{d,t} + g_d^2$$

\downarrow we need to take

• Note: $w_d \rightarrow \forall d: w_d^{t+1} \leftarrow w_d^t - \alpha g_d$ small steps

i.e. weights of a particular dimension move slowly. If $\sqrt{Z_d + \epsilon}$ which this expression depends on its gradient only!! ensures that it does not move too much

$$\text{if } \|w^{t+1} - w^t\| < \epsilon \rightarrow \text{convergence}$$

\downarrow (bottom right is not minimizing)

C3 Newton's Method

\downarrow

- Since Vf not \mathcal{C}^1 differentiable \Rightarrow \vec{s} to be precise bounded + even more turns

$$L(\vec{w} + \vec{s}) \approx L(\vec{w}) + \vec{g}(\vec{w})^\top \vec{s} + \frac{1}{2} \vec{s}^\top H(\vec{w}) \vec{s}$$

\Rightarrow faster convergence

$$H(\vec{w}) = \begin{pmatrix} \frac{\partial^2 L}{\partial w_1^2} & \frac{\partial^2 L}{\partial w_1 \partial w_2} & \dots & \frac{\partial^2 L}{\partial w_1 \partial w_n} \\ \vdots & \ddots & \ddots & \frac{\partial^2 L}{\partial w_n^2} \end{pmatrix}$$

→ Always symmetric square matrix & positive semi definite

- H needs to be twice differentiable



→ generally we don't care for finding minima

but if we still want to use it

How to
minimize
it?

$$\frac{\partial (L(\vec{w} + \vec{s}))}{\partial \vec{s}} = \underset{s}{\arg \min} L(\vec{w}) + \vec{g}(\vec{w})^T \vec{s}$$

\downarrow

$$+ \frac{1}{2} \vec{s}^T H(\vec{w})(\vec{s})$$

(That's why)

we don't use

H)

$$S_1: LHS = 0$$

S2: What should be \vec{s} ?

(Note: differential of RHS wrt S)

$$\Rightarrow 0 = 0 + \vec{g}(\vec{w}) + H(\vec{w}) \vec{s}$$

$$\Rightarrow \vec{s} = - [H(\vec{w})]^{-1} \vec{g}(\vec{w})$$

(i) Converges very fast

If Hessian approximated: Limitation of Hessian: here it's inverse which is not applying
then accurate

(ii) but if f(x) is flat in some dimensions (second derivative ≈ 0),

the inverse Hessian blows up → Alternative: Hybrid version

up → huge unstable steps

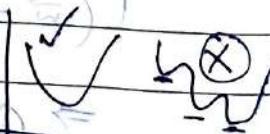
if we want to use Newton's method, we need to modify Hessian ← make it more diagonal matrix
(approximation solving method)

(here inverse become)

(i) use SGD first to get close to the optimum then switch to Newton's method since

* it's quadratic approximation works better near the minimum → most all entries 0

(3) Limitation of GD → if function not convex (multiple minima),



* quadratic approximation → using gradients & Hessian
Solve Newton's method!

- Gradient Descent ensure loss minimisation (not loss: it can go anywhere)

→ can be applied in classification as well

Topic Discussed: Stochastic G.D

(What does this expression represent?)
→ loss function

$$\cancel{f(\theta) = \min_{\theta} \frac{1}{m} \sum_{i=1}^m l_{ce}(\underbrace{\theta^T x^{(i)}, y^{(i)}}_{h(x^{(i)})})} \quad (i)$$

Apply gradient descent

$$\Rightarrow \nabla_{\theta} f(\theta) = \begin{bmatrix} \frac{\partial f(\theta)}{\partial \theta_1}, \dots, \frac{\partial f(\theta)}{\partial \theta_n} \\ \vdots \\ \frac{\partial f(\theta)}{\partial \theta_m} \end{bmatrix}$$

no. of dimension class

What's the problem with this method? (ie apply grad on (i)?)

$$\sum_{i=1}^m \nabla_{\theta} l_i(h(x^{(i)}), y^{(i)})$$

→ (we take all in gradient descent)

→ m is often large → time-consuming process

→ parallelise it ← memory consumption

→ multiple minima → get stuck

• So vanilla gradient descent not good

\downarrow what is it?

Since we want to generalise on test data

(Batch points with n dimensions)

$B \ll m$

→ Sample mini batch from $\{1, \dots, m\}$

Repeat,

Sample a mini batch $x \leftarrow \mathbb{R}^{B \times 1}$

$y \leftarrow \{1, \dots, K\}^B$
dimensions

classes

we don't get stuck

(note: B → then do gradient on it)

Changes on

every iteration

$(x^{(i)})_{i=1}^B + \{y^{(i)}\}_{i=1}^B$ Update parameters, $\theta: \theta -$

(note: minima for every batch

can be different)

? (here the local minima

is avoided)

$\alpha \left(\frac{1}{B} \sum_{i=0}^B \nabla_{\theta} l(h_i, y_i) \right)$

How to compute the gradients manually

$$\nabla_{\theta} \text{Lce} (\underbrace{\Theta^T x}_{\substack{\text{gradient} \\ \text{symbol}}}, y) = z \quad (\text{calling it})$$

Cross entropy loss

$h(\cdot)$: hypothesis function

$\max_{k=1}^K h_k(x) (= k \times m)$ from data points

$$\Rightarrow \Theta^T x \quad \begin{matrix} K \times 1 \\ m \times 1 \end{matrix}$$

Only one class or
one dimension

Calculated logits!

$$\nabla_{\theta} \text{Lce} (h, y) \leftarrow *$$

$\frac{1}{m} \sum_i h_i$

1 data point
with 4 dimensions

$$\text{Lce}(h, y) = -h_y + \log \sum_{j=1}^K \exp(h_j)$$

Method 1: Using h_i

$$\frac{\partial \text{Lce}(h, y)}{\partial h_i} = \frac{\partial}{\partial h_i} \left(-h_y + \log \sum_{j=1}^K \exp(h_j) \right)$$

$$= \frac{\partial}{\partial h_i} \left(\underbrace{-h_y}_{\text{when } i=y} + \left(\frac{\partial}{\partial h_i} \sum_{j=1}^K \exp(h_j) \right) \right)$$

It's an indicator

$$\sum_{j=1}^K \exp(h_j)$$

are not supposed to 1)

$$\text{use this product formula} = -1 \{i=y\} + \exp(h_j) \rightarrow (\text{for } j=i \text{ constant})$$

$m \gg 1$
 $x_i \rightarrow x$ divided into a sum

$$\sum_{j=1}^K \exp(h_j) \rightarrow \text{softmax}$$

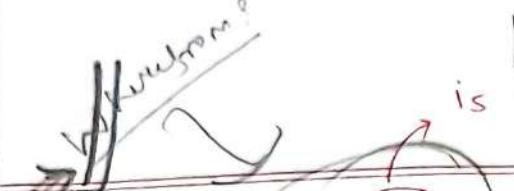
This is softmax

This gives the probability

$$- \Theta^T x, \text{softmax} = 1 \{i=y\} + \text{softmax} (1)$$

$$(\text{softmax}) = \frac{\exp(\Theta^T x_i)}{\sum_j \exp(\Theta^T x_j)}$$

NOTE: e^y usually denotes the ~~class~~^{class} vector with 1 at the true class y



$$= -e^y + z \quad (= \text{normalized } (\exp(h)))$$

$$\nabla_{\theta} l_{ce}(h, y) = z - e^y$$

* → is the simplification one

Method 2: For generalizing,

With wts $\nabla_{\theta} l_{ce}(\Theta^T x, y)$

$$= \frac{\partial}{\partial \Theta} l_{ce}(\Theta^T x, y) \quad \left\{ \begin{array}{l} \text{using chain rule} \\ (\times) \Theta^T x = \end{array} \right\}$$

instead of $\Theta \leftarrow \Theta; \text{formatting}$ \rightarrow { Here the trouble was differentiating directly with Θ }

$$= \frac{\partial}{\partial \Theta^T x} l_{ce}(\Theta^T x, y) \cdot \frac{\partial \Theta^T x}{\partial \Theta} \rightarrow (x^T \Theta)$$

$$= \frac{\partial l_{ce}(h, y)}{\partial h} \cdot x^T$$

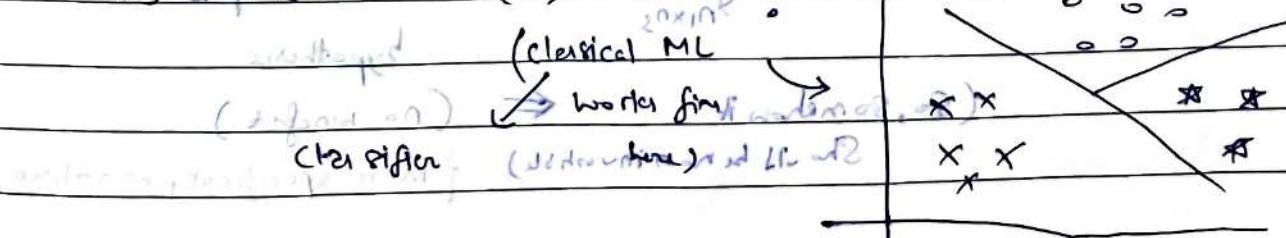
\leftarrow note: we already computed this h is the hypothesis Φ function grammar with

$$= ((z - e^y) \cdot x^T)_{K \times 1} \quad \left(\begin{array}{l} \text{Lm} \rightarrow \Theta \leftarrow n \nabla_{\theta} \\ \text{sm} \rightarrow \Theta \leftarrow \frac{1}{n} \nabla_{\theta} \end{array} \right)$$

$$= \underbrace{x}_{n \times 1} \underbrace{(z - e^y)^T}_{\text{prediction}} \underbrace{(\Phi)^T \Theta}_{\text{true}} - (y)_d \quad \left(\begin{array}{l} \text{Input} \rightarrow \text{prediction} \rightarrow \text{true} \\ \text{Input} \rightarrow \text{prediction} \rightarrow \text{true} \end{array} \right)$$

linear hypothesis

• Why do we hate $h(x) = \Theta^T x$?



but for non-linear classification,
linear classifier doesn't work

\Rightarrow we can do feature expansion

Sign

$x_1 \in \mathbb{R}, x_2 \in \mathbb{R}$ one dimension

known
 \downarrow

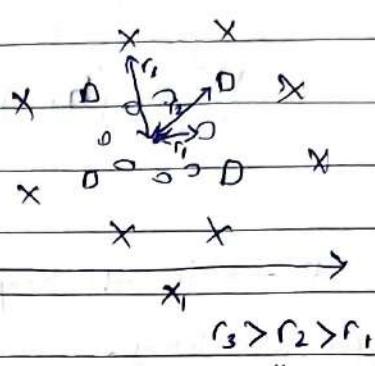
$$\text{current } \leftarrow x_1^2 + x_2^2 \text{ (introduced)} \quad x_1^2 + x_2^2$$

trick (feature expansion)

$$h(x) = \Theta^T \cdot \phi(x) \quad (\text{initially})$$

$$= \Theta^T \phi(x)$$

\downarrow (not so good)



feature expansion

can be done

That's why now we go for neural networks

• need for feature expansion: difficult to understand

Neural Networks for non-linear classification

$$T X \cdot \frac{\partial \delta / \partial \theta / \partial s}{\partial \theta} =$$

• How to come up with good $\phi(x)$?

→ kernel trick (m1) — i.e. feature expansion

→ learn by itself (m2)

• Let's assume, $h_\theta(x) = \Theta^T \phi(x)$

$$\phi(x) = w^T x$$

$$= \Theta^T w^T x^T \mathbf{1} \quad \text{where } \phi(x) = w^T x$$

$$\underbrace{n_1 n_2}_{n_1 n_3} \underbrace{n_2 n_3}_{n_1 n_3}$$

$\Rightarrow \Theta^T w =$ similar to previous hypothesis

(so, somehow it's good) \leftarrow (no benefit)

Show it is invertible) { more specifically non-linear}

Topic: Neural Networks

C1: Why do we need? \rightarrow for non linear classification

C2: (continuation of prev page)

$$h_{\theta}(w) = \Theta^T (\phi(w))$$

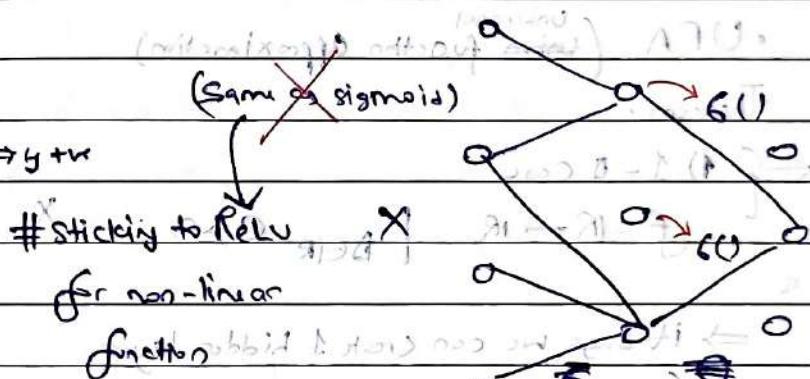
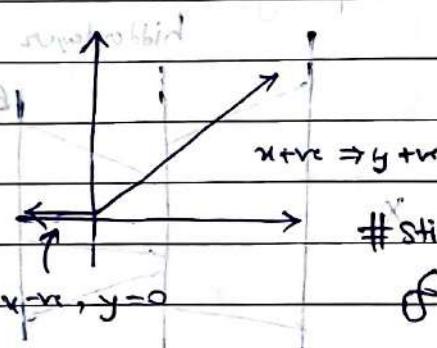
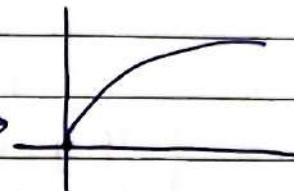
$$= \Theta^T \sigma(W^T w) \neq \tilde{\theta} \cdot x$$

\leftarrow Non linear

function (like sigmoid) \rightarrow

(This can't be broken into

first hypothesis)



finding θ : example modelling

How to optimise: $h_{\theta}(w) = \Theta^T (\phi(w))$ \rightarrow (1) non-linear
 $= \Theta^T \sigma(W^T w)$ \rightarrow (2) linear

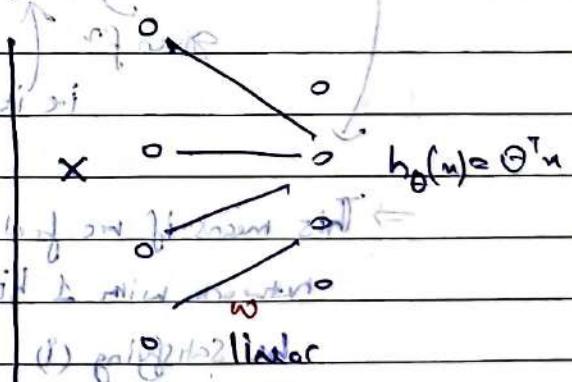
let $w \rightarrow$ be random

$\sigma \rightarrow \cosine()$

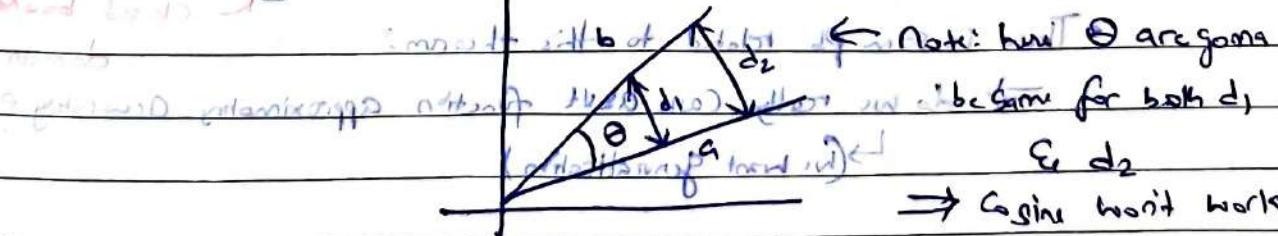
$\Theta \rightarrow$ learned

• first layer - cosine inc random w if we want to learn

• second/last layer - learned w will reflect mind structure



Random w \rightarrow



\Rightarrow cosine won't work

for all scenarios

Are the weights associated with layers?

(note the straight lines are just a boundary depiction)

$\Theta : w_1, w_2 \leftarrow$ both are parameters (learned)

$h(x) = w_2^T g(w_1^T x) \leftarrow$ 2-layer N.N (as 2 wts)

$x w_1$

\rightarrow How to find optimal value of w_1 & w_2 ?

i.e.

$$w_1^{t+1} = w_1^t + s_1 \quad \rightarrow \quad s_1 = -\alpha \vec{g}(w_1) \quad \left. \begin{array}{l} \text{How to} \\ \text{determine} \end{array} \right\}$$

$$w_2^{t+1} = w_2^t + s_2 \quad \rightarrow \quad s_2 = -\alpha \vec{g}(w_2) \quad \left. \begin{array}{l} \text{this?} \\ \text{if backprop} \end{array} \right\}$$

\rightarrow if you give me enough neurons

& the right setup, I can represent

the true $f(x)$ - even if I don't know it explicitly!

• UFA (Universal Function Approximation)

Theorem:

NOTE: $\left\{ \begin{array}{l} 1) 1-\text{D case} \\ (\text{Objective is}) \end{array} \right.$

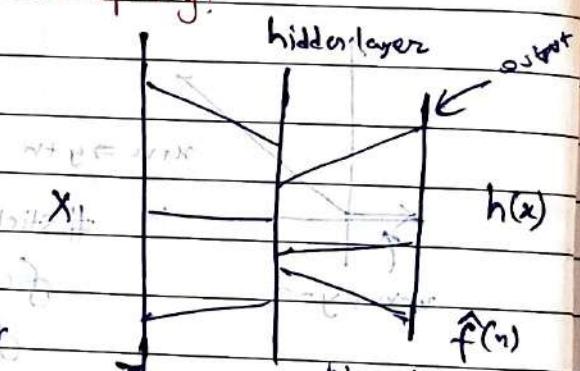
We make a

function

\Rightarrow it says we can create 1 hidden layer

approximately

$$f: \mathbb{R} \rightarrow \mathbb{R} \quad |f| \in \mathbb{R}, f' \geq 0$$



$$\max_{x \in D} |f(x) - \hat{f}(x)| \leq \epsilon$$

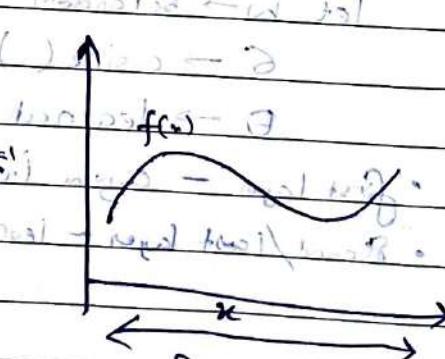
given $f(x)$

i.e. it is bounded

\Rightarrow This means if we feed x to the neural

network with 1 hidden layer, will result in

be satisfying (1)



There is some problem related to this theorem:

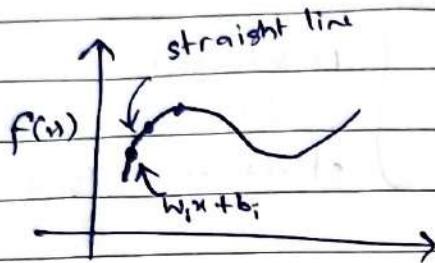
\Rightarrow do we really care about function approximating accurately?

\hookrightarrow (we want generalisation)

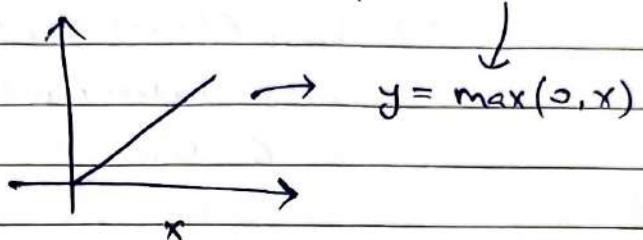
show that $w_1, w_2 \leftarrow$

show that $w_1, w_2 \leftarrow$

• Why is 1 layer hidden very powerful?



using ReLU



$$\hat{f}(u) = \sum_{i=1}^d \pm \max \{ 0, w_i u + b_i \}$$

(This can map $f(u)$ very closely)

$$y = m u + c$$

here we have to find optimal w & b

→ universal function approximator: Says given any function, then

I will add small $\frac{1}{2}$ -linear functions ϵ_i , approximate the function provided by w_i, ϵ_i by one already given (Computed using back prop)

~~• do we know $f(u)$ thus? — not exact mathematical formulation~~

~~• why would we compute if we know? — Sometimes expensive so approximation~~

~~helps to how we make it learned from the examples we have~~

• Defining loss ← computed in final layer ← (no layer specific loss)

$$h(u) = \underbrace{\sum_{i=1}^d w_i}_{\text{Input layer}} \underbrace{G(x)}_{\text{hidden layer}} \underbrace{w_i}_{\text{sum & apply}}$$

~~is weight of output of hidden neurons with no bias term~~

~~is weight of output of hidden neurons with no bias term~~

$$(I - \alpha)^T (h(u)) = ((I, \alpha) \xrightarrow{x} w_1, \dots, w_d) \xrightarrow{w_1, \dots, w_d} h(u) \xrightarrow{\text{loss}} |h - h|^2$$

~~is weight of output of hidden neurons with no bias term~~

~~is weight of output of hidden neurons with no bias term~~

~~is weight of output of hidden neurons with no bias term~~

~~is weight of output of hidden neurons with no bias term~~

activation function

is the hypothesis function

Major focus: how to do weight updates
→ back propagation

Neural Network

$$w^{t+1} = w^t + s ; s = -\alpha \nabla_{w^t} J(w)$$

$$\Rightarrow L_{ce}(G(xw_1)w_2, y)$$

w₁: random normal

G: Cosine

Not very good



• how to get $\nabla_{w^t} J(w)$

finding the right w_1 & w_2

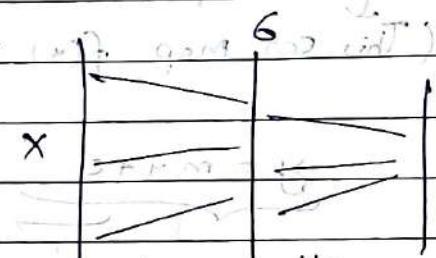
• how do we train a neural network? (a multi-layer NN)

→ What are the principles?

→ use backpropagation

(need to compute gradients w.r.t w_1 & w_2)

$$\nabla_{\{w_1, w_2\}} L_{ce}(G(xw_1)w_2, y)$$



$$\begin{aligned} \nabla_{w_2} (L_{ce}) &= \frac{\partial L_{ce}}{\partial w_2} (G(xw_1)w_2, y) \\ &= \frac{\partial L_{ce}}{\partial w_2} (G(xw_1)w_2, y) \cdot \frac{\partial G}{\partial w_2} (xw_1)w_2 \\ &= \frac{\partial L_{ce}}{\partial G} (xw_1)w_2 \cdot \frac{\partial G}{\partial w_2} (xw_1)w_2 \end{aligned}$$

Output of hidden layer

$$= \frac{\partial L_{ce}}{\partial h} (h, y) \cdot \frac{\partial G}{\partial w_1} (xw_1)$$

$$= (S - I_y) \cdot G(xw_1)$$

Softmax identity $m \times k$

$(m \times k) \times d$

Batch size of m (m samples)

• why need to do some dimension matching to represent it in form of

matrices

$$\nabla_{w_2} L_{ce}(G(xw_1)w_2, y) = G(xw_1)^T (S - I_y)$$

$d \times k$ To compute? Activation op

gradient, w_1 & w_2 w.r.t.

softmax

output

from L_2

layer

- We don't have the loss for L_1 layer but for L_2 (ie overall at the end)



$$\frac{\partial \text{loss}}{\partial w_1} = \frac{\partial \text{loss}(\sigma(xw_1)w_2, y)}{\partial w_1} = \frac{\partial \text{loss}(\sigma(xw_1)w_2, y)}{\partial (\sigma(xw_1)w_2)} \cdot \frac{\partial (\sigma(xw_1)w_2)}{\partial (\sigma(xw_1))} \cdot \frac{\partial (\sigma(xw_1))}{\partial (xw_1)}$$

$$\frac{\partial (xw_1)}{\partial w_1}$$

$$\nabla_{w_1} \text{loss}(\sigma(xw_1)w_2, y) = \frac{(S - I_y)}{m \times k} \cdot \frac{w_2}{d \times k} \cdot \frac{\sigma'(xw_1) \cdot x}{m \times d} \cdot \frac{x}{m \times n}$$

Matching the dimension:

$$= x^T \cdot \left(\begin{matrix} (S - I_y) & w_2^T \\ m \times k & k \times d \end{matrix} \right) \circ \left(\begin{matrix} \sigma'(xw_1) \\ m \times d \end{matrix} \right) \quad \text{--- (A)}$$

↑
n × m ↑
m × k ↓
m × d ↓
Output of final layer ↓
more multiplication
Weights of next layer

Gradient depending on (Input)

backward pass much more computationally expensive than forward pass: no memory sharing — as seen in (A)

- Scaling upto more layers

$$i=1 \dots L \quad (\text{L layers}) \quad (\text{for forward pass}) \quad *$$

$$\rightarrow z_{i+1} = \sigma_i(z_i w_i) \rightarrow \text{let } z_1 = x \quad \text{activation output}$$

Output for each layer

1st layer

$$\text{why?} \Rightarrow \frac{\partial \text{loss}}{\partial w_i} (z_{L+1}, y) = \frac{\partial \text{loss}}{\partial z_{L+1}} \cdot \frac{\partial z_{L+1}}{\partial z_L} \cdot \frac{\partial z_L}{\partial z_{L-1}} \cdots \frac{\partial z_2}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_i} \quad \left[\begin{matrix} \frac{\partial z_{L+1}}{\partial w_i} \\ \vdots \\ \frac{\partial z_1}{\partial w_i} \end{matrix} \right]$$

the loss taken as

z_{L+1} here we need the gradient

$$G_{i+1} = \frac{\partial \text{loss}}{\partial z_{i+1}} (z_{L+1}, y)$$

↑ flow from right to left

isn't this the wrong way?

actual output?

NO we take activation applied one)

↑ thought with why previously!!

↑ activation

What if we want to do gradient wrt z_i ?

$$\frac{\partial L(z_{L+1}, y)}{\partial z_i} = G_i = G_{i+1} \cdot \frac{\partial z_{i+1}}{\partial z_i}$$

↓
wrt activation

Gradient of i th layer is dependent on $(i+1)$
i.e. the next layer

$$= G_{i+1} \cdot \frac{\partial G_i(z_i w_i)}{\partial z_i w_i} \cdot \frac{\partial z_i w_i}{\partial z_i}$$

$$= G_{i+1} \cdot G'(z_i w_i) \cdot w_i = G_i$$

$M \times n_i$ $M \times n_{i+1}$ $n_i \times n_{i+1}$

Dimension matching

gradient from prev layer (right to left)

$$\cancel{\frac{\partial L(z_{L+1}, y)}{\partial z_i}} = (G_{i+1} \odot G'(z_i w_i)) w_i^T$$

$M \times n_i$ $M \times n_{i+1}$ $n_{i+1} \times n_i$

called incoming gradient

$$\Rightarrow \nabla_{w_i} L(z_{L+1}, y) = G_{i+1} \cdot \frac{\partial z_{i+1}}{\partial w_i}$$

$$= G_{i+1} \odot \frac{\partial G_i(z_i w_i)}{\partial z_i w_i} \cdot \frac{\partial z_i w_i}{\partial w_i}$$

$$= G_{i+1} \odot G'(z_i w_i) \cdot z_i^T$$

activation output

$$= z_i^T (G_{i+1} \odot G'(z_i w_i))$$

Vector
(also backtracks)

vector product

Jacobian

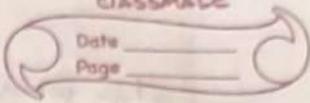
$n_i \times M$

(needs to be
stored)

$M \times n_{i+1}$
incoming backward

gradient (difficult - annoying)

NOTE: quadratic error matters here because heavy
penalty for large errors!



(*) Requirements for F.P.

Requirements for B.P.: i.e. what all req from F.P.P.

→ store activation output (memory intensive)

Multiplication req? — how many? (lower bound)

$$\Rightarrow 2 \times f.p.$$

14/08/25

Type of Differentiation

* Note:

BP

$$G_{L+1} = \nabla_{Z_{L+1}} l(z_{L+1}, y) = S - I_y$$

$G_1 \leftarrow$ same as last class

→ still gradient descent but includes back

• Gradient Descent & Backprop difference backpass & forward pass

↓
→ learning rate η is shared throughout (i.e.)

→ update only for left layer • for all layers, we can update {multiple minima}

• apply vanilla gradient descent as well but stochastic

→ single minima to optimise the memory

Next Topic: Automatic Differentiation (AD)

• Why do we need? — apply backprop in much more principle way? ??

• how to compute gradients in Software?

i) Numerical differentiation

(aka first order (linear))

$\nabla f(\theta) = \frac{f(\theta + \epsilon v) - f(\theta)}{\epsilon}$ approximation

$$\Rightarrow \frac{df(\theta)}{d\theta_i} = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon \theta_i) - f(\theta)}{\epsilon} : \text{error is } O(\epsilon)$$

↓ more exactly to do it

(this is second order (quadratic))

general differentiation

$\nabla^2 f(\theta) = \frac{f(\theta + \epsilon v) - f(\theta - \epsilon v)}{2\epsilon}$ approximation

$$\frac{df(\theta)}{d\theta_i} = \frac{f(\theta + \epsilon \theta_i) - f(\theta - \epsilon \theta_i)}{2\epsilon} : \text{more efficient}$$

as error is $O(\epsilon^2)$

• Problems: sensitive to ϵ (we can choose any ϵ)

here what we are trying to do is checking if the computed gradient is correct or not by calculating the slope!

ii) Even though it has problem, we use it for gradient checking {as in if correctly com-

$$\hat{(\theta + \epsilon\delta)}^T \nabla_{\theta} f(\theta) = f(\theta + \epsilon\delta) - f(\theta - \epsilon\delta) + O(\epsilon^2)$$

puted}

What is this? $\hat{(\theta + \epsilon\delta)}^T \nabla_{\theta} f(\theta)$

\Rightarrow it's a direction vector

We get from software it's a dot product

computed by hand

iii) Symbolic Differentiation

Numerical $f(\theta) = \prod_{i=0}^n \theta_i$

Redundant $f(\theta) = \prod_{i=0}^n \theta_i$

Form graph method

for optimisation cost: $n(n-1)$ multiplies to compute all partial gradients

$\approx n^2$

(iii) Computational graph \Rightarrow (we need to generate so to understand how computationally heavy the expression is)

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin x_2$$

\Rightarrow node: intermediate value in the computation

edge: input-output relationship

$$v_1 = x_1$$

$$v_2 = x_2$$

forward pass
 $v_3 = \ln v_1$

$$v_4 = v_1 v_2 - (v_3^2 + v_3) v_2$$

$$v_5 = \sin v_2$$

$$v_6 = v_3 - v_5$$

$$v_7 = v_4 + v_6$$

$$(i\theta\beta - \theta) v_7 - (i\theta\beta + \theta) v_7 = (\theta)^2 \beta$$

(\exists needs $\cos(\theta)$ \exists at which $\sin(\theta)$)

• now we got y

Forward mode AD

• we need to apply backward pass: (on loss) \hookrightarrow cost is $O(n)$

• how to compute the gradient?

input

$\rightarrow M_1$: from left

$$\frac{\partial v_i}{\partial u_1}$$

$$\frac{\partial v_i}{\partial u_2}$$

(we get the outputs of next layer easily)

$$\Rightarrow \dot{v}_1 = 1 \quad \left(\frac{\partial x_1=1}{\partial x_1} \right)$$

$$\dot{v}_2 = 0 \quad \left(\frac{\partial u_2=0}{\partial x_1} \right)$$

$$\dot{v}_3 = \frac{1}{v_1} = \frac{\dot{v}_1}{v_1} = 0.5 \quad \left(\frac{\partial v_3}{\partial x_1} = \frac{\partial \ln v_1}{\partial x_1} = \frac{1}{v_1} \cdot \frac{\partial v_1}{\partial x_1} \right)$$

$$\dot{v}_4 = \dot{v}_1 v_2 + \dot{v}_2 v_1 = 5$$

$$\dot{v}_5 = \cos v_2 \cdot \dot{v}_2 = 0$$

$$\dot{v}_6 = \frac{\dot{v}_1}{v_1} + \dot{v}_1 v_2 + \dot{v}_2 v_1 = \dot{v}_3 + \dot{v}_4 = 5.5$$

$$\dot{v}_7 = \dot{v}_6 - \dot{v}_5 = 5.5$$

$$\Rightarrow \frac{\partial y}{\partial u_1} = \dot{v}_7 = 5.5$$

• What is the problem with this method?

\Rightarrow if lots of input \Rightarrow huge storage $v_6 = \sqrt{V} \rightarrow$ so it is recursive AD mode

\hookrightarrow as n forward AD passes $\{$ for each input -1 pass $\}$

• What is the advantage about it? \downarrow to compute all gradients

\Rightarrow simplicity & low memory usage

\Rightarrow if $n < k$ for example $n=3, k=1000 \rightarrow$ \downarrow cost \downarrow

but for $n > k \Rightarrow$ forward mode AD (good)

problems of backpropagation as backprop is not good for gradient?

backprop output

→ M2: Reverse mode AD \leftarrow cost $\rightarrow O(K)$
 (Autograd)

• forward evaluation remains same

(for gradient computation) ↳ additionally: Compute adjoint i.e. $\frac{\partial y}{\partial v_i} = \bar{V}_i$

going from right \rightarrow left

$$\bar{V}_7 = \frac{\partial y}{\partial v_7} = 1$$

$$\begin{aligned}\bar{V}_6 &= \frac{\partial y}{\partial v_6} = \frac{\partial y}{\partial v_7} \cdot \frac{\partial v_7}{\partial v_6} && \text{use chain rule} \\ &= \bar{V}_7 \cdot \frac{\partial v_7}{\partial v_6}\end{aligned}$$

$$\begin{aligned}&\quad \dots \\ &= 1 \cdot (-1) = -1 && (\text{NOTE: } V_7 = V_6 - V_5)\end{aligned}$$

$$\begin{aligned}\bar{V}_5 &= \frac{\partial y}{\partial v_5} = \frac{\partial y}{\partial v_7} \cdot \frac{\partial v_7}{\partial v_5} && (\text{Since partial differentiation is independent}) \\ &= 1 \cdot (-1) = -1 && \text{i.e. } V_5\end{aligned}$$

$$\begin{aligned}\bar{V}_4 &= \frac{\partial y}{\partial v_4} = \frac{\partial y}{\partial v_6} \cdot \frac{\partial v_6}{\partial v_4} \\ &= \bar{V}_6 \cdot 1\end{aligned}$$

$$= 1 \cdot 1 = 1$$

$$\bar{V}_3 = \frac{\partial y}{\partial v_3} = \frac{\partial y}{\partial v_6} \cdot \frac{\partial v_6}{\partial v_3}$$

Fix (i.e. it is coming in for 2 different nodes) $\rightarrow V_2 \rightarrow V_4 \rightarrow V_6$

$$\begin{cases} \bar{V}_2 = \frac{\partial y}{\partial v_2} & \text{if it is coming in for 2 paths} \\ & \text{exist for } V_2 \end{cases}$$

Similarly for V_1 so both contributions we need to consider

$$*\overline{v}_2 = \overline{v}_s \frac{\partial v_s}{\partial v_2} + \overline{v}_q \cdot \frac{\partial v_q}{\partial v_2} \quad (\text{we say } v_s \text{ & } v_q \text{ both contribute equally to } v_2)$$

$$\Rightarrow \overline{v}_2 = v_s \cos v_2 + \overline{v}_q - v_1$$

$$= 1.716$$

$$*\overline{v}_1 = \overline{v}_q \frac{\partial v_q}{\partial v_1} + \overline{v}_2 \frac{\partial v_2}{\partial v_1}$$

$$= \overline{v}_q v_2 + \overline{v}_3 \frac{1}{v_1}$$

$$(v_1, v_2) = 5.5$$

Derivation for the multiple

How to arrive at *? (using 1 to 3 pathway case)

(1 to 3 to 1 to 0) (using no self loops diagram)



$$\overline{y} = f(v_2, v_3)$$

$$\overline{v}_1 = \frac{\partial \overline{y}}{\partial v_1} = \underbrace{\frac{\partial f(v_2, v_3)}{\partial v_2} \cdot \frac{\partial v_2}{\partial v_1} + \frac{\partial f(v_2, v_3)}{\partial v_3} \cdot \frac{\partial v_3}{\partial v_1}}$$

(if i to j = 0 then 0)

(partial adjoint can be computed for

every edge

Partial adjoint $\overline{v}_i \rightarrow j = \overline{v}_j \frac{\partial v_j}{\partial v_i}$ where $i \rightarrow j$ denotes an outgoing edge from i

$$\overline{v}_{iv} = \sum_{j \in \text{next}(i)} \overline{v}_{i \rightarrow j} \leftarrow \text{Sum all outgoing edge adjoint}$$

$$(v) \overline{v} = v \iff$$

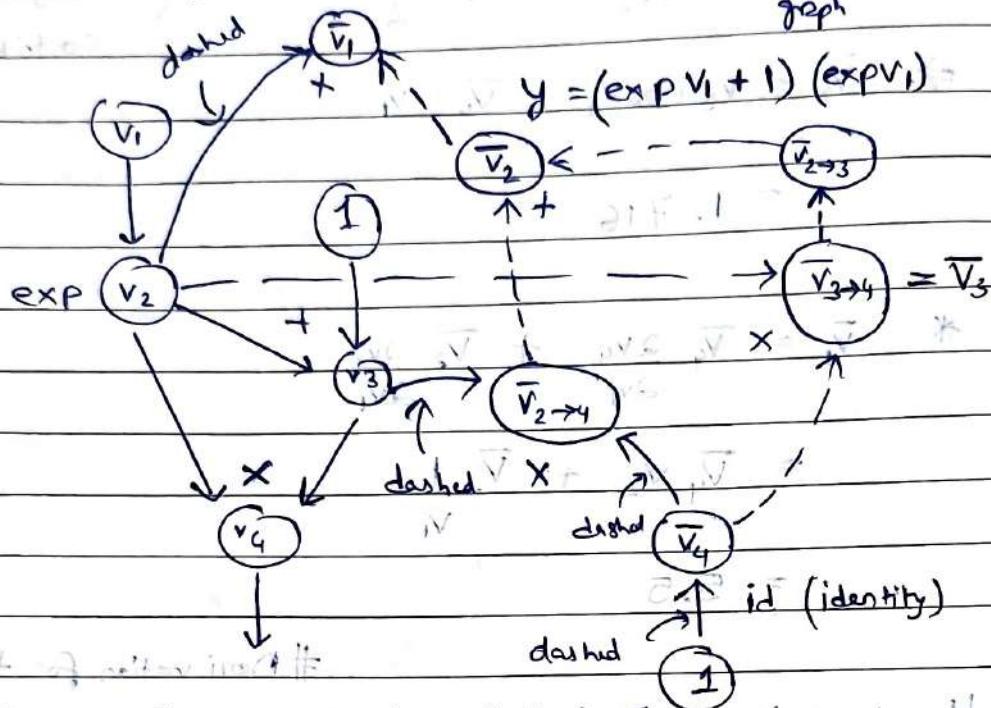
$$\overline{v} = \overline{v}$$

$$\overline{v} = v \cdot \overline{v} = (v) \overline{v} \cdot \overline{v} =$$

Running of algorithm

combined adjoint
of all precomputed
adjoints of partials
computations graph

- Reverse AD algorithm



Step 1 • If 1 comes at v_4 , sum of all partial adjoints but there are none (adjoint of $g(p=1)$)

• $v_2 \in v_3$ (for all $k \in \text{input}(v_4) \rightarrow$ we need to calc partial adjoint
 $(v_2, v_3) \bar{V}_{2 \rightarrow 3} = \bar{V}_4 \cdot \frac{\partial v_4}{\partial v_2} = \bar{V}_4 \cdot v_3$ as $v_4 = v_2 \cdot v_3$)

$\bar{V}_{3 \rightarrow 4} = \bar{V}_4 \cdot \frac{\partial v_4}{\partial v_3} = \bar{V}_4 \cdot v_2 = \bar{V}_3$

(as partial adjoint is the
read adjoint)

$i = 3$

Step 2 $\bar{V}_{2 \rightarrow 3} = \bar{V}_3 = \bar{V}_3 \frac{\partial v_3}{\partial v_2} = \bar{V}_3 \cdot \frac{\partial \exp(v_1) + 1}{\partial v_2} = \bar{V}_3 \cdot \frac{\exp(v_1)}{\partial v_2}$

$i = 2$

Step 3 $v_2 \rightarrow v_3, v_4$ (input \rightarrow)

$$\begin{aligned} \bar{V}_{1 \rightarrow 2} &= \bar{V}_2 \cdot \frac{\partial v_2}{\partial v_1} \Rightarrow \text{as } v_2 = \exp(v_1) \\ &= \bar{V}_2 \cdot \frac{\partial \exp(v_1)}{\partial v_1} = \bar{V}_2 \cdot v_1 = \bar{V}_1 \end{aligned}$$

NOTE: here we are calculating partial adjoints beforehand \Rightarrow time complexity reduced!

- no need to arrive at v_1 as it's the starting node!

NOTE: if the values change for $v_1, v_2, v_3, \dots, v_n$ then gradients are gonna change (so we may need to recompute)

Let's see if we see the computational graph, we have the exact ~~the~~ relation \rightarrow reducing the no. of unnecessary computations (\approx grade)

hence it is beneficial than

conventional backprop $\{ = \text{which one using } \eta \text{?} \}$

\rightarrow req both forward

↓ backward
no need of backward pass as already pre-computed

\rightarrow Reverse mode AD by extending computational graph

\rightarrow since backprop is implementation of reverse mode AD, what's different?

\rightarrow how I think it was, this? once the graph is plotted it's fixed (i.e. the graph is fixed once the function is written) \rightarrow on changing the input? only the numerical values flowing through those formulas change, not the chain rule itself

now is planning

variables

$0 = 0, 1 = 1$

minimum will be $\theta + \sqrt{(\theta-1)^2}, \theta = 1 \Rightarrow$

$\theta = 2$

if max of minimum is $\leftarrow 0 \Rightarrow \theta$

not existing

Here pbs is: when we do GRP, ht
update happens by looking at
stepout paths/dircn !!

classmate

Date 19/08/05

Page

Optimization ($\#$ Fully connected Network) : FCN

↓ memory hungry

$$(note: it might \leftarrow \begin{cases} Z_{i+1} = \underbrace{\sigma_i}_{\text{EIR man}}(Z_i W_i + \underbrace{1 b_i^T}_{1 \times 1 \text{ } b_i^T \times n}) & \rightarrow \text{so we use broadcasting} \\ \text{be last layer} & \text{ie instead of storing} \\ \text{is non activation} \rightarrow \text{key Qs on (FCN)} & 1 \times 1 \text{ back} \\ \text{bias} & \end{cases}$$

Q How do we choose no. of layers?
no. of neurons per layer?
⇒ we just repeat b
Matrix \otimes times without
storing it

Optimization

(i) considering a convex function, (gradient descent again discussed)

→ next, Adagrad

Nadamit method

memory and pt ←

$$\text{Loss } f(\theta) = \frac{1}{2} \theta^T P \theta + q^T \theta \quad (\text{quadratic})$$

→ keep track of past actions / directions

to avoid fig 1 later of movement

(# level set diagram
first fig)

(ii) Momentum $\rightarrow \epsilon \in [0, 1]$

→ gives better trl

Informed dirn



$$U_{t+1} = \beta U_t + (1-\beta) \nabla_{\theta} f(\theta_t)$$

(dec with time)

$$\theta_{t+1} = \theta_t - \alpha U_{t+1}$$

→ minima θ_1

→ if we go out, loss inc

Previously it was

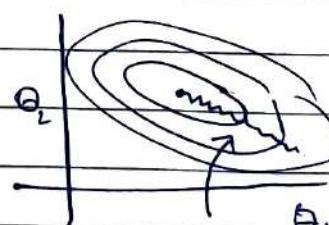
$$\text{let } U_0 = 0$$

Gradient

$$t=1, U_1 = (1-\beta) \nabla_{\theta} f(\theta_t)$$

now it's momentum θ_2

$$t=2,$$



If $\beta=0 \rightarrow$ no momentum; same like

previous case

After momentum

Pb1 → Slow start when using momentum as $v_0 = 0$

↓
for larger networks it is not good
↓
so modifying the wts update

(v_t smaller in initial iterations than in later ones)

unbiased momentum

$$\text{SOLN} \quad \Theta_{t+1} = \Theta_t - \alpha v_{t+1}$$

~~($\# = \alpha t + 1$)~~ ← unbiased factor

$$(1 - \beta^{t+1})$$

for $t = 0$, $\left(\frac{1-\beta^0}{1-\beta} \nabla \dots \right)$
(gives a big jump at start)

Nesterov Momentum

benefit: {converges faster & oscillates less}

$$\Rightarrow v_{t+1} = \beta v_t + (1 - \beta) \nabla_\theta f(\underbrace{\Theta_t - \alpha v_t}_{\text{previously only } \Theta_t})$$

$$\Theta_{t+1} = \Theta_t - \alpha v_{t+1}$$

↑ previously only Θ_t

NOTE: Momentum used for wts update

↳ till now 2-D wts

Pb 2 ↳ what about high dimensional wts?

⇒ multiple minimas, for diff wts

That's why we go for stochastic update
(i.e. batch update in epochs)

same α for Θ_1 & Θ_2 {in vanilla method}

↳ same α for all dimensions

(both tells different α for each dimensions we req.)!

↳ different α

↳ no drop

↳ Adagrad in momentum update (Adam)

(unit 1) ↳ ↳ ↳

Adam (most popular optimising method)

$$U_{t+1} = \beta_1 U_t + (1-\beta_1) \nabla_\theta f(\theta_t)$$

moving avg

$$V_{t+1} = \beta_2 V_t + (1-\beta_2) (\nabla_\theta f(\theta_t))^2$$

$$\theta_{t+1} = \theta_t - \frac{\alpha U_{t+1}}{V_{t+1}}$$

if $V_t \uparrow$, slow update

$$(\sqrt{V_{t+1}} + \epsilon)$$

else fast update

$$(\sqrt{V_{t+1}} + \epsilon) \Rightarrow \text{To avoid } 0 \text{ case}$$

optimisation: On top of this, stochastic gradient

stochastic gradient descent

(Instead of doing on small moments,

 $(\theta_t - \epsilon) + \sqrt{1-\epsilon} \cdot \text{do something with small moments})$

Initialisation of wts

• can we initialise to 0? \rightarrow no bias momentum

↳ wts the prob — backpropagation will become

pathological \rightarrow impossible to update \leftarrow ~~it's not possible to update~~

If we initialise to 0 (all gradients will be 0) — as wts are 0

• random wts \downarrow

↳ wts size distribution of wts don't change much

(adding same to original initial val.)

Easier if $\theta_0, \theta_1, \theta_2, \dots \approx 0$ \Rightarrow it's very important step: what's a good initialisation?! (just one credit goes to start with $N(0, \sigma^2 I) \ni w_i$)depends on 2 \leftarrow

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

Forward: \rightarrow how z_n mod is and?
 \rightarrow why normalisation req?

classmate

Date 25/08/23
Page

\rightarrow for stabilizing & speeding up training
 # Normalisation, Regularisation \rightarrow for preventing overfitting
 (initialisation std. $\sim \frac{1}{\sqrt{n}}$)

Why is $\frac{1}{\sqrt{n}}$ optimal? \rightarrow + \rightarrow in linear activation - $\frac{1}{\sqrt{n}}$ good but not in ReLU
 (ReLU activation considered)

(check it out) Reasoning: Why wt update should be initialized properly - as bad

• Independent R.V.: $x_i \sim N(0, 1)$, $w \sim N(0, 1/\sqrt{n})$

$$E[u_i w_i] = E[u_i] E[w_i] = 0$$

Initialisation may lead to bad result/

final value (proof)

$$\text{Var}[u_i w_i] = \text{Var}[u_i] \text{Var}[w_i] = \frac{1}{n}$$

$$\Rightarrow E[w^T u] = 0 \quad (n \times 0)$$

$$\text{Var}[w^T u] = 1 \quad (n \times \frac{1}{n})$$

$$\bullet z_n = \sum_{i=1}^n u_i w_i \quad (w^T u \rightarrow N(0, 1) \text{ by central limit theorem})$$

since, no activation non linear one,

$$\Rightarrow \text{if } z_i \sim N(0, \frac{1}{n} I) \text{ & } w \sim N(0, 1/\sqrt{n} I)$$



will have $N(0, 1)$ i.e. every layer

$$z_{i+1} = w^T z_i \sim N(0, 1/\sqrt{n} I)$$

will have

which is good to have mean 0 and normal distribution with

as it sets the variance to half!!

NOTE: if we use ReLU, then half component \leftarrow mean 0 &

add to 0 \rightarrow to compensate for loss, might distort var. 1/2

$$w \sim N(0, 2(1/\sqrt{n}))$$

NOTE: initialisation - messes up the network's training

If loss can persist throughout layers \leftarrow

(need for normalisation)

initial

→ however it's fine if you do it at initial normalisation!

follows \leftarrow

~~#Batch Norm completed) (?)~~

Why regularisation & types!!

a) Layer Norm (mean (tally))

$$\hat{z}_{i+1} = \sigma_i (w_i^T z_i + b_i)$$

→ across a dimension
we take norm

$$\hat{z}_{i+1} = \hat{z}_{i+1} - E(\hat{z}_{i+1})$$

$\frac{\text{Var}(\hat{z}_{i+1}) + \epsilon}{\sqrt{2}}$

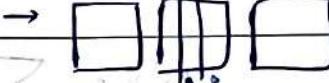
(gradient?)

- Outputs & activation not dependent on normalisation
- What's the pb here? (if we go std fully connected network)
⇒ used in transformers

⇒

(x, x)

(y, y)

↓
CNN as well

fr. So if we apply it in batch
layer since user a dimension
norm on batch, (i.e. applied) on norm
becomes 1

b) Batch Norm

- Where is it useful?

→ original data

↓

Noise added

↓ to change distribution

must and distinguishing difficult

↓ $\sqrt{I(x, y)}$ $\sim \sqrt{I} = \sqrt{I}$

→ but, model trained on out of ordered data during

training using batch norm → good result!

Problem? → unnormalized with what you see in DNN

→ Mini Batch Dependence is one characteristic of Batch Norm! When the normalisation of a single data point is calculated using the statistics

- Regularisation (mean & variance) of the specific, random minibatch it happens

→ more parameters in network - difficult to learn

→ solution space loss, if to memorise the data,

but not generalisation...)

→ Regularisation limits the complexity of the network

→ implicit
→ explicit

$$\underset{w_{1..n}}{\text{minimise}} \quad \frac{1}{m} \sum_{i=1}^m l(h_{W_{1..n}}(x^{(i)}), y^{(i)}) + \lambda \sum_{i=1}^n \|w_i\|_2^2$$

check it out

Type 1

$\rightarrow L_2$ regularisation aka weight decay (why?)

(models with)

smaller weights

are simpler &

$$w_i = w_i - \alpha \nabla_{w_i} L(h(x), y) - \alpha \lambda w_i$$

$$= (1 - \alpha \lambda) w_i - \alpha \nabla_{w_i} L(h(x), y)$$

\uparrow prevents exploding wt

\Rightarrow because of this term $(1 - \alpha \lambda)$ we say it's weight decay!

\Rightarrow as it shrinks the value at each step (as slightly less than 1)

Type 2 \rightarrow used during training only!!

\rightarrow Dropout (similar to ensemble method)

- we randomly drop some neurons from layer (???)

- randomly set some fraction of the activations at each layer to zero

$$\hat{z}_{i+1} = \sigma_i(w_i^T z_i + b_i)$$

~~\hat{z}_{i+1}~~

\rightarrow DL Basics Completed

$$(Z_{i+1})_j = \begin{cases} (\hat{z}_{i+1})_j / (1-p) & \text{with probability } 1-p \\ 0 & \text{with probability } p \end{cases}$$

For each of O neurons now with weight off by p .

\leftarrow activation

How Batch Norm Works

\rightarrow this method normalises the activations across the examples in a minibatch

* This makes the model's training sensitive to both the batch size & the other examples it is grouped with at each step!

Here focus

was CPU!!

A bit about Hardware

→ How to link DL to Hardware : Focus

• MLP

• Application - CNN in hardware

• Transformer

• DL for inference

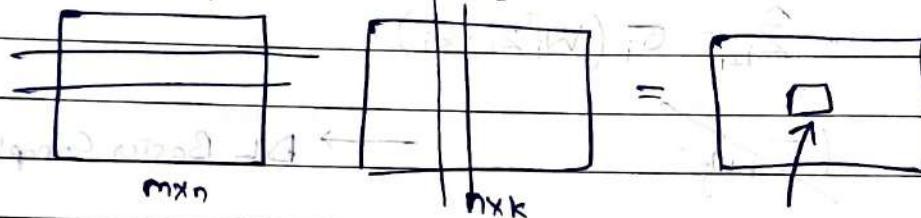
• Optimal Training Methods

→ Gain: in terms of energy

(If we use

hardware thing)

• Matrix Multiplication: Naive Way



• To get this, n computations req.

• What are the things ppl do while applying DL in hardware?

a) vectorization → (group computation better than single data)

how to add 2 arrays of length 256?

Computation)

A: 256×256 B: 256×256 To get C: 256×256

Naive: Load all 256 one by one

Optimisation: Load 4 at a time ($256/4 = 64$)

b) How to store data/matrix in memory

→ to store 256×256 directly : long/heavy

T1 → Row Major - store row by row

$$A[i,j] = A[\text{data}[i \times A.\text{shape}[1] + j]]$$

T2 → Column " " - col by col (1 col after another)

$A[i,j]$ T3 → Strides format - random position defined by strides (same offset)

$$= A[\text{data}[j \times A.\text{shape}[0]] + i]$$

$$\hookrightarrow A[i,j] = A[\text{data}[i \times A.\text{strides}[0] + j * A.\text{strides}[1]]]$$

how does this help?

- a lot of slicing in array : then helps

ex : $A[2:, :]$ ← skipping first row

- transpose

- broadcast } ← stride = 0 i.e. I can repeat sending it multiple times

Downside : as memory access not continuous

- vectorization harder

- (many, linear algebra operations may require compacting array first)

c) What other optimisations exists ?

- Parallelizing the operations (OpenMP)

(ex. loading data from 2 or different matrices at a time)

it depends on code if possible

Let in ML, since lot of matrix operation is possible

d) Matrix Vector Multiplication

→ memory heavy

→ Naive way : $O(n^2)$: this is seen when same type of memory

↑ when we do $(V \times V)$, it is used (for storing/loading
improve the placement of data in different

NOTE:

for $i = 0 \dots n$

for $j = 0 \dots n$

for $k = 0 \dots n$ } instead what we do is, we can bring A

can do is, we can bring B

but instead of B one at a time into L1 L2 " " - 7ns

(→ tile by $V \rightarrow$ P.C bring set of data)

DRAM - 200ns

Registers

L1 Cache - 0.5ns

NOTE: naive matrix multiplication

for $A \times B \rightarrow$ Load cost: $2 * \text{dram speed} * n^3$

Register cost: 3 //

classmate

Date _____
Page _____

Registers

d1) ^Tiled Matrix Multiplication

→ Block wise multiplication

→ Cost: n^3 / V_1 for A's DRAM

NOTE:

$$\bullet \text{time: } A[n/v_1][n/v_3][v_1][v_3]$$

$$B[n/v_2][n/v_3][v_2][v_3] \frac{n^3}{V_2} \text{ for B's DRAM}$$

• dram → reg time cost

= No. of iterations * amt of data brought

• Load cost

= time cost

* Speed!!

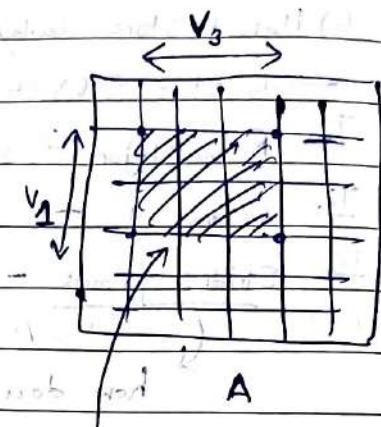
Registers cost: $V_1 \times V_3$ for A's register cost

Increase $V_2 \times V_3$ || B's || || register

Fine!! $V_1 \times V_2$ || C's || || register

for B: $V_2 \times V_3$

for C: $V_1 \times V_2$
stored



d2) Cache-line-aware tiling

→ Use of Cache previously not used

→ Loads from DRAM can be reduced

→ by using cache

Idea is:

from A → for $\frac{n}{b_1}$ times bring $b_1 \times n$

" B → for $\frac{n}{b_2}$ times bring $b_2 \times n$

constraint: $b_1 \times n + b_2 \times n \leq \text{L1 cache size}$

(as the registers part

for doing the

third for loop

reg) \rightarrow like $V_1 \times V_2$ at a time computed!!

HT-add
cost

L1 speed (instead of DRAM now)

$$\left(\frac{n^3}{V_1} + \frac{n^3}{V_2} \right) * L1 \text{ speed}$$

DRAM load
cost

$$\left\{ \begin{array}{l} \text{DRAM speed} \rightarrow \left(\frac{n^2 + n^3}{b_1} \right) * \text{DRAM speed} \\ (\text{extra for read } \frac{n^3}{b_1}) \end{array} \right.$$

01/09/05

GPU for DL

- CPU & GPU together coexist

- Why GPU?

→ on CPU — we optimise using memory

but more? — parallelise multiplication 4 times

at most

{if only 4 cores}

!!

but ~~one~~ (like give) 4 controllers

→ which not req as single operation

→ 1 controller enough

hence GPU ~~short to~~ { → more cores req !! }

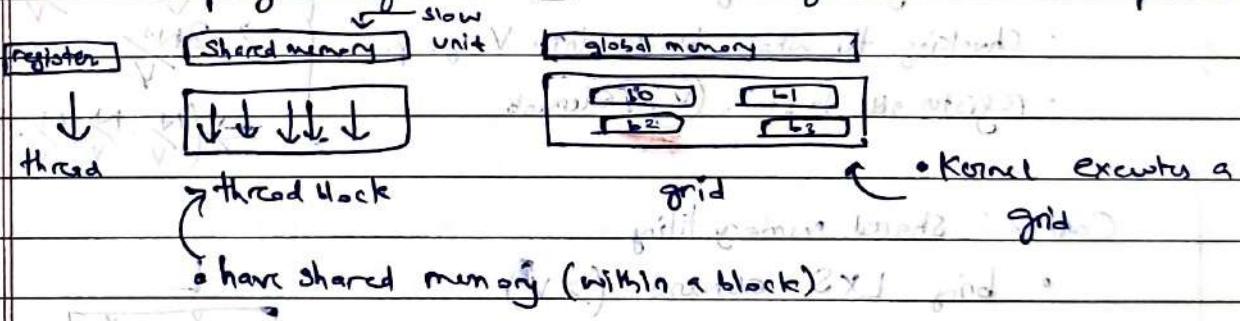
parallel

- massive computing units

- bare min instruction

→ Garnishing GPU hardware with cache is difficult

→ GPU programming mode: SIMD — single instruction multiple threads



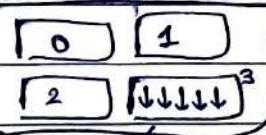
• have shared memory (within a block)

• executes same code, but can take different paths

- When we launch a kernel what exactly happens?

→ Thread Identified by blockid, gridid

(for global
id)



{blocks dim \times blockDim
+ thread id}

(5x3+4) \rightarrow in the reference
blocks

threadId in the blocks

grid block

↓

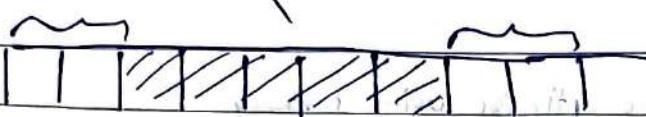
V.P.D size buf at

Use of hierarchy's memory to optimise task in GPU

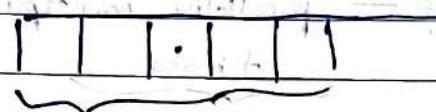
→ like we did in CPU

Example window sum \approx (1D CNN) ↪ NOTE: the idea is multiple threads executing

Input:



Output:



the same

⇒ each block will

be computed at

the same time

by looking

if I want to compute only these 5

⇒ 2 from left, & 2 from right only req

at rad left

⇒ not whole input array

& right!!

⇒ threads per block + 2 * radius

matrix multiplication on GPU

here if no tiling $\Rightarrow N^3$

but if tiling $\Rightarrow V^3$

Case 1: register tiling

• chunking the whole block into $V \times V$ groups

• register able to store $V \times V$ elements

$$\Rightarrow (N \cdot N \cdot V \cdot V)$$

Case 2: shared memory tiling

• bring $L \times S \times V$ elements ($> V$) into shared mem L

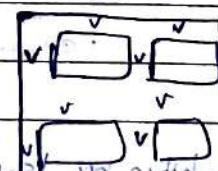
• NOTE: $L \cdot V = 0$ (reg)

global \rightarrow shared

what about S? \leftarrow (not reg)

$$\text{copy} = 2 * N^2 / L$$

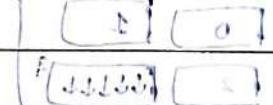
$$\text{Shared} \rightarrow \text{register} = 2N^3 / V$$



• for each core → warp level optimisation also exists ↑

↳ (skip slide 18)

↓ bl fast +



↓ bl fast

↓ bl

Multiprocuring Stream

ability to use thread & register

↓ width \rightarrow Thread width \leftarrow more data brought in register but less threads

thus we have auto tune

to find nice L & V

blocks & shared memory

↓ global memory

more data brought in memory than available blocks

- Mathematical constraints on N.N with the given hardware constraint

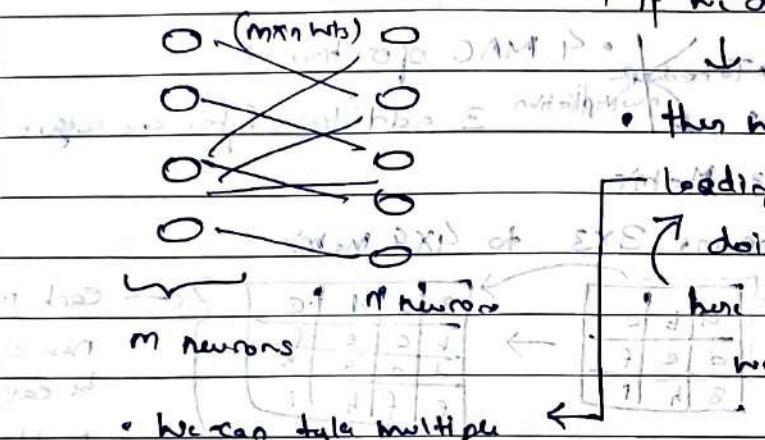
Solving these constraints
in mathematical form $\leftarrow \begin{cases} \text{whatever we discussed before} \\ \rightarrow \text{memory hierarchy} \end{cases}$

Specific to CNN

- $M \times N$ weights - huge storage & computation req \hookrightarrow fully connected

fully connected

\Rightarrow if we don't want to load all $M \times N$



- then we do (ig) 1D-CNN

loading (Rxs) wts at a time &

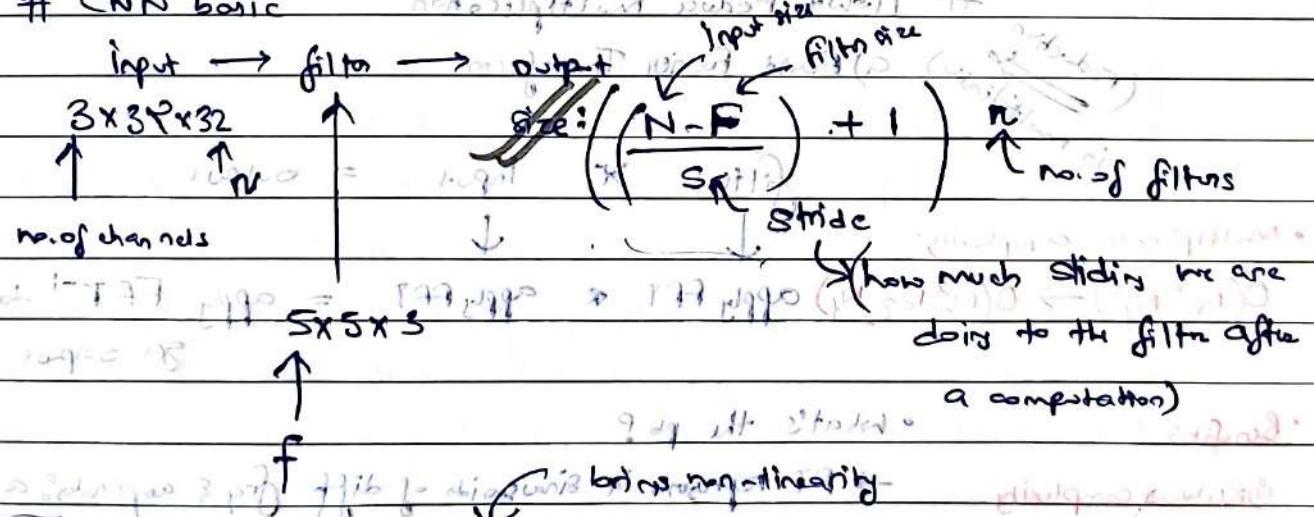
doing 'partial sum (something)'
here w.r.t instead of processing in 1D

we process in 2D

- We can take multiple

- Wts (filter or kernel) it's a function also : \rightarrow input presented in 2D
- \therefore wts = Rxs at a time

CNN basic



- What kind of activation functions?

- How large is the output? (sum of all elements)

\rightarrow After padding / with zeroizing, the input size becomes $\left(\frac{N-F+2P}{S} + 1 \right)$

$$S=1$$

$$F=2$$

$$9 \times 2 \rightarrow 4$$

Why should we do it?

\rightarrow how many lots of zeros added extra on each side

NOTE! in this method, if we look closely, we are reducing the size of inputs but no. of computations, it increases.

$$\Rightarrow N \times R \times S \times C \times (E \times F) \quad \left\{ \begin{array}{l} \text{for computation!!} \\ \text{Output} \end{array} \right. \quad O(n^7)$$

↑ batches ↑

• this is known as Convolution

→ How to optimise this?

↳ Computation?

• 4 MAC operations +

• To reduce multiplication

3 addition? for one output

Not efficient
(Same no. of multi)

Ideally

To split Matrix

Does it solve
anything?

→ another form

of representation?

Convert 3x3 to 4x4 matrix

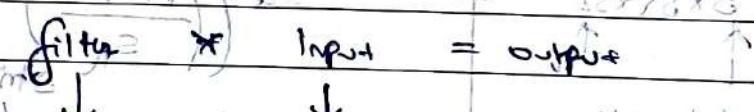
$$\begin{matrix} a & b & c \\ d & e & f \\ g & h & i \end{matrix} \rightarrow \begin{matrix} a & b & d & e \\ b & c & e & f \\ d & e & g & h \\ e & f & h & i \end{matrix}$$

each row
now can
be easily
multiplied

→ converted matrix becomes huge, size!

How to reduce multiplication

(Reduction in multiplication) a) Fast Fourier Transform



Multiplication complexity:

$$O(N_s^2 N_f^2) \rightarrow O(N_s^2 \log_2 N_f) \text{ apply FFT} * \text{apply FFT} = \text{apply } FFT^{-1} \text{ to filter}$$

Benefits:

• What's the pb?

(i) Reduced complexity

→ FT represents in sinusoids of diff freq & amplitudes. a

(ii) Larger convolution kernel,

↓ Signal processing of activation for better order

Better the performance due to which memory requirement (bigger than R & S)

→ though computation reduced, → energy hungry too!

↓ will be next

more details coming

this was no

Improvement → Can we have few more addition than multiplication? (Yes)
 Count on

FFT

$$\text{ex. } (a+bi)(c+di) = (ac+bd) + (bc+ad)i$$

4. MUL

3 add

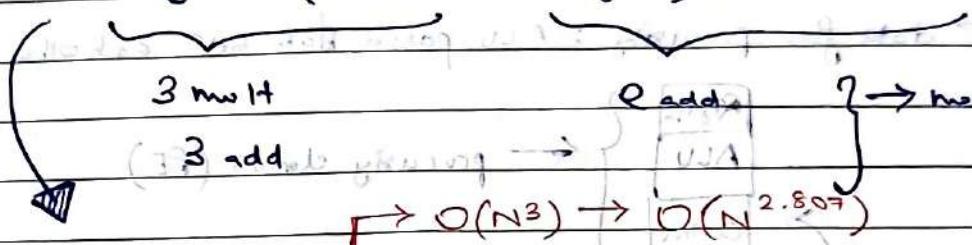
$$k_1 = c(a+b)$$

$$k_2 = a(d-c)$$

$$k_3 = b(c+d)$$

$$\text{Real part} = k_1 - k_3$$

$$\text{Imaginary part} = k_1 + k_2$$



• Reduces multiplies ← ~~#~~ Strassen's Algorithm — not much red but yes!

from 8 matrix mults ~~#~~ Winograd - 1D - $F(2,3)$ — targets convolution instead

• numerical stability loss
• requires more memory

* output size ↑ filter size ↓ of matrix mult! Hply

can be extended to 2D

⇒ (more add &

subtract & errors

accuracy)

↑ moment to moment p. error

↓ b b b

↓ b b b

09/09/25

GMAM

What is a convolution?

~~# Accelerator Architectures~~

• Why MACC Memory hungry?

→ 4. memory access req, per MAC op

→ MAC not liked in ML mostly because of this

→ Add nearby memory? helps in reducing the pk

→ Also parallelisation

wt?

• ASIC - has spatial architecture for DCE

→ Give all ALUs each individual controlling register can access

→ Nearby ALU can access others via register



Network on Chip

• SIMD/SIMT Architecture (In GPU mostly)

→ less communication betw ALUs

→ single controller (global one)

→ 2 tradeoffs in GPU

- (i) memory & blocks } discussed previously
- (ii) registers & threads }

• Spatial Architecture

→ data flow processing : ALU passes data betw each other



processing element (PE)

Creates a hierarchy of memory !!

- own memory } ← Small } ← So how to avoid re-access from DRAM?
- mostly PE memory }
- register file
- DRAM

32/64/128

→ #flops is of bandwidth?

- kind of how to utilize this small

own memory & mostly PE memory

→ generally first filter then activation

filter forward → filter J in height for 3AM -

filter output right! { but we can run : (i) filter + activation both
parallelizing, or (ii) only filter

(iii) Only activation

→ 3 kinds of data stored for a single convolution op

weights, input, and output

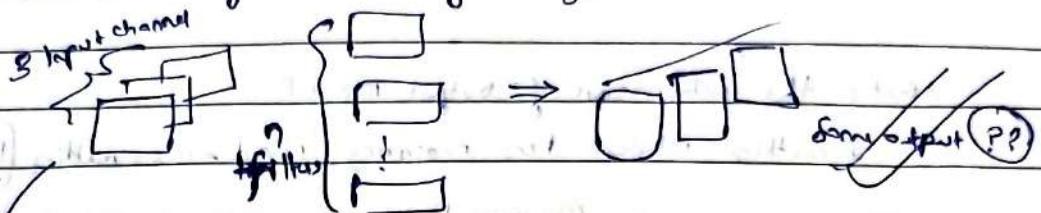
(i) Activation input

(ii) Activation output

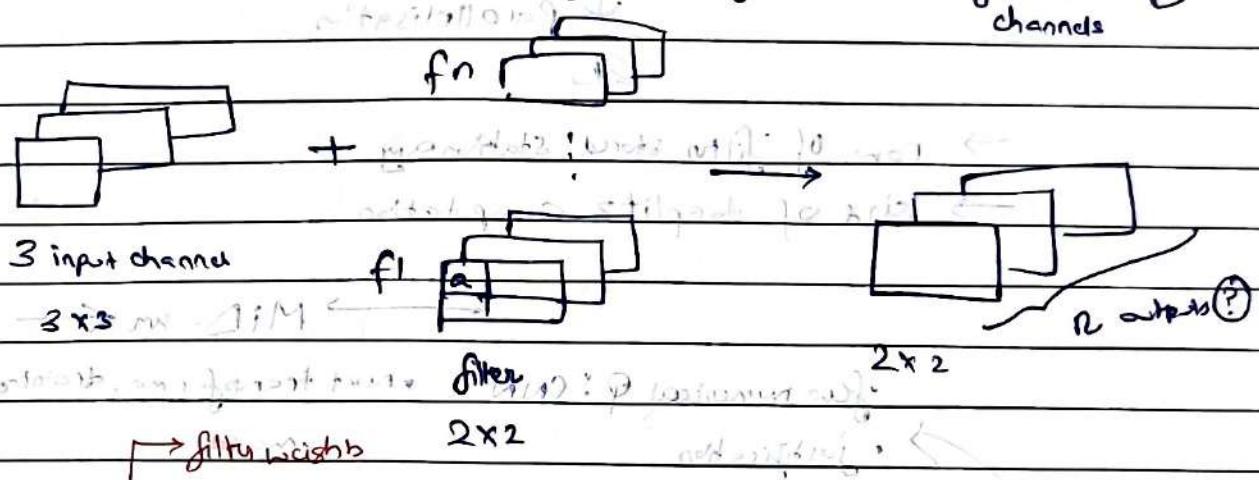
- Why are there multiple filters?

→ Shouldn't filter be constant w.r.t no. of channels?

→ how is no. of channel & filter affecting?



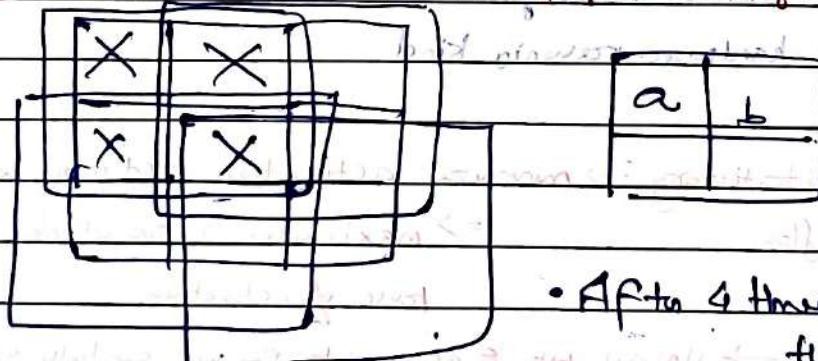
{ depth of filter = no. of channels ?
no. of filter = no. of output channels



- In Weight Stationary Design
- first compute for small Θ with all the possible matrices

- minimise energy the first input image

of reading weights \Rightarrow maximum convolutional result & a filter's use of weights



- After 4 times \rightarrow change/update the register value to b

- i.e. r times we update & r is filter dimension ($n \times m$) part of
- for each of those Θ we need to bring activation

- minimize the energy of reading & writing partial sums
- in output stationary design → maximise local accumulation of partial sums
- partial sum stored in register
 ↳ not weights / activation

What's the advantage of output row?

→ something related to memory read or something (?)

maximise reuse & accumulation at the RF

- Row Stationary ↑ level for all types of data (weights, activation, partial sums)
- can be done using 1D
 ↓ parallelisation

2D

→ row of filter stored stationary

→ kind of trapezoidal computation



few numerical Q: CNN • next transformer, training loop

• justification

• how to train CNN

• by hand (Floyd Warshall algorithm)

• optimisation: (minimising error)

• More on basis of DL

• hardware research kind

• input stationary: → minimize activation read energy consumption

dataflow

⇒ maximise convolutional & fmap

Used for sparse

reuse of activation

CNN strategy → Unicast Wts & accumulate partial sums spatially across

→ PE array

and it's implemented with 7x7

2x2 (max) window

and it's implemented with 3x3 window

with reuse

MidSEM revision

(i) Note: $y \in \{1, \dots, k\}^m = \begin{bmatrix} y(1) \\ \vdots \\ y(m) \end{bmatrix}_{m \times 1}$ where $y(1)$ is $[1, k]$

(ii) $\text{ML Basis} > L1 - v2.pdf > pg 24 - \text{optimisation}$
 {check back}

(iii) Backpropagation: Forward & Backward Pass

(iv) Forward Pass

initialise: $Z_1 = X$

iterate: $Z_{i+1} = \sigma_i(Z_i w_i) \quad i=1 \dots L$

(i) Backward Pass

initialise: $G_{L+1} = \nabla_{Z_{L+1}} \ell(Z_{L+1}, y) = S - I_y$

Iterate: $G_i = (G_{i+1} \circ \sigma'_i(Z_i w_i)) w_i^\top, i=L \dots 1$

$$\nabla_{w_i} \ell(Z_{K+1}, y) = Z_i^\top (G_{i+1} \circ \sigma'(Z_i w_i))$$

(v) Practice Questions of Automatic Differentiation (AD)

↳ Repeat the slide

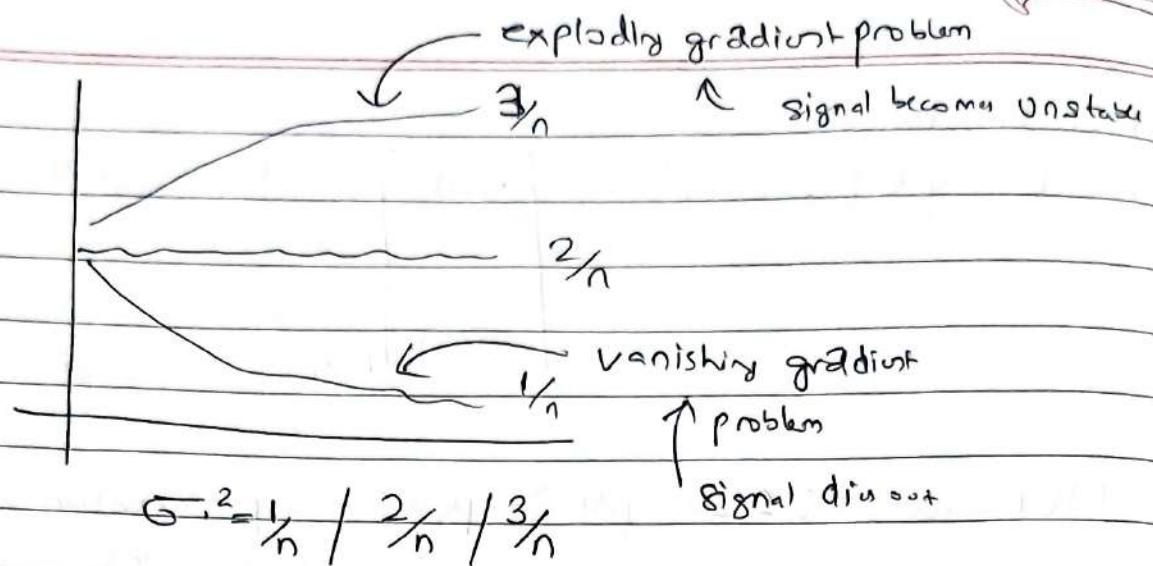
(vi) In the initialisation case,

$$w_i \sim N(0, \sigma^2 I) \quad \{ \text{randomly initializing} \}$$

The choice of variance σ^2 affect 2 quantities:

a) The norm of the forward activations, Z_i

b) The norm of the gradients $\nabla_{w_i} \ell(h_\theta(x), y)$



Reasoning for RELU_g (why $\frac{2}{n}$)

$$\text{Var}(y) = \text{Var}\left(\sum_{i=1}^n w_i n_i\right)$$

$$\text{as ID} \quad = \sum_{i=1}^n \text{Var}(w_i n_i)$$

$$I - S = (B - S) \quad \Rightarrow \quad I = B + S \\ = n \cdot \text{Var}(w n)$$

$$= n \left(E[w^2] - [E(w)]^2 \right)$$

$$= n \left(E[w^2] E[n^2] - (E[w] E[n])^2 \right)$$

$$\therefore E[w] = 0$$

$$= n \left(E[w^2] E[n^2] \right) \xrightarrow{*} \text{Var}(w) = E(w^2) \quad *$$

$$= n \text{Var}(w) E(n^2) \xrightarrow{*} E(n^4) = E(m \sim (0, \rho)^2)$$

$$= n \text{Var}(w) E(\rho^2)$$

$$= \int_0^\infty \max(0, p)^2 f(p) dp$$

$$\text{Var}(y) = n \text{Var}(w) \text{Var}(\rho)$$

$$= \int_{-\infty}^0 0^2 f(p) dp + \int_0^\infty p^2 f(p) dp$$

$$= \frac{1}{2} E[\rho^2] \quad \because E[\rho] = \int_{-\infty}^{\infty}$$

Final output $\text{var}(y) = \text{var}(\rho)$ which is input

$$\Rightarrow \frac{1}{n} = \frac{\text{Var}(\omega)}{n}$$

$$\boxed{\frac{1}{n} = \text{Var}(\omega)}$$

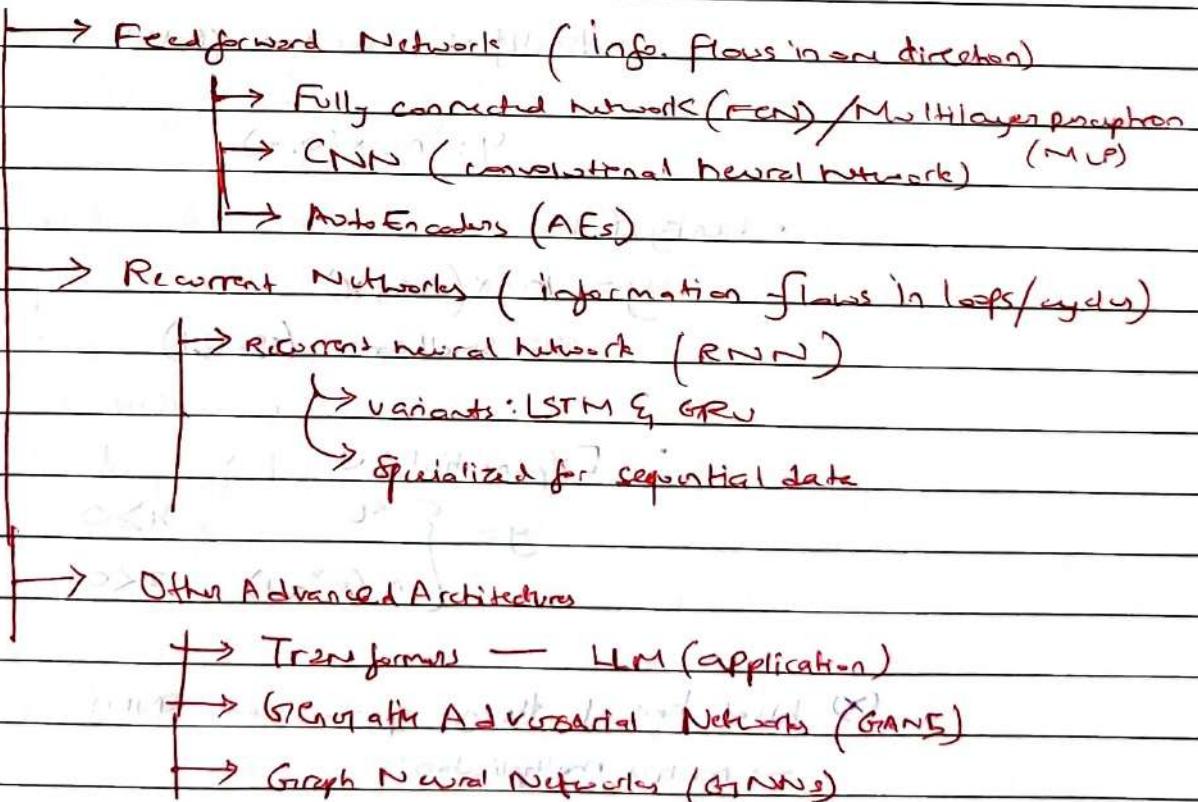
(v) Questions of Normalisation

- (i) Layer norm
 - (ii) Batch norm
 - (iii) Minibatch dependence
- } Which one is efficient & why?
any numerical question?

(vi) NOTE:

- depth of filter == depth of input map
- Many filters can be present
- for each filter we get a map \Rightarrow which determines the depth of the o/p

(vii) ANN/NN



AI

→ ML

→ Supervised Learning

→ Unsupervised //

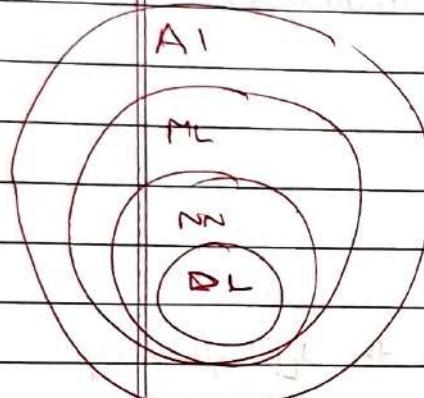
→ Reinforcement //

→ Deep // - utilizes ANN

→ ANN

→ RNN

→ Advanced Architectures



(ix) Activation Functions

- Sigmoid - $y = \frac{1}{1+e^{-x}}$

- Hyperbolic Tangent - $y = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

(x) Rectified Linear Unit (ReLU)

- $y = \max(0, x)$

• Leaky ReLU

- $y = \max(\alpha x, x)$

(where $\alpha = \text{small const. (e.g. 0.1)}$)

• Exponential LU

$$y = \begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

(and why) LU - confident

(xi) What kind of theory governs CNN?
→ Matrix multiplication?

$$h(x) = w \cdot x$$

classmate

Date 22/09/25
Page _____

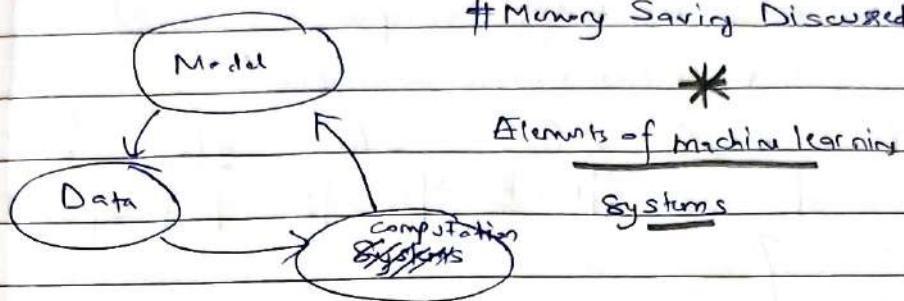
- Training Large Models - req for transforms

→ 2 things ^{Included:} memory saving

+ parallel operations

Sticking to GPU
(Optimization)

Memory Saving Discussed



- In computational graph, we mention the loss-grad, linear grad etc which (I think for backpropagation)
- but in inference not req

- What consumes memory?
 - i) weights {model} + momentum for Adam etc high order optimisation
 - ii) optimisation states {more in CNN? \Rightarrow huge memory req}
 - iii) intermediate activation values - requiring backprop

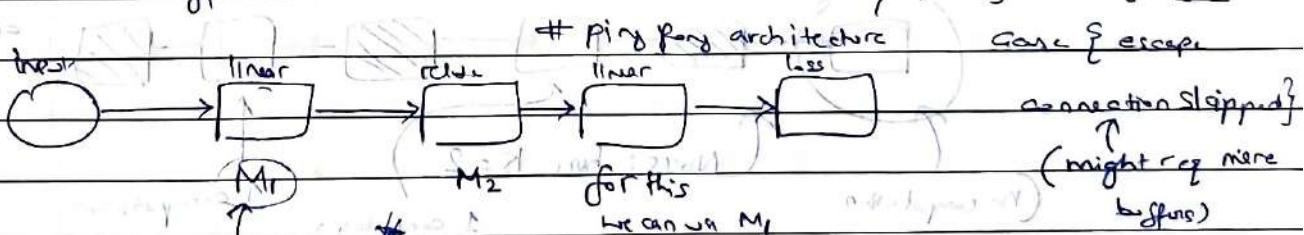
- What is used during inference? → weights

→ model weights is already optimised states {no}

handled by converting nn to CNN {small filter size considered at a time}

- How many memory buffers required for inference?

2 buffers



- Storing h_i, w already therein model (no explicit storing)

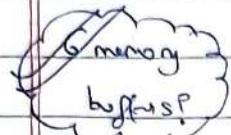
↑ both new hidden layer ~~XXXXXX~~ ↓ Model saved in memory

↳ going with first \Leftarrow 1st buffer now end

↓ Output

↳ all the time \Leftarrow output node

NOTE: To find out updates, we need to require store the outputs of each layer : $O(N)$ memory req



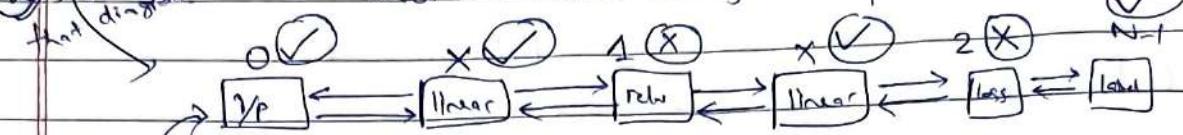
||

Sublinear Memory
Cost Theory

Instead we can do : Checkpointing some

- do we really need to do forward pass at runtime?
- Obj: which can reduce memory consumption

why for diagonal
that doesn't



Anyways need to store

- bit more forward pass to save memory

→ Don't store this but we can do forward

pass after reaching there & get the values
as prev. layer info. stored

• reduction in 2 memory buffers i.e. $6 - 2 = 4$

• Not just 1 method - there can be many checkpointing techniques

it can be represented in form of mathematical exp

→ so now &
happens!!

$$\text{Memory cost} = O(N) + O(K)$$

Using pipelining
architecture

Checkpoint cost

→ how many we are skipping through (Note: $K = N$; $K = \sqrt{N}$)

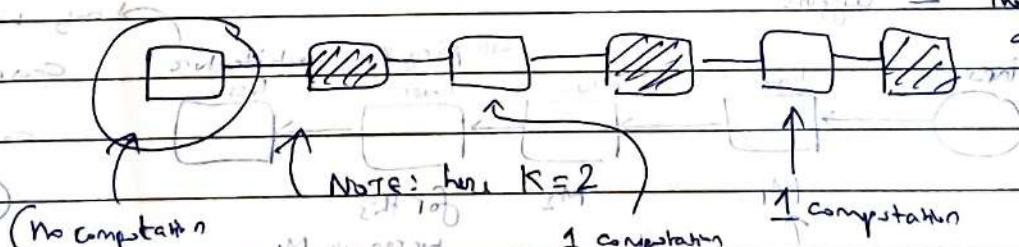
Note: N should be reversible

Sliding
generally

(for computing
how much

memory
are using)

In transformer, we represent layers in form of blocks



Note: here $K=2$

1 computation

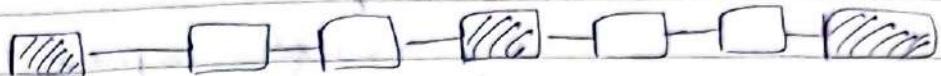
(req) (initialization) means almost parallel to forward pass.

Programmable hardware # No. of recomputations can be more than K

but we are sliding by $K \Rightarrow K-1$ of them aren't computed

When doing feed
pass again!!

So we need all of them in
memory!!



↑

is it possible in nn?

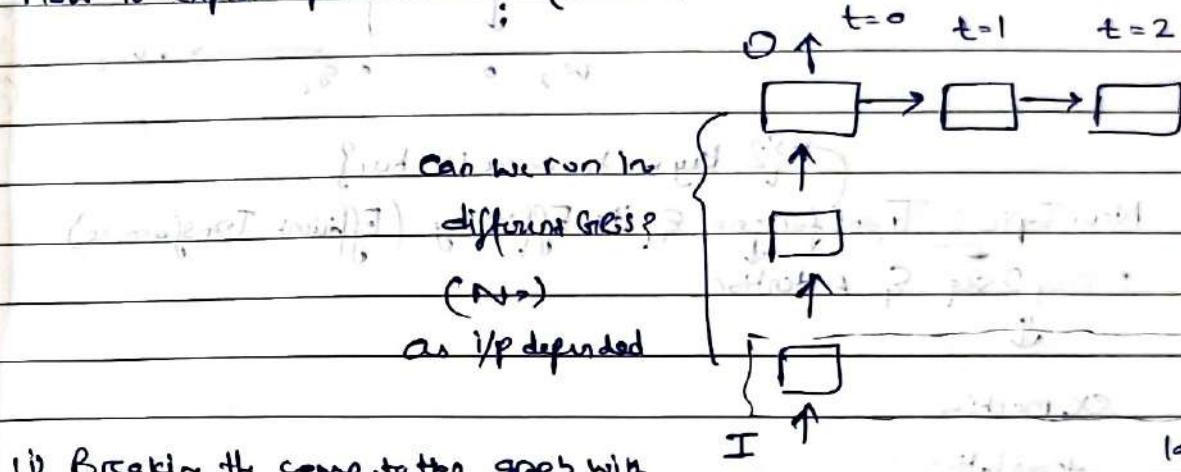
⇒ more memory req

What is KMP? ③

06/10/25

Topic: Training Large Models

How to exploit parallelism? (in GPU)



(i) Breaking the computation graph into

- Send/receive commands into 2 halves

(ii) In 1st cycle we cannot but somehow what we can do is when find pass happening → we can kind of do backward pass in parallel

M2: Data parallel training: Split the data among different GPUs, then sum it up (the partials)

→ How to do sum up of tensors → All reduction operation
→ How to replace the replica of models? ↓

→ Collect all G_i
→ Sum up all G_i = G₀
→ pass G₀ to all replicas so they can update by itself

Alt? ←
→ do global
↪ update &

{ What's the full backprop? }

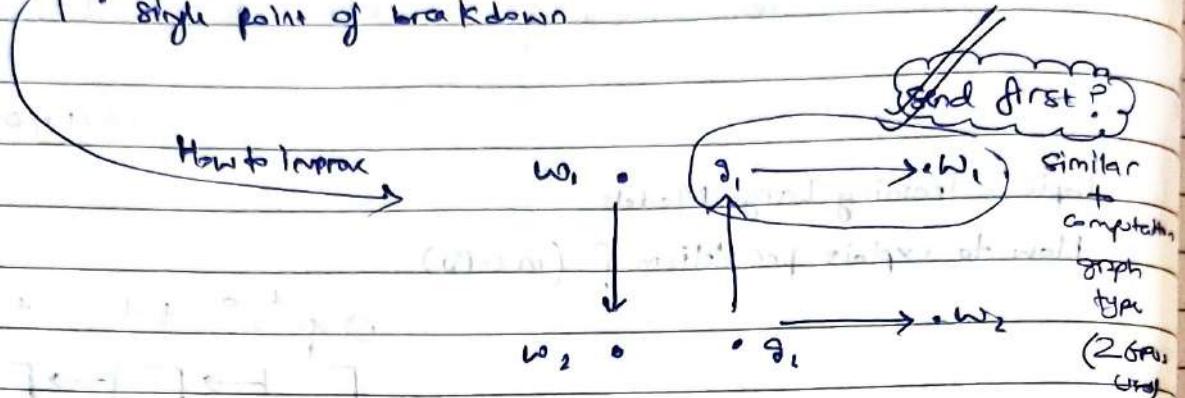
- each GPU will run nicely!
- All models will remain same but

then pass: Params serve abstracto! no!

Parameter Server Model - benefit?

- one server failure no problem → as in one replica failure
- fallbacks?

- privacy problem - as g_i shared
- single point of breakdown



→ {2 keywords come in here?}

New Topic: Transform & its Efficiency (Efficient Transformer)

- Seq2seq & Attention

↓

ex. machine

translation

- ⇒ encoder-decoder framework
- uses source representation
 - gets input → generate target seq
 - builds a sort of representation

embedding

• processes forward + previous prediction

to generate probability distribution

- simple model: 2 RNN

→ only 1 vector to represent whole next utter. for diff next utt

↓ input: that's given to decoder at t+1, etc.

→ may not be perfect representation of output

→ if seq is long → suffer from vanishing gradient

problem (gradient flow is faded)

→ EMA

input = old + new

(new new old old old)

problem

but max weight third element MA

problem

! can affect next element: every nth

#Attention

- Importance vector : like out of whole - some given proportion

no single vector only passed

- for current computation, we have attention computation which kind of determine off't of all input vectors which is important !!
or state ,

$$\begin{aligned} & \text{q}_k \text{ values 'bet' 0 to 1} \\ & \text{if } \sum q_k = 1 \end{aligned}$$

Applying softmax

To make everything differential

(like kind of give wts to all states)

How to compute score?

→ Dot product (s_i, h_t)

→ Intuitively it gives 1 for the decoder

07/10/25

Transformers : { has ~~cross~~ attention in all encoders; decoding one

↑ Why?

more 3D (3D) + what

Reg 1: parameters should not increase with input length

↳ [Sharing wts] !!

Reg 2: (----) should handle text of different lengths (i.e. independent of text length)

Q: Query, Key, Value, b: 3. Important terms in transformers.

Reg 3: Connection bet'n words → strength of these connections depends on words themselves

Req for design of transformer: like basic turns up with many of .

→ using only 3 wts: W_q, W_k, W_v {independent of input {i.e. shared}}

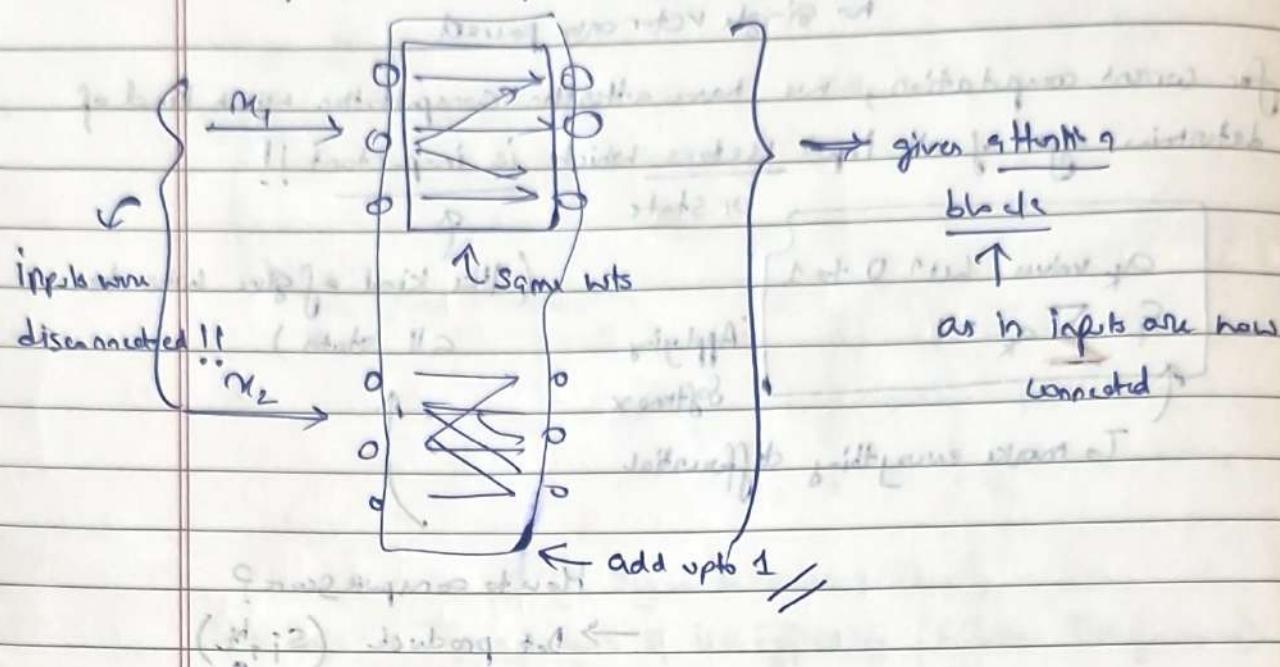
each word get's its own weight matrix \rightarrow length is not a problem

$$W_q \times \boxed{?} = \text{Query } q(P, P) \text{ where } P \text{ is not } \rightarrow$$

this is done through dot product → take attention then

• Attention as routing

a) Input $\rightarrow m_1, m_2, m_3$



b) Sparsity in wt matrix

c) hyper network - attention wts computed in somewhere else

d) How to compute attention wts - compute similarity betn query &

key words returns top n attention terms and then average them

Takes up $O(n^2)$ space

Can explore next class - optimisation of space

so scale it up!

Encoder + Position Encoding (P.P.) - What, When, How, what ...

Discusses signal (using sinusoids)

(F) Multi Head Attention: sort of doing training we split the

data & train \rightarrow output is combined but while taking care with

all data to reduce multiplications

\rightarrow Why do we need this?

\rightarrow to process this input who require a feed forward Neural Network

Networks are given - $\{W_i\}$ of depth d & E you give \leftarrow

Residual Connection: helps in training large data

\hookrightarrow sort of add the previous data!!

and nothing else \leftarrow taking the Complete Transformation (Pg 90) - covered!!

↳ like what kind of customisations

- How to use transformers for NLP — We use Tokenish
→ converting text to numerical value (31 - 111)

Next Lecture: Efficient Transformer // — How to improve the Architecture

13/10/25

Efficient Transformer

- Self Attention Lit Computation \rightarrow taken $O(n^2)$
- How to optimize? transform algorithm
- MHA: Multi Head Attention: most efficient: In Transformer steps
- Do we need to pay attention to all words? \rightarrow (No)

\Rightarrow fix the viewpoint \rightarrow learnt

To do this: apply Gaussian scale: left centred

SM1: attention are hardcoded

• What's the problem?

↑ in mind at word

M2: local attention & keeping window: K previous tokens

(99) step # longformer - Sliding window: 3 in right & 3 on left
{ handcrafted solution }

\hookrightarrow receptive view field is limited

: inputs \rightarrow 3 similar windows conv ps

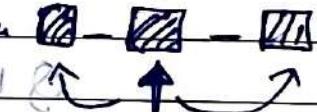
(ii) Global + Sliding windows

\rightarrow Not all global attention

but some of them +

Sliding windows

(ii) do dilation: i.e. instead of 3 concen-

\hookrightarrow they do it like 

\hookrightarrow leads to increase in field of view

• but performance is issue

• no. of computation steps !!

(02/10 - 16/10)

Inference

↳ reduce Matrix size = low rank ~~optimal~~ matrix

d - original embedding dimension

Note: Keys & values of size $N \times d$
 \Downarrow reduce to $k \times d$ ↳ What pb with it? \Rightarrow the sequence length is actually small for inference on this model!!if done on large \leftarrow (go neural)

• it fails

What to do? \rightarrow previously only this

• Can we care about data & architecture both?

 \Rightarrow based on diff ^{tokens} seq pos change the sliding window 'k'↳ # Routing attention? tokens learn / lead to attention
of tokens that belong to same cluster //

how to form it?

↳ longer K \Rightarrow longer clustering! same centroid for both keys

↳ better design & make more parallel - need Dp + Q question (??)

but still big with memory

so better optimization?

↳ Limitation of this design:

↳ Slower - build on k centroids shift to second (??)



matrix locality lost

↳ Where do we need large context windows? \rightarrow wordwin to \Rightarrow Question answering
just 23 surrounding token

!! need contiguous form

(Slide - 45/50)

Model Deployment

- not backprop exactly - but automatic differentiation
- Model trained - but how to use in different applications?

↓

→ as in different environment

data preprocessing
post "

- AT guided compilation - can it help the model optimisation for diff devices?

Learned Models

represented by neural simpler than AD one as it's inference

Represented in computational graphs, & lots of them

learnt by inference engine — like coreML, TF Lite

- How to represent it?

→ inference done using ping pong architecture

2 buffer

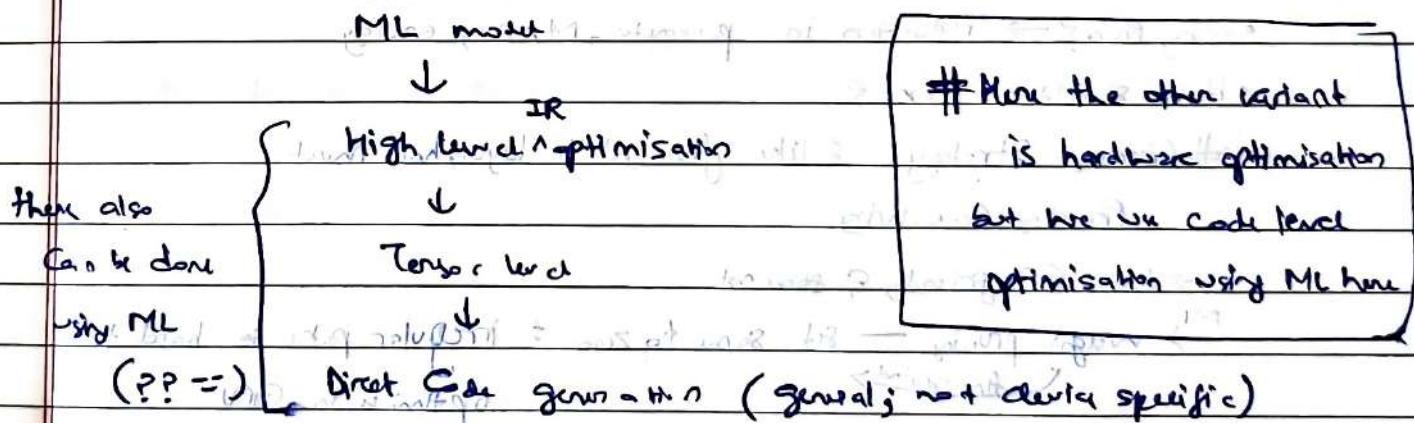
By Apple

By Google

They use libraries for it

→ but need a lot of customization & burden

Machine Learning Compilation



So for this one, we try to use the ML to search for all possible config

Pruning (#ML Perf)

Pruning

Instead of keeping all non-zero layers, remove some of them.

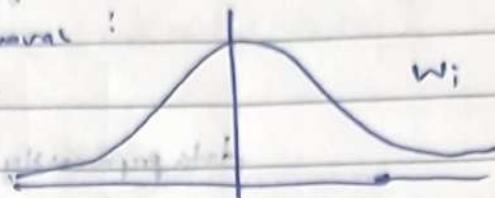
How to figure out which one to remove?

(→ as in weight removal)

S1

Blind Removal

{either remove layers or remove activations}
don't work → i.e. removing neurons with $w=0$



Note: only deal with additions

optimization req!!

→ extreme pruning +

fine tuning helps!!

• Hashing $\xrightarrow{\text{why?}}$ repeated multiplication of same type

S2 \downarrow to do multiplication

as it's expensive but addition not!!

S3 Iterative

pruning & fine tuning

$$\begin{array}{c} \downarrow \\ 1 \\ \uparrow \\ 5 \end{array} \quad \begin{array}{c} \downarrow \\ 7 \\ \uparrow \\ 3 \end{array} \quad \begin{array}{c} \downarrow \\ 11 \\ \uparrow \\ 7 \end{array} \quad \begin{array}{c} \downarrow \\ 34 = () \end{array}$$

$$\begin{array}{c} \downarrow \\ 1 \\ \uparrow \\ 5 \end{array} \quad \begin{array}{c} \downarrow \\ 7 \\ \uparrow \\ 3 \end{array} \quad \begin{array}{c} \downarrow \\ 11 \\ \uparrow \\ 7 \end{array} \quad \begin{array}{c} \downarrow \\ 34 = () \end{array}$$

number of activations go into final result

ratio $15/10/25$, mid 2M #

• Why Prune? → reduction in parameters, MACs, energy

• How should we prune?

Pruning Strategy: like from which layers, how much?

Pruning Granularity

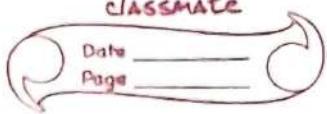
M1 → Some: GPU friendly & some not

\rightarrow weight pruning — set some to zero: irregular parts are hard to handle (fine grain) \rightarrow optimization in GPU

M2 → Coarse grained — set some rows to zero

at most 1000, no diff. of 0.03

some advice: No row zeros etc



- What can we prune?

- Input & of channel

- Kernel + (.....)

~~# what is kernel?~~

~~# channel?~~

~~# is pruning & pooling group~~

(Benefit)

Fine grained Pruning : flexible pruning indices

larger compression ratios

speed on custom hardware (not on GPU)

Pattern based Pruning : $\{(N:M)\}$ i.e. in every M block, N of them will be pruned?

ex. 2:4

X	.	X	.	X	X

← instead of storing all cells

simply store the non-zero cells

with corresponding indices

↓
more compact

Compressed matrix

Channel based Pruning : (reduces channel numbers in every layer)

↓ latency is preserved either uniformly or non-uniformly

latency is preserved either uniformly or non-uniformly

(??)

↓ 2M 3 and 4 channels

synopsis → we know there is a method

• How do we determine which cells to prune? → to do it, but how

→ equal among

do we select the cells?

M1 → remove less important one

→ how to we determine?

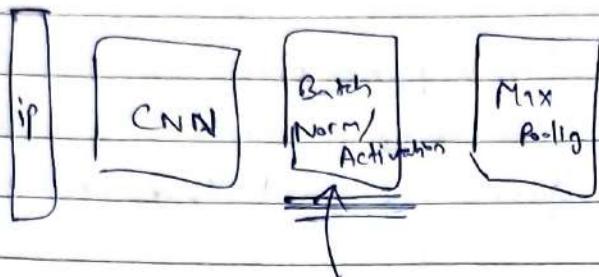
M2 → heuristic pruning criteria : L1, L2 norm & composite importance for

↓ now!!

• Can the algorithm learn it & do?

• Scaling factor associated with filter

- In batch normalization \rightarrow we kind of learn scaling factor
 \rightarrow that is only reutilized
 for the filter or



- Utilise scaling factor from here

Second Order based Pruning

minimising the loss after pruning is done! using taylor's series

Percentage of Zero Based Pruning

(here neuron some context involved!!)

- Heuristic methods: doge calculation (formula)

Neural network method: can use backprop to find the

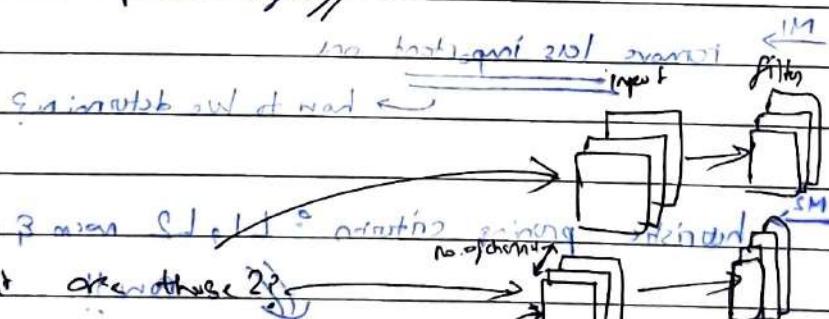
wights?

something: local loss & MSE

method 2: regression based pruning

and iteration \rightarrow some reduction in loss for the this has planned

pruned layers,



What are these?

what is kernel nothing but $W \times h$

what is stride nothing but?

Principle way to Prune :
 ↗ easy for specific hardware

• Pruning Ratio for each layer

M1 - uniform prune : latency high

M2 - AMC : latency low

↓ ↗ how fast we determine the result

• use RL

• as in come up with different sets (L_1, L_2, L_3, L_4, L_5)

• use RL to find the best set

M^{2-a} if we do grid search \Rightarrow sensitivity per layer

(Sensitivity analysis)



Q is RL & grid search sens? (NO)

→ Drawback: Analysing Accuracy & pruning for only a specific layer at a time & trying to make inference

Threshold T may not be achieved,

M^{2-b} & we go for RL method : overall goal (known)



per layer we need to determine

AMC: AutoML for Model Computation

• What is error limit? \Rightarrow 10% diff

NOTE: AutoML gives smaller model \downarrow despite giving same accuracy!

1x1 CNN - hard to prune
as only 1 unit

2x3 CNN - representation

N2-c : NotAdapt : doesn't use RL (non RL based)

• try few combinations pruning

↓
fine tune \Rightarrow accuracy high

Select K

• is the accuracy determined for that layer or the whole model?

How to improve performance of pruned model

→ on pruned → sparsity high

⇒ reduce learning rate : $\frac{1}{100}$ or $\frac{1}{10}$ of original

⇒

Regularisation

↳ penalise non-zero params

encourage smaller ||

• When L1, L2 better. P. will have sparsity of 20%

and both have about 20%

↳ apply both of them in case of finetuning

Now model is pruned in size

↳ how to optimise for hardware constraints

(3 methods) to target

(i) sparse weights → no need for their multiplication

(ii) sparse activation

(iii) quantization → making it fit in 8 bits

• Storing sparse matrices

SIM : hardware accelerator

↳ CSR format (M × N) of 2M + 3MN

↳ store lists of rows & cols (for non-zero)

↳ strip to efficient matrix patterns using SIMD instructions

reduce no. of computations (slide 86)

openmp parallel - multi EXE

multithreaded

multicore - SIMD EXE

(fixed 19 ms) 19 x 19 binary multiplication : 2 - 5 M

strong cache misses due to

elements being stored at diff. addresses ← interwoven

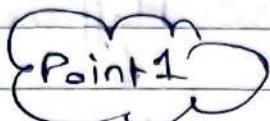
→ 20 loops / 64

4 to 18

28/10/23

Quantization

a) 32 bit $\xrightarrow[\text{cannot be done directly so how to?}]{} 4 \text{ bit}$



b) 32 floating point no

0 — zero / subnormal

1 - 254 — normal

255 — infinity / NAN

c) Why learning rate when kept high, bits turned to be NAN?

↓
how can we exceed a 32 bit FP range?

d) importance of w, b, x, Δ, ℒ : differ } \Rightarrow exponent &

exponent — controls range

fraction — precision

fraction size differs

in neural networks w, b \rightarrow exponent less req. on range

-1 to 4 } \Rightarrow precision

for input data \rightarrow more range req

e) Input data continuous $\xrightarrow{\text{transformed to discrete set}}$

through quantization

K-Means

Quantization ensures !!

\rightarrow if Point1 holds so centroids &

old no we get only !!

\rightarrow if centroids go very low after finetuning \rightarrow do cluster reassignment!

Rounding off

or generalization

but affects accuracy

• Pruning first \rightarrow Quantization next

• NOTE: Huffman Coding

\rightarrow one kind of bits representation
 \rightarrow motivation: all weights do not come with same frequency

• How do we know linear quantized multiplication is efficient?

$$Y = \underbrace{WX}_{\text{FP}}$$

{ as all

parameters are binary \rightarrow no multiplication in quantized process

Complex process \rightarrow as all elements already quantized & already quantized

• NOTE: $Z_W = 0$ {empirical observation}

{ conventionally?

31/10/23

• Quantizers { Unif. \downarrow \rightarrow tensorial nature of

for whole channel 1 (not good)

(i) have scaling factor per channel

(ii) per channel

(iii) per group

y: pertensor

~~Types~~

$$S = \frac{f_{\max} - f_{\min}}{2^N - 1}$$

type: vs Quant

Sig: per vector

b) Quantization error: if measures how good quantization is

\rightarrow it should be low

c) EFMO - Exponential 8 Mantissa 0/1

$M \times 4 \rightarrow SIM2$: Signed 1 bit, Mantissa 2 bits

(S1) WHAT SHOULD BE THE SCALING FACTOR?

$\rightarrow \gamma_{\min} \& \gamma_{\max}$: unknown

so using statistical observation on training data, we figure it out for $\gamma_{\min} \& \gamma_{\max}$

- needn't be extreme vals
- How to determine? — Minimise MP'sq error
or minimise KL divergence
- then follow with clipping: it helps

(S2) Rounding Schemes

(i) Stochastic Rounding

$\hookrightarrow W_x$: Activation output

\hookrightarrow here we try to learn in what range we should

~~round~~ quantized by learning from activation output

- MobileNet V.1/2 — efficient modules (like no redundancy)

\hookrightarrow addition of weights

\hookrightarrow Post Quant — Training doesn't work well

\Rightarrow So hardware specific we need to determine!

\Rightarrow During training we learn to quantize

\hookrightarrow we have both W_f & quantized wts during training

\hookrightarrow we use quantized wts in testing

- Why do we care about binary & ternary quantization?

\hookrightarrow because we want to avoid Multiplication for float

addition, subtraction, multiplication operations!!

(float) with fixed number of bits \hookrightarrow γ \rightarrow \uparrow

(MPB) more room \uparrow

{last topic}

Offline KD → memory hungry

- Logits - output before softmax applied
 - Knowledge Distillation = using bigger model to train tiny models - so softmax output similar
- (to make ←)

smaller networks

as good
as large
model

Student Teacher model

Temperature: helps in controlling softmax output

already trained

Matching of Student & Teacher networks: logits

(i) weights should match when Student networks projected to same dimensions

(ii) distillation loss: some pb!!

match logits & output dimensions same, to match channels

How to find similar projections? - 1x1 convolution

MLP for matching & learning

{ dimensions not same}

on projection - another loss introduced as later
are learning to match the wts(iii) ~~match intermediate features / activation output~~(iv) Matching gradient wrt input // $\frac{\partial L}{\partial w}$

helps in generating attention maps

softmax result needed to deg. DNN

(v) Matching sparsity patterns: (0s)

(vi) Matching relational information

→ 2 ways: relation within the model (s_i, s_j)" across samples (s_i, s_{i+1})

self

Offline KD

(i) iterative classification & distillation !!

(ii) Pb ? - req some step

May not work in all situations : not better accuracy

Online KD

(i) mutual learning

(ii) no pretrained Teacher

} ← will not size of tiny same
as big?
→ then whole is tiny com-
ing into picture?

Combined Self + Online

- Early exit networks : previous stage layers inpond after the future layers trained \Rightarrow future layers can be dumped

- Limitations : (i)

Previous Distillation for classification

↓ (done) !!

how to do for object detection

(3 HWs) { (i) object classification - foreground & background
that need to be taken care of } (ii) regression - bounding boxes
(iii) final output

how to match this output for
Student & teacher?

$L_b(R_s, R_t, y)$ interpretation !!

at 11:30 AM ←

(Pg 44)

• If teacher bad \Rightarrow student not follow teacher
~~if~~

• If teacher good \Rightarrow upper bound for student is teacher's

(Pg 55) - KD for LLMs (transformer)

o can we do KD of attentions?

(Pg 61) - network augmentation? - for small network (not good)

↳ like dropout, cut etc

↳ so they do NetAug !!

NetAug : During training increase the network!!
 ↴
 only works for small networks
 (not big ones)!!

→ ENDSEM ←

- a) How to modify Transformer architecture to handle images instead of text?

Step 1: Input Representation

↪ Transformer for text

Text → Tokens → Embedding Layer

for Images. Patch Based Approach as in Vision Transformer

Image → Non overlapping → Flatter patch

→ patches → convert to vector = patch vector

→ Project it to Transformer's embedding dimension

↓
Consider them as input tokens

↪ Pixel Based Approach

Pixel as tokens → pass it to

& RGB Values

embedding layer

as features

- Hybrid Approach

Use CNN to extract initial feature maps



flatten
↓

Embed these features as tokens

Step 2: Attention Mechanism

- Transformer for text

Self attention computes (i) pair wise interactions between tokens

(ii) capturing dependencies regardless of their position in the sequence

(iii) modulated by positional encoding

- for images

- Global Self Attention as in Vision Transformer

- (i) Apply self attention on seq of patch embeddings

- (ii) Each patch attends to all other patches capturing global spatial relationships

- Local Sparse Attention as in Swin Transformer

- (i) Restrict attention to nearby patches or pixels to model local spatial dependencies similar to convolutions

- (ii) Shifted windows allow gradual cross region mixing

Step 3: Positional Encoding

- in case of Texts

(i) It is added to tokens embeddings to encode the order of tokens in a 1D sequence

- in case of Image

(ii) 1D Positional Encoding

assign a positional embedding to each patch based on its index in the flattened Sequence

(iii) 2D Positional Encoding

encode (x, y) of each patch or pixel in the 2D grid using sinusoidal functions or learned embeddings

(iv) Relative Positional encoding

encode relative distances between patches (e.g. Manhattan distance) in the attention mechanism

b) Computing Attention output for first token
~~what about second or third?~~
 (Similar manner)

Given:

$$Q = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}$$

$$K = \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$V = \begin{bmatrix} 2 & 0 \\ 0 & 2 \\ 1 & 1 \end{bmatrix}$$

first token query = $q_1 = [1, 0]$

1. Dot products $q_1 \cdot k_j \quad [(q_1 \cdot k^T)]$

Score vector = $q_1 \cdot k_j$

$$\begin{aligned} &= \left[(1, 0) \cdot (1, 1), (1, 0) \cdot (0, 1), (1, 0) \cdot (1, 0) \right] \\ &= [1, 0, 1] \end{aligned}$$

2. Scale by $\sqrt{d_k} = 8$

$$[1, 0, 1] = [0.125, 0, 0.125]$$

3. Softmax

$$e^{0.125}, e^0 = 1, e^{-0.125}$$

$$\text{Sum} = e^{0.125} + e^0 + e^{-0.125} = S$$

$$w_1 = \frac{e^{0.125}}{S}$$

$$w_2 = \frac{e^0}{S}$$

$$w_3 = \frac{e^{-0.125}}{S}$$

4. Weighted sum of V

$$v_1 = [2, 0] \quad v_2 = [0, 2] \quad v_3 = [1, 1]$$

$$\text{Output: } w_1 v_1 + w_2 v_2 + w_3 v_3$$

$$\text{N-comp: } 2w_1 + w_2 \cdot 0 + w_3 = 2w_1 + w_3$$

$$y-11 : 2w_2 + w_3$$

\Rightarrow attention output: $[2, y]$

has keys/values

classmate
Date _____
Page _____

Cross attention

c) Suppose in one encoder-decoder attention head, all query vectors for a given output position in the decoder are zero, while encoder keys and values are arbitrary.

a) What is the resulting attention weight vector?

(i) Attention Scores : $q \cdot k_i = 0$

for every key as query vector
is a zero vector

(ii) Scaled Scores are all 0

(iii) after softmax,

$$\text{Weights} = \left[\frac{1}{N}, \frac{1}{N}, \dots, \frac{1}{N} \right]$$

\Rightarrow Uniform distribution over all encoder tokens

b) What is the Cross attention output for that position?

Cross attention output :

$$= \sum_{i=1}^N w_i v_i$$

Since every weight $w_i = \frac{1}{N}$

$$\text{Output} = \frac{1}{N} \sum_{i=1}^N v_i$$

\Rightarrow The output is the mean of all encoder Value vectors

d) # Light Weight Transformer

$$Q = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \\ 0 & 0 \end{bmatrix} \quad K = \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 1 & 0 \\ 0 & 0 \end{bmatrix} \quad V = \begin{bmatrix} 2 & 0 \\ 0 & 2 \\ 1 & 1 \\ 0 & 0 \end{bmatrix}$$

(i) first token: $q_1 = [1, 0]$

$$\sqrt{d_K} = \sqrt{8}$$

(ii) Score : $q_1 \cdot k_j$

with $k_1: (1, 0) \cdot (1, 1) = 1$

" $k_2: (0, 1) \cdot (1, 0) = 0$

" $k_3: (1, 0) \cdot (1, 0) = 1$

" $k_4: (1, 0) \cdot (0, 0) = 0$

$$[1, 0, 1, 0]$$

(iii) scale it : $\left[\frac{1}{\sqrt{8}}, 0, \frac{1}{\sqrt{8}}, 0 \right]$

(iv) apply softmax:

$$a_1 = e^{1/\sqrt{8}}$$

$$a_2 = e^0 = 1$$

$$a_3 = e^{-1/\sqrt{8}} = a_1$$

$$a_4 = e^0 = 1$$

$$S = a_1 + a_2 + a_3 + a_4$$

(v) weights $\Rightarrow w_1 = a_1/S, w_2 = a_2/S, w_3 = a_3/S, w_4 = a_4/S$

$$w_4 = a_4/S$$

(vi) Weighted sum of V

$$H\text{-component: } \frac{2a_1}{s} + \frac{a_2}{s} = \frac{3a_1}{s}$$

$$Y\text{-component: } \frac{2a_2}{s} + \frac{a_2}{s}$$

e) Difference b/w ~~that~~ Attention with ^{more}_{fewer} heads, smaller d_k v/s ~~heads~~ heads.

larger d_k

on factors:

- (i) expressivity (ability to capture distinct relationships)
- (ii) memory usage
- (iii) training stability

Case 1:

expressivity

→ 16 → 32

a) More heads, Smaller d_k

- Pros:
- high diversity of attention patterns
 - Each head can focus on a different relationship (syntax, negation, long range etc)
 - better for tasks needing many linguistic cues (translation, Q A)

- Cons:
- Smaller $d_k \Rightarrow$ each head has less capacity
 - per head representation is weaker

(# B reads over depth) (many simple specialists)

4 → 128

b) Fewer heads, larger d_k

Pros: → each head is more expressive → can model complex patterns

→ stronger per head feature extraction

Cons: → fewer heads → less diversity

→ might miss some specialized patterns

(# Depth over breadth) few powerful specialists

Case 2:

memory usage:

(a) More heads

→ more attention score matrices → slightly higher memory

→ must store $16 \times n \times n$ attention matrices

→ more intermediate tensors per head

→ more softmax + projection buffers

(b) Less heads

→ fewer attention score matrices → lower memory

→ only $4 \times n \times n$ score matrices

→ fewer intermediate activations

NOTE: Difference is small because total dimension and n^2 dominate memory.

Case 3: Training Stability

a) More heads

- Pros: → More heads: more diverse gradient signals
 → each head learns different patterns → reduces overfitting
 → diversity smooths the optimisation landscape
 → can improve stability

Cons: → very small d_k can cause unstable attention scores

- Scores become very noisy
 - Softmax may fluctuate noise
- risk of head collapse (weak heads) due to low capacity

Key Idea: Stability from diversity, instability from tiny d_k

b) Fewer Heads

Pros: → Larger d_k → lower variance in

$$\frac{QK^T}{\sqrt{d_k}}$$

→ Softmax becomes smoother → more stable gradients

→ more expressive heads handle deeper patterns
 clearly

Cons : → Less diversity : fewer gradient directions
 → May reduce robustness since fewer heads share the learning load

Key Idea : Stability from stronger heads, but less diversity

f) Layer Normalization → what?

→ applied after self attention and feed forward networks in each transformer layer to stabilize training

why?

Why Layer Norm preferred over Batch Norm in Transformers ?

→ due to its robustness in handling variable sequence lengths

→ ability to stabilize training in deep networks

→ computationally efficient for parallel processing

When Batch Norm can be used ?

→ In specific scenarios such as an NLP Task plus

→ with fixed length sequences, larger batch sizes (trained on a homogeneous dataset with consistent statistical properties across samples)

Case 1: Sequence Length Variability

Layer Norm

(i) normalizes within a single token's features (across embedding dimension)

(ii) independent of sequence length

(iii) padding doesn't affect normalization

(iv) stable for NLP because each token gets normalized separately

Idea : Token wise normalization; length invariant

Batch Norm

- (i) normalizes across the batch (across tokens &nd Exmaples)
- (ii) requires consistent batch structure \rightarrow assumes homogenous input size
- (iii) variable length sequences cause issues: padding Skews mean/variance
- (iv) statistics become unreliable for small batches or lots of padding
- (v) can cause unstable training for NLP

Idea: sensitive to padding, & sequence length

Case 2: Training Dynamics

Layer Norm

- (i) normalizes each token independently \rightarrow stable for variable length NLP seq
- (ii) keeps activations and gradients consistent across layers \rightarrow crucial for deep transformers!
- (iii) works well with residual connections, preventing gradient explosion/vanishing

Batch Norm

- (i) depends on batch statistics, which are noisy for small or variable batches in NLP
- (ii) padding affects mean/variance \rightarrow leads to unstable normalization
- (iii) inference uses running averages \rightarrow can mismatch training statistics

Case 3: Computational Efficiency

(i) Layer Norm

(i) computes mean/variance per token over feature dimension \rightarrow cost $O(d)$

(ii) for a sequence of n tokens \rightarrow cost $O(n \cdot d)$ (independent of batch size)

(iii) no reliance on batch size \rightarrow simpler, faster and no cross-device sync needed

(iv) very efficient for large scale NLP with varying batch sizes

Batch Norm

(i) computes statistics across the batch \rightarrow cost $(N \cdot d)$ per layer

(ii) less efficient for small batches and memory heavy for large batches

(iii) requires synchronization across devices in distributed training \rightarrow extra overhead

(iv) computationally heavier and less suited for NLP transformers

g)