Numerical Representation.   32-bit sheem
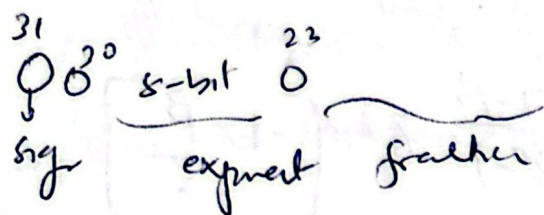
ys
) $2^{N+1}$, by 2's complement.  (Eg# $2^{9999}$ ⌐ $512 \times (2)^{333}$ → signfga post $= 2.9999$
                                                                    most s.value

Floating point numbers



31  Q $8^{20}$ 8-bit $0^{23}$                    o    $85.612 = 8 \times w^1 + 5 \times w^0 + 6 \times 6^{-1}, 1\bar{e}^2, 2\bar{e}^3$
sig    expnent  fralher

→ Modern 67/u  heue 16-bit flouty number, less precision


Quatnatun deluils

→ in  NN  linecr leeyer  have two matrices weights& Bcase,
   represented by floating point No , we aim to represent
   it by integers.   ($Y = X W + B$) int32

→ we aim t find reversible mapphing of float ⇌ int
   such that the performcam ce remain same
                [-127-127]



$2^4 = 32bit$-flut   $2^3 = 8bit$

244    127    244. 162 -

origin        φ        D-φ
                      loss-scme
                      of o

**Types:** Symetric vs Asymetric

| Symetric | Asymetric |
|---|---|
| $0$ <br> $\downarrow$ mapped to <br> $0$ alway <br><br> $[-127, +127]$ <br> $8$ bit <br> if input is symetric | $0$ <br> $\downarrow$ <br> diffret <br><br> $[0, 255]$ <br> if input is Asymetric. |

## Asymmetic Quantizatei

$$x_q = clamp\left(\left[\frac{x}{s}\right] + z ; 0 : 2^{n-1}\right), \quad s = \frac{a - \beta}{2^n - 1} \quad \overset{big \quad small}{\underset{\nearrow}{}} \quad \overset{round \; off}{\underset{\nearrow}{}} \quad z = \left[-1 \times \frac{\beta}{s}\right]$$

de quantize:

$$x_f = S(x_q - z)$$

## Symmetic Quantization

$$x_q = clamp\left(\left[\frac{x_f}{s}\right] ; \underset{-117}{-(2^{n-1}-1)} , \underset{127}{2^{n-1}-1}\right), \quad s = \frac{abs(a)}{2^{n-1}-1}$$

$$x_f = S x_q$$

---

→ inputs dynamic quantizater on fly. ($\alpha, \beta$ calculet)

→ Calibration: $y =$ output of quantized matrices so how can we convert it floating point when we have quantized so, we do it by calibration. (we run inference on the model using few inputs and observes the typical output (max, min) to calculate the scale & the zeros. (Post training Quantization)

OK, why?,

en adapting to specific task, LLMs have "low-intrinsic dimensity"
nd can still learn efficiently despite random projection
→ a smaller subspace. That's the hypothesis here.

( W contain many parameters that convey same information
as other, we can get rid of them without decreasing
performance. This kind of matrice = (rank-deficient)

→ (Rank = independent vectors)

→ Quantization.

The issue.

- Most LLM's have billions of parameter, if each parameter
is 32 bit then for llama 2 (7B) = $\frac{7 \times 10^9 \times 32}{8 \times 10^9}$ = 28GB disk
is require to store parameters & load in RAM (for inference)

- Just like human, PC are slow in computing floating point nos.

Solution = Quantization.

- it aims to reduce amount of bits required to represent
each parameter (usually by converting float to int)
9bit almost (10GB - 19B) depending on type of
quantization. (⊗ round or truncate )

→) also speed up computation (int) are fast to work with.

Advantages

→ less consumption when loading model.
→ less inference time due to same data type.
→ less energy consumption

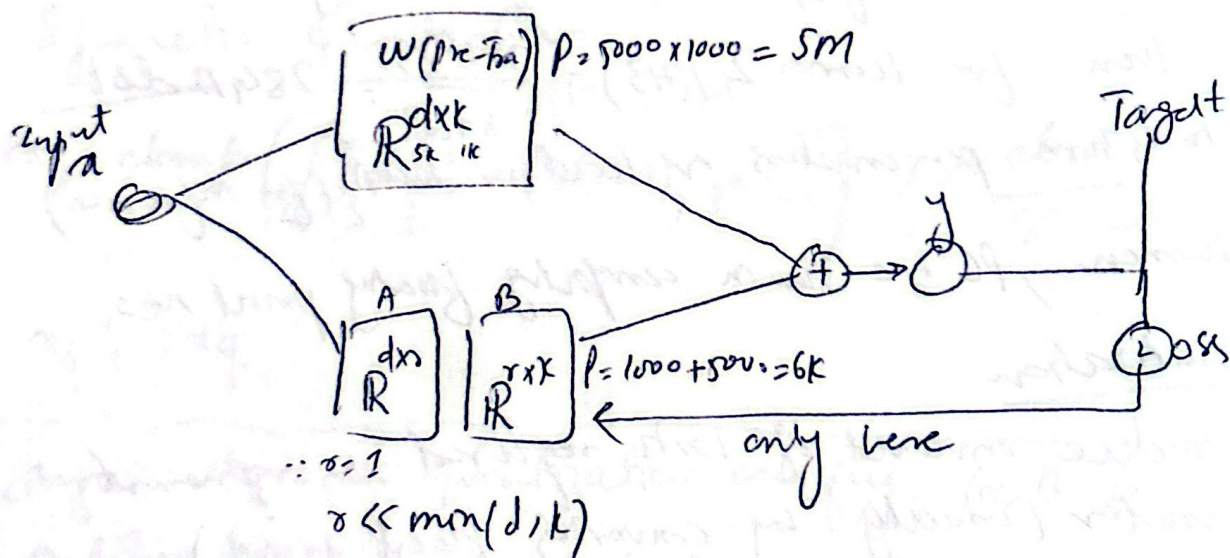# LORA :- low-Rank-Adaption of LLM's

→ To fine tuned LLM (pretrained) on specific task.

→ By freezing the dates.

## Problems in Full-Finetuing.

- computationaly expensive
- storage requirements for checkpoints are expensive.
- If we have multiple find tuned, we have to load all weights

---

## Fine-tuning

- we want to produce 2 matrix that multiply produce same result as $w$

- $(B, A)$ is lower representation of $w$, although dimension are same but we loose some information and structure that are



$$W(\text{Pre-Tra}) \quad P = 5000 \times 1000 = 5M$$

$$\mathbb{R}^{d \times k}$$
$$\mathbb{R}_{5k \cdot 1k}$$

input $a$

A
$$\mathbb{R}^{d \times s}$$

B
$$\mathbb{R}^{r \times k} \quad P = 1000 + 500 = 6K$$

$\therefore r = 1$

$r \ll \min(d, k)$

Target

$y$

$\oplus$

Loss

only here

## Benefits

→ less Parameters to train and store

→ less Parameter ⇒ less storage requirement

→ fast Backpropogation.

→ Faster switching b/w finetuned model

How do we choose $\alpha, \beta$ ??

Min-Max :

$a = max(V)$

$\beta = min(V)$

sensitive to outliers

Percentile :- set the range to the percentile of the distribution of $V$, to reduce sensitivity to outliers.

MSE: Choose $(\alpha, \beta)$ such that MSE is min b/w original & quantized

we usually use GridSearch

Cross-Entropy. when values in tensor are nA important but the distribution is. (Eg avg should chang, and other shouldn't also) so, we choose $\alpha, \beta$ the way cross entropy b/w $V$ and $\hat{V}$ is min
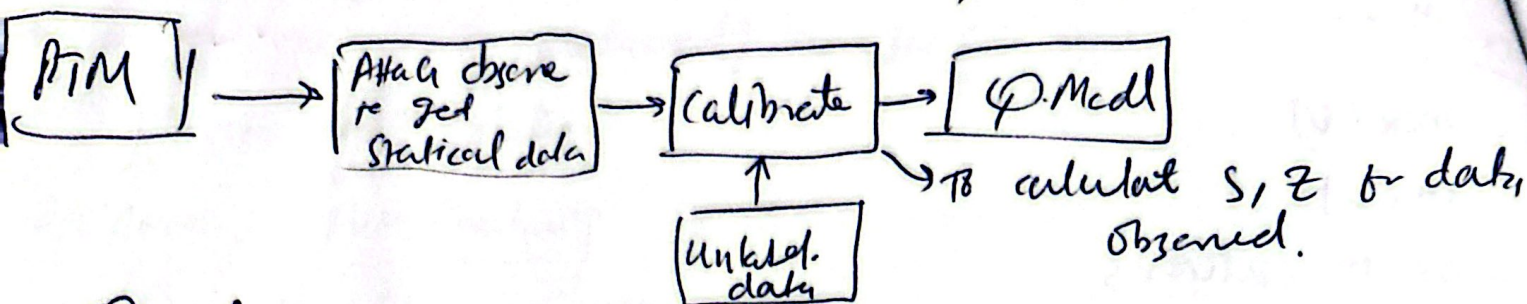
$$\underset{\alpha, \beta}{arg\ min}\ CrossEntropy\left(\left(softmax(V), softmax(\hat{V})\right)\right)$$
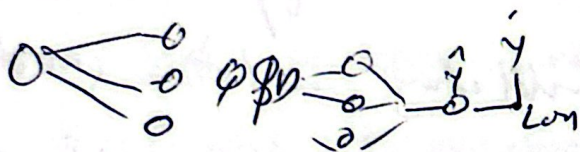
Quantization Granularity.

In CNNs,
we have many kernels of different size, so we loose some their quantization range if we used same $\alpha, \beta$ so we do CHANNEL WISE quantization, we calculat $\alpha, \beta$ for each channel.

# Post-Training Quantization :- (PTQ)

PTM → Attach Observer to get Statical data → Calibrate → Q.Model

↑ Unlabel. data

→ To calculate $S, Z$ for data observed.

## Quantization - Aware Training:

To make model more robust to quantization,
we do fake quantization/dequantization b/w layers while
training to introduce some error (quantization), so it can will
more robust to this error



Train further quantized based on observer connected during traing AD

## QAT: Gradient. (BP should also calculate gradient w. rit operation)

→ Gradient of Quantization (Operations) are non-differentiable.

→ So Backpropagation Algorithm calculate by approximation
of STE (Strat-through Estimator).

$$\alpha \,\text{---}\, \beta = (\text{Gradient} = 1) \quad b/w \; A \& B$$

except if (Gradient = 0)

why it works?? effects?? on loss

next:-
GPTQ, AWQ

The goal of QAT is help model reach local minima which is
more wide so that if weight after moves little the loss doesn't
more a lot.