

15/04/24

# How does Quantization Work

<https://newsletter.maartengrootendorst.com/p/a-visual-guide-to-quantization%C2%A7the-era-of-bit-llms-bitnet>

→ IEEE 754 standard of representing numbers

→ FP 16 bit :- (float 16 bit)

Sign  
(-1<sup>n</sup>)

2<sup>4</sup> | 2<sup>3</sup> | 2<sup>2</sup> | 2<sup>1</sup> | 2<sup>0</sup>

Exponent

1 bit

5 bits

2<sup>-1</sup> | 2<sup>-2</sup> | 2<sup>-3</sup> | . | . | . | . | . | 2<sup>-n</sup>

Fraction (mantissa)  
(Significant)

10 bits

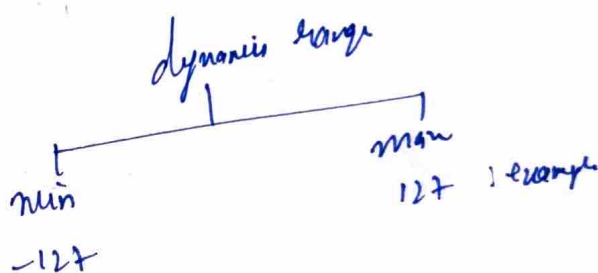
n = (0, 1)

→ higher the bit value, higher the precision

→ Dynamic Range :- Interval of representable numbers a given representation

→ Precision :- Distance between two neighboring values

→ low distance = high precision  
Inversely Proportional.



$$\text{memory} = \frac{\text{nr-bits}}{8} \times \text{nr-params.}$$

→ most models are in FP32 → Full Precision  
70B Param model need 280 280 GB of memory

70B model

32 Bit ⇒  $\frac{32}{8} \times 70B = 280 \text{ GB}$

64 Bit ⇒  $\frac{64}{8} \times 70B = 560 \text{ GB}$

16 Bit ⇒  $\frac{16}{8} \times 70B = 140 \text{ GB}$

→ If Precision ↓ Accuracy of model ↓

→ but we should, focus on reducing bits while maintaining accuracy

→ int8 → [-128 to 127]  
8 bit

→ FP32 → [-3.4e38 to 3.4e38]  
32 bits

Symmetric Quantization =  
not symmetric around the zero.

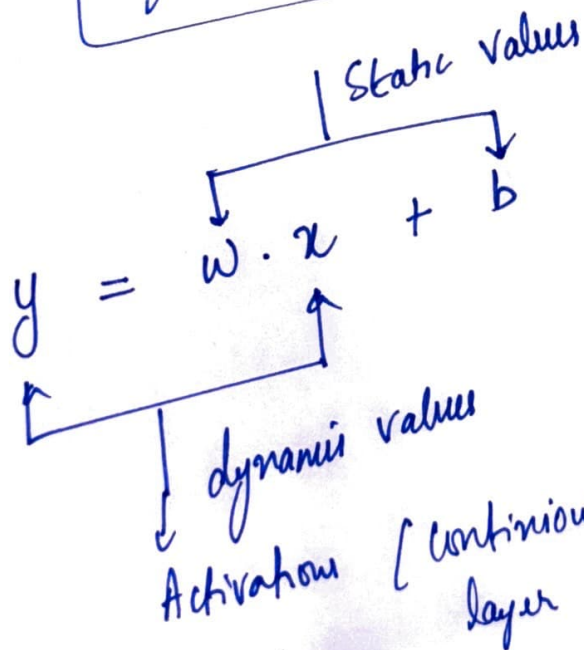
Asymmetric Quant =  
Vice versa

→ Absolute Maximum Quantization :- (abs max) Quantization form of symmetric

→ Scaling factor  $\Rightarrow S = \frac{2^{b-1} - 1}{\alpha}$

b = # of bytes that we want to quantize to  
 $\alpha$  = highest absolute value

$$x_{\text{quantized}} = \text{round}(S \cdot x)$$



## d) Calibrating Quantization

Methods :-

- ① (PTQ) Post Training Quantization :-  
Quantization after training
- ② (QAT) Quantization Aware Training :-  
" during training / fine tuning.

PTQ :- done after training. Ver. Symmetric or Asymmetric quantization.  
For Activation ( $x, y$ )  $\rightarrow$  As we don't know their values or it requires inference of the model.

It requires  $\left. \begin{array}{l} \text{① Dynamic Quantization} \\ \text{② Static Quantization} \end{array} \right\}$  for Activation ( $x, y$ )