

Quantization : Conversion from higher memory format to lower memory format.

Weights
of model

32 bits	32 bits	
	32 bits	

full precision 32
bits

We can convert this 32 bits into 8 bits.

int 32 \rightarrow int 8 $>$ Help w easier inference

How to perform Quantization

① Symmetric Quantization

② Asymmetric Quantization

(all the data is evenly distributed)
Symmetric Uint 8 :

$$\begin{array}{ccc} [0 \dots 1000] & \longrightarrow & [0 \dots 255] \\ x_{\min} & & q_{\min} \\ x_{\max} & & q_{\max} \end{array}$$

Min/Max Scalar

$$\begin{array}{ccc} 0 & & 1000 \\ \downarrow & & \downarrow \\ 0 & & 255 \end{array}$$

$$\text{Scale} = \frac{x_{\max} - x_{\min}}{q_{\max} - q_{\min}} = 3.92$$

convert 250 \rightarrow 8 bit

$$\Rightarrow \left(\frac{250}{3.92} \right) \Rightarrow 64$$

Asymmetric Quantization (non uniform distribution of data)

$$[-20.0 \dots \dots \dots 1000.0]$$



$$[0 \dots \dots \dots 255]$$

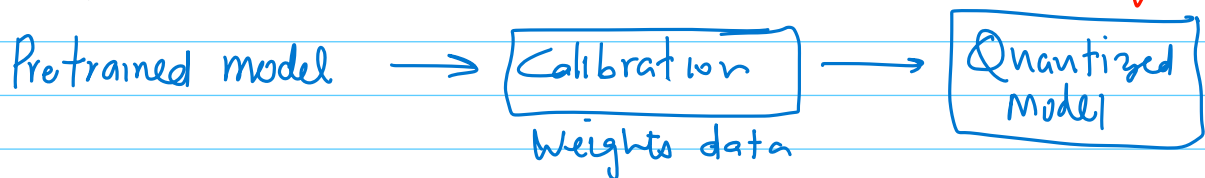
$$\text{Scale} = \frac{1000 + 20}{255} = 4.0 \text{ Scale factor}$$

Conversion :

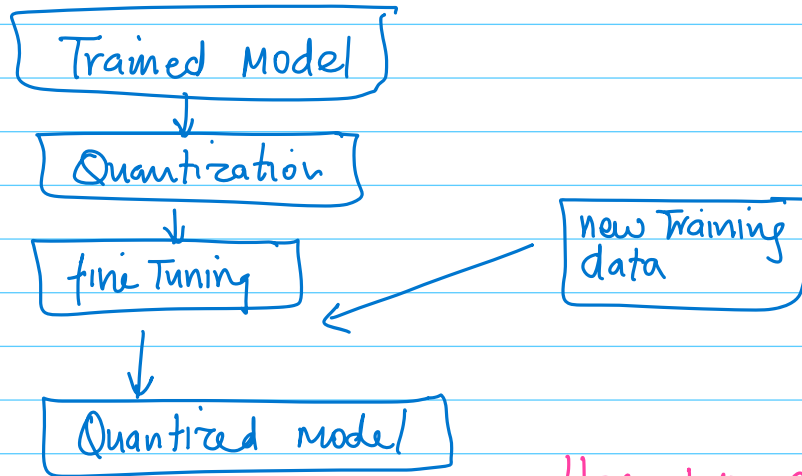
$$\frac{-20}{4} = -5.0 + \boxed{5.0}$$

add same # in
+ve sign (zero factor)

① Post Training Quantization (we don't use this technique)



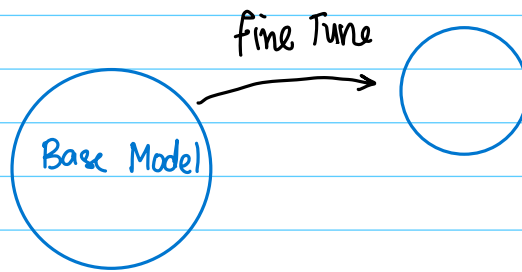
② Quantization Aware training (QAT) (This is what we use)



Here we don't lose any accuracy after this

LoRA, QLoRA intuition

Low Rank Adaptation of LLMs

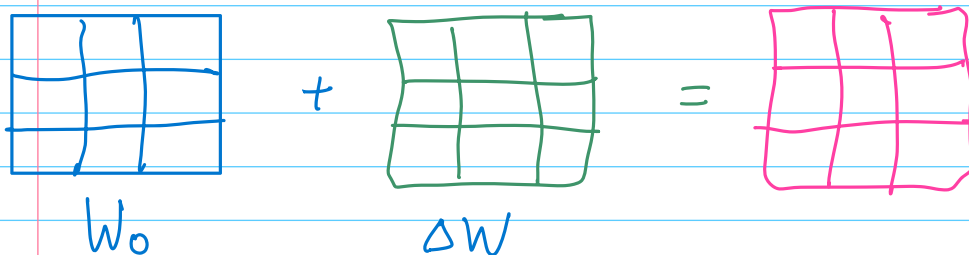


- ① full Parameter fine tuning : Training all weights
- ② Domain fine tuning
- ③ Specific task fine tuning

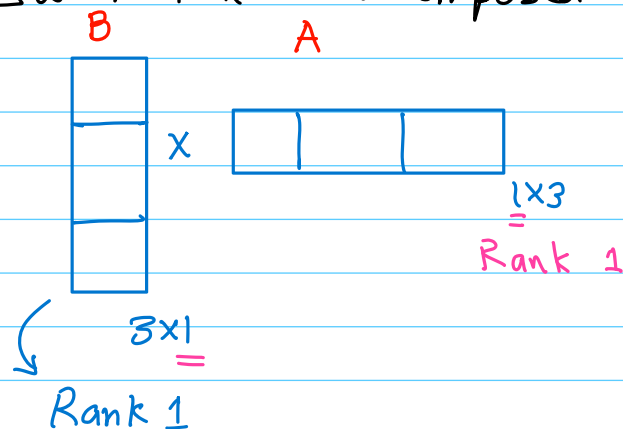
What does LoRA do?

Decomposed matrices

$$W_0 + \Delta W = W_0 + \underbrace{BA}_{\text{fine tuned weights}}$$



This ΔW matrix is decomposed like:



With Increasing $R \uparrow$ the parameters \uparrow
and we use higher ranks when model wants to
learn complex things.

QLoRA \rightarrow Quantized LoRA. \Rightarrow 16bit
↓
4bit