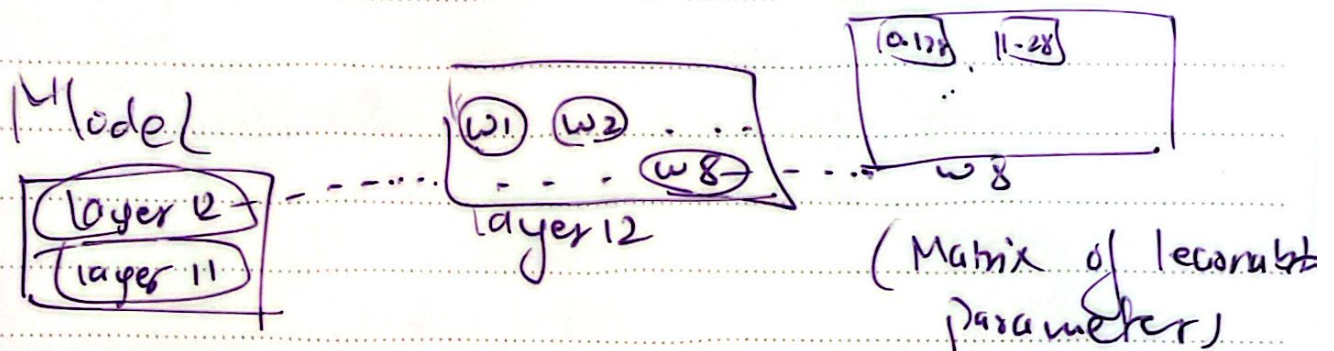


Lesson 3 Loading Models in different Datatypes



→ It is possible to inspect each parameter, datatype

→ casting using torch to load data as float16 instead of float32, saves almost half the memory.

→ look markdown for code L3

Lesson 4 Linear Quantization Theory

Linear Quantization

- ① quantize extreme values
→ highest value to int8 max
lowest value to int8 min



191.6	-13.5	720.6
92.14	295.5	-18.4
0	614.6	245.5

127
128

Rest are mapped with Linear mapping

some math formula

final small matrix \rightarrow $s = \text{scale}$, $j = \text{zero point}$

can get original with formula
some error

pip install transformers

pip install quantize
torch

look at Mark down

24

pythia is in FP32, 400 million param

$\rightarrow 400 \times 10^6$ param

32 bits = 4 bytes

$\rightarrow 8 \text{ bits} = 1 \text{ byte}$

$= 400 \times 10^6 \times 4 \text{ bytes}$

$= 1600 \times 10^6 \text{ bytes}$

$= 1600 \text{ megabytes}$

$\approx 1.6 \text{ Gbytes}$



LLM Quantization

- Many papers
- Some require calibration step
- Some ready to go
 - LL.M. INT8 (only 8 bit)
 - Linear quantization
 - QLoRA (only 4 bit)
 - HQQ (up to 2-bit)

Llama 27B

- 28 GB storage in FP32 (32-bit precision)
- ~4 GB in 4 bit precision in GGUF format

→ The Bloke

Performance Degradation

flipping face

Open LLM leaderboard

