

Quantization fundamental

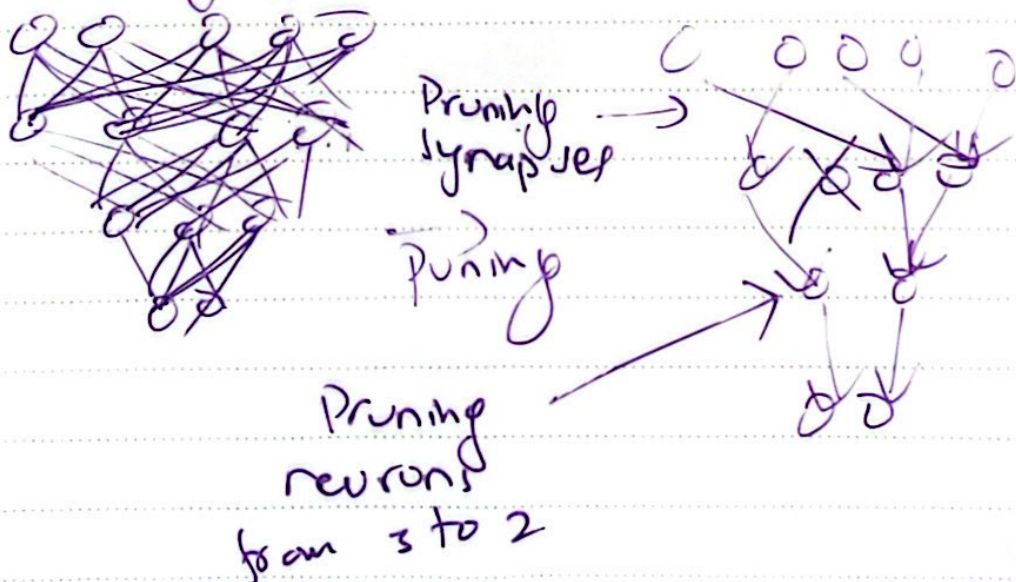
→ GPT → 80 GB (in 2022)

→ In 2023, 70 GB → 280 GB required parameters

→ Soln Model compression

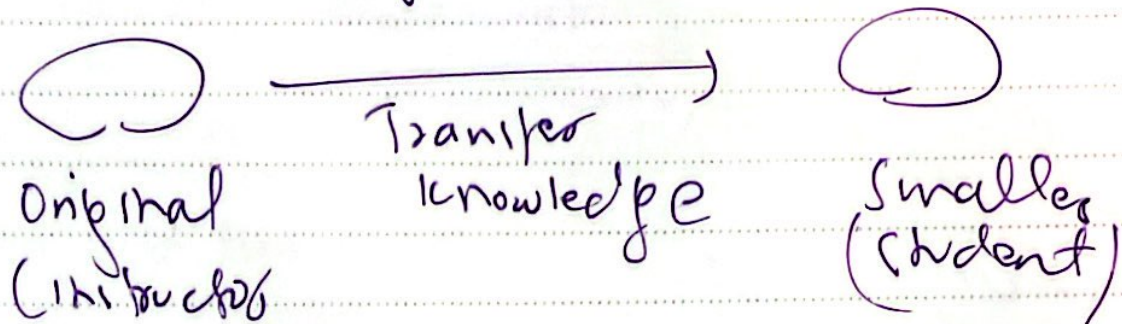
① Pruning :

→ Remove certain layers or which are not an absolute necessity



③ Knowledge Distillation

→ train smaller model (student) from original (~~the~~ instructor)



→ needs high compute power

④ Quantization

→ the parameters of model in lower precision

13.5	14.3	8.5
-4.7	-3.2	-6.4
-0.4	1.3	3.2

(FP32 to INT8)
to INT8

13	4	8
-5	-3	-6
0	1	4

FP32: 4 bytes for each

1 byte each

total 36 bytes

→ total
total: 9 bytes



→ quantization error, challenge is to lower this error

→ Quant Model

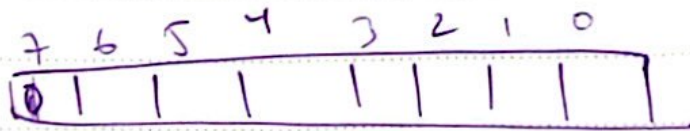
→ from transformer model quantization.

Data Types

① Unsigned int8

→ 0 to $2^n - 1$

→ 8 bit $-10 - 2^8 \div 2 = 255$



$$\begin{array}{ccccccc} \downarrow & & & & & & \downarrow \\ 2^7 & + & & + & + & \dots & 2^0 = 0 \end{array}$$

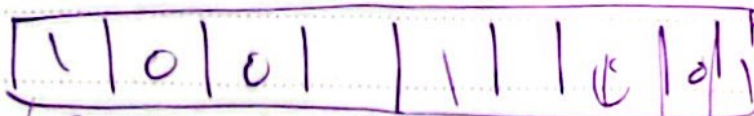
② Signed int8

→ also represents negative

→ $[-2^{n-1}, 2^{n-1} - 1]$
range



8 bit



last bit \swarrow 0 1 2 3 4 5 6 7

when set 7

it will negative

for 8 bit

$$2^{-7} = -128$$

2's complement

Then add this for rest

in our case

$$-128 + 0 + 0 + 0 + 2^3 + 0 + 0 + 2^0$$

$$= -119$$

$$0 \quad 0 \quad 1 \quad 0 = 2$$

+

$$1 \quad 1 \quad 1 \quad 0 = -2$$

=

$$0 \quad 0 \quad 0 \quad 0 = 0$$

Binary addition carry over.



In torch

torch.int8

int16

torch.uint8

int32

int64

torch.iinfo (torch.uint8)

③ Floating point:

→ 3 components

① sign → positive & negative
1 bit

② Exponent (range)
→ impact the representable
range of number

③ Fraction (precision): impact on
precision of the number,

FP32 BF:16

→ sign: 1 bit

Exponent (range) 8 bit

Fraction : 7 bit

Total: 16 bit

→ very low & high

→ different formulas for decoding

Positive & negative values



torch.float16
bfloat16
float32
float64

fp32 \rightarrow Precision better maximum $\sim 10^{+32}$

fp16 better $\sim 10^{+4}$

Bf16 good $\sim 10^{38}$

value = $1/3$ \rightarrow gonna bit 64 bit

t = torch.tensor(value, dtype=torch.float64)
(tensor node)

Downcasting

\rightarrow convert higher data type to lower
 \rightarrow usually loss of data.

d = torch.randn(1000, dtype=torch.float32)

d[:5]

\rightarrow d.to(dtype=torch.float16)



`torch.dot(d, ...)`

Matrix Multiplication

- Result will be different, logs of data

Pros: Reduced memory footprint

- > Efficient GPU memory
- large batch size
- increase in speed

Cons: less precision.

Use case:

Mixed Precision training

-> computation in lower precision
(FP16 / BF16 / FP8)

-> store & update weights in higher
precision FP32

