

NLP方向总结-显存优化



量化Quantization

量化是通过从一个数据类型到另一个数据类型的“四舍五入”来完成的, 会导致信息损失, 是一种有损的压缩, 看似微小的错误在通过模型的各层传播时往往会积累和增长, 并导致性能下降

模型量化压缩技术

BFloat16是一种浮点格式, 它用于表示16位二进制浮点数。它的名称来自于“Brain Floating Point format”, 因为它是由谷歌公司的Brain团队开发的, 与传统的浮点格式 (如IEEE 754) 相比, BFloat16使用更少的位数来表示一个浮点数, 这在机器学习和人工智能等计算密集型应用中非常有用, 因为它可以减少内存使用和加速计算。

BFloat16的格式是1位符号位、8位指数和7位尾数。它可以表示的范围和精度与IEEE 754单精度浮点数 (32位) 相似, 但是它不能表示与IEEE 754双精度浮点数 (64位) 相同的精度和范围。BFloat16广泛应用于机器学习框架和芯片设计中, 例如TensorFlow和Google的TPU芯片。

可以表示的数字范围比较小, 但是相对于FP16而言精度更高。在TensorFlow和Google的TPU芯片等机器学习框架和芯片中广泛使用。可以用于进行混合精度计算, 其中权重和梯度使用FP32或更高精度的浮点数, 而激活使用BFloat16, 从而提高计算速度和减少内存使用。

使用16位二进制格式来表示实数, 其中6位用于指数, 10位用于尾数, 1位用于符号。可以表示的数字范围和精度相对于BFloat16而言略低。在深度学习等计算密集型任务中, 使用FP16可以大大减少内存使用, 加快计算速度。可以用于进行混合精度计算, 从而提高计算速度和减少内存使用。

BFloat16相对于FP16而言精度更高, 但是支持程度较低, 而FP16则具有更广泛的支持程度, 但是精度略低。相对于FP16, 失去了3比特的精度。

FP16

FP16是16位浮点格式, FP32是32位浮点格式。使用16位二进制格式来表示实数。可以表示的数字范围比较小, 但是精度相对较低。在深度学习等计算密集型任务中, 使用FP16可以大大减少内存使用, 加快计算速度。硬件支持程度有限, 一些较旧的GPU可能不支持FP16。

FP32

使用32位二进制格式来表示实数。可以表示的数字范围更广, 精度更高。在大多数应用中被广泛使用, 因为它提供了足够的精度和数字范围。硬件支持程度广泛, 几乎所有现代CPU和GPU都支持FP32。

其它

科学计数法表示一个实数, 包含一个基数 (在计算机中通常为2) 和指数, 例如: 1.23×10^{-4} 表示为 12300, 其中基数为10, 指数为4。在计算机中, 指数表示实数的幂次, 尾数表示实数的小数部分。FP16中的5位指数可以表示-15到+16之间的指数, 其中一部分用于表示正的指数, 一部分用于表示负的指数。10位尾数可以表示小数点后的10位二进制制, 用于表示实数的小数部分。例如, 考虑一个实数1.25, 它可以写成 $1.25 = 1.01 \times 2^0$, 其中指数为0, 尾数为0.01 (二进制下为0.25)。在FP16中, 这个实数可以写成以下二进制表示:

符号位: 0 (正数)
指数位: 01110 (14, 表示0+15)
尾数位: 0100000000 (512/1024, 表示0.5+0.25=0.75)
将这些位组合起来, 得到以下16位二进制制:
0011101001000000。 指数, 尾数

如果你使用AdamW优化器, 每个参数需要8个字节, 10亿个参数需要8GB的GPU内存

GPU设备上加载一个模型, 每十亿个参数在float32精度下要花费4GB, float16要花费2GB, int8要花费1GB

如果模型是1760亿参数, 用半精度推理, 那么需要的GPU内存是1760亿 * 2byte = 352GB

FP32被称为全精度 (4字节), 而BF16和FP16被称为半精度 (2字节)。在此基础上, int8 (INT8) 数据类型由8字节表示组成, 可以存储 2^8 个不同的值 (对于有符号的整数, 在[0, 255]或[-128, 127]之间)

BLOOM-1760亿参数需要8个80GB的A100进行推理, 如果微调, 需要72个

使用LLM.int8()的BLOOM-176B比fp16版本慢了大约15%到23% Huggingface的LLM.int8()

