

Quantization Abhinav Bhatele, Daniel Nichols



Announcements

- Interim report for the project is due on April 22
- Sign up for presentation slot (link on Piazza)
- Send data requirements (if needed) to course email by end of day today





Floating Point Numbers

- Computers cannot store arbitrary precision
- Need approximations
- IEEE 754 floating point standard







Do we need all 64 bits?

- FP64 yields ~16 decimals of precision
 - <u>https://en.wikipedia.org/wiki/IEEE 754#Basic and interchange formats</u>
- Deep learning does not need all this precision
 - We often use fp32, fp16, bf16, or some combination for training





Do we need all 64 bits?

- FP64 yields ~16 decimals of precision
 - https://en.wikipedia.org/wiki/IEEE 754#Basic and interchange formats
- Deep learning does not need all this precision
 - We often use fp32, fp16, bf16, or some combination for training
- We usually need even less bits for inference!



Quantization

- Map existing weights/activations from higher to lower precision
 - i.e. fp32 -> fp16
- Why quantize?
 - Save memory
 - Potentially faster compute
 - Potentially lower energy consumption
- When not to quantize
 - Lack of native hardware support (if performance matters)
 - Operations sensitive to precision







Quantization: An Example

- fp32 can represent a much wider range of values
- How do we map them to int8?







• Suppose we want to quantize a model with N parameters from fp32 to 8 bit integers

An fp32 float can be written as an affine transformation of an 8 bit integer





Quantization: An Example

- fp32 can represent a much wider range of values
- How do we map them to int8?

Any other issues or improvements?



Abhinav Bhatele, Daniel Nichols (CMSC828G)

 $\hat{x} = \left| \begin{array}{c} x \\ - \\ \mu \end{array} + \eta \right|$

• Suppose we want to quantize a model with N parameters from fp32 to 8 bit integers

An fp32 float can be written as an affine transformation of an 8 bit integer





Types of Quantization

- Post-training quantization
 - Quantize the weights of a pre-trained model
 - Simple, less compute-intensive
- Quantization-aware Training
 - Adapt training process so that quantization is easier (i.e. ensure weights all in particular range)
 - Better quantizations at the cost of compute and training complexity







- Layer-wise quantization
- Efficient solution to optimization problem









Abhinav Bhatele, Daniel Nichols (CMSC828G)

$\operatorname{argmin}_{\widehat{\mathbf{W}}} ||\mathbf{W}\mathbf{X} - \widehat{\mathbf{W}}\mathbf{X}||_2^2.$

Weight Matrix / Block

AWQ: Activation-Aware Weight Quantization

- Quantize weights of an LLM
- Not all weights are equally important
- Based on previous paper "SmoothQuant: Accurate and Efficient Post-Training"
 - Quantization for Large Language Models"







Performance Motivation







Generation is extremely memory bound



Abhinav Bhatele, Daniel Nichols (CMSC828G)

Quantization can significantly decrease the memory footprint



Why not quantize activations and weights?

- For lower arithmetic intensities weight-only quantization has a higher ceiling
- For larger scales quantizing both may be beneficial







Naive Quantization Methods Hurt Model Performance

-1.8	2.1	3.4
-0.1	0.3	-1.2
1.7	0.5	-0.6
-1.9	-1.0	-3.9

W





Potential Weight-Aware Solution

	W	
-1.8	2.1	3.4
-0.1	0.3	-1.2
1.7	0.5	-0.6
-1.9	-1.0	-3.9

Use LI or L2 norm to find most important weights



Abhinav Bhatele, Daniel Nichols (CMSC828G)



This also doesn't help much in practice

Use Activations to Determine Important Weights

	W	
-1.8	2.1	3.4
-0.1	0.3	-1.2
1.7	0.5	-0.6
-1.9	-1.0	-3.9







Abhinav Bhatele, Daniel Nichols (CMSC828G)

This would be really difficult to implement and optimize!



Observation: Scaling Weights

• What happens when we scale particular weights before quantization?

y = wx,

Linear quantization

 $Q(\boldsymbol{w}) = \Delta \cdot \text{Round}$





Abhinav Bhatele, Daniel Nichols (CMSC828G)

$$y = Q(\boldsymbol{w})\boldsymbol{x}$$

Forward pass and forward pass with quantization

$$\left(\frac{\boldsymbol{w}}{\Delta}\right), \quad \Delta = \frac{\max\left(|\boldsymbol{w}|\right)}{2^{N-1}}$$

$$\cdot s) \cdot \left(\frac{x}{s}\right)$$

What happens when we scale by s>l?

Observation: Scaling Weights

• What happens when we scale particular weights before quantization?

$$y = \boldsymbol{w} \boldsymbol{x}, \quad y = Q(\boldsymbol{w}) \boldsymbol{x}$$

 $Q(\boldsymbol{w}) = \Delta \cdot \text{Round} \left(\frac{\boldsymbol{w}}{\Delta}\right)$
 $\Delta = \frac{\max\left(|\boldsymbol{w}|\right)}{2^{N-1}}$
 $Q(\boldsymbol{w} \cdot \boldsymbol{s}) \cdot \left(\frac{\boldsymbol{x}}{s}\right)$

Observation 2: Scaling a single or small subset of w values has little impact on Δ . $\Delta' \approx \Delta$



Abhinav Bhatele, Daniel Nichols (CMSC828G)

tion I: Expected error rounding is ~0.25; condent of w and s

$$= \Delta' \cdot \operatorname{Round} \left(\frac{ws}{\Delta'}\right) \cdot \frac{x}{s}$$

Observation 3: x and Δ are fp16 and have no quantization error

Observation: Scaling Weights

• What happens when we scale particular weights before quantization?

Err

$$y = wx, \quad y = Q(w)x$$

$$Q(w) = \Delta \cdot \text{Round}\left(\frac{w}{\Delta}\right)$$

$$\Delta = \frac{\max\left(|w|\right)}{2^{N-1}}$$
Based on observation 3
we can represent the
quantization errors
$$\text{Err}\left(Q(w \cdot s)\left(\frac{x}{s}\right)\right) = \Delta' \cdot \text{RoundErr}\left(\frac{ws}{\Delta'}\right)$$

Scaled term has lower relative error for s > /!



$$\frac{\left(Q(w \cdot s)\left(\frac{x}{s}\right)\right)}{\operatorname{rr}\left(Q(w)x\right)} = \frac{\Delta'}{\Delta \cdot s} \approx$$



But Does It Work?

• Try different s in practice





Abhinav Bhatele, Daniel Nichols (CMSC828G)

///////////////////////////////////////		11 1		
As s gets lar assumption	rger the Δ' on breaks d	$\sim \Delta$ own		
s = 1.25	s = 1.5	s = 2	s = 4	
2.8%	4.4%	8.2%	21.2%	
1.005	1.013	1.038	1.213	The relative e
0.804	0.676	0.519	0.303	does keep ge smaller
12.87	12.48	11.92	12.36	

In practice, some *s*^{*} gives the best modeling performance



Use Activations to Determine Important Weights

_	W	
-1.8	2.1	3.4
-0.1	0.3	-1.2
1.7	0.5	-0.6
-1.9	-1.0	-3.9







Abhinav Bhatele, Daniel Nichols (CMSC828G)

Use search/sweep over calibration dataset to find optimal scaling parameters

AWQ Results

AWQ is performs better or near 1%	ĺ
fp16, but is much faster and more	ŀ
memory efficient	ļ

OPT	$(\mathbf{PPL}\downarrow)$	1.3B

FP16	14.62
RTN	119.47
1% FP16	16.91
s = 2	18.63
AWO	16.32



Abhinav Bhatele, Daniel Nichols (CMSC828G)

6.7B	13B	30B
10.86	10.13	9.56
23.54	46.04	18.80
11.39	10.43	9.85
11.92	10.80	10.32
11.39	10.56	9.77
	6.7B 10.86 23.54 11.39 11.92 11.39	6.7B13B10.8610.1323.5446.0411.3910.4311.9210.8011.3910.56

A fixed value of *s* gives decent results, but it's better to select dynamically

AWQ Comparison with Other Methods

AWQ yields better perplexity than other SotA approaches (GPTQ*)

PPL↓			Llama-2		LLaMA			
		7B	7B 13B 70		7B	13B 30B		65B
FP16	_	5.47	4.88	3.32	5.68	5.09	4.10	3.53
	RTN	6.66	5.52	3.98	7.01	5.88	4.88	4.24
INT3	GPTQ	6.43	5.48	3.88	8.81	5.66	4.88	4.17
g128	GPTQ-R	6.42	5.41	3.86	6.53	5.64	4.74	4.21
	AWQ	6.24	5.32	3.74	6.35	5.52	4.61	3.95
	RTN	5.73	4.98	3.46	5.96	5.25	4.23	3.67
INT4	GPTQ	5.69	4.98	3.42	6.22	5.23	4.24	3.66
g128	GPTQ-R	5.63	4.99	3.43	5.83	5.20	4.22	3.66
	AWQ	5.60	4.97	3.41	5.78	5.19	4.21	3.62



AWQ Downstream Tasks

Math and coding downstream tasks

MBPP (7B)	pass@1	pass@10	GSM8K	7B	13B	70B
FP16	38.53	49.77	FP16	13.87	26.16	56.41
RTN	37.51	48.49	RTN	11.07	21.23	53.98
GPTQ	31.97	44.75	GPTQ	12.13	24.26	56.03
AWQ	40.64	49.25	AWQ	13.57	25.25	56.40



AWQ performance is close to fp16 performance

AWQ Runtime Results

TinyChat benchmark for edge devices







		AWÇ	is faste mei	er and sa mory	ves on			
			11/	/				
39	Durs (FP16)	38	ours (AWQ	, W4A16)	$\begin{array}{c}60\\45\\20\end{array}$	MO	60	52
	21 FP16 OOM	11 12	FP16 OOM ⁹	7 9 22	30 00 15 HI 0	FP16 OC FP16 OC	FP16 OC	
na-2 3)	Llama-2 (13B)	MPT (7B)	MPT (30B)	Falcon (7B)	0 Llama-2 (7B)	Llama-2 (13B)	MPT (7B)	Falcon (7B)
	(b) Jetson C	orin mobil	e GPU		(c) RT	X 4070 laj	ptop G	PU

Binarization

- Special form of quantization
- Map weights to {-1,1} or {-1,0,1}
- Train using int8 activations and BitLinear layers

 $\widetilde{W} = \operatorname{Roup}$



Abhinav Bhatele, Daniel Nichols (CMSC828G)

• 1.58-bit models, "The Era of 1-bit LLMs: All Large Language Models are in 1.58 Bits"

What is large potential optimization we can make assuming we have ternary weights?

$$\operatorname{ndClip}(\frac{W}{\gamma+\epsilon}, -1, 1),$$

RoundClip(x, a, b) = max(a, min(b, round(x))),

$$\frac{1}{nm}\sum_{ij}|W_{ij}|.$$



Potential Savings with Binarized Models

Binarized models are faster and take up less memory





Potential Savings with Binarized Models

Using exclusively int8 additions is extremely energy efficient









UNIVERSITY OF MARYLAND