

Kronecker Products and GPT2

Ayoub @ Uni Passau

March 6, 2024

Abstract

GPT2 checkpoints from HuggingFace reach a approximate loss of 3.11 on OpenWebText. In this work, we factorize the main parameters of GPT2, as Kronecker Products leading to a compression factor of $x\%$, and ask the following question: What is the fastest way to reach back to the 3.11 loss using the Kronecker factors. And how to optimally manipulate training and distillation of the small compressed network. We particularly investigate freezing (for both training and distillation), and smooth transitioning (explained later on).

- **freezing** the pre-trained weights, helps boost accuracy.
- One by One distillation: drop the weights slowly and distill through the network.

Methodology:

1. Train GPT2 to a certain level. – 3.11.
2. Decompose the MLP weights.
3. Try a bunch of methods and see which one guarantees getting back to the original loss as fast as possible.

Contents

1	Resutls	2
2	Introduction and Setup	2
3	Training setups	2
4	Kronecker Decomposition	2
4.1	The 95M model	2
4.2	Different decomposition schemes:	3
5	Impact of down-scaling: How small can we go	3
5.1	Initialization	4
5.2	Questions	4

1 Results

Results so far on wikitext103, wikitext2 and lambada.

# Params	Model	Datasets		
		wikitext-103	wikitext-2	Lambada
124M	GPT2	29.16	24.67	45.28
82M	DistilGPT2	44.53	36.48	76.00
81M	KronyPT-81M-1350	41.98	34.99	-
81M	KronyPT-81M-3950	-	-	64.92
81M	TQCompressedGPT2	40.28	32.25	64.72
81M	KnGPT-2 (Huawei)	40.97	32.81	67.62

Table 1: Comparison of different models on various datasets.

2 Introduction and Setup

A lot of work has been done on the compression of LLMs using factorization methods, including Kronecker Products. Most notably for encoder based based, and rarely applied to decoder based methods ([1], [2]).

We study how far we can push the compression ratio, and then impact of down-scaling on the quality of the model. We also investigate adding multiple Kronecker factors.

3 Training setups

Here we discuss the main. After many I am going for the main, mostly

4 Kronecker Decomposition

In this section, we study different Kronecker decomposition setups, and the percentage of compression it would lead to. So far we only decompose the weights `c_fc.weight` and `c_proj.weight` (each has $2.3M$ in the original GPT2-small architecture.). Each transformer layer (there are 12 in total) has two of these weights. They count to $56.6M$ in total (45% of GPT2 $124M$). Hence any significant reduction to these matrices would lead to a remarkable compression ratio of the whole model. We choose not to compress the other weights, namely, attention weights and embedding matrix for various reasons that we will expose later on.

4.1 The 95M model

The most basic strategy is to divide one of the dimensions of each W by 2, this would lead to a 95M model. The parameter $p_1 = \text{c_fc.weight}$ (resp.

$p_2 = \text{c_proj.weight}$) has a shape of (3072, 768) (resp. (768, 3072)). We first try the following decomposition: $p_1 = \underbrace{W_{11}}_{(3072, 384)} \otimes \underbrace{W_{12}}_{(1, 2)}$ and $p_2 = \underbrace{W_{21}}_{(384, 3072)} \otimes \underbrace{W_{22}}_{(2, 1)}$

This decomposition would lead to to reduction of $28M$ (50%). The new network would have approx $95M$. Our goal is to eventually reach the $82M$ mark, similar to DistilGPT2, and other Factorized models (inserts other refs here).

4.2 Different decomposition schemes:

It is reasonable to aim for a decomposition that guarantees the maximum rank we can get. Since the Rank of Kronecker products is multiplicative, i.e., $\text{rank}(A \otimes B) = \text{rank}(A) \cdot \text{rank}(B)$ ¹, we can easily compute the rank of each possible decomposition. In our case, we have $W \in R^{(m, n)}$ where $m = 3072$ and $n = 768$. Hence, for each layer of GPT2, we aim to find the “optimal” $A \in R^{(m_1, n_1)}$, and $B \in R^{(m_2, n_2)}$, i.e.,:

$$W \approx A \otimes B, \quad m = m_1 m_2, n = n_1 n_2$$

W.l.o.g, for each decomposition (A, B) the maximum rank we can reach is $\min(m_1, n_1) \cdot \min(m_2, n_2)$. And each of the reduced decompositions would have exactly $m_1 n_1 + m_2 n_2$ parameters. Hence, theoretically, the maximum rank we can get is 384 of a (3072, 768) matrix. The following table summarizes all possible combinations, alongside the reduction it would lead to per layer, and the total number of parameters in GPT2, for only those decompositions of rank 384.

5 Impact of down-scaling: How small can we go

In this section, we study the impact of down-scaling on the compressed models, we train 4 different architectures, with variant dimensions:

- 67M: diverges // try with higher batch size (the least we can get)
- 81M: Super.
- 81M: with a Variant scheme
- 95M: duh duh duh. (the highest we can get)

Keep in mind:

- You can probably stabilize training using various tricks. It’s just an endless loop. We are not going to play the What-IF game. One single config, that’s it.

¹Link to proof: <https://math.stackexchange.com/questions/707091/elementary-proof-for-textrank-left-a-otimes-b-right-textrank-a-cdot>

- One could try to have a mixed strategy:
 - higher levels of compression in the early layers.
 - higher levels of compression in the late layers.
 - Every odd layer
 - In the middle
-
-

5.1 Initialization

Here we test a new initialization trick based on pruning for the 95M models. And compare it to the VL method. Pruning shows a much better

5.2 The fuck is Van Loan

XX

5.3 Pruning based Initialization

XX

5.4 Questions

- When we plug new KP matrices to the pre-trained model, is it useful to freeze the other weights that are already pre-trained? Or using different learning rates schemes maybe?
- When don't don't freeze the weights. Does the network rely on other pre-trained methods?
 - ****The idea that I'm challenging here****: A lot of work on NNs compression using factorized methods, claim that the network only need to be (post-)trained on small
 - But, what the fail to mention or elaborate on at least, is that not the weights matrices of the network are decomposed.
 - And with all the residual connection that are present in most attention networks, one could suspect, that maybe the weights that were not decomposed are taking over...
 - A good remedy for this is to freeze the original weights during the post-training, and only allow the new dropped in (factorized matrices) to be updated with backprop.

- This strategy could also be implemented in distillation. We refer to this method as ****forced distillation.****
- When we freeze, we investigate how useful is it to distill matrices one by one, rather than drop in the matrices all at once. And also investigate if the order of dropping has any significance, bottom up or top bottom...

References

- [1] M. Tahaei, E. Charlaix, V. Nia, A. Ghodsi, and M. Rezagholizadeh, “Kroneckerbert: Significant compression of pre-trained language models through kronecker decomposition and knowledge distillation,” in *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2116–2127, 2022.
- [2] A. Edalati, M. Tahaei, A. Rashid, V. P. Nia, J. J. Clark, and M. Rezagholizadeh, “Kronecker decomposition for gpt compression,” *arXiv preprint arXiv:2110.08152*, 2021.