

Geometry-aware training of factorized layers in tensor Tucker format

Emanuele Zangrando¹, Steffen Schotthöfer², Gianluca Ceruti³, Jonas Kusch⁴, Francesco Tudisco^{1,5}

¹ Gran Sasso Science Institute

² Oak Ridge National Laboratory

³ University of Innsbruck

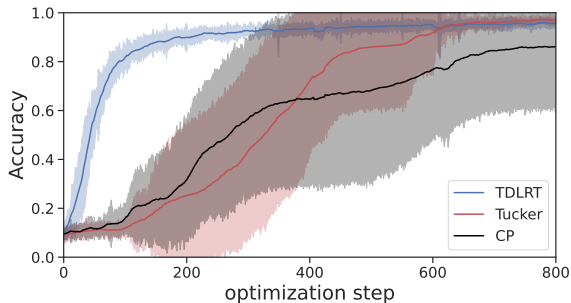
⁴ Norwegian University of Life Sciences

⁵ University of Edinburgh



Problem

- Train or fine-tune neural network with tensor layers (e.g. convolutions) in a low-rank memory efficient format
- train networks in a rank-adaptive fashion
- Avoid instabilities of factorized gradient descent



Original dynamics

$W \in \mathbb{R}^{n_1 \times \dots \times n_d}$, full-training consists in discretizing the gradient flow

$$\dot{W} = -\nabla L(W(t))$$

$$\text{space: } O(\prod_{i=1}^d n_i)$$

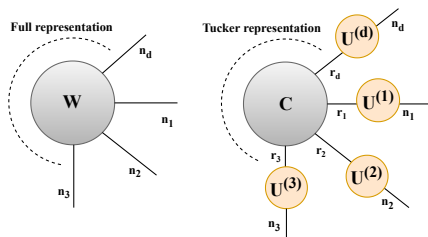
Tucker decomposed layer:

$$W_{i_1, \dots, i_d} = C^{j_1, \dots, j_d} U_{j_1, i_1}^{(1)} \dots U_{j_d, i_d}^{(d)}$$

Block-dynamics:

$$\dot{U}^{(i)} = -\nabla_{U^{(i)}} L(C \times_j U^{(j)})$$

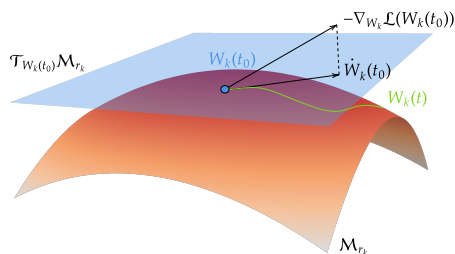
$$\dot{C} = -\nabla_C L(C \times_j U^{(j)})$$



Low-rank dynamics approximation

$$\mathcal{M}_\rho := \{W \in \mathbb{R}^{n_1 \times \dots \times n_d} \mid \text{rank}(\text{Mat}_i(W)) = r_i\}, \quad \rho = (r_1, \dots, r_d)$$

$$\dot{W} = -P_{T_W \mathcal{M}_\rho} \nabla L(W)$$



Best 1st order local approximation of the original dynamics

$$\dot{U}^{(i)} = -(I - U^{(i)} U^{(i)\top}) \text{Mat}_i(\nabla L(W) \times_{j \neq i} U^{(j)\top}) \text{Mat}_i(C)^\dagger$$

$$\dot{C} = -\nabla L(W) \times_{j=1}^d U^{(j)\top}$$

$$\text{space: } O(\prod_{i=1}^d n_i r_i + \prod_{i=1}^d r_i)$$

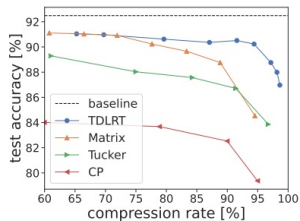
Results

- Theoretical guarantees of approximation of the original problem and descent of the loss, together with convergence to critical points in the stochastic setting

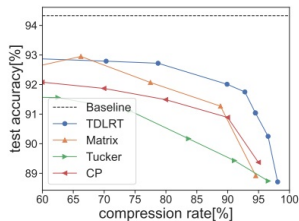
$$\begin{aligned}\|W(t) - C(t) \times_j U_j(t)\| &\leq c_1 \varepsilon + c_2 \lambda + c_3 \tau / \lambda \\ \liminf_{t \rightarrow \infty} \mathbb{E} \|\nabla L(W(t+1)) \times_j U_j(t) U_j(t)^\top\|^2 &= 0\end{aligned}$$

- Efficient and robust integration of the projected system to optimize the neural net (constants c_i do not depend on the singular values!)
- Variety of different experiments ranging from low-rank training from scratch to fine-tuning pretrained models

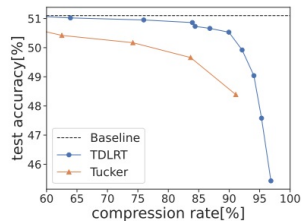
Experiments



(b) VGG16 Cifar10



(c) ResNet18 Cifar10



(d) ResNet18 Tiny-Imagenet

Figure: Compressed training from scratch

Experiments

Table 2: Fine-tuning performance metrics on Deberta V3 Glue benchmark (left) and on Stable diffusion Dreambooth (right).

GLUE	LoRA	TDRIT(Ours)	method	loss	# params
# params	1.33M (rank 8)	0.9M ($\tau = 0.15$)	LoRA ($r = 8$)	0.260	5 M
CoLa (Corr.)	0.6759	0.7065	LoRA ($r = 5$)	0.269	3 M
MRPC (Acc.)	0.8971	0.9052	LoRA ($r = 3$)	0.274	1.8 M
QQP (Acc.)	0.9131	0.9215	Ours ($\tau = 0.02$)	0.2635	1.8 M
RTE (Acc.)	0.8535	0.8713	Ours ($\tau = 0.1$)	0.272	1.5 M
SST2 (Acc.)	0.9484	0.9594			

Figure: Fine-tuning with low-rank tensor and matrix adapters

Poster

We are happy to welcome you at our poster session for further discussion

Friday 13 December, 4:30-7:30 pm local time

