

Tensorising Neural Networks

Alexander Novikov, Dmitry Podoprikin, Anton Osokin, Dmitry Vetrov

UBC MLRG 2020 Winter Term 1

Presented by Gursimran Singh (simar)

Why tensors?

Many objects in machine learning can be treated as tensors:

- ▶ Data cubes (RGB images, videos, different shapes/orientations)
- ▶ Weight matrices can be treated as tensors, both in Conv-layers and fully-connected layers

Using tensor decompositions we can compress data!

Motivation

Problem: Neural network is too large to fit into memory.

Approaches:

- Distributed neural network
 - distribute parameters
 - challenge: training [Elastic Averaging SGD (NIPS'15)]
- Model compression
 - reduce required space

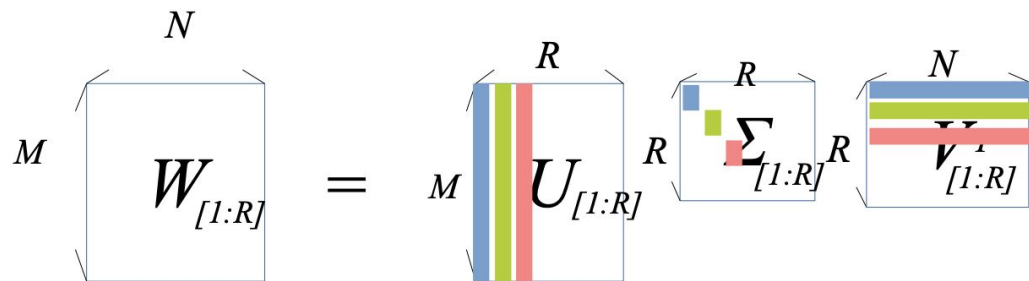
Problem Formulation

- Given $M \times N$ weight matrix W of a fully-connected layer

$$o(x; \theta) = f(W^T x + b)$$

- Goal
 - reduce space complexity
- Requirement
 - compact with back-propagation

Naive Method: Low-rank SVD



SVD Decomposition

$$W(i, j) = \sum_{r=1}^R U(i, r) \Sigma(r, r) V(j, r)^T$$

Naive Method: Low-rank SVD

The diagram illustrates the decomposition of matrix W into three matrices: U , Σ , and V^T . Matrix W is $M \times N$. Matrix U is $M \times R$. Matrix Σ is $R \times R$. Matrix V^T is $R \times N$. The decomposition is shown as $W = (U \Sigma V^T)$. Below this, matrix A is $M \times R$ and matrix B is $R \times N$.

$$W = \sum_{r=1}^R A_r B_r^T$$

space complexity

$$O(M \times N) \rightarrow O(R(M + N))$$

Naive Method: Low-rank SVD

By low-rank SVD $W = \sum_{r=1}^R A_r B_r^T$

Instead of updating W , $\frac{\partial E}{\partial W}$

update components

$$\frac{\partial E}{\partial A_r} \quad \frac{\partial E}{\partial B_r}$$

Naive Method: Low-rank SVD

- To integrated with back-propagation

$$o = f(W^T x + b)$$

– have calculate 3 gradients:

- objective wrt parameter $\frac{\partial E}{\partial o}$

- output wrt input $\frac{\partial o}{\partial x}$

$$W = \sum_{r=1}^R A_r B_r^T$$

No change in these

- output wrt parameter $\frac{\partial o}{\partial A_r} = f'(W^T x + b) B_r \mathbf{1}^T x$

$$\frac{\partial o}{\partial B_r} = f'(W^T x + b) \mathbf{1} A_r^T x$$

Simple equations to compute gradients

TT-Decomposition: Two ideas to do better

Low-rank SVD works but can we do better than that?

Two key Ideas:

recursively applying low-rank SVD

$$1) \quad r(M + N) \leq MN$$

Two ideas to do better

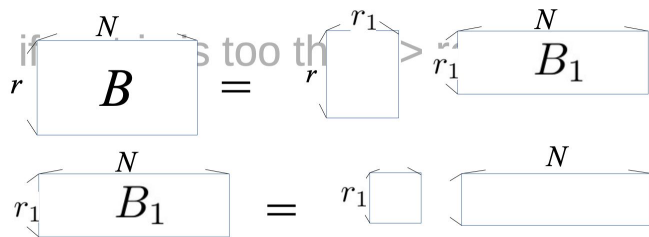
Low-rank SVD works but can we do better than that?

Two key Ideas:

recursively applying low-rank SVD

1)

$$r(M + N) \leq MN$$



Two ideas to do better

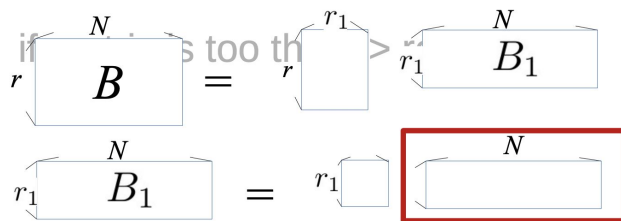
Low-rank SVD works but can we do better than that?

Two key Ideas:

recursively applying low-rank SVD

1)

$$r(M + N) \leq MN$$



gets thinner and thinner

Two ideas to do better

Low-rank SVD works but can we do better than that?

Two key Ideas:

recursively applying low-rank SVD

1)

$$r(M + N) \leq MN$$

2) if matrix is too thin => reshape

$$\frac{M}{m} + mN \leq M + N$$

Two ideas to do better

Low-rank SVD works but can we do better than that?

Two key Ideas:

recursively applying low-rank SVD

1)

$$r(M + N) \leq MN$$

2) if matrix is too thin => reshape

$$\frac{M}{m} + mN \leq M + N$$

Given two matrices X,Y with fixed total elements $XY = C$

Min when $X \sim Y$; func $\rightarrow C/X+X$

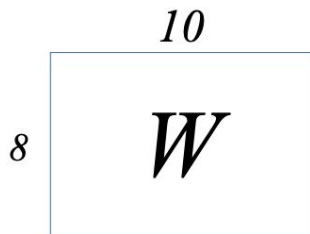
Eg - $50+2 = 52$

$25+4 = 29$ (less than 52)

Tensor-Train Decomposition

Combination of the two ideas we discussed

If we want to TT-decompose a matrix



First, need to reshape it into a

tensor $\mathcal{W} : 2 \times 2 \times 2 \times 2 \times 5$

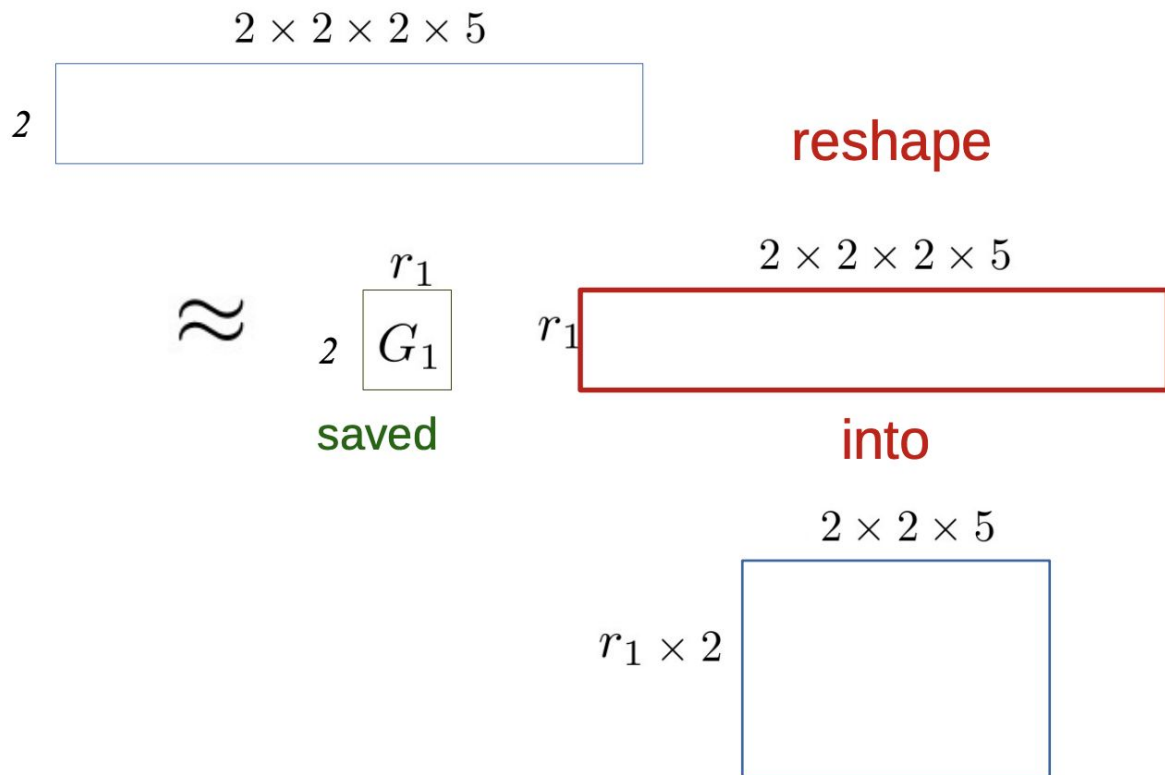
Tensor-Train Decomposition

unfold \mathcal{W} by 1st dimension

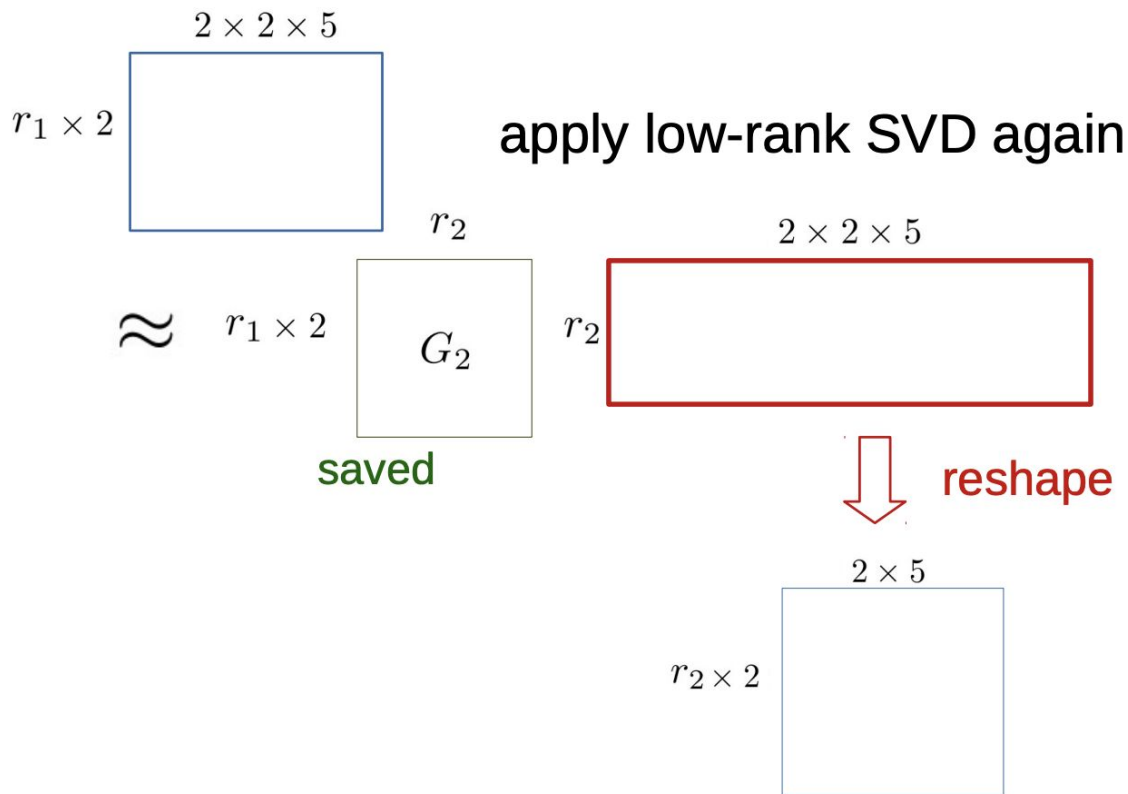
$$\begin{array}{c} 2 \times 2 \times 2 \times 5 \\ \boxed{\phantom{\text{matrix}}} \end{array} \quad \text{and apply low-rank SVD}$$
$$\approx \begin{array}{c} r_1 \\ \boxed{G_1} \end{array} \begin{array}{c} 2 \times 2 \times 2 \times 5 \\ \boxed{\phantom{\text{matrix}}} \end{array}$$

The diagram illustrates the unfolding of a 4-mode tensor \mathcal{W} (with dimensions $2 \times 2 \times 2 \times 5$) into a matrix of size $2 \times (2 \times 2 \times 2 \times 5)$. This matrix is then decomposed via low-rank SVD into a core tensor G_1 of size $r_1 \times 2$ and a residual tensor of size $r_1 \times (2 \times 2 \times 2 \times 5)$.

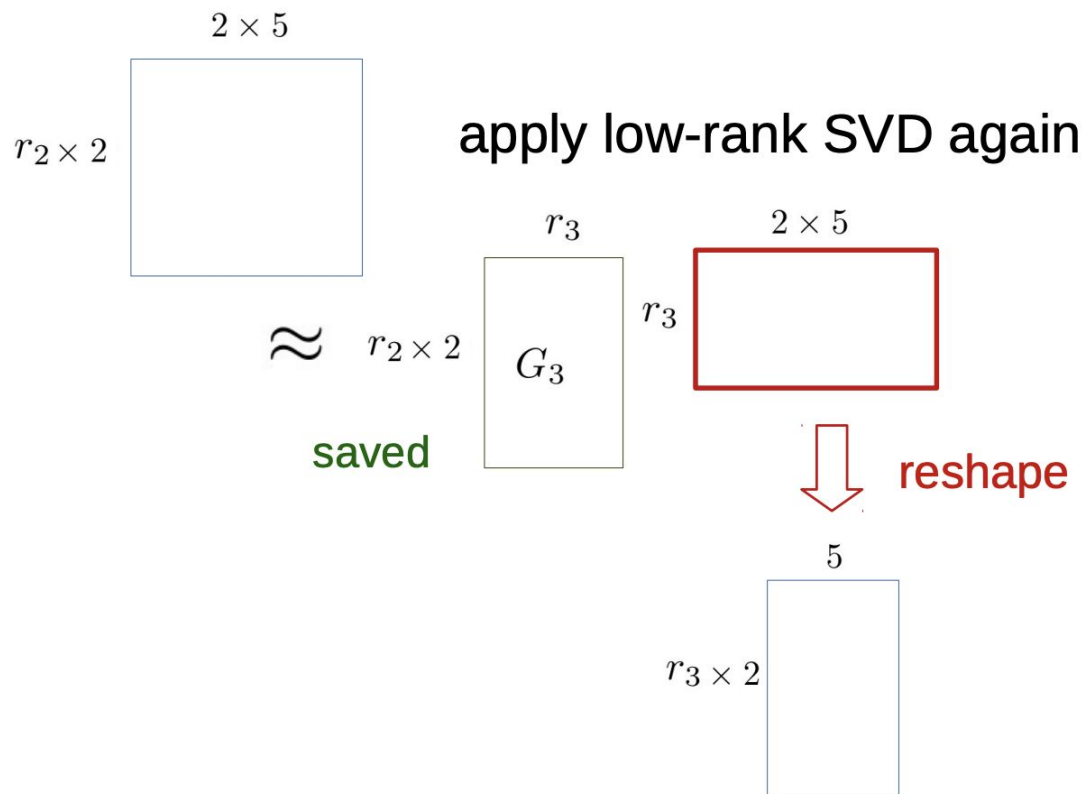
Tensor-Train Decomposition



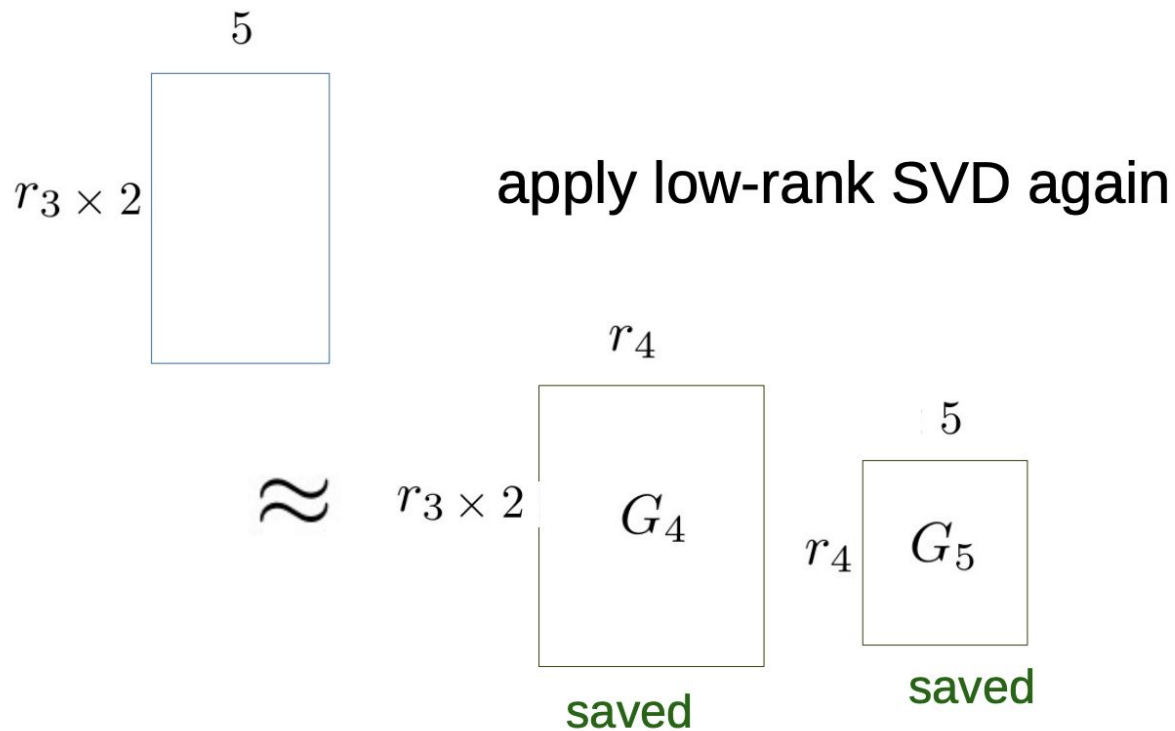
Tensor-Train Decomposition



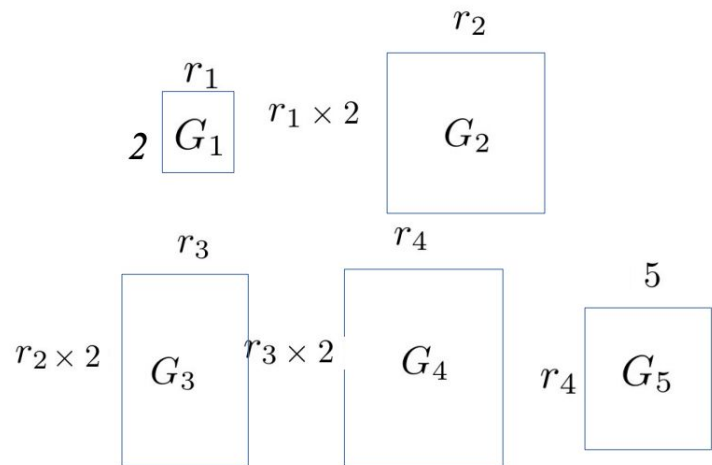
Tensor-Train Decomposition



Tensor-Train Decomposition



Tensor-Train Decomposition



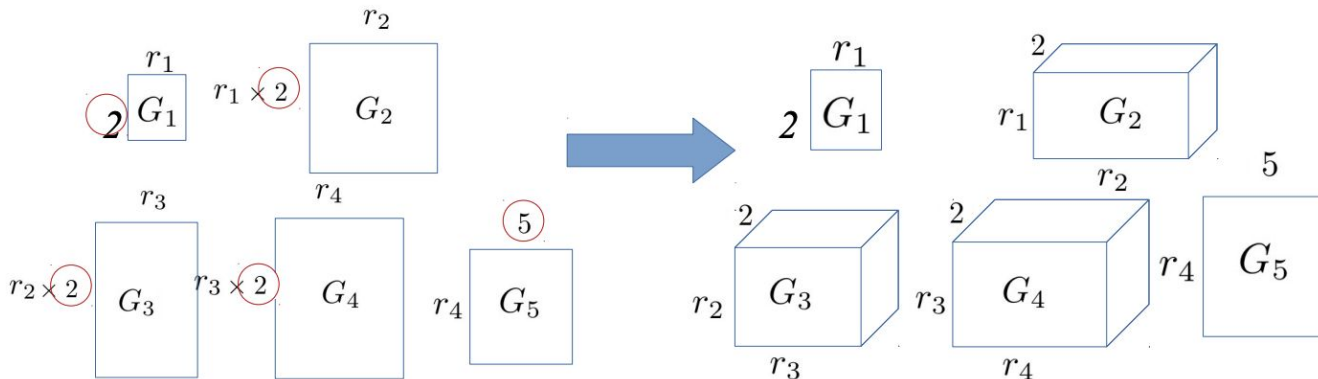
can approximate



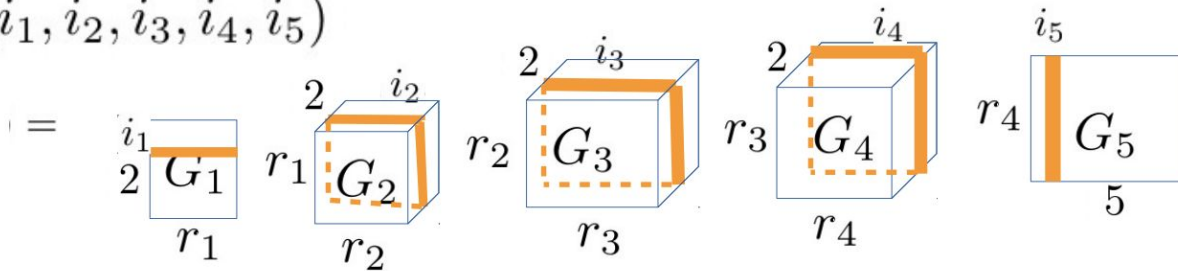
$$\mathcal{W} : 2 \times 2 \times 2 \times 2 \times 5$$

Tensor-Train Decomposition

fold into core tensors



$$\mathcal{W}(i_1, i_2, i_3, i_4, i_5)$$



Tensor-Train Decomposition

$$\mathcal{W}(i_1, i_2, i_3, i_4, i_5)$$

The diagram illustrates the decomposition of a 5D tensor \mathcal{W} into five cores G_1, G_2, G_3, G_4, G_5 . The tensor is represented as a 5D box with indices i_1, i_2, i_3, i_4, i_5 . The decomposition is shown as a sequence of cores connected by ranks r_1, r_2, r_3, r_4 . G_1 is a 2x2 matrix with index i_1 and rank r_1 . G_2 is a 2x2x2 cube with indices i_2 and i_3 , and ranks r_1 and r_2 . G_3 is a 2x2x2 cube with index i_3 and rank r_2 . G_4 is a 2x2x2 cube with index i_4 and rank r_3 . G_5 is a 2x5 matrix with index i_5 and rank r_4 .

$$\mathcal{A}(\mathbf{i}) = G_1[i_1]G_2[i_2] \cdots G_d[i_d]$$

Tensor-Train format

TT-rank = SVD decomposition rank

$$(r_1, r_2, r_3, r_4, r_5)$$

TT-SVD algorithm for TT-Decomposition

Suppose, we want to approximate:

$$A(i_1, \dots, i_d) \approx G_1(i_1)G_2(i_2)G_3(i_3)G_4(i_4)$$

1. A_1 is an $n_1 \times (n_2 n_3 n_4)$ reshape of A .
2. $U_1, S_1, V_1 = \text{SVD}(A_1)$, U_1 is $n_1 \times r_1$ — first core
3. $A_2 = S_1 V_1^*$, A_2 is $r_1 \times (n_2 n_3 n_4)$.
Reshape it into a $(r_1 n_2) \times (n_3 n_4)$ matrix
4. Compute its SVD:
 $U_2, S_2, V_2 = \text{SVD}(A_2)$,
 U_2 is $(r_1 n_2) \times r_2$ — second core, V_2 is $r_2 \times (n_3 n_4)$
5. $A_3 = S_2 V_2^*$,
6. Compute its SVD:
 $U_3 S_3 V_3 = \text{SVD}(A_3)$, U_3 is $(r_2 n_3) \times r_3$, V_3 is
 $r_3 \times n_4$

Properties of TT-decomposition

- Has been shown that for an arbitrary tensor A a TT-representation exists but is not unique.
 - It's natural to seek a representation with the lowest ranks
- TT-representation is very efficient in terms of memory if ranks are small

$$\sum_{k=1}^d n_k r_{k-1} r_k \quad \text{vs} \quad \prod_{k=1}^d n_k$$


- Efficient rounding to prevent explosion of TT-Ranks
- Efficiently perform several types of operations on tensors
 - Addition/ multiplication of constant
 - Summation and the entrywise product of tensors (results TT-tensors with more rank)
 - Global characteristics - sum of all elements and the Frobenius norm
 - Sum two TT-matrices
 - Matrix-by-vector (matrix-by-matrix) product

Properties of TT-decomposition

Operation	Output rank	Complexity
$\mathbf{A} \cdot \text{const}$	r_A	$O(dr_A)$
$\mathbf{A} + \text{const}$	$r_A + 1$	$O(dnr_A^2)$
$\mathbf{A} + \mathbf{B}$	$r_A + r_B$	$O(dn(r_A + r_B)^2)$
$\mathbf{A} \odot \mathbf{B}$	$r_A r_B$	$O(dnr_A^2 r_B^2)$
$\text{sum}(\mathbf{A})$	—	$O(dnr_A^2)$
...		

Tensor network diagrams

N-index tensor = shape with N lines

$$T^{s_1 s_2 s_3 \cdots s_N} = \text{Diagram of a tensor with } N \text{ lines labeled } s_1, s_2, s_3, s_4, \dots, s_N$$


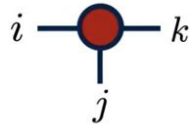
Low-order tensor examples



v_j



M_{ij}

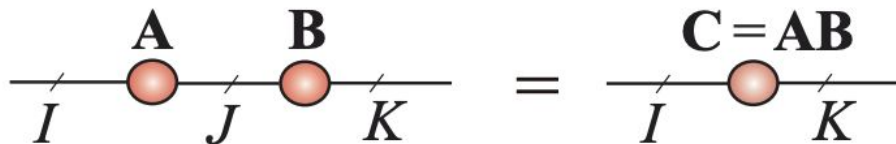


T_{ijk}

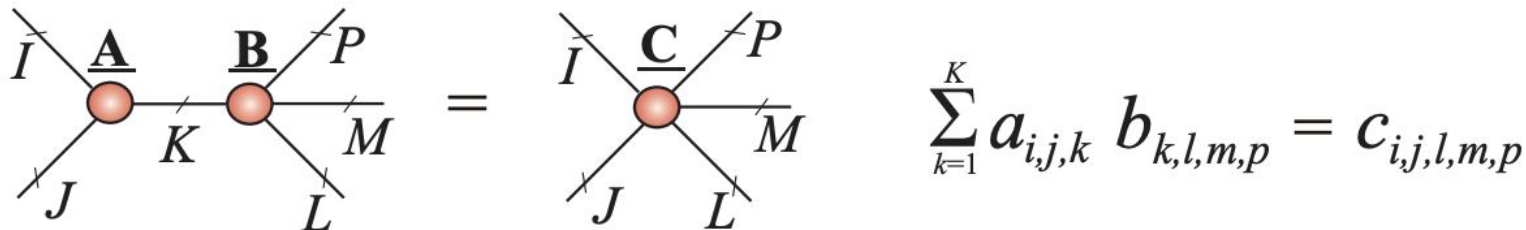
- Matrix-vector multiplication



- Matrix-matrix multiplication



- Tensor contraction



Formal definition: TT Decomposition of Tensor

Tensor \mathbf{A} can be decomposed to TT format as:

$$\mathbf{A}(i_1, i_2, \dots, i_d) = \mathbf{G}_1[i_1] \mathbf{G}_2[i_2] \dots \mathbf{G}_d[i_d]$$

Where:

$$\mathbf{G}_k[i_k] \in \mathbb{R}^{r_{k-1} \times r_k} \quad , \quad r_0 = r_d = 1$$

- TT-cores: \mathbf{G}_k
- TT-ranks: r_k
- TT max rank $r = \max r_k \quad , \quad k = 0, \dots, d$

Compression:

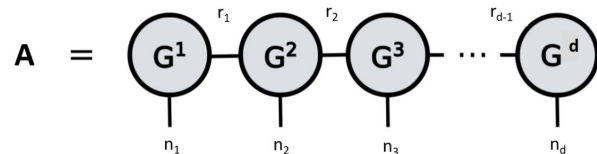
$$O(n^d) \rightarrow O(ndr^2)$$

$$\sum_{k=1}^d n_k r_{k-1} r_k = O(ndr^2)$$

Formal definition: TT Decomposition of Tensor

Tensor \mathbf{A} can be decomposed to TT format as:

$$\mathbf{A}(i_1, i_2, \dots, i_d) = \mathbf{G}_1[i_1] \mathbf{G}_2[i_2] \dots \mathbf{G}_d[i_d]$$



Where:

$$\mathbf{G}_k[i_k] \in \mathbb{R}^{r_{k-1} \times r_k}, \quad r_0 = r_d = 1$$

- TT-cores: \mathbf{G}_k
- TT-ranks: r_k
- TT max rank $r = \max r_k$, $k = 0, \dots, d$

Compression:

$$O(n^d) \rightarrow O(ndr^2)$$

$$\sum_{k=1}^d n_k r_{k-1} r_k = O(ndr^2)$$

Formal definition: TT-Vector

Consider a vector $\mathbf{b} \in \mathbb{R}^N$

Where: $N = \prod_{k=1}^d n_k$

We can establish a bijection:

$\mu : l \in \{1, \dots, N\} \mapsto (\mu_1(l), \dots, \mu_d(l))$ Where:

Where:

$\mu_k(l) \in \{1, \dots, n_k\}$

We can represent it using a tensor B:

$$B \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_d}$$

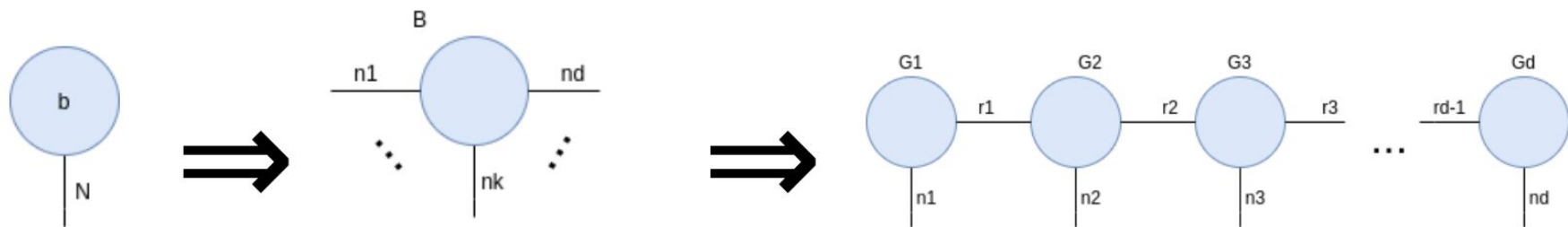
$$B((\mu_1(l), \dots, \mu_d(l))) = b_l$$

Step1: Convert the large matrix into a tensor

Step2: Decompose into TT-representation to get a TT-vector

Formal definition: TT-Vector network diagram

$$b(l) = B((\mu_1(l), \dots, \mu_d(l))) = \underbrace{G_1[\mu_1(l)]}_{1 \times r_1} \underbrace{G_2[\mu_2(l)]}_{r_1 \times r_2} \dots \underbrace{G_d[\mu_d(l)]}_{r_{d-1} \times 1}$$



Formal definition: TT-Matrix

Consider a matrix A :

Where: $A \in \mathbb{R}^{M \times N}$

And: $M = \prod_{k=1}^d m_k$, $N = \prod_{k=1}^d n_k$

We can establish the bijections:

$$\nu : t \in \{1, \dots, M\} \mapsto (\nu_1(t), \dots, \nu_d(t))$$

And:

$$\mu : l \in \{1, \dots, N\} \mapsto (\mu_1(l), \dots, \mu_d(l))$$

Where:

$$\underline{\nu}(t) = (\nu_1(t), \dots, \nu_d(t)) \text{ and } \underline{\mu}(l) = (\mu_1(l), \dots, \mu_d(l))$$

Formal definition: TT-Matrix

We can represent using a tensor W :

$$W \in \mathbb{R}^{m_1 n_1 \times m_2 n_2 \times \dots \times m_d n_d}$$

Where:

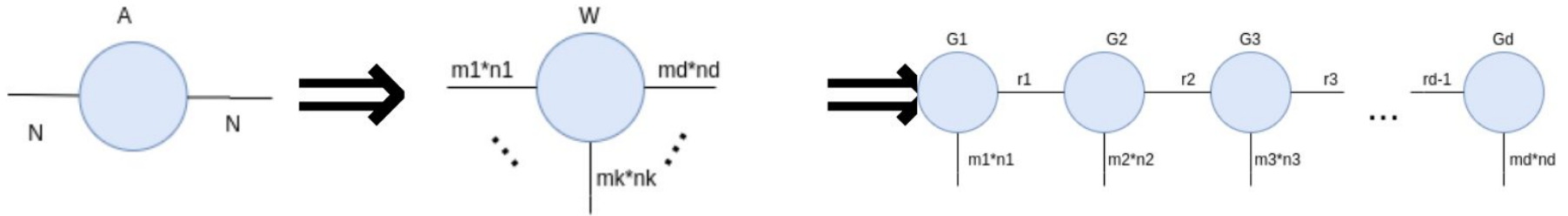
$$W(t, \ell) = \mathcal{W}((\nu_1(t), \mu_1(\ell)), \dots, (\nu_d(t), \mu_d(\ell)))$$

Cores: $\mathbf{G}_k[\nu_k(t), \mu_k(\ell)], k = 1, \dots, d,$

Index: $(\nu_k(t), \mu_k(\ell))$

Formal definition: TT-Matrix network diagram

$$A(t, l) = W((\nu_1(t), \mu_1(l)), \dots, (\nu_d(t), \mu_d(l))) = \underbrace{G_1[\nu_1(t), \mu_1(l)]}_{1 \times r_1} \underbrace{G_2[\nu_2(t), \mu_2(l)]}_{r_1 \times r_2} \dots \underbrace{G_d[\nu_d(t), \mu_d(l)]}_{r_{d-1} \times 1}$$



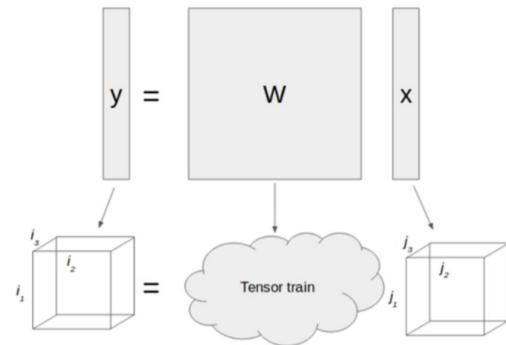
TensorNet:

TT-layer is a fully-connected layer with the weight matrix stored in the TT-format.

- A neural network with **one or more TT-layers as TensorNet**.

FC-layer:

$$\underbrace{\mathbf{y}}_M = \underbrace{\mathbf{W}}_{M \times N} \underbrace{\mathbf{x}}_N + \underbrace{\mathbf{b}}_M$$



TT-Layer

$$\mathcal{Y}(i_1, \dots, i_d) = \sum_{j_1, \dots, j_d} \mathbf{G}_1[i_1, j_1] \dots \mathbf{G}_d[i_d, j_d] \mathcal{X}(j_1, \dots, j_d) + \mathcal{B}(i_1, \dots, i_d).$$

A TT-layer transforms a d -dimensional tensor \mathcal{X} (formed from the corresponding vector \mathbf{x}) to the d -dimensional tensor \mathcal{Y} (which correspond to the output vector \mathbf{y}). We assume that the weight matrix \mathbf{W} is represented in the TT-format with the cores $\mathbf{G}_k[i_k, j_k]$.

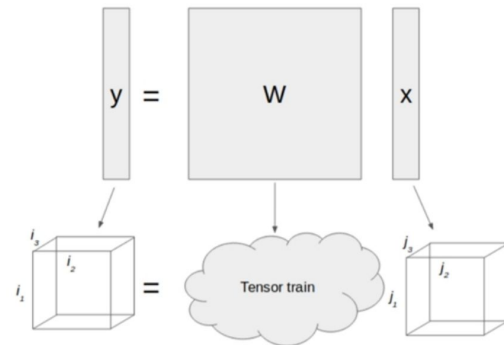
TensorNet:

TT-layer is a fully-connected layer with the weight matrix stored in the TT-format.

- A neural network with **one or more TT-layers as TensorNet**.

FC-layer:

$$\underbrace{\mathbf{y}}_M = \underbrace{\mathbf{W}}_{M \times N} \underbrace{\mathbf{x}}_N + \underbrace{\mathbf{b}}_M$$



TT-Layer

$$\mathcal{Y}(i_1, \dots, i_d) = \sum_{j_1, \dots, j_d} \mathbf{G}_1[i_1, j_1] \dots \mathbf{G}_d[i_d, j_d] \mathcal{X}(j_1, \dots, j_d) + \mathcal{B}(i_1, \dots, i_d).$$

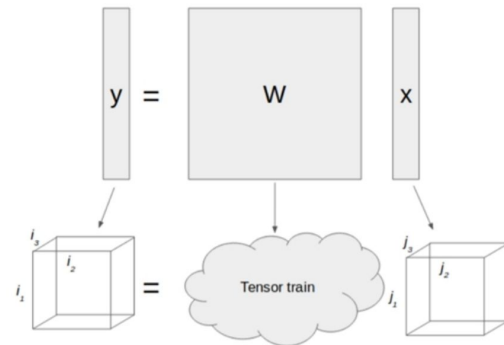
TensorNet:

TT-layer is a fully-connected layer with the weight matrix stored in the TT-format.

- A neural network with **one or more TT-layers as TensorNet**.

FC-layer:

$$\underbrace{\mathbf{y}}_M = \underbrace{\mathbf{W}}_{M \times N} \underbrace{\mathbf{x}}_N + \underbrace{\mathbf{b}}_M$$



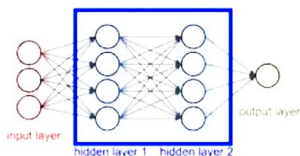
TT-Layer

$$\mathcal{Y}(i_1, \dots, i_d) = \sum_{j_1, \dots, j_d} \mathbf{G}_1[i_1, j_1] \dots \mathbf{G}_d[i_d, j_d] \mathcal{X}(j_1, \dots, j_d) + \mathcal{B}(i_1, \dots, i_d).$$

Forward pass: $O(MN) \rightarrow O(dr^2 m \max\{M, N\})$

No of cores
Maximal rank
 $\max(m_k)$

TensorNet network diagram



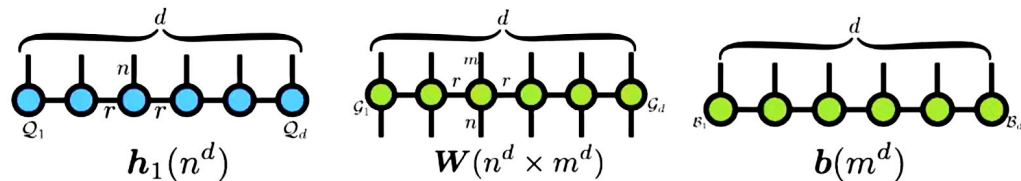
$$\mathbf{h}_2 = \sigma(\mathbf{W}^T \mathbf{h}_1 + \mathbf{b})$$

Hidden layer

Weights

Hidden layer

TT
Representation



$$\mathbf{h}_2 = \sigma \left(\begin{array}{c} \text{Tensor } \mathbf{W} \\ \text{Tensor } \mathbf{h}_1 \end{array} + \text{Tensor } \mathbf{b} \right)$$

↓ TN Contraction

$$= \sigma \left(\begin{array}{c} \text{Contracted } \mathbf{W} \text{ and } \mathbf{h}_1 \\ \text{Tensor } \mathbf{b} \end{array} \right)$$

$$= \sigma \left(\text{Final contracted tensor} \right)$$

Backpropagation

$$\frac{\partial L}{\partial \mathbf{x}} = \mathbf{W}^\top \frac{\partial L}{\partial \mathbf{y}}, \quad \boxed{\frac{\partial L}{\partial \mathbf{W}} = \frac{\partial L}{\partial \mathbf{y}} \mathbf{x}^\top}, \quad \frac{\partial L}{\partial \mathbf{b}} = \frac{\partial L}{\partial \mathbf{y}}.$$

layer. To compute the gradient of the loss function w.r.t. the bias vector \mathbf{b} and w.r.t. the input vector \mathbf{x} one can use equations (6). The latter can be applied using the matrix-by-vector product (where the matrix is in the TT-format) with the complexity of $O(dr^2n \max\{m, n\}^d) = O(dr^2n \max\{M, N\})$.

$$\underbrace{\frac{\partial L}{\partial \mathbf{G}_k[\tilde{i}_k, \tilde{j}_k]}}_{r_{k-1} \times r_k} = \sum_i \frac{\partial L}{\partial \mathcal{Y}(i)} \frac{\partial \mathcal{Y}(i)}{\partial \mathbf{G}_k[\tilde{i}_k, \tilde{j}_k]}. \quad O(M r_{k-1} r_k)$$

Backpropagation

$$\mathcal{Y}(\mathbf{i}) = \sum_j G_1[i_1, j_1] \cdots G_k[i_k, j_k] \cdots G_d[i_d, j_d] \mathcal{X}(\mathbf{j}) + \mathcal{B}(\mathbf{i})$$

a^T X b

$$\frac{\partial \mathcal{Y}(\mathbf{i})}{\partial G_k[i_k, j_k]}$$

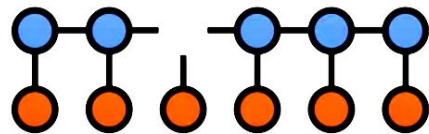
calculus rule

$$\frac{\partial(a^T X b)}{\partial X} = ab^T$$

Backpropagation

$$\mathcal{Y}(\mathbf{i}) = \sum_{\mathbf{j}} \underbrace{G_1[i_1, j_1] \cdots G_k[i_k, j_k]}_{a^T} \underbrace{X}_{\mathbf{X}} \underbrace{G_d[i_d, j_d]}_b \mathcal{X}(\mathbf{j}) + \mathcal{B}(\mathbf{i})$$

$$\frac{\partial \mathcal{Y}(\mathbf{i})}{\partial G_k[i_k, j_k]} = \sum_{\mathbf{j} \setminus j_k} (G_1[i_1, j_1] \cdots)^T (\cdots G_d[i_d, j_d])^T \mathcal{X}(\mathbf{j})$$



Backpropagation

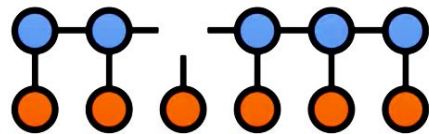
$$\mathcal{Y}(\mathbf{i}) = \sum_{\mathbf{j}} \underbrace{G_1[i_1, j_1] \cdots G_k[i_k, j_k]}_{a^T} \underbrace{X}_{\mathbf{X}} \underbrace{G_d[i_d, j_d]}_b \mathcal{X}(\mathbf{j}) + \mathcal{B}(\mathbf{i})$$

$$\frac{\partial \mathcal{Y}(\mathbf{i})}{\partial G_k[i_k, j_k]} = \sum_{\mathbf{j} \setminus j_k} (G_1[i_1, j_1] \cdots)^T (\cdots G_d[i_d, j_d])^T \mathcal{X}(\mathbf{j})$$

$O(d \bar{r}^2 m \max\{M, N\})$

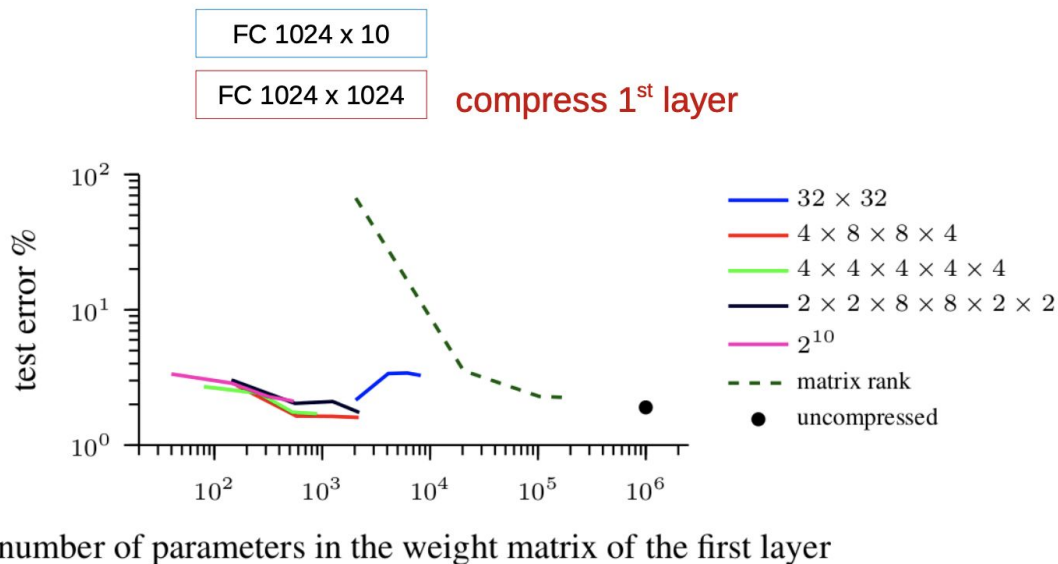
backward pass

$$O(MN) \rightarrow O(d^2 r^4 m \max\{M, N\})$$



Experimental results: MNIST

- Small: MNIST



TT-Layers provide much better flexibility than the matrix rank keeping same compression level
TT-layers with too small number of values for each tensor dimension and with too few dimensions
perform worse than their more balanced counterparts

Experimental results: ImageNet

Architecture	TT-layers compr.	vgg-16 compr.	vgg-19 compr.	vgg-16 top 1	vgg-16 top 5	vgg-19 top 1	vgg-19 top 5
FC FC FC	1	1	1	30.9	11.2	29.0	10.1
TT4 FC FC	50972	3.9	3.5	31.2	11.2	29.8	10.4
TT2 FC FC	194622	3.9	3.5	31.5	11.5	30.4	10.9
TT1 FC FC	713614	3.9	3.5	33.3	12.8	31.9	11.8
TT4 TT4 FC	37732	7.4	6	32.2	12.3	31.6	11.7
MR1 FC FC	3521	3.9	3.5	99.5	97.6	99.8	99
MR5 FC FC	704	3.9	3.5	81.7	53.9	79.1	52.4
MR50 FC FC	70	3.7	3.4	36.7	14.9	34.5	15.8

Table 2: Substituting the fully-connected layers with the TT-layers in vgg-16 and vgg-19 networks on the ImageNet dataset. FC stands for a fully-connected layer; TT \square stands for a TT-layer with all the TT-ranks equal “ \square ”; MR \square stands for a fully-connected layer with the matrix rank restricted to “ \square ”. We report the compression rate of the TT-layers matrices and of the whole network in the second, third and fourth columns.

Great compression factor of 194622 with 0.3 accuracy drop

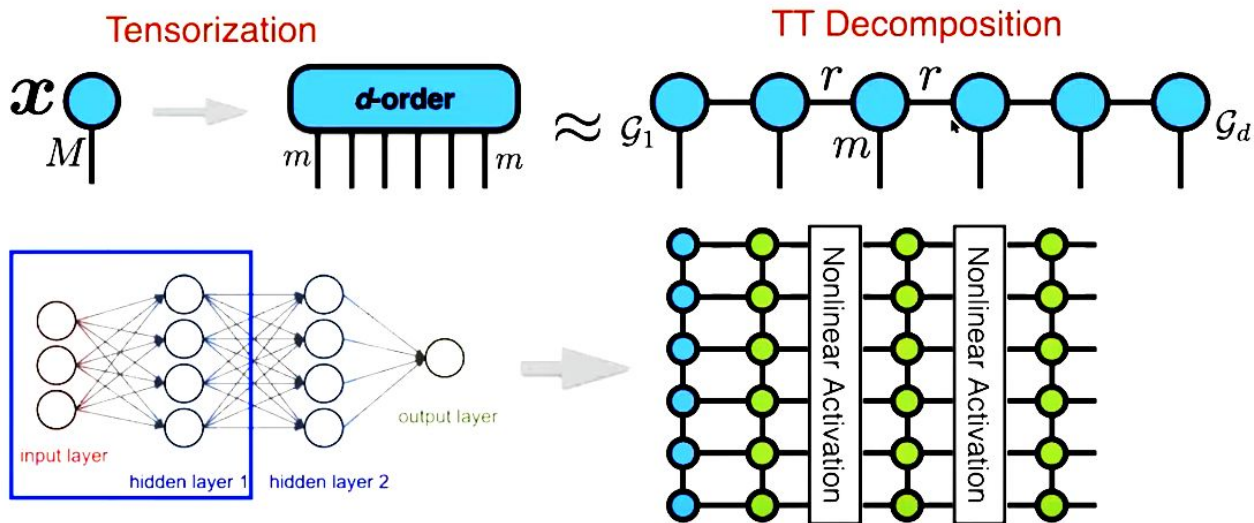
Experimental results: ImageNet

Type	1 im. time (ms)	100 im. time (ms)
CPU fully-connected layer	16.1	97.2
CPU TT-layer	1.2	94.7
GPU fully-connected layer	2.7	33
GPU TT-layer	1.9	12.9

Table 3: Inference time for a 25088×4096 fully-connected layer and its corresponding TT-layer with all the TT-ranks equal 4. The memory usage for feeding forward one image is 392MB for the fully-connected layer and 0.766MB for the TT-layer.

TT-Layer has better inference time in comparison to FC-Layer

Challenges



- ▶ Input data may not admit low-rank TT approximation (small r)
- ▶ Nonlinear activation destroy TT format

Similar works

Lebedev V. et al. Speeding-up convolutional neural networks using fine-tuned cp-decomposition arXiv:1412.6553.

8.5x speedup with 1% accuracy drop

Recent example: Yang, Yinchong, Denis Krompass, and Volker Tresp. "Tensor-Train Recurrent Neural Networks for Video Classification." arXiv:1707.01786

3000 parameters in TT-LSTM vs 71,884,800 in LSTM

Accuracy is better due to additional regularisation

Thanks

References

[1] Tensorizing Neural Network; NIPS 2015 slides [[link](#)]

[2] Alexander Novikov, Dmitry Podoprikin, Anton Osokin, Dmitry Vetrov, Tensorizing Neural Networks; NIPS 2015

[3] Slides by Moussa Traore Mehraveh Javan [[link](#)]

[4] Tensor Train in machine learning Slides [[link](#)]

[5] More slides -> [[Link1](#)], [[Link2](#)]

[6] Lecture Notes [[Link](#)]