

# Tensorizing Neural Networks

Alexander Novikov<sup>1</sup>   Dmitry Podoprikin<sup>1</sup>   Anton Osokin<sup>2</sup>  
Dmitry Vetrov<sup>1,3</sup>

<sup>1</sup>Skolkovo Institute of Science and Technology, Russia   <sup>2</sup>SIERRA, INRIA, France  
<sup>3</sup>Higher School of Economics, Russia

September 4, 2015

# Machine learning

Lets consider a classification or regression problem.



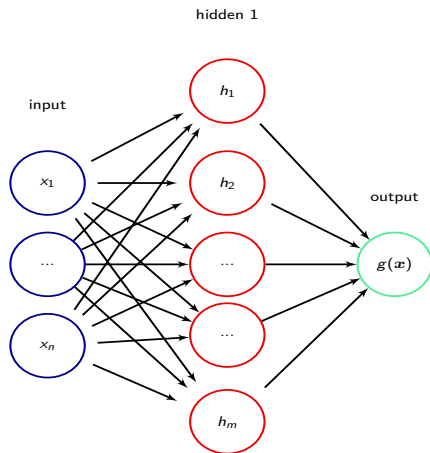
# Neural networks

Lets use a composite function:

$$g(\mathbf{x}) = f(\mathbf{W}_2 \cdot \underbrace{f(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)}_{h \in \mathbb{R}^m} + \mathbf{b}_2)$$

$$h_k = f((\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)_k)$$

$$f(x) = \frac{1}{1 + \exp(-x)}$$



$$g(\mathbf{x}) = f(\mathbf{W}_2 \cdot f(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2)$$

To teach the neural network minimize its' error on the training set:

$$\mathbf{W}_1^*, \mathbf{b}_1^*, \mathbf{W}_2^*, \mathbf{b}_2^* = \arg \min_{\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2} \sum_q (y_q - g(\mathbf{x}_q))^2$$

We will compress fully-connected layers.

Why do we care about memory?

- State-of-the-art deep networks doesn't fit to mobile devices;
- Up to 95% percent of the parameters are in the fully connected layers;
- A shallow network with a huge fully connected layer can achieve almost the same accuracy, as an ensemble of deep CNNs (Ba and Caruana 2014).

- 1 Neural networks
- 2 Matrix formats**
- 3 TensorNet
- 4 Experiments
- 5 Future work

# Compact formats for a matrix

- Constant:  $W(k, \ell) = a$ .
  - Too restrictive
- Zip archive.
  - You need to uncompress the matrix before the multiplication, RAM requirement is not reduced
- Sparse matrix (most of the elements are zero).
  - How to choose which elements are zero?
  - A little bit hard to implement efficiently (especially on a GPU)

# Matrix rank decomposition

Lets consider an  $M \times N$  matrix  $\mathbf{W}$  with the rank equals  $r$ . We can use  $(M + N)r$  parameters instead of  $MN$ :

$$\underbrace{\mathbf{W}}_{M \times N} = \underbrace{\mathbf{A}}_{M \times r} \underbrace{\mathbf{B}}_{r \times N}$$

It works, but we want to do better.



Tensor Train (TT) decomposition (Oseledets 2011):

- A compact representation for vectors, matrices and tensors;
- Allows for efficient application of linear algebra operations.

Tensor  $\mathbf{A}$  is said to be in the *TT-format*, if

$$\mathbf{A}(i_1, \dots, i_d) = \mathbf{G}_1[i_1] \mathbf{G}_2[i_2] \cdots \mathbf{G}_d[i_d], \quad i_k \in \{1, \dots, n\},$$

where  $\mathbf{G}_k[i_k]$  — is a matrix of size  $r_{k-1} \times r_k$ ,  $r_0 = r_d = 1$ .

Notation & terminology:

- $\mathbf{G}_k$  — *TT-cores*;
- $r_k$  — *TT-ranks*;
- $r = \max_{k=0, \dots, d} r_k$  — *the maximal TT-rank*.

The TT-format uses  $O(ndr^2)$  memory to store  $n^d$  elements. **Efficient only if the TT-rank is small.**

## TT-format: example

$$\mathbf{A}(i_1, i_2, i_3) = i_1 + i_2 + i_3,$$

$$i_1 \in \{1, 2, 3\}, i_2 \in \{1, 2, 3, 4\}, i_3 \in \{1, 2, 3, 4, 5\}.$$

$$\mathbf{A}(i_1, i_2, i_3) = G_1[i_1]G_2[i_2]G_3[i_3],$$

$$A(i_1, i_2, i_3) = i_1 + i_2 + i_3,$$

$$i_1 \in \{1, 2, 3\}, i_2 \in \{1, 2, 3, 4\}, i_3 \in \{1, 2, 3, 4, 5\}.$$

$$A(i_1, i_2, i_3) = G_1[i_1]G_2[i_2]G_3[i_3],$$

$$G_1[i_1] = \begin{bmatrix} i_1 & 1 \end{bmatrix} \quad G_2[i_2] = \begin{bmatrix} 1 & 0 \\ i_2 & 1 \end{bmatrix} \quad G_3[i_3] = \begin{bmatrix} 1 \\ i_3 \end{bmatrix}$$

Lets check:

$$\begin{aligned} A(i_1, i_2, i_3) &= \begin{bmatrix} i_1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ i_2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ i_3 \end{bmatrix} = \\ &= \begin{bmatrix} i_1 + i_2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ i_3 \end{bmatrix} = i_1 + i_2 + i_3. \end{aligned}$$

$$\mathbf{A}(i_1, i_2, i_3) = i_1 + i_2 + i_3,$$

$$i_1 \in \{1, 2, 3\}, i_2 \in \{1, 2, 3, 4\}, i_3 \in \{1, 2, 3, 4, 5\}.$$

$$\mathbf{A}(i_1, i_2, i_3) = G_1[i_1]G_2[i_2]G_3[i_3],$$

$$G_1 = \left( \begin{bmatrix} 1 & 1 \end{bmatrix}, \begin{bmatrix} 2 & 1 \end{bmatrix}, \begin{bmatrix} 3 & 1 \end{bmatrix} \right)$$

$$G_2 = \left( \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 2 & 1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 3 & 1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 4 & 1 \end{bmatrix} \right)$$

$$G_3 = \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 \\ 3 \end{bmatrix}, \begin{bmatrix} 1 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 \\ 5 \end{bmatrix} \right)$$

The tensor has  $3 \cdot 4 \cdot 5 = 60$  elements.

The TT-format use 32 parameters to describe it.

# Mapping example: a vector

How to store a vector  $\mathbf{b}$  in the TT-format?

Build a mapping from a vector  $\mathbf{b}$  indices to tensor's elements:

$$t \leftrightarrow \mathbf{i} = (i_1, \dots, i_d)$$

Example (Matlab reshape):

$$B(1, 1, 1) = \mathbf{b}(t(1, 1, 1)) = \mathbf{b}(1),$$

$$B(2, 1, 1) = \mathbf{b}(t(2, 1, 1)) = \mathbf{b}(2),$$

$\dots,$

$$B(2, 3, 3) = \mathbf{b}(t(2, 3, 3)) = \mathbf{b}(18).$$

Now lets use the TT-format for the tensor  $\mathbf{B}$ .

# Matrices in the TT-format

Build a mapping from row / column indices of matrix  $\mathbf{W} = [\mathbf{W}(t, \ell)]$  to vectors  $\mathbf{i}$  and  $\mathbf{j}$ :  $t \leftrightarrow \mathbf{i} = (i_1, \dots, i_d)$  and  $\ell \leftrightarrow \mathbf{j} = (j_1, \dots, j_d)$ .

TT-format for matrix  $\mathbf{W}$ :

$$\mathbf{W}(i_1, \dots, i_d; j_1, \dots, j_d) = \mathbf{W}(t(\mathbf{i}), \ell(\mathbf{j})) = \underbrace{\mathbf{G}_1[i_1, j_1]}_{1 \times r} \underbrace{\mathbf{G}_2[i_2, j_2]}_{r \times r} \dots \underbrace{\mathbf{G}_d[i_d, j_d]}_{r \times 1}$$

The TT-format exists for any matrix  $\mathbf{W}$  and uses  $O(dmnr^2)$  memory to store  $m^d n^d$  elements. **Efficient only if the TT-rank is small.**

- 1 Neural networks
- 2 Matrix formats
- 3 TensorNet**
- 4 Experiments
- 5 Future work



Input is a  $N$ -dimensional vector  $\mathbf{x}$ , output is a  $M$ -dimensional vector  $\mathbf{h}$ :

$$\mathbf{h} = \mathbf{W}\mathbf{x} + \mathbf{b}.$$

$\mathbf{W}$  is represented in the TT-format:

$$h(i_1, \dots, i_d) = \sum_{j_1, \dots, j_d} \mathbf{G}_1[i_1, j_1] \dots \mathbf{G}_d[i_d, j_d] \mathbf{x}(j_1, \dots, j_d) + \mathbf{b}(i_1, \dots, i_d)$$

The parameters are the vector  $\mathbf{b}$  and the TT-cores  $\{\mathbf{G}_k\}_{k=1}^d$

# Tensor Train layer learning

$$g(\mathbf{x}) = \sum_{j_1, \dots, j_d} \mathbf{G}_1[i_1, j_1] \dots \mathbf{G}_d[i_d, j_d] \mathbf{x}(j_1, \dots, j_d) + \mathbf{b}(i_1, \dots, i_d)$$

Instead of optimizing over all possible matrices  $\mathbf{W}$ :

$$\mathbf{W}^*, \mathbf{b}^* = \arg \min_{\mathbf{W}, \mathbf{b}} \sum_q (y_q - g(\mathbf{x}_q))^2,$$

we optimize over the matrices representable in the TT-format with rank  $r$ :

$$\mathbf{G}_1^*, \dots, \mathbf{G}_d^*, \mathbf{b}^* = \arg \min_{\mathbf{G}_1, \dots, \mathbf{G}_d, \mathbf{b}} \sum_q (y_q - g(\mathbf{x}_q))^2.$$

# Tensor Train layer: the Jacobian

The Jacobian of the linear transformation:

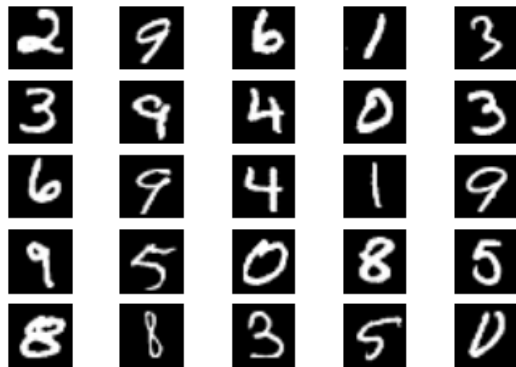
$$\mathbf{h}(\mathbf{i}) = \sum_j \mathbf{G}_1[i_1, j_1] \dots \mathbf{G}_k[i_k, j_k] \dots \mathbf{G}_d[i_d, j_d] \mathbf{x}(j) + \mathbf{b}(\mathbf{i}).$$

$$\begin{aligned} \frac{\partial \mathbf{h}(\mathbf{i})}{\partial \mathbf{G}_k[i_k, j_k]} &= \sum_{j^{\setminus k}} \overbrace{\mathbf{G}_1[i_1, j_1] \dots \mathbf{G}_k[i_k, j_k]}^{r \times 1} \overbrace{\dots \mathbf{G}_d[i_d, j_d]}^{1 \times r} \mathbf{x}(j) = \\ &\sum_{j^{\setminus k, d}} \mathbf{G}_1[i_1, j_1] \dots \mathbf{G}_k[i_k, j_k] \dots \mathbf{G}_{d-1}[i_{d-1}, j_{d-1}] \\ &\underbrace{\sum_{j_d} \mathbf{G}_d[i_d, j_d] \mathbf{x}(j)}_{r \times mn^{d-1}} \end{aligned}$$

The complexity is  $O(d^2 r^4 m \max\{M, N\})$ .

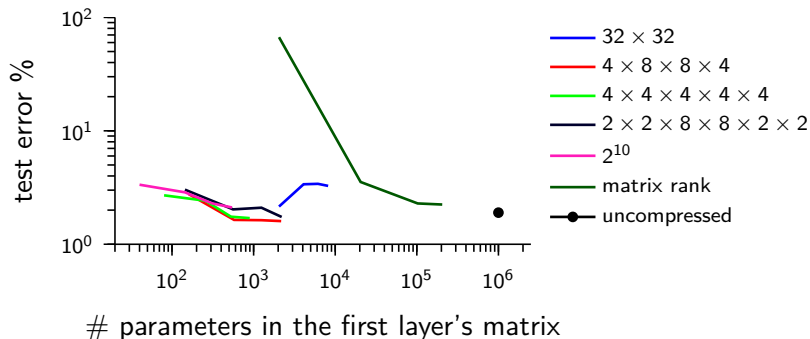
- 1 Neural networks
- 2 Matrix formats
- 3 TensorNet
- 4 Experiments**
- 5 Future work

## Handwritten digits recognition



# Mnist cont'd

A two layered neural network with  $1024 \times 1024$  and  $1024 \times 10$  matrices.

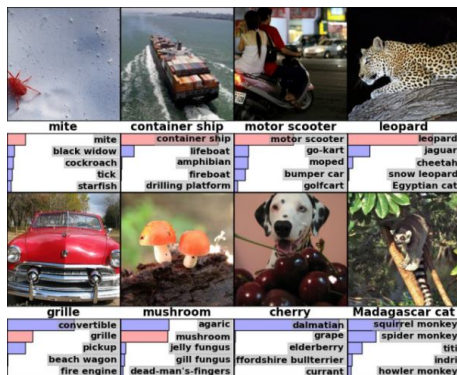


Type	1 im. time (ms)	100 im. time (ms)
CPU fully-connected layer	16.09	97.2
CPU TT-layer	1.24	94.7
GPU fully-connected layer	2.7	33
GPU TT-layer	1.92	12.86

25 088 × 4 096 layer, the fully-connected layer uses 392MB and the TT-layer uses 0.766MB.

# Image classification

## Cifar & ImageNet



We used a  $262\,144 \times 4\,096$  TT-matrix to outperform other non-convolutional neural networks on Cifar.



Architecture	Matrices compr.	Network compr.	Error
FC FC FC	1	1	11.2
TT4 FC FC	50 972	3.9	11.2
TT2 FC FC	194 622	3.9	11.5
TT1 FC FC	713 614	3.9	12.8
TT4 TT4 FC	37 732	7.4	12.3
LR1 FC FC	3 521	3.9	97.6
LR5 FC FC	704	3.9	53.9
LR50 FC FC	70	3.7	14.9

A 3-layered network, FC stands for the traditional layer; TT $\square$  stands for the TT-layer with all the TT-ranks equal " $\square$ "; LR $\square$  stands for the rank decomposition with the rank equal " $\square$ ".

- Get rid of the convolutions.
  - Convolution matrix can be represented in the TT-format with small ranks
- Fully tensorial network.
  - Use the TT-format for all the intermediate calculations
  - Allows billions of hidden neurons
- Riemannian optimization.
  - Highly efficient for the optimization with the ranks constraints

- Ba, Jimmy and Rich Caruana (2014). “Do Deep Nets Really Need to be Deep?” In: *Advances in Neural Information Processing Systems 27*. Ed. by Z. Ghahramani et al. Curran Associates, Inc., pp. 2654–2662.
- Oseledets, I. V. (2011). “Tensor-Train Decomposition”. In: *SIAM J. Scientific Computing* 33.5, pp. 2295–2317.