

Tensorizing Neural Networks

Alexander Novikov

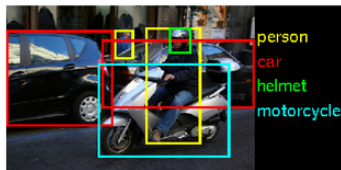
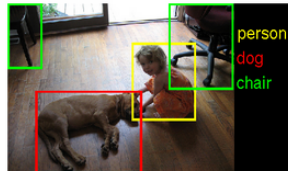
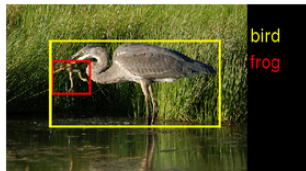
Dmitry Podoprikin

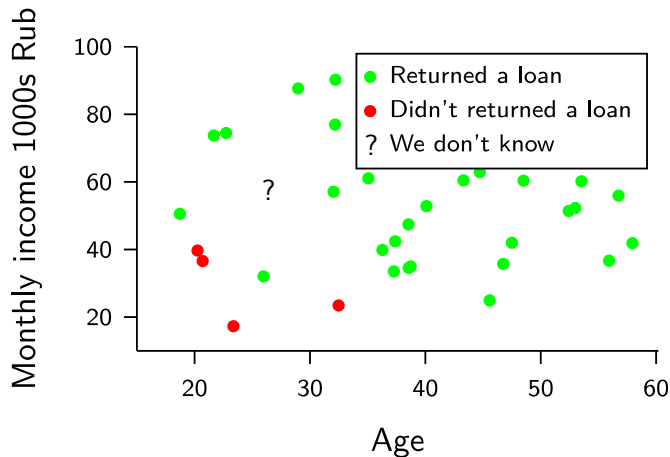
Anton Osokin

Dmitry Vetrov

August 24, 2015

Machine learning









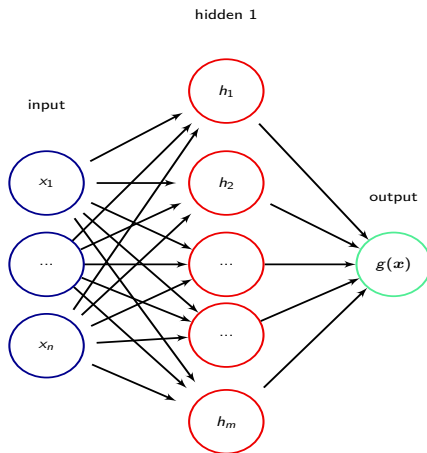
Neural networks

Lets use a composite function:

$$g(x) = f(W_2 \cdot \underbrace{f(W_1 x + b_1)}_{h \in \mathbb{R}^m} + b_2)$$

$$h_k = f((W_1 x + b_1)_k)$$

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



$$g(x) = f(W_2 \cdot f(W_1 x + b_1) + b_2)$$

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

$$W_1^*, b_1^*, W_2^*, b_2^* = \arg \min_{W_1, b_1, W_2, b_2} \sum_q (y_q - g(x_q))^2$$

You always want to use larger layers and to use more of them.

- 1 It's hard to train too many layers (too deep networks);
- 2 For large layers matrix \mathbf{W} doesn't fit to memory (especially on mobile devices).

1 Neural networks

2 Matrix formats

3 TensorNet

4 Experiments

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 even more.

Tensor Train (TT) decomposition (Oseledets 2011):

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

Mapping example: a vector

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):

$$B(1, 1, 1) = b(t(1, 1, 1)) = b(1),$$

$$B(2, 1, 1) = b(t(2, 1, 1)) = b(2),$$

$\dots,$

$$B(2, 3, 3) = b(t(2, 3, 3)) = b(18).$$

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} :

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

Notation & terminology:

- $\mathbf{W} \in \mathbb{R}^{M \times N}$, $M = m^d$, $N = n^d$;
- $i_k \in \{1, \dots, m\}$, $j_k \in \{1, \dots, n\}$;
- \mathbf{G}_k — TT-cores;
- r — TT-rank;

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

Tensor Train layer

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:

$$h(i_1, \dots, i_d) = \sum_{j_1, \dots, j_d} G_1[i_1, j_1] \dots G_d[i_d, j_d] x(j_1, \dots, j_d) + b(i_1, \dots, i_d)$$

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:

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

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

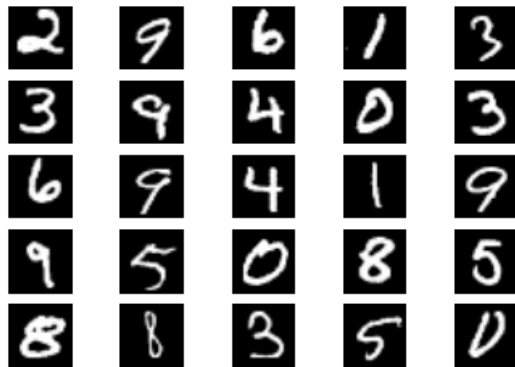
1 Neural networks

2 Matrix formats

3 TensorNet

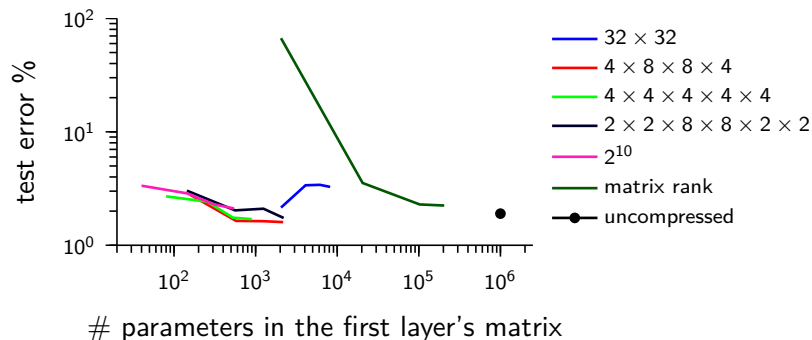
4 Experiments

Handwritten digits recognition



Mnist cont'd

A two layered neural network with 1024×1024 and 1024×10 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 262144×4096 TT-matrix to outperform other plain (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 ".

Oseledets, I. V. (2011). “Tensor-Train Decomposition”. In: *SIAM J. Scientific Computing* 33.5, pp. 2295–2317.