

TensorNet: putting neural networks on a Tensor Train

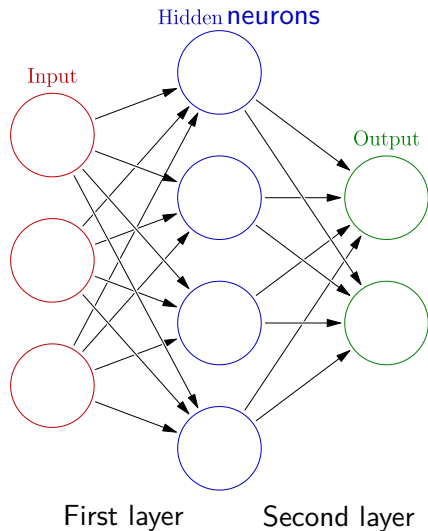
A. Novikov D. Podoprikin A. Osokin D. Vetrov

March 6, 2015

Outline

- 1 Neural networks
- 2 Tensor train
- 3 TensorNet
- 4 Experiments

Neural networks



Why do we care about memory?

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

Matrix rank decomposition

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

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

Drawbacks of the rank decomposition

Problems:

- 1 The low rank compression rate is limited (we want more);
- 2 There is no practical way to train low rank shallow networks.

1 Neural networks

2 Tensor train

3 TensorNet

4 Experiments

Tensor Train (TT) decomposition:

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

Mapping example: vector

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

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

Example (Matlab reshape):

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

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

$\dots,$

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

Matrices in the TT-format

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

TT-format for matrix \mathbf{W} :

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

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;

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

1 Neural networks

2 Tensor train

3 TensorNet

4 Experiments

Tensor Train layer: feedforward

Input is a $N \times 1$ vector \mathbf{x} , output is a $M \times 1$ vector \mathbf{y} :

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

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

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

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

Backpropagation

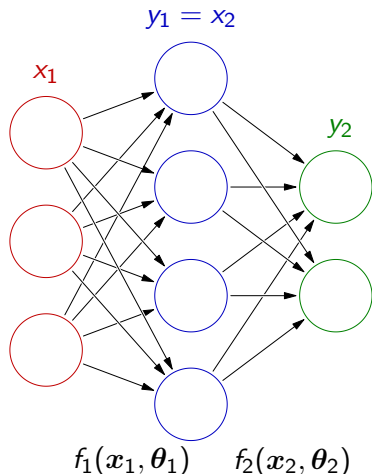
$$L = \frac{1}{2} \sum_{s=1}^S \|y_2^s - y^s\|_2^2$$

$$\frac{\partial L}{\partial y_2} = \sum_{s=1}^S (y_2^s - y^s)$$

$$\frac{\partial L}{\partial x_2} = \sum_i \frac{\partial L}{\partial y_2(i)} \frac{\partial y_2(i)}{\partial x_2}$$

$$= \frac{\partial f_2}{\partial x_2} \frac{\partial L}{\partial y_2}$$

$$\frac{\partial L}{\partial \theta_2} = \frac{\partial f_2}{\partial \theta_2} \frac{\partial L}{\partial y_2}$$

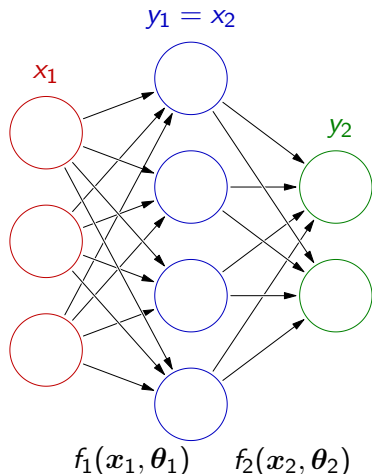


Backpropagation cont'd

From each layer we need only this:

$$\frac{\partial L}{\partial \mathbf{x}_k} = \mathbf{g}_k^{\mathbf{x}}(\mathbf{x}_k, \frac{\partial L}{\partial \mathbf{y}_k})$$

$$\frac{\partial L}{\partial \boldsymbol{\theta}_k} = \mathbf{g}_k^{\boldsymbol{\theta}}(\mathbf{x}_k, \frac{\partial L}{\partial \mathbf{y}_k})$$



Tensor Train layer: backpropagation

Input: vectors $\frac{\partial L}{\partial \mathbf{y}} \in \mathbb{R}^M$ and $\mathbf{x} \in \mathbb{R}^N$.

Output: $\frac{\partial L}{\partial \mathbf{x}}$, $\frac{\partial L}{\partial \mathbf{b}}$ and $\frac{\partial L}{\partial \mathbf{G}_k[i_k j_k]}$.

Tensor Train layer: backpropagation

Input: vectors $\frac{\partial L}{\partial \mathbf{y}} \in \mathbb{R}^M$ and $\mathbf{x} \in \mathbb{R}^N$.

Output: $\frac{\partial L}{\partial \mathbf{x}}$, $\frac{\partial L}{\partial \mathbf{b}}$ and $\frac{\partial L}{\partial \mathbf{G}_k[i_k, j_k]}$.

$$\frac{\partial L}{\partial \mathbf{x}} = \mathbf{W}^\top \frac{\partial L}{\partial \mathbf{y}},$$

$$\frac{\partial L}{\partial \mathbf{b}} = \frac{\partial L}{\partial \mathbf{y}}.$$

$$\underbrace{\frac{\partial L}{\partial \mathbf{G}_k[i_k, j_k]}}_{r \times r} = \sum_{i \setminus k} \frac{\partial L}{\partial \mathbf{y}(i)} \frac{\partial \mathbf{y}(i)}{\partial \mathbf{G}_k[i_k, j_k]}$$

Tensor Train layer: Jacobian

We want to differentiate the following expression:

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

Tensor Train layer: Jacobian

We want to differentiate the following expression:

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

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

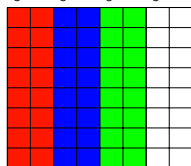
$$W(i_1, i_2, i_3; j_1, j_2, j_3) = \underbrace{G_1[i_1, j_1]}_{\in \mathbb{R}} \underbrace{G_2[i_2, j_2]}_{\in \mathbb{R}} \underbrace{G_3[i_3, j_3]}_{\in \mathbb{R}}$$

$$W \in \mathbb{R}^{64 \times 64}$$

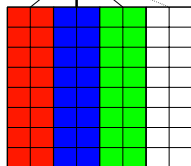
Input x and output y are re-shaped to $4 \times 4 \times 4$ tensor.

To vanish all dashed line weights, set $G_3[i_3 = 2, j_3 = 4] = 0$.

$$i_3 = 1 \quad i_3 = 2 \quad i_3 = 3 \quad i_3 = 4$$



Hidden units



Input image

$$j_3 = 1 \quad j_3 = 2 \quad j_3 = 3 \quad j_3 = 4$$

1 Neural networks

2 Tensor train

3 TensorNet

4 Experiments

Mnist dataset, two layered neural network. The input 28×28 image is reshaped to $2 \times 2 \times 7 \times 2 \times 2 \times 7$ tensor.

Network	Error	Neurons	I layer params	II layer
Baseline	2.34%	500	400 000	5 000
TensorNet (one TT-layer)	2.26%	46 656	420	466 560
TensorNet (two TT-layer)	2.07%	15 625	3 360	1 380
TensorNet	2.6%	15 625	350	156 250
TensorNet (random order)	3.5%	15 625	350	156 250

Mnist dataset, two layered neural network. The input 28×28 image is reshaped to $4 \times 7 \times 4 \times 7$ tensor.

Network	Error	Neurons	I layer params	II layer
Baseline	2.34%	500	400 000	5 000
TensorNet (one TT-layer)	1.68%	4 096	1 760	40 960

Deep convolutional neural network for CIFAR-10 (image classification).
We compressed the last two fully connected layers x11 (the error increased from 23.25% to 23.74%).

- 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.