

Computationally efficient methods for partition function estimation in graphical models in TT-format

D. Vetrov

Lomonosov Moscow State University

MIT

Tensor — multidimensional array (function of its indices)

$$\mathbf{A}(\mathbf{x}) = \mathbf{A}(x_1, \dots, x_n), \quad x_i \in \{1, \dots, d_i\}$$

Terminology:

- n — *dimensionality*;
- x_i — *indices*.

A tensor \mathbf{A} is said to be represented in the **TT-format** if

$$\mathbf{A}(x_1, \dots, x_n) = G_1^{\mathbf{A}}[x_1] G_2^{\mathbf{A}}[x_2] \cdots G_n^{\mathbf{A}}[x_n],$$

where $G_i^{\mathbf{A}}[x_i]$ is a matrix of size $r_{i-1}(\mathbf{A}) \times r_i(\mathbf{A})$, $r_0(\mathbf{A}) = r_n(\mathbf{A}) = 1$.

Terminology:

- $G_i^{\mathbf{A}}$ — **TT-cores**;
- $r_i(\mathbf{A})$ — **TT-ranks**;
- $r(\mathbf{A}) = \max_{i=0, \dots, n} r_i(\mathbf{A})$ — **maximal TT-rank**.

The TT-format uses $O(ndr^2(\mathbf{A}))$ memory to store $O(d^n)$ elements ($d = \max_{i=1, \dots, n} d_i$).

Operations on tensors in the TT-format

Operation	Output rank
$\mathbf{C} = c \cdot \mathbf{A}$	$r(\mathbf{C}) = r(\mathbf{A})$
$\mathbf{C} = \mathbf{A} + c$	$r(\mathbf{C}) = r(\mathbf{A}) + 1$
$\mathbf{C} = \mathbf{A} + \mathbf{B}$	$r(\mathbf{C}) \leq r(\mathbf{A}) + r(\mathbf{B})$
$\mathbf{C} = \mathbf{A} \odot \mathbf{B}$	$r(\mathbf{C}) \leq r(\mathbf{A}) r(\mathbf{B})$
$c = M\mathbf{b}$	$r(c) \leq r(M) r(\mathbf{b})$
$\mathbf{C} = \text{round}(\mathbf{A}, \varepsilon)$	$r(\mathbf{C}) \leq r(\mathbf{A})$
sum \mathbf{A}	–
$\ \mathbf{A}\ _F$	–

Suppose that you have a TT-decomposition $\mathbf{A}(x) = G_1^{\mathbf{A}}[x_1] \cdots G_n^{\mathbf{A}}[x_n]$ with non-optimal TT-ranks.

The TT-rounding procedure (Oseledets 2011)

$$\widehat{\mathbf{A}} = \text{round}(\mathbf{A}, \varepsilon), \quad \varepsilon \geq 0$$

finds a tensor $\widehat{\mathbf{A}}$:

- 1 $\|\mathbf{A} - \widehat{\mathbf{A}}\|_F \leq \varepsilon \|\mathbf{A}\|_F$;
- 2 its TT-ranks are minimal among all tensors \mathbf{B} :
 $\|\mathbf{A} - \mathbf{B}\|_F \leq \frac{\varepsilon}{\sqrt{n-1}} \|\mathbf{A}\|_F$.

Here $\|\mathbf{A}\|_F = \sqrt{\sum_{x_1, \dots, x_n} \mathbf{A}^2(x_1, \dots, x_n)}$.

Finding a TT-representation of a tensor

- There are analytical formulae for some special cases;

Example

$$\mathbf{A}(x_1, \dots, x_n) = \sum_{i=1}^n f_i(x_i),$$

$$G_i^{\mathbf{A}}[x_i] = \begin{bmatrix} 1 & 0 \\ f_i(x_i) & 1 \end{bmatrix}, \quad i = 2, \dots, n-1,$$

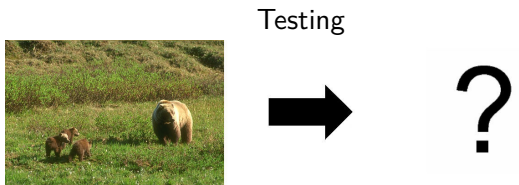
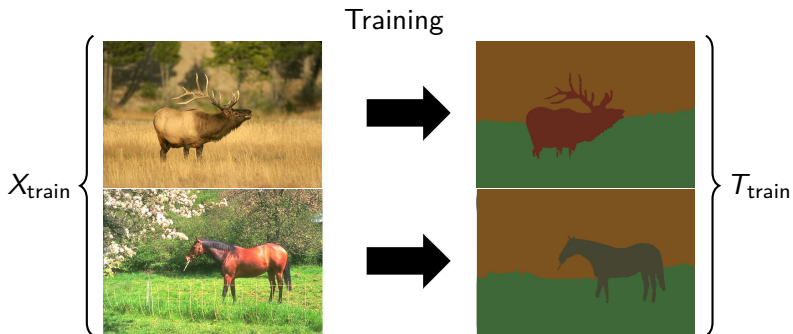
$$G_1^{\mathbf{A}}[x_1] = \begin{bmatrix} f_1(x_1) & 1 \end{bmatrix}, \quad G_n^{\mathbf{A}}[x_n] = \begin{bmatrix} 1 \\ f_n(x_n) \end{bmatrix}.$$

Finding a TT-representation of a tensor

- There are analytical formulae for some special cases;
- **TT-SVD**: finds an exact TT-representation for a tensor but suitable only for low dimensionality n ;
- **AMEn-cross**: builds a TT-approximation of a tensor by using only a small fraction of its elements; suitable for high dimensionality n but doesn't have strong theoretical guarantees.

- 1 Tensor train
- 2 Motivation example**
- 3 Results
- 4 Experimental evaluation
- 5 High Order Potentials
- 6 Conclusion

Semantic segmentation



Probabilistic approach

- Define a **probabilistic model** on the set of all possible labellings.
- Let $p(T|X, W)$ measure the probability of the labelling T given the image X and the parameters of the model W .
- The goal is to find the labelling T^* that maximizes $p(T|X, W)$:

$$T^* = \arg \max_T p(T|X, W).$$

This is called the **maximum a posteriori (MAP) inference**.

- We will use **Markov random fields (MRFs)** to define the probabilistic model $p(T|X, W)$.

The MAP-inference problem now corresponds to the following problem:

$$\max_T p(T|X, W) = \max_T \frac{1}{Z(X, W)} \prod_{c \in \mathcal{C}} \Psi_c(T_c; X, W).$$

Further we demonstrate how one can address such a problem using the **Tensor-Train (TT) framework**.

We will assume that

- the parameters of the model W are already chosen;
- we are performing the MAP-inference for the concrete image X .

So, to simplify notation, we won't explicitly write X, W any more:

$$\max_T p(T) = \max_T \frac{1}{Z} \prod_{c \in \mathcal{C}} \Psi_c(T_c)$$

MAP-inference & energy minimization

The MAP-inference problem

$$\max_{\mathbf{x}} P(\mathbf{x}) = \max_{\mathbf{x}} \frac{1}{Z} \prod_{\ell=1}^m \Psi_{\ell}(\mathbf{x}^{\ell})$$

is equivalent to the following problem:

$$\min_{\mathbf{x}} \sum_{\ell=1}^m [-\ln \Psi_{\ell}(\mathbf{x}^{\ell})].$$

MAP-inference & energy minimization

The MAP-inference problem

$$\max_{\mathbf{x}} P(\mathbf{x}) = \max_{\mathbf{x}} \frac{1}{Z} \prod_{\ell=1}^m \Psi_{\ell}(\mathbf{x}^{\ell})$$

is equivalent to the following problem:

$$\min_{\mathbf{x}} \sum_{\ell=1}^m [-\ln \Psi_{\ell}(\mathbf{x}^{\ell})].$$

Terminology:

- The terms $\Theta_{\ell}(\mathbf{x}^{\ell}) = -\ln \Psi_{\ell}(\mathbf{x}^{\ell})$ are called MRF **potentials**.
- Their sum $E(\mathbf{x}) = \sum_{\ell=1}^m \Theta_{\ell}(\mathbf{x}^{\ell})$ is called MRF **energy**.

So, the MAP-inference is equivalent to **energy minimization**:

$$\min_{\mathbf{x}} E(\mathbf{x}) = \min_{\mathbf{x}} \sum_{\ell=1}^m \Theta_{\ell}(\mathbf{x}^{\ell}).$$

Maximum likelihood estimation:

$$\begin{aligned}W^* &= \arg \max_W \mathbf{P}(T_{\text{train}} | X_{\text{train}}, W) \\&= \arg \max_W \frac{1}{Z(X_{\text{train}}, W)} \prod_{\ell=1}^m \Psi_{\ell}(T_{\text{train}}; X_{\text{train}}, W) \\&= \arg \max_W \left(\sum_{\ell=1}^m \log \Psi_{\ell}(T_{\text{train}}; X_{\text{train}}, W) - \log Z(X_{\text{train}}, W) \right)\end{aligned}$$

We need to solve the following optimization task:

$$\max_W \left(\sum_{\ell=1}^m \log \Psi_{\ell}(T_{\text{train}}; X_{\text{train}}, W) - \log Z(X_{\text{train}}, W) \right)$$
$$Z(X_{\text{train}}, W) = \sum_T \prod_{\ell=1}^m \Psi_{\ell}(T; X_{\text{train}}, W)$$

To use optimization methods we need to compute **the value** and **the gradient** of the objective function.

It is easy to compute the gradient of the log-partition function if we know the **marginal distributions**¹:

$$\mathbf{P}(t_i) = \sum_{T \setminus t_i} \mathbf{P}(T)$$

¹S. Nowozin and C. Lampert (2010). “Structured Learning and Prediction in Computer Vision”. In: *Foundations and Trends in Computer Graphics and Vision* 6.3–4, pp. 185–365.

Assume that we fixed the parameters W and the training data $X_{\text{train}}, T_{\text{train}}$.

We will focus on three problems (hereinafter we use \mathbf{x} to denote the variables of the model):

- MAP-inference: $\min_{\mathbf{x}} \mathbf{E}(\mathbf{x})$
- The partition function estimation: $Z = \sum_{\mathbf{x}} \hat{\mathbf{P}}(\mathbf{x})$;
- The marginal distributions estimation: $\mathbf{P}(x_i) = \frac{1}{Z} \sum_{\mathbf{x} \setminus x_i} \hat{\mathbf{P}}(\mathbf{x})$.

- 1 Tensor train
- 2 Motivation example
- 3 Results**
- 4 Experimental evaluation
- 5 High Order Potentials
- 6 Conclusion

- The energy $E(\mathbf{x})$ can be considered as an n -dimensional tensor:

$$E(\mathbf{x}) = E(x_1, \dots, x_n).$$

- Then energy minimization corresponds to **finding the minimal element in the tensor** $E(\mathbf{x})$.
- If the energy $E(\mathbf{x})$ were represented in the TT-format, we could use a special algorithm from the TT-framework to find the minimal element.
- **How to convert the energy tensor into the TT-format?**
AMEn-algorithm?

- The energy $E(\mathbf{x})$ can be considered as an n -dimensional tensor:

$$E(\mathbf{x}) = E(x_1, \dots, x_n).$$

- Then energy minimization corresponds to **finding the minimal element in the tensor $E(\mathbf{x})$** .
- If the energy $E(\mathbf{x})$ were represented in the TT-format, we could use a special algorithm from the TT-framework to find the minimal element.
- **How to convert the energy tensor into the TT-format?**
AMEn-algorithm? Possible, but there is also a much better way!

The idea of the algorithm

- Let's try to take into account the structure of the energy tensor E .

Recall:
$$E(\mathbf{x}) = \sum_{\ell=1}^m \Theta_{\ell}(\mathbf{x}^{\ell}).$$

- Each potential $\Theta_{\ell}(\mathbf{x}^{\ell})$ can be considered as an n -dimensional tensor $\Theta_{\ell}(\mathbf{x})$ if we add **inessential variables** $\mathbf{x} \setminus \mathbf{x}^{\ell}$ for non-existing dimensions: $\Theta_{\ell}(\mathbf{x}) \equiv \Theta_{\ell}(\mathbf{x}^{\ell})$.

- The energy $E(\mathbf{x})$ can be expressed as a sum of the tensors $\Theta_{\ell}(\mathbf{x})$:

$$E(\mathbf{x}) = \sum_{\ell=1}^m \Theta_{\ell}(\mathbf{x}).$$

- If the tensors Θ_{ℓ} were represented in the TT-format, we could exploit the **summation operation** on tensors in the TT-format to build the TT-representation for the tensor E .
- How to find the TT-decomposition for each tensor Θ_{ℓ} ?**

Converting potentials into the TT-Format

- As opposed to the energy $E(\mathbf{x})$, each potential $\Theta_\ell(\mathbf{x}^\ell)$ depends only on part of the all variables and is **usually of low dimensionality**.
- To compute the TT-decomposition of the tensor $\Theta_\ell(\mathbf{x}^\ell)$, we can use the **TT-SVD algorithm**.
- All that remains is to add the inessential variables $\mathbf{x} \setminus \mathbf{x}^\ell$ to $\Theta_\ell(\mathbf{x}^\ell)$ so as to make it n -dimensional.
- These inessential variables can be added **constructively**:
 - Let $\mathbf{x} = (x_1, x_2, x_3, x_4, x_5)$, $\mathbf{x}^\ell = (x_1, x_2, x_4)$.
 - Suppose that after TT-SVD we have:

$$\Theta_\ell(x_1, x_2, x_4) = G_1[x_1]G_2[x_2]G_4[x_4].$$

- To introduce x_3, x_5 , we need to **define the missing cores** $G_3[x_3], G_5[x_5]$.
- Define them as identity matrices:

$$\Theta_\ell(x_1, x_2, x_3, x_4, x_5) = G_1[x_1]G_2[x_2] \underbrace{I}_{\equiv G_3[x_3]} G_4[x_4] \underbrace{I}_{\equiv G_5[x_5]} .$$

- The maximal TT-rank hasn't increased!

The algorithm & its theoretical guarantees

- 1 Compute the TT-decomposition for each individual potential $\Theta_\ell(\mathbf{x}^\ell)$.
- 2 Add the inessential variables to each $\Theta_\ell(\mathbf{x}^\ell)$ to obtain $\Theta_\ell(\mathbf{x})$.
- 3 Use the TT-summation to build $E(\mathbf{x})$: $E(\mathbf{x}) = \sum_{\ell=1}^m \Theta_\ell(\mathbf{x})$.

The algorithm & its theoretical guarantees

- 1 Compute the TT-decomposition for each individual potential $\Theta_\ell(\mathbf{x}^\ell)$.
- 2 Add the inessential variables to each $\Theta_\ell(\mathbf{x}^\ell)$ to obtain $\Theta_\ell(\mathbf{x})$.
- 3 Use the TT-summation to build $\mathbf{E}(\mathbf{x})$: $\mathbf{E}(\mathbf{x}) = \sum_{\ell=1}^m \Theta_\ell(\mathbf{x})$.

Theorem

The maximal TT-rank of the tensor \mathbf{E} constructed by the algorithm is polynomially bounded:

$$r(\mathbf{E}) \leq d^{\frac{p}{2}} m,$$

where

- *d is the number of values that each variable can take;*
- *m is the total number of potentials;*
- *p is the maximal order of a potential (i.e. the maximal $|\mathbf{x}^\ell|$).*

Consider $d = 2$, $p = 2$. Then $r(\mathbf{E}) \leq 2m$ (linear dependence on m).

The TT-format for probability tensor

We can find the TT-representation of $\hat{\mathbf{P}}$

$$\hat{\mathbf{P}} = \underset{\ell=1}{\overset{m}{\odot}} \Psi_{\ell}.$$

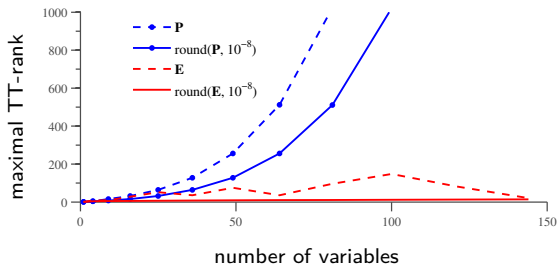
However, the TT-ranks of $\hat{\mathbf{P}}$ are exponential.

The TT-format for probability tensor

We can find the TT-representation of $\hat{\mathbf{P}}$

$$\hat{\mathbf{P}} = \bigodot_{\ell=1}^m \Psi_{\ell}.$$

However, the TT-ranks of $\hat{\mathbf{P}}$ are exponential.



The algorithm motivation

- The TT-ranks of $\hat{\mathbf{P}}$ grow exponentially;
- We must approximate the partition function Z **without explicitly building** the TT- representation of $\hat{\mathbf{P}}$.

The partition function estimation

Recall the definition of the partition function: $Z = \sum_{\mathbf{x}} \widehat{\mathbf{P}}(\mathbf{x})$.

Let us transform $\widehat{\mathbf{P}}(\mathbf{x})$.

Kronecker product property:

$$a \cdot b = a \otimes b, \quad a \in \mathbb{R}, b \in \mathbb{R}.$$

Applying this property we get

$$\widehat{\mathbf{P}}(\mathbf{x}) = \prod_{\ell=1}^m \underbrace{\Psi_{\ell}(\mathbf{x})}_{\in \mathbb{R}} = \bigotimes_{\ell=1}^m \Psi_{\ell}(\mathbf{x}) = \bigotimes_{\ell=1}^m \left(G_1^{\ell}[x_1] \cdots G_n^{\ell}[x_n] \right).$$

The partition function estimation

Recall the definition of the partition function: $Z = \sum_{\mathbf{x}} \widehat{\mathbf{P}}(\mathbf{x})$.

Let us transform $\widehat{\mathbf{P}}(\mathbf{x})$.

Kronecker product property:

$$a \cdot b = a \otimes b, \quad a \in \mathbb{R}, b \in \mathbb{R}.$$

Applying this property we get

$$\widehat{\mathbf{P}}(\mathbf{x}) = \prod_{\ell=1}^m \underbrace{\Psi_{\ell}(\mathbf{x})}_{\in \mathbb{R}} = \bigotimes_{\ell=1}^m \Psi_{\ell}(\mathbf{x}) = \bigotimes_{\ell=1}^m \left(G_1^{\ell}[x_1] \cdots G_n^{\ell}[x_n] \right).$$

Mixed product property:

$$AC \otimes BD = (A \otimes B)(C \otimes D).$$

Hence:

$$\widehat{\mathbf{P}}(\mathbf{x}) = \left(G_1^1[x_1] \otimes \cdots \otimes G_1^m[x_1] \right) \cdots \left(G_n^1[x_n] \otimes \cdots \otimes G_n^m[x_n] \right).$$

The partition function estimation cont'd

$$\hat{\mathbf{P}}(\mathbf{x}) = \left(G_1^1[x_1] \otimes \cdots \otimes G_1^m[x_1] \right) \cdots \left(G_n^1[x_n] \otimes \cdots \otimes G_n^m[x_n] \right).$$

Denote: $A_i[x_i] = G_i^1[x_i] \otimes \cdots \otimes G_i^m[x_i]$.

The partition function estimation cont'd

$$\hat{\mathbf{P}}(\mathbf{x}) = \left(G_1^1[x_1] \otimes \cdots \otimes G_1^m[x_1] \right) \cdots \left(G_n^1[x_n] \otimes \cdots \otimes G_n^m[x_n] \right).$$

Denote: $A_i[x_i] = G_i^1[x_i] \otimes \cdots \otimes G_i^m[x_i]$.

Finally,

$$\begin{aligned} Z &= \sum_{\mathbf{x}} \hat{\mathbf{P}}(\mathbf{x}) = \sum_{x_1, \dots, x_n} A_1[x_1] \cdots A_n[x_n] \\ &= \left(\sum_{x_1} A_1[x_1] \right) \cdots \left(\sum_{x_n} A_n[x_n] \right) = B_1 \cdots B_n, \end{aligned}$$

where $B_i = \sum_{x_i=1}^d A_i[x_i]$.

The algorithm

$$Z = B_1 \cdots B_n,$$

The algorithm:

- 1: Find TT-cores $G_1^\ell, \dots, G_n^\ell$ for Ψ_ℓ
- 2: Initialize $\mathbf{f}_{n+1} := \mathbf{1}$
- 3: **for** $i := n$ **downto** 1 **do**
- 4: Construct TT-matrix $A_i[x_i] := \bigotimes_{\ell=1}^m G_i^\ell[x_i]$
- 5: $B_i := \sum_{x_i=1}^d A_i[x_i]$
- 6: $\mathbf{f}_i := \text{round}(B_i \mathbf{f}_{i+1}, \varepsilon)$
- 7: **end for**
- 8: $\tilde{Z} := \mathbf{f}_1$

Our approach can be generalized to compute the marginal distributions:

$$\begin{aligned}\hat{p}_i(x_i) &= \sum_{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n} \hat{\mathbf{P}}(\mathbf{x}) = \sum_{\mathbf{x} \setminus x_i} A_1[x_1] \dots A_n[x_n] = \\ &= B_1 \dots B_{i-1} A_i[x_i] B_{i+1} \dots B_n,\end{aligned}$$

where

$$B_i = \sum_{x_i=1}^d A_i[x_i].$$

We can explicitly normalize the obtained distribution

$$p_i(x_i) = \frac{\hat{p}_i(x_i)}{\sum_{x=1}^d \hat{p}_i(x)}$$

- 1 Tensor train
- 2 Motivation example
- 3 Results
- 4 Experimental evaluation**
- 5 High Order Potentials
- 6 Conclusion

MAP-inference

The **TT-method** for the MAP-inference:

- 1 Convert the energy into the TT-format;
- 2 Find the minimal element in the energy tensor.

We compare the TT-method with the popular **TRW-S algorithm** on several real-world **image segmentation** problems from the OpenGM database.

Problem	Variables	Labels	TRW-S	TT	Time (sec)
gm6	320	3	45.03	43.11	637
gm29	212	3	56.81	56.21	224
gm66	198	3	75.19	74.92	172
gm105	237	3	67.81	67.71	230
gm32	100	7	150.50	289.29	257
gm70	122	7	121.78	163.60	399
gm85	143	7	168.30	228.40	1 912
gm192	99	7	114.51	174.78	180

Partition function

Ising model:

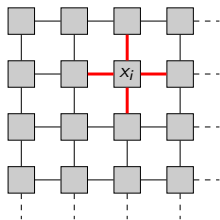
$$\widehat{\mathbf{P}}(\mathbf{x}) = \prod_{i=1}^n \exp\left(-\frac{1}{T} h_i x_i\right) \prod_{(i,j) \in \mathcal{E}} \exp\left(-\frac{1}{T} c_{ij} x_i x_j\right)$$

where $x_i \in \{-1, 1\}$.

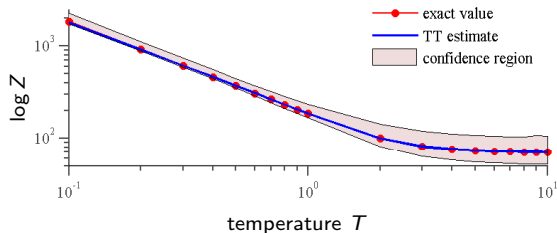
Notation:

- *Temperature* T ;
- *Unary coefficients* h_i ;
- *Pairwise coefficients* c_{ij} .

We use 4-connected grid of size 10×10 .



Theoretical guaranties

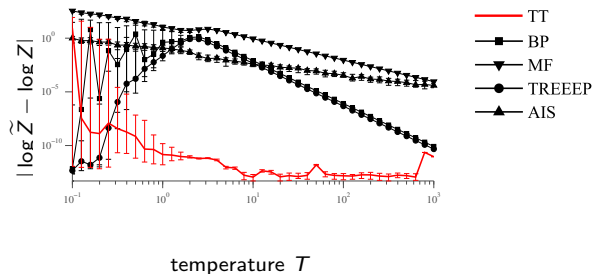


Confidence bounds on the computation of the partition function Z ($c_{ij} = 1$).

Methods we compare with:

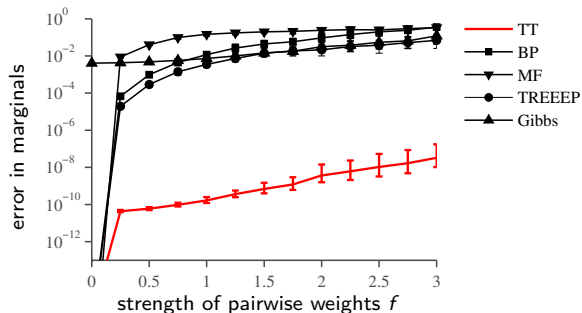
- Belief Propagation (BP, Kschischang, Frey, and Loeliger 2001 Kschischang et al., 2001);
- Tree Expectation Propagation (TREEEP, Minka and Qi 2004 Minka & Qi, 2004);
- Mean Field method (MF, Wainwright and Jordan 2008 Wainwright & Jordan, 2008);
- Annealed importance sampling method (AIS, Neal 2001 Neal, 2001);
- Gibbs sampling (Gibbs, Wainwright and Jordan 2008 Wainwright & Jordan, 2008).

Comparison



Comparison on Ising model (all pairwise weights are equal $c_{ij} = 1$) against state-of-the-art algorithms.

Marginal distributions



Ising models, $T = 1$, $c_{ij} \sim U[-f, f]$.

- 1 Tensor train
- 2 Motivation example
- 3 Results
- 4 Experimental evaluation
- 5 High Order Potentials**
- 6 Conclusion

High Order Potentials

- Sometimes it is convenient to use **potentials of high order**, i.e. those which depend on many variables. E.g., the potential

$$\Theta_{\ell}(\mathbf{x}) = \underbrace{\left[\sum_{i=1}^n x_i \leq a \right]}_{\text{indicator function}},$$

which depends on all the variables, could be used to specify some preference on the minimal value of the area of foreground in the problem of segmenting an image into background/foreground.

- We can't use the TT-SVD algorithm any more to convert such potentials into the TT-format!
- However, for some of these potentials we can explicitly construct the TT-representation, i.e. we can derive **analytical formulae** for the corresponding TT-cores.
- Such TT-representations will be of low TT-rank!

Sparse Potential

- Consider a so-called sparse potential:

$$\Theta_\ell(x_{i_1}, \dots, x_{i_w}) = [x_{i_1} = \beta_1] \dots [x_{i_w} = \beta_w].$$

It always equals zero with the exception of only one configuration.

- Such a potential admits a TT-representation

$$\Theta_\ell(x_{i_1}, \dots, x_{i_w}) = G_{i_1}[x_{i_1}] \dots G_{i_w}[x_{i_w}]$$

with the following TT-cores:

$$G_{i_v}[x_{i_v}] = [x_{i_v} = \beta_v], \quad v = 1, \dots, w.$$

- In this case each TT-core is simply a number (1-by-1 matrix) for every concrete value of x_{i_v} . Hence, the **maximal TT-rank equals 1**.
- A more general sparse potential which differs from zero on $s > 1$ configurations can be obtained as a sum of several potentials of the above type. Thus, the TT-rank of a general sparse potential is bounded above by s .

- Consider the potential

$$\Theta_\ell(\mathbf{x}) = \left[\sum_{i=1}^n x_i \leq a \right],$$

where $x_i \in \{0, 1\}$ and $a \in \mathbb{Z}_+$.

- This potential can be analytically represented in the TT-format with **the maximal TT-rank equal to $a + 1$** :

$$G_i[x_i] = (S_a)^{x_i}, \quad (i = 2, \dots, n-1),$$
$$G_1[x_1] = \underbrace{[0 \dots 0 \ 1 \dots 1]}_{x_1}, \quad G_n[x_n] = (S_a)^{x_n} \underbrace{[0 \dots 0 \ 1]}_a^T,$$

where $S_a = \underbrace{\begin{bmatrix} O & I_a \\ O & O \end{bmatrix}}_{(a+1) \times (a+1)}.$

Area Potential cont'd

Key property of S_a : $\underbrace{[0 \dots 0 \ 1 \dots 1]}_k S_a = \underbrace{[0 \dots 0 \ 1 \dots 1]}_{k+1}$.

Consider, e.g., that $a = 3$. In this case

$$S_a = \left[\begin{array}{c|cccc} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ \hline 0 & 0 & 0 & 0 \end{array} \right].$$

- $[1111]S_a = [0111]$ (the sum of all rows);
- $[0111]S_a = [0011]$ (the sum of rows 2, 3, 4);
- and so on.

Then

$$\begin{aligned} \Theta_\ell(\mathbf{x}) &= G_1[x_1]G_2[x_2]G_3[x_3] \dots G_n[x_n] \\ &= \underbrace{[0 \dots 0 \ 1 \dots 1]}_{x_1} (S_a)^{x_2} (S_a)^{x_3} \dots (S_a)^{x_n} \underbrace{[0 \dots 0 \ 1]}^T \end{aligned}$$

- 1 Tensor train
- 2 Motivation example
- 3 Results
- 4 Experimental evaluation
- 5 High Order Potentials
- 6 Conclusion**

Our main contributions are:

- We have proposed an algorithm that converts MRF energy into the TT-format exactly.
- We have derived an upper bound on the TT-ranks of the energy tensor constructed by the proposed algorithm.
- We have proposed an algorithm for estimating the partition function and the marginal distributions. The key feature of the algorithm is that it does not explicitly construct the TT-representation of the unnormalized joint distribution.
- We have bounded the errors of the partition function estimation.

- Analytical formulae of TT-representations for other types of high order potentials.
- Better theoretical guaranties for the partition function estimation.
- Better algorithm for finding the minimal element in a tensor represented in the TT-format.

- Kschischang, F. R., B. J. Frey, and H.-A. Loeliger (2001). “Factor Graphs and the Sum-Product Algorithm”. In: *IEEE Transactions on Information Theory* 47.2, pp. 498–519.
- Minka, T. and Y. Qi (2004). “Tree-structured Approximations by Expectation Propagation”. In: *Advances in Neural Information Processing Systems 16 (NIPS)*, pp. 193–200.
- Neal, R. (2001). “Annealed importance sampling”. In: *Statistics and Computing* 11, pp. 125–139.
- Nowozin, S. and C. Lampert (2010). “Structured Learning and Prediction in Computer Vision”. In: *Foundations and Trends in Computer Graphics and Vision* 6.3–4, pp. 185–365.
- Oseledets, I. V. (2011). “Tensor-Train Decomposition”. In: *SIAM J. Scientific Computing* 33.5, pp. 2295–2317.

Wainwright, M. J. and M. I. Jordan (2008). “Graphical models, exponential families, and variational inference”. In: *Foundations and Trends in Machine Learning* 1.1–2, pp. 1–305.