

---

# Matrix Product Operator Restricted Boltzmann Machines

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 A restricted Boltzmann machine (RBM) learns a probabilistic distribution over its  
2 input samples and has numerous uses like dimensionality reduction, classification  
3 and generative modeling. Conventional RBMs accept vectorized data that dismisses  
4 potentially important structural information in the original tensor (multi-way) input.  
5 Matrix-variate and tensor-variate RBMs, named MvRBM and TvRBM, have been  
6 proposed but are all restrictive by construction. This work presents the matrix  
7 product operator RBM (MPORBM) that utilizes a tensor network generalization  
8 of Mv/TvRBM, preserves input formats in both the visible and hidden layers,  
9 and results in higher expressive power. A novel training algorithm integrating  
10 contrastive divergence and an alternating optimization procedure is also developed.

## 11 1 Introduction

12 A restricted Boltzmann machine (RBM) [1] is a probabilistic model that employs a layer of hidden  
13 variables to achieve highly expressive marginal distributions. RBMs are an unsupervised learn-  
14 ing technique and have been extensively explored and applied in various fields [2–4]. However,  
15 conventional RBMs are designed for vector data and cannot directly deal with matrices and higher-  
16 dimensional data, which are common in real-life. The traditional approach to apply an RBM on  
17 high-dimensional data is through vectorization of the data which leads to two drawbacks. First, the  
18 spatial information in the original data is lost, thus weakening the model’s capability to represent these  
19 structural correlations. Second, the fully connected visible and hidden units may cause overfitting  
20 since the intrinsic low-rank property of many real-life data is disregarded.

21 Researchers have been motivated to develop corresponding multiway RBMs [5, 3]. However, those  
22 works are both aiming to capture the interaction among different vector inputs and are hence not  
23 directly applicable to matrix and tensor data. The first RBM designed for tensor inputs is described in  
24 [6], where the visible layer is represented as a tensor but the hidden layer is still a vector. Furthermore,  
25 the connection between the visible and hidden layers is described by a canonical polyadic (CP) tensor  
26 decomposition [7], which constrains the model representation capability [8]. Another RBM related  
27 model that utilizes tensor input is the matrix-variate RBM (MvRBM) [8]. The visible and hidden  
28 layers in an MvRBM are both matrices. Nonetheless, to limit the number of parameters, an MvRBM  
29 models the connection between the visible and hidden layers through two separate matrices, which  
30 restricts the ability of the model to capture correlations between different data modes.

31 All these issues have motivated this work. Specifically, we propose a matrix product operator (MPO)  
32 restricted Boltzmann machine (MPORBM) where both the visible and hidden layers are in tensor  
33 forms. Moreover, MPORBM utilizes a general and powerful tensor network, namely an MPO, to  
34 connect the tensorial visible and hidden layers. By doing so, an MPORBM achieves a more powerful  
35 model representation capacity than MvRBM and at the same time greatly reduces the number of  
36 model parameters compared to a standard RBM.

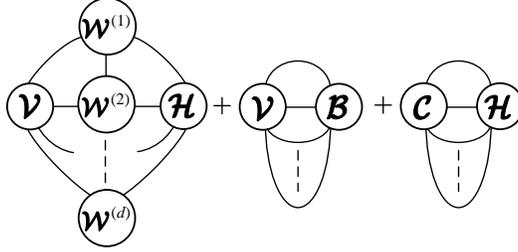


Figure 1: Negative energy functions ( $-E$ ) of the MPORBM.

## 37 2 Method

38 In an MPORBM, both the visible layer  $\mathcal{V} \in \mathbb{R}^{I_1 \times \dots \times I_d}$  and the hidden layer  $\mathcal{H} \in \mathbb{R}^{J_1 \times \dots \times J_d}$  are  
 39  $d$ -way tensors. As a result, the weight matrix  $\mathcal{W}$  is now a  $2d$ -way tensor  $\mathcal{W} \in \mathbb{R}^{I_1 \times \dots \times I_d \times J_1 \times \dots \times J_d}$ ,  
 40 which is represented by an MPO instead in order to lift the curse of dimensionality. Per definition,  
 41 the corresponding MPO decomposition represents each entry of  $\mathcal{W}$  as

$$\mathcal{W}(i_1, \dots, i_d, j_1, \dots, j_d) = \sum_{r_1, r_2, \dots, r_d}^{R_1, R_2, \dots, R_d} \mathcal{W}^{(1)}(r_1, i_1, j_1, r_2) \cdots \mathcal{W}^{(d)}(r_d, i_d, j_d, r_1). \quad (1)$$

42 The “building blocks” of the MPO are the 4-way tensors  $\mathcal{W}^{(1)}, \dots, \mathcal{W}^{(d)}$ , also called the MPO-cores.  
 43 The dimensions  $R_1, \dots, R_d$  of the summation indices  $r_1, \dots, r_d$  are called the MPO-ranks. With both  
 44 the visible and hidden layers being tensors, it is therefore also required that the bias vectors  $\mathbf{b}, \mathbf{c}$  are  
 45 tensors  $\mathcal{B} \in \mathbb{R}^{I_1 \times \dots \times I_d}$ ,  $\mathcal{C} \in \mathbb{R}^{J_1 \times \dots \times J_d}$ , respectively. A tensor network diagram representation of  
 46 the negative energy function of the MPORBM is shown in Figure 1, where each tensor is represented by  
 47 by a node in the network and the number of edges connected to a node represents the order of the  
 48 corresponding tensor. The vertical edges between the different MPO-cores  $\mathcal{W}^{(1)}, \dots, \mathcal{W}^{(d)}$  represent  
 49 the summations in (1) and are the key ingredients in being able to express generic weight tensors  $\mathcal{W}$ .  
 50 The storage complexity of an MPORBM with uniform ranks and dimensions is  $O(dIJR^2)$ , which is  
 51 linear on the order  $d$  and therefore removes the curse of dimensionality. The MvRBM model can be  
 52 interpreted as a very specific case of an MPORBM where there are only 2 MPO-cores without any  
 53 vertical edge, which limits the expressive power. The corresponding conditional distribution over the  
 54 hidden or visible layer involves the summation of the weight MPO with either the hidden or visible  
 55 layer tensors into a  $d$ -way tensor, which is then added elementwise with the corresponding bias tensor.  
 56 The final step in the computation of the conditional probability is an elementwise application of the  
 57 logistic sigmoid function on the resulting tensor.

58 Let  $\Theta = \{\mathcal{B}, \mathcal{C}, \mathcal{W}^{(1)}, \mathcal{W}^{(2)}, \dots, \mathcal{W}^{(d)}\}$  denote the model parameters. The model learning task is  
 59 then formulated into maximizing the training data likelihood:

$$\mathcal{L}(\mathcal{V}; \Theta) = p(\mathcal{V}; \Theta) = \sum_{\mathcal{H}} p(\mathcal{V}, \mathcal{H}; \Theta) \quad (2)$$

60 with respect to the model parameter  $\Theta$ . Similar to the standard RBM [1], the expression of the  
 61 gradient of the log-likelihood is:

$$\frac{\partial}{\partial \Theta} \log \mathcal{L}(\mathcal{V}; \Theta) = -\mathbb{E}_{\mathcal{H}|\mathcal{V}} \left[ \frac{\partial E(\mathcal{V}, \mathcal{H})}{\partial \Theta} \right] + \mathbb{E}_{\mathcal{V}, \mathcal{H}} \left[ \frac{\partial E(\mathcal{V}, \mathcal{H})}{\partial \Theta} \right] \quad (3)$$

62 We mainly use the contrastive divergence (CD) procedure to train the MPORBM model. First, a  
 63 Gibbs chain is initialized with one particular training sample  $\mathcal{V}_{(0)} = \mathcal{X}_{train}$ , followed by  $K$  times  
 64 Gibbs sampling which results in the chain  $\{(\mathcal{V}_{(0)}, \mathcal{H}_{(0)}), (\mathcal{V}_{(1)}, \mathcal{H}_{(1)}), \dots, (\mathcal{V}_{(K)}, \mathcal{H}_{(K)})\}$ . The  
 65 model expectation is then approximated by  $\{\mathcal{V}_{(K)}\}$ . The derivative of  $\log \mathcal{L}(\mathcal{V}; \Theta)$  with respect to  
 66 the  $k$ -th MPO-core  $\mathcal{W}^{(k)}$  can be computed by removing  $\mathcal{W}^{(k)}$  from two tensor network diagrams  
 67 (one diagram with  $\mathcal{V}_{(0)}, \mathcal{H}_{(0)}$  and one with  $\mathcal{V}_{(K)}, \mathcal{H}_{(K)}$ ), taking the elementwise difference and  
 68 summing over all edges. The derivatives of the log-likelihood with respect to the bias tensors  $\mathcal{B}, \mathcal{C}$  are

$$\frac{\partial}{\partial \mathcal{B}} \log \mathcal{L}(\mathcal{V}; \Theta) = \mathcal{V}_{(0)} - \mathcal{V}_{(K)}, \quad \frac{\partial}{\partial \mathcal{C}} \log \mathcal{L}(\mathcal{V}; \Theta) = \mathcal{H}_{(0)} - \mathcal{H}_{(K)}.$$

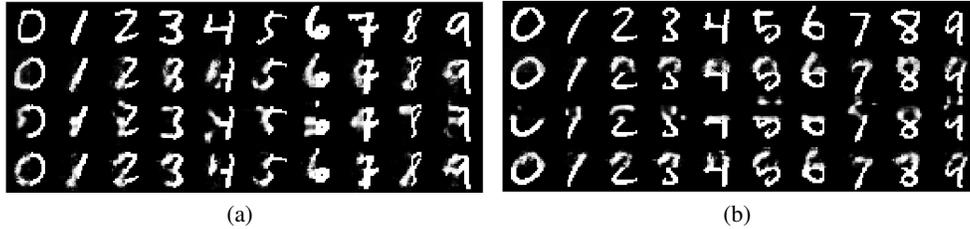


Figure 2: Image completion results when given only the (a) right half; and (b) bottom half. Top row: original binarized images; 2nd row: RBM completion; 3rd row: MvRBM completion; 4th row: MPORBM completion.

69 Instead of updating all MPO-cores simultaneously with one batch of input training data, we employ  
 70 the alternating optimization procedure (AOP). This involves updating only one MPO-core at a time  
 71 while keeping the others unchanged using the same batch of input training data. We name this  
 72 parameter learning algorithm CD-AOP. The superiority of AOP over simultaneously updating all  
 73 MPO-cores, which we call CD-SU henceforth, will be demonstrated through numerical experiments.

### 74 3 Experiments

75 In the first experiment, we demonstrate the superior data classification accuracy of MPORBM using  
 76 standard datasets, namely, the Binary Alphadigits, normalized DrivFace, Arcene and COIL-100  
 77 datasets. The vectorized sample sizes of these datasets vary from 320 to 49152. We assume a binary  
 78 input in our RBM setting, so for non-binary datasets a multi-bit vector is used to represent each value  
 79 in the original data. The trained RBM models were employed to extract features from the hidden  
 80 layer. These features were then utilized to train a  $K$  Nearest Neighbor ( $K$ -NN) classifier with  $K = 1$   
 81 for all experiments. Table 1 lists the resulting classification errors. The restrictive expressive power  
 82 of the weight matrix in MvRBM explains why it has the worst classification performance for all  
 83 datasets. The worse performance of the standard RBM may be caused by overfitting due to the small  
 84 training sample size. For COIL-100 dataset, the standard RBM fails to learn the large number of  
 85 parameters in the full weight matrix due to out-of-memory errors. We need to mention that CD-AOP  
 86 algorithm achieves a surprisingly 0% test error because of the small test sample number. Moreover,  
 87 the CD-AOP algorithm shows a higher classification accuracy than CD-SU, which indicates that the  
 88 alternating updating scheme is more suitable for the proposed MPORBM model.

Table 1: Classification errors of different RBM models.

Datasets	RBM	MvRBM	MPORBM CD-SU	MPORBM CD-AOP
Alphadigits	28.10%	31.20%	28.10%	26.90%
DrivFace	24.20%	15.48%	9.68%	8.06%
Arcene	45.00%	34.00%	32.00%	27.00%
COIL-100	—	6.82%	6.82%	0.00%

89 Finally, we show that an MPORBM is good at generative modeling exemplified by image completion.  
 90 We tested this generative task on the binarized MNIST dataset: one half of the image was provided to  
 91 the trained RBM models to complete the other half. Figure 2 shows the completed images of different  
 92 RBM models when given the same randomly selected right and bottom halves, respectively. It is  
 93 clear that MvRBM is not able to complete the image well, which further confirms the efficacy of the  
 94 MPO generalization.

### 95 4 Conclusion

96 The MPORBM is proposed, which preserves the tensorial nature of the input data and utilizes a matrix  
 97 product operator (MPO) to represent the weight matrix. The MPORBM generalizes all existing RBM  
 98 models to tensor inputs and has better storage complexity since the number of parameters grows only  
 99 linearly with the order of the tensor. Experiments have verified the superiority of MPORBM over  
 100 traditional counterparts.

101 **References**

- 102 [1] Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural computation*,  
103 14(8):1771–1800, 2002.
- 104 [2] Hugo Larochelle and Yoshua Bengio. Classification using discriminative restricted Boltzmann machines. In  
105 *Proceedings of the 25th international conference on Machine learning*, pages 536–543. ACM, 2008.
- 106 [3] Alex Krizhevsky, Geoffrey Hinton, et al. Factored 3-way restricted Boltzmann machines for modeling  
107 natural images. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and*  
108 *Statistics*, pages 621–628, 2010.
- 109 [4] Abdel-rahman Mohamed and Geoffrey Hinton. Phone recognition using restricted Boltzmann machines.  
110 In *Acoustics Speech and Signal Processing (ICASSP), 2010 IEEE International Conference on*, pages  
111 4354–4357. IEEE, 2010.
- 112 [5] Graham W Taylor and Geoffrey E Hinton. Factored conditional restricted Boltzmann machines for modeling  
113 motion style. In *Proceedings of the 26th annual international conference on machine learning*, pages  
114 1025–1032. ACM, 2009.
- 115 [6] Tu Dinh Nguyen, Truyen Tran, Dinh Q Phung, et al. Tensor-Variate Restricted Boltzmann Machines. In  
116 *AAAI*, pages 2887–2893, 2015.
- 117 [7] Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. *SIAM review*, 51(3):455–500,  
118 2009.
- 119 [8] Guanglei Qi, Yanfeng Sun, Junbin Gao, et al. Matrix variate restricted Boltzmann machine. In *Neural*  
120 *Networks (IJCNN), 2016 International Joint Conference on*, pages 389–395. IEEE, 2016.