

VEST: Very Sparse Tucker Factorization of Large-Scale Tensors

MoonJeong Park¹, Jun-Gi Jang², and Lee Sael(✉)²

¹ Daegu Gyeongbuk Institute of Science and Technology, Daegu, Korea

² Department of Computer Science and Engineering, Seoul National University, Seoul, Korea
saelllee@snu.ac.kr

Abstract. Given a large tensor, how can we decompose it to sparse core tensor and factor matrices such that it is easier to interpret the results? How can we do this without reducing the accuracy? Existing approaches either output dense results or give low accuracy. In this paper, we propose VEST, a tensor factorization method for partially observable data to output a very sparse core tensor and factor matrices. VEST performs initial decomposition, determines unimportant entries in the decomposition results, removes the unimportant entries, and updates the remaining entries. To determine unimportant entries, we define and use entry-wise ‘responsibility’ for the decomposed results. The entries are updated iteratively using a carefully derived coordinate descent rule in parallel for scalable computation. VEST also includes an auto-search algorithm to give a good trade-off between sparsity and accuracy. Extensive experiments show that our method VEST is at least 2.2 times sparser and at least 2.8 times more accurate compared to competitors. Moreover, VEST is scalable in terms of dimensionality, number of observable entries, and number of threads. Thanks to VEST, we successfully interpret the decomposition result of real-world tensor data based on the sparsity pattern of the factor matrices.

Keywords: Scalable tensor factorization · Tucker · Interpretability · Sparsity

1 Introduction

How can we factorize a large tensor to sparse core tensor and factor matrices such that outputs are easier to interpret? How can we do this without sacrificing accuracy? A tensor is a powerful tool for representing multi-modal data. Tensor factorization outputs a core tensor and factor matrices which reveal the latent relation of the data. Tensor factorization can also be viewed as a tool for multi-linear regression problem where only the target values, i.e., values of input tensor entries, are known. In this perspective, the columns of factor matrices act as latent components, their values as latent feature values, and the cells of the core tensor as weights of the relations between the latent components [8]. Sparse tensor factorization aims to output sparse core tensor and factor matrices. As sparse linear regression model enhances its interpretability [19], sparse factor matrices and a core improve interpretability.

There are two widely used approaches for sparse factorization. The first approach adds an L_1 norm as sparse regularizer to the factorization objective function [14,12].

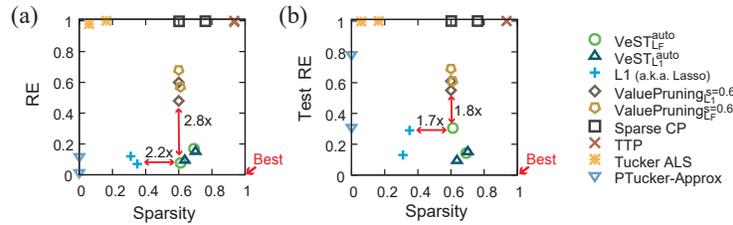


Fig. 1: Sparsity and accuracy of VEST and competitors on Yelp-s and AmazonFood-s. VEST generates sparse and accurate results that generalize well on unseen data: points are located in the bottom right areas of both RE (reconstruction error) and Test RE plots. However, the sparsity is sensitive to lambda values. The second approach removes elements with small values from the core tensor or the factor matrices [21,2,22]. However, removing such elements does not necessarily leads to small reconstruction error, and thus value-based pruning sacrifices accuracy.

Other application-specific approaches include utilizing domain-specific knowledge as sparsity constraints [11], constructing factor matrix from sparse input sampling [10,13], using smoothing matrices [17,7], and learning sparse dictionary for image data [18,23,5]. These methods are based on several strong assumptions. The first approach assumes that prior classification of each mode, e.g., gene sets in omics data, are known. The second approach assumes that input tensors are already sparse and interpretable, e.g., network data. The third and fourth approaches assume values in input tensor are smooth, and unsmoothing them does not affect the overall meaning of the data, e.g., image data. However, these strong assumptions limit their use in general tensor factorizations.

In this paper, we propose VEST (VERY Sparse Tucker factorization), a scalable and accurate Tucker factorization method to generate sparse factors and a core tensor for large-scale partially observable input tensor. VEST outputs very sparse factorization results by carefully determining the importance of elements of factors and the core, and pruning unimportant ones. VEST guarantees that the sparsity non-decreases in the update process by carefully derived update rules. Often, increasing sparsity too much degrades accuracy. VEST gives an algorithm to automatically determine a reasonable sparsity which offers a good balance with regard to accuracy. The very sparse result of VEST helps interpreting the result of Tucker factorization and easily revealing the relations of dimensions in multilinear regression.

Our main contributions are as follows:

- **Algorithm.** We propose VEST, an efficient Tucker factorization method for partially observable data, which produces very sparse outputs for better interpretability without loss of accuracy. VEST also provides an algorithm to automatically determine a sparsity which gives a good trade-off with regard to accuracy.
- **Theory.** We carefully derive parallelizable coordinate decent update rules for core tensor and factor matrices, and prove their correctness. We also analyze the time and space complexity of VEST.
- **Performance.** VEST provides better sparsity, accuracy, and scalability compared to other methods (see Figure 1).

The codes and datasets used in this paper are available at <http://github.com/leesael/VeST>.

2 Preliminaries and Related Works

We introduce concepts of tensor and its operations, Tucker factorization, and the standard algorithm for Tucker. Table 1 lists the symbols used.

Table 1: Table of symbols and definitions.

\mathcal{X}	input tensor ($\in \mathbb{R}^{I_1 \times \dots \times I_N}$)	\mathcal{G}	core tensor ($\in \mathbb{R}^{J_1 \times \dots \times J_N}$)
N	order of \mathcal{X}	I_n, J_n	dimensionality of the n th mode of \mathcal{X} and \mathcal{G} , respectively
$\mathbf{A}^{(n)}$	n th factor matrix ($\in \mathbb{R}^{I_n \times J_n}$)	$a_{i_n j_n}^{(n)}$	(i_n, j_n) th element of $\mathbf{A}^{(n)}$
Ω	set of observable entries of \mathcal{X}	$ \Omega $	number of observable entries of \mathcal{X}
$\Omega_{i_n}^{(n)}$	set of observable entries whose n th mode index is i_n	λ	regularization parameter for core and factor matrices
$\ \mathcal{X}\ _F$	Frobenius norm of tensor \mathcal{X}	$\ \mathcal{X}\ _1$	sum of absolute values of tensor \mathcal{X}
α	an entry (i_1, \dots, i_N) of input tensor \mathcal{X}	β	an element (j_1, \dots, j_N) of core tensor \mathcal{G}
$\alpha_{i_n=i}$	an entry $(i_1, \dots, i_n = i, \dots, i_N)$ of input tensor \mathcal{X}	$\beta_{j_n=j}$	an element $(j_1, \dots, j_n = j, \dots, j_N)$ of core tensor \mathcal{G}

2.1 Tensor and its Operations

Tensor is multi-dimensional array that contains numbers. An ‘order’ or ‘mode’ is the number of tensor dimensions, where a 1st-order tensor represents a vector and a 2nd-order tensor represents a matrix. We denote vectors by boldface lowercase letters (e.g., \mathbf{a}), matrices by boldface capital letters (e.g., \mathbf{A}), and three or higher order tensors by boldface Euler script letters (e.g., \mathcal{X}). An entry of a 3rd-order tensor can be expressed with three indices. For example, the (i_1, i_2, i_3) th entry of a 3rd-order tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ is denoted by $x_{i_1 i_2 i_3}$, where index i_n spans from 1 to I_n .

The size of a tensor is often evaluated by the Frobenius norm. Given an N -order tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$, the *Frobenius norm* of \mathcal{X} is $\|\mathcal{X}\|_F = \sqrt{\sum_{\forall \alpha \in \mathcal{X}} \mathcal{X}_\alpha^2}$, where $\alpha = (i_1, \dots, i_N)$ is an index to an entry of input tensor \mathcal{X} . Tensor decomposition often involves matricization of a tensor, and product between a tensor and a matrix. The *mode- n matricization* of a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ is denoted as $\mathbf{X}_{(n)}$ and the mapping from an entry (i_1, \dots, i_N) of \mathcal{X} to an entry (i_n, j) of $\mathbf{X}_{(n)}$ is given by $j = 1 + \sum_{k=1, k \neq n}^N [(i_k - 1) \prod_{m=1, m \neq n}^{k-1} I_m]$. Also, the *n -mode product* of a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ with a matrix $\mathbf{U} \in \mathbb{R}^{J \times I_n}$ is denoted by $\mathcal{X} \times_n \mathbf{U} \in \mathbb{R}^{I_1 \times \dots \times I_{n-1} \times J \times I_{n+1} \times \dots \times I_N}$. Entry-wise, we have $(\mathcal{X} \times_n \mathbf{U})_{i_1 \dots i_{n-1} j i_{n+1} \dots i_N} = \sum_{i_n=1}^{I_n} (\mathcal{X}_{\alpha_{i_n=i}} u_{ji})$.

2.2 Tucker Factorization

Our proposed method VEST is built on top of Tucker factorization, one of the most popular tensor factorization methods. Given an N th-order tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$, Tucker factorization approximates \mathcal{X} by a core tensor $\mathcal{G} \in \mathbb{R}^{J_1 \times \dots \times J_N}$ and factor matrices $\{\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times J_n} | n = 1 \dots N\}$ by minimizing the full reconstruction error: $\min_{\mathcal{G}, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}} \|\mathcal{X} - \mathcal{G} \times_1 \mathbf{A}^{(1)} \dots \times_N \mathbf{A}^{(N)}\|_F$.

Figure 2 illustrates a Tucker factorization result for a 3rd-order tensor. Typically, a core tensor \mathcal{G} is assumed to be smaller and denser than the input tensor \mathcal{X} . Each factor

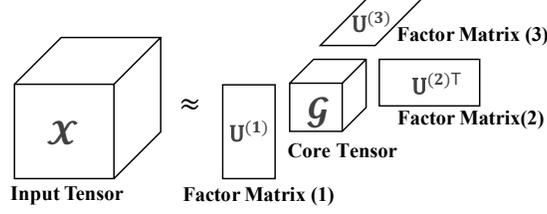


Fig. 2: Tucker factorization of a 3-way tensor.

matrix $\mathbf{A}^{(n)}$ represents the latent features of the object related to the n th mode of \mathcal{X} , and each element of a core tensor \mathcal{G} indicates the weights of the relations composed of columns of factor matrices.

However, in real-world, data are often incomplete with some missing entries. To accommodate for the missing data, a partially observable Tucker factorization is needed. Given a tensor $\mathcal{X} (\in \mathbb{R}^{I_1 \times \dots \times I_N})$ with observable entries Ω , the goal of *partially observable Tucker factorization* of \mathcal{X} is to find factor matrices $\mathbf{A}^{(n)} (\in \mathbb{R}^{I_n \times J_n}, n = 1, \dots, N)$ and a core tensor $\mathcal{G} (\in \mathbb{R}^{J_1 \times \dots \times J_N})$, which minimize the following loss:

$$L_F(\mathcal{G}, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}) = \sum_{\forall \alpha \in \Omega} \left(\mathbf{x}_\alpha - \sum_{\forall \beta \in \mathcal{G}} \mathcal{G}_\beta \prod_{n=1}^N a_{i_n j_n}^{(n)} \right)^2 + \lambda (\|\mathcal{G}\|_F^2 + \sum_{n=1}^N \|\mathbf{A}^{(n)}\|_F^2), \quad (1)$$

where α is an observable entry (i_1, \dots, i_N) of input tensor \mathcal{X} , β is an element (j_1, \dots, j_N) of core tensor \mathcal{G} , and $\lambda > 0$ is a regularization parameter. Note that the reconstruction error in Eq. (1) depends only on the observable entries of \mathcal{X} , and L_F regularization is used in Eq. (1) to prevent overfitting.

Tucker factorization often results in dense core and factor matrices. One of the approaches for sparsifying results is by including a sparsity constraint in the form of L_1 norm, a.k.a., Lasso, into the objective function. Given a tensor \mathcal{X} with observable entries Ω , the goal of *partially observable Tucker factorization via sparse regularizer* of \mathcal{X} is to find factor matrices and a core tensor that minimize the following loss:

$$L_1(\mathcal{G}, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}) = \sum_{\forall \alpha \in \Omega} \left(\mathbf{x}_\alpha - \sum_{\forall \beta \in \mathcal{G}} \mathcal{G}_\beta \prod_{n=1}^N a_{i_n j_n}^{(n)} \right)^2 + \lambda (\|\mathcal{G}\|_1 + \sum_{n=1}^N \|\mathbf{A}^{(n)}\|_1), \quad (2)$$

which changed the L_2 regularization term of Eq. (1) to L_1 . Again the reconstruction error in Eq. (2) depends only on the observable entries of \mathcal{X} , and L_1 regularization is used to enforce sparsity. Another approach for sparsifying results is by pruning. That is, a *partially observable Tucker factorization via minimal element value pruning* of \mathcal{X} is to optimize on either Eq. (1) or Eq. (2), and sets the smallest s ratio of the elements to zero in the core and factor matrices.

Evaluation of tensor decomposition and the prediction of the missing entry values (a.k.a., tensor completion) involves reconstruction. Given core tensor \mathcal{G} and factor matrices $\mathbf{A}^{(n)}$, the *reconstruction* of the original tensor \mathcal{X} is defined as $\mathcal{X} \approx \mathcal{G} \times_1 \mathbf{A}^{(1)} \dots \times_N \mathbf{A}^{(N)}$.

2.3 Tucker ALS Algorithm

A widely used technique for minimizing the loss functions Eq. (1) and Eq. (2) in a standard tensor factorization is alternating least squares (ALS) [8], which updates a factor matrix or a core tensor while keeping all others fixed.

Algorithm 1: Tucker-ALS for Fully Observable Tensors (HOOI)

Input : Tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$, and core tensor dimensionality J_1, \dots, J_N .
Output: Factor matrices $\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times J_n}$ ($n = 1, \dots, N$), and core tensor $\mathcal{G} \in \mathbb{R}^{J_1 \times J_2 \times \dots \times J_N}$.

- 1 initialize all factor matrices $\mathbf{A}^{(n)}$
- 2 **repeat**
- 3 **for** $n = 1 \dots N$ **do**
- 4 $\mathcal{Y} \leftarrow \mathcal{X} \times_1 \mathbf{A}^{(1)\top} \dots \times_{n-1} \mathbf{A}^{(n-1)\top} \times_{n+1} \mathbf{A}^{(n+1)\top} \dots \times_N \mathbf{A}^{(N)\top}$
- 5 $\mathbf{A}^{(n)} \leftarrow J_n$ leading left singular vectors of $\mathcal{Y}_{(n)}$
- 6 **until** reconstruction error converges or exceeds maximum iteration;
- 7 $\mathcal{G} \leftarrow \mathcal{X} \times_1 \mathbf{A}^{(1)\top} \dots \times_N \mathbf{A}^{(N)\top}$

Algorithm 1 describes a vanilla Tucker factorization algorithm based on ALS, which is called the *higher-order orthogonal iteration* (HOOI) (see [8] for details) that works on fully observable tensor. Notice that Algorithm 1 assumes missing entries of \mathcal{X} as zeros during the update process (lines 4-5). However, setting missing values to zero enforces Tucker ALS to factorize the original tensor such that missing values become zero when reconstructed. Note that the missing values are often nonzero values that are unknown. Thus, setting missing values to zero inserts false information into the factorization which results in higher reconstruction error as well as higher generalization error. Moreover, Algorithm 1 computes SVD (singular vector decomposition) given $\mathcal{Y}_{(n)}$, which often results in dense matrices, thus tensor-ALS results in overall dense core tensor and factor matrices. Also, Algorithm 1 requires storing a full-dense matrix $\mathcal{Y}_{(n)}$, and the amount of memory needed for storing $\mathcal{Y}_{(n)}$ is $O(I_n \prod_{m \neq n} J_m)$. The required memory grows rapidly when the order, the mode dimensionality, or the rank of a tensor increase, and ultimately causes *intermediate data explosion* [6].

In summary, the vanilla Tucker-ALS algorithm results in high generalization error in the presence of missing data, results in dense and thus hard-to-interpret core tensor and factor matrices, and cannot be applied to large data. Therefore, Algorithm 1 needs to be revised to focus only on observed entries, make sparse outputs, and be scaled for large-scale tensors at the same time.

3 Proposed Method

In this section, we propose VEST (Very Sparse Tucker factorization), a method for partially observed large scale tensor, that results in very sparse core tensor and factor matrices. Sparse results of tensor factorization increase interpretability and provide a scheme for better compression. To maximize sparsity without losing accuracy, VEST iteratively updates core tensor and factor matrices, and prunes unimportant elements from the core tensor and factor matrices. However, there are several challenges in designing an efficient update and pruning rules.

- **Evaluating importance of elements.** Vital elements of core tensor and factor matrices should not be pruned. How can we evaluate their importance?
- **Automatically determining the sparsity.** There is a trade-off relationship between sparsity and accuracy. How can we automatically determine an appropriate sparsity which gives a good balance with regards to accuracy?
- **Updating factors while guaranteeing non-decreasing sparsity.** The update process of factors and the core in the regular Tucker-ALS does not guarantee that the sparsity improves over the update process. How can we guarantee that update rules improve the sparsity?

We have the following main ideas to address the above challenges which we describe in detail in later subsections.

- **Design responsibility indicator** to evaluate contribution of each element on the accuracy (Section 3.2).
- **Design auto-search algorithm VEST_{*}^{auto}** to find a good sparsity that resides near the maximum sparsity just before the reconstruction error shoots up (Section 3.3).
- **Design element-wise update rules** to independently update each element of factor matrices and the core tensor. Element-wise update rules guarantee that the sparsity non-decreases by keeping pruned elements to zeros (Sections 3.4 and 3.5).

3.1 Overview

Algorithm 2: VEST: Very Sparse Tucker Factorization

Input : Tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$, core tensor dimensionality J_1, \dots, J_N , and target sparsity s (if manual-mode).

Output: Sparse factor matrices $A^{(n)} \in \mathbb{R}^{I_n \times J_n}$ ($n = 1, \dots, N$) and sparse core tensor $\mathcal{G} \in \mathbb{R}^{J_1 \times J_2 \times \dots \times J_N}$.

- 1 randomly initialize $\mathbf{A}^{(n)}$ ($n = 1, \dots, N$) and \mathcal{G} ; set $pr = \text{INIT_PR}$, $iterN = 0$.
- 2 **repeat**
- 3 update unpruned elements of $A^{(n)}$ ($n = 1, \dots, N$) ▷ Algorithm 4
- 4 update unpruned elements of \mathcal{G} ▷ Eq.(12) or (14)
- 5 compute RE using observable entries Ω ▷ Eq.(4)
- 6 **if** *should_prune()* **then**
- 7 prune $pr = \min(\text{INIT_PR} * iterN, \text{MAX_PR})$ ratio of elements e in $A^{(n)}$ and \mathcal{G} based
 on $Resp(e)$ of the elements ▷ Algorithm 3
- 8 **until** RE converges or $iterN++$ exceeds maximum iteration;
- 9 **for** $n = 1 \dots N$ **do**
- 10 $\mathbf{U}^{(n)} \mathbf{B}^{(n)} \leftarrow \mathbf{A}^{(n)}$, and set $\mathbf{A}^{(n)} \leftarrow \mathbf{U}^{(n)}$
- 11 $\mathcal{G} \leftarrow \mathcal{G} \times_n \mathbf{B}^{(n)}$

VEST is a scalable Tucker factorization method that results in very sparse core tensor and factor matrices for partially observable data (see Algorithm 2). First, VEST initializes all elements of the core tensor and factor matrices with random real values between 0 and 1 (line 1). Next, VEST iteratively updates the core tensor and factor matrices while pruning their elements (lines 3-7). In lines 3-4, VEST updates unpruned

elements of the core tensor and factor matrices by element-wise update rules (Section 3.4), guaranteeing that the sparsity non-decreases. Then VEST prunes unimportant elements in the core tensor and factor matrices (lines 6-7). Importance of each element e is evaluated by responsibility $Resp(e)$ which indicates how largely the element contributes to the accuracy (Section 3.2). $should_prune()$ function determines when to stop pruning: if desired sparsity s is achieved (in the manual version $VEST_{*}^{man}$) or the reconstruction error shows a rapid increase (in the automatic version $VEST_{*}^{auto}$). Motivated from simulated annealing, we gradually increase the pruning ratio as iterations proceed (line 7); this enables to explore larger search space in the beginning, while reducing the extent of the search to reduce to a minimum in the later iterations. The iterations proceed until the reconstruction error converges or the maximum iteration is reached. Finally, VEST standardizes all columns of factor matrices such that their norm is equal to one, and updates core tensor accordingly (lines 9-11). Specifically, $A^{(n)}$ is decomposed to $\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times J_n}$ and $\mathbf{B}^{(n)}$ where columns of $\mathbf{U}^{(n)}$ are unit vectors and $\mathbf{B}^{(n)}$ is a diagonal matrix whose (i, i) th element is the norm of $A^{(n)}$'s i th column. The core tensors are updated to maintain the same reconstruction error [9].

3.2 Evaluating Importance of Elements by Responsibility

VEST calculates the *responsibility* of each element which represents its contribution to the overall reconstruction accuracy over the observable elements of the input tensor to determine and prune unimportant elements. The intuition is that reconstruction error increases significantly when a vital element of the core tensor or factor matrices is set to zero, i.e., pruned. On the other hand, if the reconstruction error after pruning is similar or even smaller to that before pruning, the pruned element is insignificant. Formally, the responsibility is defined as follows.

Definition 1 (Responsibility). *Responsibility of an element e in a factor matrix ($e = a_{ij}^{(n)}$) or core tensor ($e = \mathcal{G}_\gamma$) is given by*

$$Resp(e) = \frac{RE(e) - RE}{RE}, \quad (3)$$

where

$$RE = \sqrt{\sum_{\forall \alpha=(i_1, \dots, i_N) \in \Omega} \left(\mathbf{x}_\alpha - \sum_{\forall \beta=(j_1, \dots, j_N) \in \mathcal{G}} \mathcal{G}_\beta \prod_{n=1}^N a_{i_n j_n}^{(n)} \right)^2} / \|\mathcal{X}\|_F \quad (4)$$

is the normalized reconstruction error over the observable entries Ω of the original tensor \mathcal{X} , and $RE(e)$ is residual reconstruction error defined when the element e is set to zero (Eq. (5) and (6)). \square

Definition 2 (Residual reconstruction error). *The residual reconstruction error $RE(\mathcal{G}_\gamma)$ for (j_1, \dots, j_N) th element γ in core tensor \mathcal{G} is as follows:*

$$(RE(\mathcal{G}_\gamma))^2 = \frac{\sum_{\forall \alpha=(i_1, \dots, i_N) \in \Omega} \left(\mathbf{x}_\alpha - \sum_{\forall \beta=(j_1, \dots, j_N) \neq \gamma \in \mathcal{G}} \mathcal{G}_\beta \prod_{n=1}^N a_{i_n j_n}^{(n)} \right)^2}{\|\mathcal{X}\|_F^2} \quad (5)$$

The residual reconstruction error $RE(a_{i,j}^{(n)})$ for an (i,j) th element $a_{i,j}^{(n)}$ in a factor matrix $\mathbf{A}^{(n)}$ is as follows:

$$(RE(a_{i,j}^{(n)}))^2 = RE^2 + \frac{\sum_{\forall \alpha \in \Omega_{i_n}^{(n)}} (2 \cdot (\mathbf{X}_\alpha - B(\alpha)) + B_{j_n=j}(\alpha)) \cdot (B_{j_n=j}(\alpha))}{\|\mathbf{X}\|_F^2}, \quad (6)$$

where $B(\alpha)$ is the entry-wise reconstruction defined as

$$\mathbf{X}_{\alpha=(i_1, \dots, i_N)} \approx B(\alpha) = \sum_{\forall \beta=(j_1, \dots, j_N) \in \mathcal{S}} \mathfrak{G}_\beta \prod_{n=1}^N a_{i_n j_n}^{(n)}, \quad (7)$$

and $B_{j_n=j}(\alpha)$ and $B_{j_n \neq j}(\alpha)$ are the partial reconstruction functions defined as

$$B_{j_n=j}(\alpha) = \sum_{\forall \beta_{j_n=j}} \mathfrak{G}_{\beta_{j_n=j}} \prod_{n=1}^N a_{i_n j_n}^{(n)}, \quad B_{j_n \neq j}(\alpha) = \sum_{\forall \beta_{j_n \neq j}} \mathfrak{G}_{\beta_{j_n \neq j}} \prod_{n=1}^N a_{i_n j_n}^{(n)}. \quad (8)$$

□

The proof of correctness for the derivation of the Eq.(6) is provided in the supplementary material [16]. Note that both definitions are derived from the element-wise reformulation of the reconstruction error.

3.3 Pruning

After calculation of responsibility values, VEST prunes core tensor and factor matrices. Pruning is performed iteratively, each time after the core tensor and factor matrices are updated. The process of pruning an element consists of setting the value of the element to zero and marking the element as pruned in a marking table. The marked elements are excluded from the update step. To prune elements with low responsibility values, VEST sorts elements of the core tensor and each factor matrix, respectively, by the responsibility $Resp(e)$ in ascending order. Then, VEST prunes smallest $pr|\mathcal{S}|$ elements from core tensor and smallest $pr|A^{(n)}|$ from each factor matrix, where pr is the pruning rate of the current iteration. VEST starts with a small pruning rate pr (INIT_PR) and slowly increases pr until maximum pruning rate (MAX_PR) is reached. The default values of INIT_PR and MAX_PR are set to 0.01 and 0.1, respectively.

To determine when to stop pruning, VEST provides two different algorithms, the manual version $VEST_*^{man}$ and the automatic version $VEST_*^{auto}$; the subscript * denotes whether $L1$ or L_F regularization is used.

- $VEST_*^{man}$ takes a target sparsity s as an input from the user and stops pruning when the total sparsity reaches s . That is, $VEST_*^{man}$ enables users to decide on the lower bound of the final sparsity.
- $VEST_*^{auto}$ determines the final sparsity automatically. $VEST_*^{auto}$ does this by tracking changes in the reconstruction error and determines to stop pruning when an elbow point of the reconstruction error curve is reached. The elbow point is estimated as the point when the second derivative of the RE curve, estimated as $(RE_t + RE_{t-2} - 2 * RE_{t-1})/pr_t$ where RE_t and pr_t are RE and pruning rate at t^{th} iteration, respectively, exceeds a small threshold (0.05 used).

3.4 Element-Wise Update Rules

VEST updates elements of the core tensor and factor matrices based on a coordinate decent approach in parallel. It enables VEST to update the core tensor and factor matrices without changing the value of the pruned elements. VEST checks the marking table which indicates whether elements have been pruned, and updates only the un-pruned elements. The update of an element is performed with observable tensor entries and fixed values of other elements in the factor matrices and the core tensor. The update rules for the core tensor and factor matrices are derived by setting the partial derivative of the loss function to zero and solving for each element. In previous works [15,11], this approach has been shown and proven to converge faster with higher accuracy than existing approaches. Advantages of our update rules are that 1) accuracy is high and convergence is faster [15], 2) parallelization and selective updates are possible because all the elements are independently updated, and 3) the size of intermediate data is small, making the algorithm scalable.

Element-wise update rules with L_F regularization The update rule for an element $a_{i_n j_n}^{(n)}$ of factor matrix $A^{(n)}$ is derived by setting the partial derivative of loss function (Eq. (1)) with regard to $a_{i_n j_n}^{(n)}$ to zero.

Lemma 1 (Update rule for factor matrices with L_F regularization).

$$a_{i_n j_n}^{(n)} \leftarrow \arg \min_{a_{i_n j_n}^{(n)}} \mathcal{L}_F(\mathcal{G}, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}) = \frac{\left(\sum_{\forall \alpha \in \Omega_{i_n}^{(n)}} \mathbf{x}_\alpha \delta_\alpha^{(n)}(j_n) \right) - \left(\sum_{\forall t \neq j_n} \mathbf{v}_{i_n j_n}^{(n)}(t) \cdot a_{i_n t}^{(n)} \right)}{\mathbf{v}_{i_n j_n}^{(n)}(j_n) + \lambda}, \quad (9)$$

where $\mathbf{v}_{i_n j_n}^{(n)}$ is a length J_n vector whose j th element is

$$\mathbf{v}_{i_n j_n}^{(n)}(j) = \sum_{\forall \alpha \in \Omega_{i_n}^{(n)}} \delta_\alpha^{(n)}(j) \delta_\alpha^{(n)}(j_n), \quad (10)$$

$\delta_\alpha^{(n)}$ is a length J_n vector whose j th element is

$$\delta_\alpha^{(n)}(j) = \sum_{\forall \beta_{j_n=j} \in \mathcal{S}} \mathcal{G}_{\beta_{j_n=j}} \prod_{k \neq n} a_{i_k j_k}^{(k)}, \quad (11)$$

$\Omega_{i_n}^{(n)}$ is the subset of Ω whose index of n th mode is i_n , and λ is a regularization parameter. \square

The derivation of the core tensor update rule is similar to that of the factor matrix. The update rule for the β^{th} element \mathcal{G}_β of the core tensor \mathcal{G} is given as follows.

Lemma 2 (Update rule for core tensor with L_F regularization).

$$\mathcal{G}_\beta \leftarrow \frac{\sum_{\forall \alpha \in \Omega} (\mathbf{x}_\alpha - \sum_{\forall \gamma \neq \beta} \mathcal{G}_\gamma \prod_{n=1}^N a_{i_n j_n}^{(n)}) \cdot \prod_{n=1}^N a_{i_n j_n}^{(n)}}{\lambda + \sum_{\forall \alpha \in \Omega} \left(\prod_{n=1}^N a_{i_n j_n}^{(n)} \right)^2} \quad (12)$$

□

Element-wise update rules with L_1 regularization For an element $a_{i_n j_n}^{(n)}$ of factor matrix $A^{(n)}$, the element-wise update rule with L_1 regularization is provided in the following Lemmas.

Lemma 3 (Update rule for factor matrix with L_1 regularization).

$$\arg \min_{a_{i_n j_n}^{(n)}} L_1(\mathcal{G}, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}) = \begin{cases} (\lambda - g_{fm})/d_{fm} & \text{if } g_{fm} > \lambda \\ -(\lambda + g_{fm})/d_{fm} & \text{if } g_{fm} < -\lambda \\ 0 & \text{otherwise} \end{cases} \quad (13)$$

where $g_{fm} = 2(\sum_{\forall \alpha \in \Omega_{i_n}^{(n)}} \mathcal{X}_\alpha \delta_\alpha^{(n)}(j_n)) - (\sum_{\forall t \neq j_n} \mathbf{v}_{i_n j_n}^{(n)}(t) \cdot a_{i_n t}^{(n)})$, $d_{fm} = 2\mathbf{v}_{i_n j_n}^{(n)}(j_n)$, and $\mathbf{v}_{i_n j_n}^{(n)}$, $\delta_\alpha^{(n)}$, $\Omega_{i_n}^{(n)}$, and λ follow the same specification provided in Lemma 1. □

For an element \mathcal{G}_β of core tensor, the element-wise update rule with L_1 regularization is as follows:

Lemma 4 (Update rule for core tensor with L_1 regularization).

$$\arg \min_{\mathcal{G}_\beta} L_1(\mathcal{G}, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}) = \begin{cases} (\lambda - g_c)/d_c & \text{if } g > \lambda \\ -(\lambda + g_c)/d_c & \text{if } g < -\lambda \\ 0 & \text{otherwise} \end{cases} \quad (14)$$

where $g_c = -2 \sum_{\forall \alpha \in \Omega} (\mathcal{X}_\alpha - \sum_{\forall \gamma \neq \beta} \mathcal{G}_\gamma \prod_{n=1}^N a_{i_n j_n}^{(n)}) \cdot \prod_{n=1}^N a_{i_n j_n}^{(n)}$, and $d_c = 2 \sum_{\forall \alpha \in \Omega} (\prod_{n=1}^N a_{i_n j_n}^{(n)})^2$. □

The proofs of Lemmas 1 to 4 are provided in the supplementary material [16].

3.5 Parallel Update Algorithms

Responsibility calculation and factor matrices updates are performed in parallel. Algorithm 3 describes the pruning process where responsibility values of the core tensor and factor matrices are calculated in parallel for each observable entries of the input tensor. Note that the use of $B(\alpha)$ in line 5 enabled fast computing of $Resp(\mathcal{G}_\beta)$ in line 6; for a given β in line 4, computing line 5 requires $O(|\Omega|)$ rather than $O(|\Omega||\mathcal{G}|)$ since there is no need to compute $\sum_{\forall \beta \neq \gamma \in \mathcal{G}} \mathcal{G}_\gamma \prod_{n=1}^N a_{i_n j_n}^{(n)}$ in Eq. (5) from scratch.

The element-wise update of factor matrix $A^{(n)}$ is performed in parallel for each rows of factor matrices using either the L_F or L_1 regularization (see Algorithm 4). Elements of the core tensor are dependent on each other and thus cannot be updated in parallel. However, considering that typical size $|\mathcal{G}|$ of the core tensor is small, the core tensor updates are a minor burden in the computational process.

Algorithm 3: Parallel Pruning

Input : Tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$, factor matrices $\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times J_n}$ ($n = 1, \dots, N$), core tensor $\mathcal{G} \in J_1 \times \dots \times J_N$, and pruning ratio pr .

Output: Pruned $\mathbf{A}^{(n)}$ ($n = 1, \dots, N$) and \mathcal{G}

- 1 **for** $\alpha = \forall(i_1, \dots, i_N) \in \Omega$ ▷ **in parallel**
- 2 calculate $B(\alpha) = \sum_{\forall \beta=(j_1, \dots, j_N) \in \mathcal{G}} \mathcal{G}_\beta \prod_{n=1}^N a_{i_n j_n}^{(n)}$
- 3 calculate $\mathcal{X}_\alpha - B(\alpha)$ ▷ Eq.(4)
- 4 **for** $\beta = \forall(j_1, \dots, j_N) \in \mathcal{G}$ ▷ **in parallel**
- 5 calculate $\sum_{\forall \alpha \in \Omega} (\mathcal{X}_\alpha - B(\alpha) + \mathcal{G}_\beta \prod_{n=1}^N a_{i_n j_n}^{(n)})$
- 6 calculate $Resp(\mathcal{G}_\beta)$ ▷ Eq.(3), (5)
- 7 sort core tensor elements by $Resp(\mathcal{G}_\beta)$ values in an ascending order
- 8 **for** $i_n = 1 \dots I_n$ **do**
- 9 **for** $j_n = 1 \dots J_n$ **do** ▷ **in parallel**
- 10 **for** $\alpha = \forall(i_1, \dots, i_N) \in \Omega_{i_n}^{(n)}$ **do**
- 11 calculate $(2(\mathcal{X}_\alpha - B(\alpha)) + B_{j_n=j}(\alpha)) \cdot B_{j_n=j}(\alpha)$
- 12 calculate $Resp(a_{i_n j_n}^{(n)})$ ▷ Eq.(3), (6)
- 13 sort factor matrix elements by $Resp(a_{i_n j_n}^{(n)})$ values in an ascending order
- 14 prune smallest $pr|\mathcal{G}|$ and $pr|A^{(n)}|$ elements of \mathcal{G} and $A^{(n)}$ ($n = 1, \dots, N$), respectively.

Algorithm 4: Parallel Element-Wise Factor Matrix Update

Input : Tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$, factor matrices $A^{(n)} \in \mathbb{R}^{I_n \times J_n}$ ($n = 1, \dots, N$), and core tensor $\mathcal{G} \in J_1 \times \dots \times J_N$.

Output: Updated factor matrices $A^{(n)} \in \mathbb{R}^{I_n \times J_n}$ ($n = 1, \dots, N$)

- 1 **for** $n = 1 \dots N$ **do** ▷ n th factor matrix
- 2 **for** $i_n = 1 \dots I_n$ **do** ▷ **in parallel**
- 3 **for** $j_n = 1 \dots J_n$ **do**
- 4 **if** $a_{i_n j_n}^{(n)}$ is pruned **then**
- 5 continue
- 6 **for** $\alpha = \forall(i_1, \dots, i_N) \in \Omega_{i_n}^{(n)}$ **do**
- 7 **for** $\beta = \forall(j_1, \dots, j_N) \in \mathcal{G}$ **do** ▷ compute δ
- 8 $\delta_\alpha^{(n)}(j_n) \leftarrow \delta_\alpha^{(n)}(j_n) + \mathcal{G}_\beta \prod_{\forall k \neq n} a_{i_k j_k}^{(k)}$
- 9 accumulate $\mathcal{X}_\alpha \delta_\alpha^{(n)}(j_n)$, and update $\mathbf{v}_{i_n j_n}^{(n)}$ ▷ Eq.(10), (11)
- 10 update $a_{i_n j_n}^{(n)}$ using Eq. (9) for L_F (use Eq. (13) for L_1)

Lemma 5 (Complexity of VEST per iteration).

The time complexity per iteration of VEST is $O(N^2 J |\mathcal{G}| |\Omega| / T + NIJ \log(IJ) + |\mathcal{G}| \log |\mathcal{G}|)$, and the memory complexity is $O(TJ + J^N + NIJ)$, where T is the number of threads, N is the order, I is the dimensionality of input tensor, and J is the

dimensionality of the core tensor assuming the dimensionalities are equal for all modes. \square

The full proof of Lemma 5 is in the supplementary material [16].

4 Experiments

We conduct experiments to answer the following questions.

1. **Performance comparison (Section 4.2).** How accurately and sparsely does VEST decompose a given tensor compared to other methods?
2. **Sparsity and accuracy (Section 4.3).** Does VEST successfully prune redundant information in the decomposition without hurting the accuracy? Does $\text{VEST}_*^{\text{auto}}$ find reasonable sparsity and accuracy trade-off point?
3. **Data scalability (Section 4.4).** How scalable is VEST?
4. **Interpretable discoveries (Section 4.5).** How interpretable are the VEST results for discoveries on real-world tensors?

4.1 Experimental Settings

Datasets. We used three real-world datasets and synthetic datasets as summarized in Table 2. The real-world datasets are MovieLens³, Yelp⁴, and AmazonFood⁵. MovieLens is a 4th order tensor of movie ratings containing (user, movie, year, hour). Yelp is a 3rd order tensor of business services rating data containing (user, business, year-month). AmazonFood is a 3rd order tensor of food review scores from Amazon containing (product, user, year-month). To compare with other methods, we used subsets Yelp-s and AmazonFood-s of 3rd order tensors which are made denser than their originals. The density of Yelp-s and AmazonFood-s are 0.01 and 0.02, respectively. We also generated synthetic random tensors of various sizes and orders to test data scalability.

Table 2: Summary of datasets and hyperparameters used.

Name	Order	Dimensionality	Ranks	$ \Omega $	$ \Omega _{\text{test}}$
MovieLens	4	$138K \times 27K \times 21 \times 24$	$6 \times 6 \times 2 \times 2$	18M	2M
Yelp	3	$71K \times 16K \times 108$	$10 \times 10 \times 10$	301K	33K
AmazonFood	3	$74K \times 256K \times 143$	$9 \times 9 \times 14$	511K	57K
Yelp-s	3	$50 \times 50 \times 10$	$5 \times 5 \times 5$	235	32
AmazonFood-s	3	$50 \times 50 \times 10$	$5 \times 5 \times 5$	444	51
Synthetic	$3 - 10$	$10^3 - 10^8$	$3 \times \dots \times 3$	$10^3 - 10^7$	-

Environment. VEST was written in C++ with OPENMP [4] and ARMADILLO [20] libraries for parallelization. Methods L1 (Lasso) and Value Pruning were run on VEST framework with the difference just in the pruning approaches. We used the codes provided by the authors for TTP [21] (R) and Sparse CP [2] (Matlab). Tucker-ALS was

³ <https://grouplens.org/datasets/movielens/>

⁴ http://www.yelp.com/dataset_challenge/

⁵ <http://snap.stanford.edu/data/web-FineFoods.html>

performed via Tensor Toolbox for Matlab [3]. All experiments were done on a single machine equipped with an Intel Xeon E5-2630 v4 2.2GHz CPU (10 cores/20 threads) and 512GB memory. All reported measures are averages of five runs, unless otherwise stated.

Competitors. We compared VEST with the following methods.

- L1 (Lasso): A Tucker factorization method with lasso sparsity constraint implemented as $\text{VEST}_{L_1}^{man}$ with sparsity $s = 0$.
- Value Pruning: A Tucker factorization method with value pruning at the last step implemented as VEST_*^{man} with sparsity $s = 0$ followed by value pruning with ratio 0.6 for L_F and L_1 losses.
- TTP [21]: A tensor decomposition method that results in sparse components.
- Sparse CP [2]: CP decomposition method with lasso penalty.
- Tucker-ALS [8]: Conventional Tucker factorization method (HOOD).
- PTucker-Approx [15]: Tucker decomposition method for partially observable tensor with iterative element value pruning.

4.2 Performance Comparison

We compared the accuracy of $\text{VEST}_{L_F}^{auto}$ and $\text{VEST}_{L_1}^{auto}$ with those of the competitors on datasets Yelp-s and AmazonFood-s (Table 2). The comparison was performed on smaller and denser datasets of order three and not on original real-world datasets due to limitations of the TTR and Sparse CP.

We measured and compared normalized reconstruction errors (RE) over observable entries in input tensors. As shown in Fig. 1(a), $\text{VEST}_{L_F}^{auto}$ and $\text{VEST}_{L_1}^{auto}$ decomposed a given tensor with at least 2.8 times lower RE compared to other methods at a similar sparsity. Fig. 1(a) also shows that at a similar RE value, $\text{VEST}_{L_F}^*$ and $\text{VEST}_{L_1}^*$ output at least 2.2 times more sparse factor matrices and core tensor compared to other methods, where the sparsity is measured as the ratio of number of nonzero values in \mathcal{G} and $\mathbf{A}^{(n)}$ over $|\mathcal{G}| + |\sum_{n=1}^N \mathbf{A}^{(n)}|$.

To answer how well VEST predicts missing entries compared to other methods, we measured REs of the reconstructed missing values (Test RE). After learning the factor matrices and the core tensor using 90% of the observed entries, we calculated the Test REs on the remaining 10% of the observable entries. Fig 1(b) shows that VEST predicts missing entries at least 1.8 times more accurately compared to others, in addition to providing at least 1.7 times sparser results.

4.3 Sparsity and Accuracy

We tested the sparsity and accuracy of outputs of VEST on three full size real-world datasets: MovieLens, Yelp, and AmazonFood, in Figure 3. First, we investigated how the sparsity affects the normalized reconstruction error (RE). Note that the REs are not affected much until the sparsities are above 0.6; this shows that there is redundant information in the decomposition results, and VEST successfully finds and removes such redundancies to get compact outputs. Second, we investigated how VEST_*^{auto} automatically finds a desired sparsity which gives a good tradeoff with regard to accuracy.

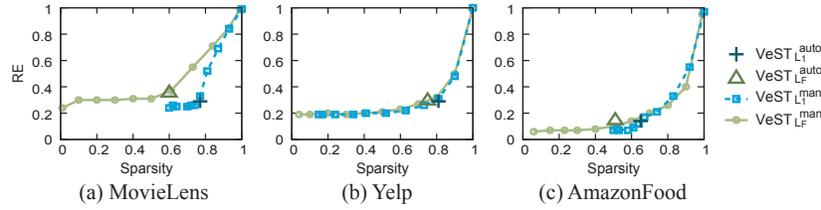


Fig. 3: Sparsity against reconstruction error (RE) of VEST with varying target sparsity s . Note that 1) VEST successfully removes redundant information in the decomposition without hurting the accuracy, and 2) $\text{VEST}_{*}^{\text{auto}}$ automatically finds the sparsity at the elbow point of the RE curve.

Note that in all the datasets, $\text{VEST}_{*}^{\text{auto}}$ successfully finds sparsity points at the elbows of the RE curves, resulting in reasonable sparsity and RE trade-offs. Similar trend was observed for the sparsity and the test RE (see the supplementary material[16]).

4.4 Data Scalability

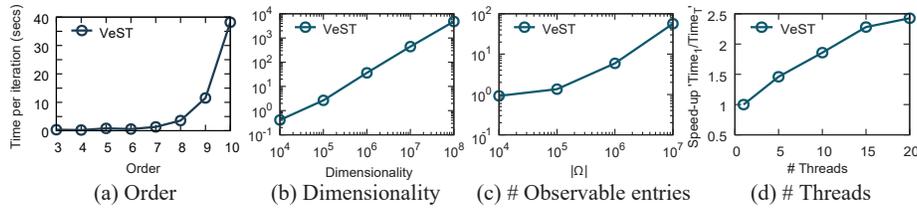


Fig. 4: Scalability of VEST. Data scalability of VEST varying (a) tensor order, (b) tensor dimensionality, (c) number of observable entries, and (d) number of threads.

We evaluated the scalability of VEST by generating synthetic tensors varying the order, the dimensionality, and the number of observable entries, and measuring the running time (see Figure 4). For convenience, the dimensionality of each mode in the input tensor, as well as the dimensionality of core tensor, were set equal, i.e., $I_1 = I_2 = \dots = I_N$ and $J_1 = J_2 = \dots = J_N = 3$, respectively.

Order. Data scalability on the order of input tensor is tested on synthetic tensors of varying orders from 3 to 10. For each input tensor, the dimensionality and the number of observable entries were fixed to $|\Omega| = 10^3$. Figure 4 (a) shows that VEST scales quadratically with regard to order, as discussed in Section 3.5.

Dimensionality. Data scalability on the dimensionality was tested on input tensors of varying dimensionalities from 10^3 to 10^8 with each mode having equal dimension. The order was set to three, and $|\Omega|$ was set equal to the dimensionality. Figure 4 (b) shows that VEST has near-linear scalability in terms of the dimensionality.

Number of Observable Entries. Data scalability on the number of observable entries was tested on input tensors by varying the number of observable entries from 10^3 to 10^7 . The order was set to three, and the dimensionality was fixed to 10^3 . Figure 4 (c) shows that VEST has near-linear scalability in term of the number of observable entries.

Effectiveness of Parallelization. We evaluated the parallelization scalability of VEST by increasing the number of threads from 1 to 20 and measuring $\text{Time}_1/\text{Time}_T$ where

$Time_T$ is the running time per iteration using T threads. Figure 4(d) shows near-linear scalability of VEST in terms of the number of threads used.

4.5 Discovery

We evaluated interpretability of VEST by investigating the factorization results of MovieLens dataset and visually showing that the sparse results enhance interpretability. It is difficult to analyze dense results generated by vanilla methods without post-processing. In contrast to existing methods, we can easily identify interesting factors generated by VEST based on sparsity of each row of a factor matrix.

Discovery of Greatest Movies. We found that a few rows of the movie-associated factor matrix are fully dense although the goal of VEST is to generate sparse results. Such rows corresponded to popular movies rated by diverse users. Figure 5 shows popular movies which have the largest number of non-zero and the sums of values. According to Empire magazine [1], 14 out of 20 movies we found were included in the 100 greatest movies. The remaining six movies, including ‘Sixth Sense’, ‘Kill Bill’, and ‘Fifth Element’, that were not in the 100 list were also very famous.

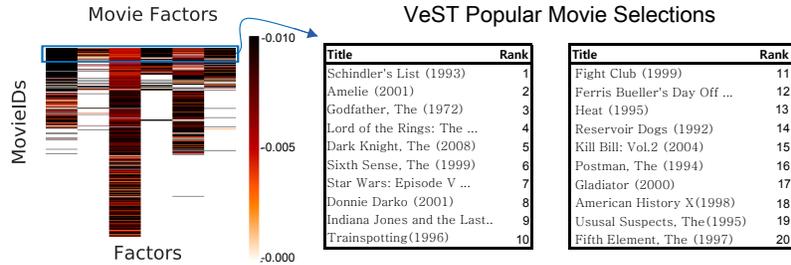


Fig. 5: Popular movies discovered by VEST. 14 of 20 movies we find are included in the 100 greatest movies introduced in Empire magazine. The remaining six movies (Sixth Sense, Kill Bill, Ferris Bueller’s Day Off, Postman, American History X, and Fifth Element) are also famous movies and were rated by various users.

5 Conclusion

We proposed VEST, a very-sparse Tucker factorization method for sparse and partially observable tensors. By deriving the element-wise partial differential equations, determining the importance of elements by responsibilities, and parallel distribution of computational work, VEST successfully offers very sparse and accurate results that are applicable for large partially observable tensors. VEST generates at least 2.2 times more sparse results compared to other methods for partially observable tensors and at least 2.8 times accurate results compared to other sparse factorization methods. VEST also shows near linear scalability regarding tensor dimensionality, number of observable entries, and number of threads. Thanks to the increased sparsity that leads to improved interpretability by VEST, we were able to discover interesting patterns related to the greatest movies in the factor matrix of a real-world movie rating tensor data. Future works include better initialization for Tucker factorization, integration of prior knowledge, and effective visualization of tensor results.

References

1. The 100 greatest movies (2018), <https://www.empireonline.com/movies/features/best-movies/>
2. Allen, G.: Sparse higher-order principal components analysis. In: *Artificial Intelligence and Statistics*. pp. 27–36 (2012)
3. Bader, B.W., Kolda, T.G., et al.: Tensor toolbox for matlab v. 3.0, version 00
4. Dagum, L., Menon, R.: Openmp: An industry-standard api for shared-memory programming. *IEEE Comput. Sci. Eng.* **5**(1), 46–55 (Jan 1998)
5. Jiang, F., Liu, X.y., Lu, H., Shen, R.: Efficient Multi-Dimensional Tensor Sparse Coding Using t-Linear Combination. In: *AAAI 2018*. pp. 3326–3333 (2018)
6. Kang, U., Papalexakis, E.E., Harpale, A., Faloutsos, C.: Gigatensor: scaling tensor analysis up by 100 times - algorithms and discoveries. In: *KDD*. pp. 316–324 (2012)
7. Kim, Y.D., Choi, S.: Nonnegative tucker decomposition. In: *CVPR'07*. pp. 1–8. IEEE (2007)
8. Kolda, T.G., Bader, B.W.: Tensor decompositions and applications. *SIAM review* **51**(3), 455–500 (2009)
9. Kolda, T.G.: *Multilinear operators for higher-order decompositions*, vol. 2. United States. Department of Energy (2006)
10. Lee, J., Choi, D., Sael, L.: CTD: Fast, accurate, and interpretable method for static and dynamic tensor decompositions. *PLoS One* **13**(7), e0200579 (2018)
11. Lee, J., Oh, S., Sael, L.: GIFT: Guided and Interpretable Factorization for Tensors with an application to large-scale multi-platform cancer analysis. *Bioinformatics* **34**(24), 4151–4158 (2018)
12. Madrid-Padilla, O.H., Scott, J.: Tensor decomposition with generalized lasso penalties. *Journal of Computational and Graphical Statistics* **26**(3), 537–546 (2017)
13. Mahoney, M.W., Maggioni, M., Drineas, P.: Tensor-CUR Decompositions for Tensor-Based Data. *SIAM J. Matrix Anal. Appl.* **30**(3), 957–987 (jan 2008)
14. Mørup, M., Hansen, L.K., Arnfred, S.M.: Algorithms for sparse nonnegative tucker decompositions. *Neural computation* **20**(8), 2112–2131 (2008)
15. Oh, S., Park, N., Sael, L., Kang, U.: Scalable Tucker factorization for sparse tensors - algorithms and discoveries. In: *ICDE*. IEEE Computer Society, Paris, France (2018)
16. Park, M., Jang, J.G., Lee, S.: Supplementary material of VeST (2019), <http://github.com/leesael/VeST/paper/supp-material.pdf>
17. Pascual-Montano, A., Carazo, J.M., Kochi, K., Lehmann, D., Pascual-Marqui, R.D.: Nonsmooth nonnegative matrix factorization (nsNMF). *IEEE TPAMIT* **28**(3), 403–415 (2006)
18. Qi, N., Shi, Y., Sun, X., Yin, B.: TenSR: Multi-dimensional Tensor Sparse Representation. 2016 IEEE CVPR pp. 5916–5925 (2016)
19. Ribeiro, M.T., Singh, S., Guestrin, C.: "why should I trust you?": Explaining the predictions of any classifier. In: *ACM SIGKDD'16*. pp. 1135–1144 (2016)
20. Sanderson, C., Curtin, R.: Armadillo: a template-based c++ library for linear algebra. *Journal of Open Source Software* (2016)
21. Sun, W.W., Lu, J., Liu, H., Cheng, G.: Provable sparse tensor decomposition. *Journal of the Royal Statistical Society: Series B* **79**(3), 899–916 (jun 2017)
22. Yi, S., Lai, Z., He, Z., ming Cheung, Y., Liu, Y.: Joint sparse principal component analysis. *Pattern Recognition* **61**(2), 524–536 (2017)
23. Zemin, Z., Aeron, S.: Denoising and completion of 3D data via multidimensional dictionary learning. pp. 2371–2377 (2016)