

International Conference on Intelligent Computing, Communication & Convergence
(ICCC-2015)

Conference Organized by Interscience Institute of Management and Technology,
Bhubaneswar, Odisha, India

Lowering Data Dimensionality in Big Data For The Benefit of Precision Agriculture

Sabarina K ¹, Priya N²

¹PG Student, Jerusalem College of Engineering, Chennai

²Assistant Professor, Jerusalem College of Engineering, Chennai

Abstract

Predictive analytics can be used to make smarter decisions in farming by collecting real-time data on weather, soil and air quality, crop maturity and even equipment and labor costs and availability. This is known as precision agriculture. Big data is expected to play an important role in precision agriculture for managing real-time data analysis with massive streaming data. The data analysis efficiency and throughput would be a challenge with the massive increase in size of big data. The unstructured streaming data received from different agricultural sources would contain multiple dimensions and not the entire content is needed for performing analysis. The core data which is small but that alone enough to represent the entire content should be extracted. This paper explains how to systematically reduce the size of big data by applying a tensor based feature reduction model. The data decomposition and core value extraction is done with the help of IHOSVD algorithm. This way it reduces the overall file size by eliminating unwanted data dimensions. The time involved in data analysis and CPU usage will be significantly reduced when dimensionality reduced data is used in place of raw (unprocessed) data.

Keywords: Big data; Precision agriculture; Dimensionality reduction; HOSVD.

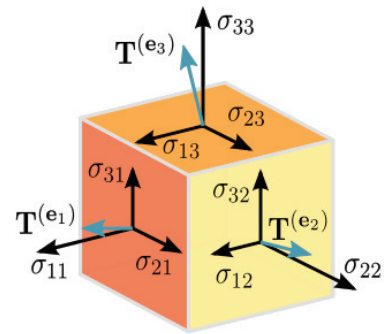
1. Introduction

The current ratio of agricultural land to world population is sufficient to meet the food requirements of the world and hence we are not much concerned about agricultural productivity. However the global population is expected to grow exponentially in the coming years wherein the total area of agricultural land is getting reduced day by day.

Hence it is very essential to switch from conventional agricultural methods to modern precision agriculture for increasing the agricultural productivity and maintaining safer food supply. The precision agriculture functions together with information technology to make a difference in the agriculture methodologies. It is expected that the precision agriculture would help us to reach sustainable agriculture where the food demand and supply will always be maintained at safer level. The Data Mining and Big Data technologies are advancing in agriculture, changing how various agricultural activities are sequenced. Making predictions about crops based on the real-time data gathered from the farms will change the current scenario in agriculture [10]. Improvising the speed of data analysis by reducing the data dimensions will help the farmers getting the tips on time. Agricultural predictions done using data mining techniques may not be effective when the data dimensions are high. The query accuracy and efficiency degrade rapidly as the dimension increases. Dimensionality reduction is an effective approach for reducing the dimensions and size of data.

The most popular approaches that are involved in data dimensionality reduction include Principal Component Analysis (PCA) [2], Incremental Singular Value Decomposition (SVD) [7], and Dynamic Tensor Analysis (DTA) [4]. These methods work fine for low dimensional data but suffer from serious performance issues when being applied on high-dimension data and fail to extract the core data sets from streaming big data. This paper explains an incremental dimensionality reduction method for reducing data dimensions and extracting the core set of data from massive big data.

Recursive Incremental HOSVD Method: The high dimensional streaming data is converted into tensors and then the tensors are unfolded into matrices. Later the IHOSVD algorithm is recursively applied to extract the core tensor by decomposing the matrices [1]. The IHOSVD algorithm is explained in a later section of this paper.



2. Preliminaries

2.1 Tensors

Tensors are mathematical elements that describe physical properties similar to vectors and scalars. They help to establish mathematical relationship between vectors, scalars, and other tensors. Typical examples of such relations include the cross product, the dot product, and linear maps. Vectors and scalars are also individually considered as tensors. Tensors are normally represented using multi-dimensional array of numerical values. The order of a tensor is determined based on the number of dimensions needed to represent it in an array. Scalars are viewed as 0th-order tensors as they are single numbers. Vectors are viewed as 1st-order tensors as they can be represented using one dimensional arrays. A linear map is considered as a 2nd-order tensor as it can be represented by a 2-dimensional array. The array dimensionality should not be confused with the dimension of the underlying vector space while determining the tensor order.

2.2 Singular Value Decomposition (SVD)

The singular value decomposition (SVD) is a mathematical approach used for splitting a complex matrix into simple matrices by applying factorization. This approach is extensively used in data analytics and signal processing applications. The singular value decomposition of a complex matrix R is represented in factorization form as follows [1] [2] [8],

$$R = U_1 \Sigma U_2^*$$

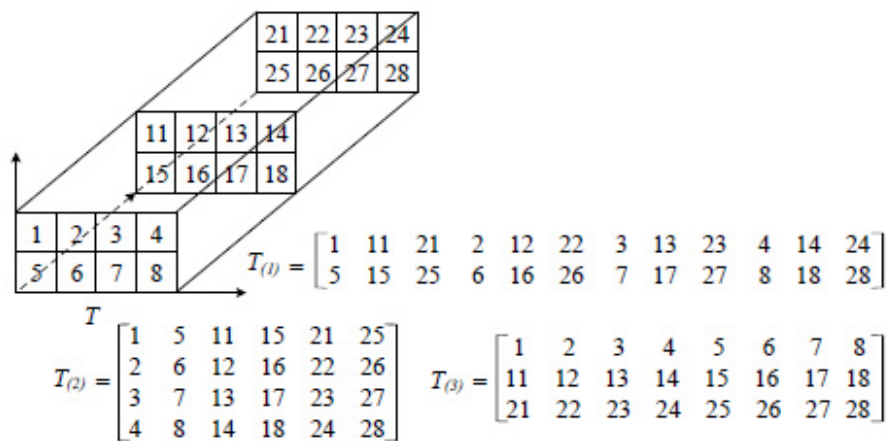
Where U_1 is an $m \times m$ real or complex unitary matrix, Σ is an $m \times n$ rectangular diagonal matrix with non-negative real numbers on the diagonal, and U_2^* (the conjugate transpose of U_2 , or simply the transpose of U_2 if U_2 is real) is an $n \times n$ real or complex unitary matrix. The diagonal entries $\Sigma_{i,i}$ of Σ are known as the singular values of R . The m columns of U_1 and the n columns of U_2 are called the left-singular vectors and right-singular vectors of R , respectively.

2.3 Tensor unfolding

Given a P -order tensor $T \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_P}$, the tensor unfolding [1]

$T_{(p)} \in \mathbb{R}^{I_p \times (I_{p+1} I_{p+2} \dots I_P I_1 I_2 \dots I_{p-1})}$ contains the element $t_{i_1 i_2 \dots i_p i_{p+1} \dots i_P}$ at the position with row number i_p and column number that is equal to $(i_{p+1} - 1) I_{p+2} \dots I_P I_1 \dots I_{p-1} + (i_{p+2} - 1) I_{p+3} \dots I_P I_1 \dots I_{p-1} + \dots + (i_2 - 1) I_3 I_4 \dots I_{p-1} + \dots + i_{p-1}$.

Consider a three-order tensor $T \in \mathbb{R}^{2 \times 4 \times 3}$, Fig. 2 shows the three unfolded matrices $T_{(1)}$, $T_{(2)}$ and $T_{(3)}$



2.4 Tensor Order and Tensor Dimension

The tensor order and tensor dimension are two key elements used in dimensionality reduction methodologies. We should understand the difference between them before using them in dimensionality reduction. Tensor $T \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_P}$ has P orders, and order i ($1 \leq i \leq P$) has I_i dimensions. A P -order tensor can be unfolded to P matrices [1] [6] [8]. For the mode- i unfolded matrix $T(i)$, the number of rows is equal to I_i , while the number of columns is equal to $\prod_{1 \leq j \leq P; j \neq i} I_j$. It will be very difficult to establish relationship among the data records when the data contains redundant, repetitive and unwanted data. Hence it is essential to identify the core data from the input dataset after removing the redundancy and unwanted data. Normally the number of tensor orders remains the same while the dimensionality is significantly reduced during the dimensionality reduction process.

2.5 Unified Tensors

Big data are heterogeneous in nature and it is composed of unstructured data d_u , semi-structured data d_{semi} and structured data d_s . Due to the requirement of processing all types of heterogeneous data, a unified data tensorization operation is performed using the following equation [1]

$$f : (d_u \cup d_{semi} \cup d_s) \rightarrow \underbrace{T_u \cup T_{semi} \cup T_s}_T$$

3. Incremental High Order Singular Value Decomposition

This paper propose an IHOSVD method for incremental dimensionality reduction on streaming data. The real-time data obtained from sensors placed in an agricultural field can be considered for this data reduction exercise [9]. The IHOSVD method consists of three algorithms that are used for recursive matrix singular value decomposition and incremental tensor decomposition [1]. The three algorithms are discussed below,

Algorithm 1 is a recursive algorithm with recursive function given in the below equation,

$$f(M_i, C_i) = \begin{cases} \text{svd}(M_1), & i = 1 \\ \text{mix}(f(M_{i-1}, C_{i-1}), C_i), & i > 1 \end{cases}$$

During the running process, function f will call itself over and over again to decompose matrices M_i and C_i . Each successive call reduces the size of matrix and moves closer to a solution until matrix M_1 is reached finally, the recursion stops and the function can exit.

Algorithm 1: Recursive singular value decomposition on matrices, $(U, \Sigma, V) = \text{R-MSvd}(M_i, C_i)$.

Input:

Initial matrix M_i .

Incremental matrix C_i .

Output:

Decomposition results U, S, V of matrix $[M_i C_i]$.

Algorithm 2: Merging the incremental matrix with decomposition results of previous matrix, $(U, \Sigma, V) = \text{mix}(M_{i-1}, C_{i-1}, U_j, \Sigma_j, V_j)$.

Input:

Initial matrix M_{i-1} and incremental matrix C_{i-1} .

Decomposition results U_j, Σ_j, V_j of matrix M .

Output:

New decomposition results U, Σ, V .

Algorithm 3: Incremental singular value decomposition on tensors,

$(S, [U, \Sigma, V]_{\text{new}}) = I - T \text{ Svd}(\chi, T, [U, \Sigma, V]_{\text{initial}})$.

Input:

New tensor $\chi \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_P}$

Previous tensor $T \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_P}$

Previous unfolded matrices SVD results $[U, \Sigma, V]_{\text{initial}}$.

Output:

New truncated SVD results $[U, \Sigma, V]_{\text{new}}$.

New core tensor S .

The consolidated steps comprising all the three algorithms are listed below,

1. Extend tensor χ and tensor T to identical dimensionality.
2. Unfold new tensor χ to matrices $\chi(1), \dots, \chi(P)$.
3. Decompose the matrix $\chi(1)$ into $U(1), S(1), V(1)$ by applying SVD
4. Execute step 10 to 14 recursively after incrementally passing matrices $\chi(2), \dots, \chi(P)$ together with decomposition results of previous matrix
5. Receive orthogonal bases U, S, V for the tensor χ
6. Truncate the new orthogonal bases.
7. Combine new tensor χ with initial tensor T .
8. Obtain new core tensor S with n-mode product
9. Return S , and $[U, S, V]_{\text{new}}$.
10. Project $\chi(2)$ on the orthogonal space spanned by $U(1)$
11. Compute matrix H which is orthogonal to $U(1)$.
12. Obtain the unitary orthogonal basis J from matrix H
13. Determine the new decomposition results $U(2), S(2), V(2)$ by applying SVD on $[U(1) J]$
14. Return the new decomposed results $[U(2) S(2) V(2)]$ (Note: This decomposed results are created for consolidated matrices $\chi(1)$ & $\chi(2)$)

4. Conclusion

Big data and predictive analytics is expected to make a difference in the agricultural industry by providing on-time tips to the farmers based on real-time data gathered from different agricultural sources. Considering the volume and heterogeneous nature of future data, a well-defined mechanism is needed to reduce the memory and time involved in data analysis. Reducing the file size by extracting core value from massive big data is one of the approach to improve the data analysis performance.

References:

- [1] Liwei Kuang, Fei Hao, Laurence T. Yang, Man Lin, Changqing Luo, and Geyong Min “A Tensor-Based Approach for Big Data Representation and Dimensionality Reduction” *IEEE Transaction on Emerging Topics in Computing*-2014
- [2] L. De Lathauwer, B. De Moor, and J. Vandewalle, “A Multilinear Singular Value Decomposition,” *SIAM Journal on Matrix Analysis and Applications*, vol. 21, no. 4, pp. 1253–1278, 2000.
- [3] I. Horrocks, P. F. Patel-Schneider, and F. Van Harmelen, “From SHIQ and RDF to OWL: The Making of a Web Ontology Language,” *Web Semantics: Science, Services and Agents on the World Wide Web*, vol. 1, no. 1, pp. 7–26, 2003.
- [4] J. Sun, D. Tao, and C. Faloutsos, “Beyond Streams and Graphs: Dynamic Tensor Analysis,” in *Proc. of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2006, pp. 374–383.
- [5] Georg Rub, Rudolf Kruse, Martin Schneider, and Peter Wagner, “Data Mining with Neural Networks for Wheat Yield Prediction,” *Springer-Verlag Berlin Heidelberg* 2008
- [6] P. Symeonidis, A. Nanopoulos, and Y. Manolopoulos, “Tag Recommendations Based on Tensor Dimensionality Reduction,” in *Proc. of the 2008 ACM Conference on Recommender Systems*. ACM, 2008, pp. 43–50.
- [7] Q. Li, X. Shi, and D. Schonfeld, “A General Framework for Robust HOSVD-Based Indexing and Retrieval with High-Order Tensor Data,” in *Proc. of the 36th IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2011, pp. 873–876
- [8] M. Kim and K. S. Candan, “Approximate Tensor Decomposition within a Tensor-Relational Algebraic Framework,” in *Proc. of the 20th ACM International Conference on Information and Knowledge Management*. ACM, 2011, pp. 1737–1742
- [9] Hideya Ochiai, Member, IEEE, Hiroki Ishizuka, Yuya Kawakami, and Hiroshi Esaki, Member, IEEE “A DTN-Based Sensor Data Gathering for Agricultural Applications,” *IEEE sensors journal*, vol. 11, no. 11, November 2011
- [10] Gopala Krishna Moorthy .K, Dr.C.Yaashuwanth, Venkatesh.K, “A Wireless Remote Monitoring Of Agriculture Using Zigbee”, *International Journal of Engineering and Innovative Technology (IJEIT)* Volume 2, Issue 8, February 2013
- [11] Xindong Wu, Fellow, IEEE, Xingquan Zhu, Gong-Qing Wu and Wei Ding, Senior Member, IEEE “Data Mining with Big Data”-*IEEE Transaction on Knowledge And Data Engineering* -2014

- [12] Lin Gu, Student Member, IEEE, Deze Zeng, Member, IEEE, and Song Guo, Senior Member, IEEE "Cost Minimization for Big Data Processing in Geo-Distributed Data Centers" IEEE Transactions on Emerging Topics in Computing-2014
- [13] Chunxiao Jiang, Member, IEEE, Yan Chen, Member, IEEE, and K.J Ray Liu, Fellow, IEEE "Graphical Evolutionary Game for Information Diffusion Over Social Networks"-IEEE journal of Selected Topics in Signal Processing-2014
- [14] Sangeeta Bansal, Dr. Ajay Rana, "Transitioning from Relational Databases to Big Data", International Journal of Advanced Research in Computer Science and Software Engineering, Volume 4, Issue 1, January 2014