

1 Encoding and Decoding Techniques for Distributed Data Storage 2 Systems

3 Ms.nbarasi A¹ and Dr. K.Vivekanandan²

4 ¹ Bharathiar University,Coimbatore,India.

5 *Received: 12 June 2011 Accepted: 8 July 2011 Published: 20 July 2011*

7 Abstract

8 Dimensionality reduction is the conversion of high-dimensional data into a meaningful
9 representation of reduced data. Preferably, the reduced representation has a dimensionality
10 that corresponds to the essential dimensionality of the data. The essential dimensionality of
11 data is the minimum number of parameters needed to account for the observed properties of
12 the data [4]. Dimensionality reduction is important in many domains, since it facilitates
13 classification, visualization, and compression of high-dimensional data, by helpful the curse of
14 dimensionality and other undesired properties of high-dimensional spaces [5]. Dimension
15 reduction can be beneficial not only for reasons of computational efficiency but also because it
16 can improve the accuracy of the analysis. In this research area, it significantly reduces the
17 storage spaces.

18 *Index terms*— Dimensionality reduction, high-dimension and storage.

20 1 INTRODUCTION

21 he analysis and mining of large volumes of transaction data for making business decisions.

22 Today it has become a key success factor than ever for vendors to understand their customers and their buying
23 patterns. If they don't they will lose them. In order to gain competitive advantage it is necessary to understand
24 the relationships that prevail across the data items among millions of transactions. The amount of data currently
25 available for studying the buying pattern is extensive and increasing rapidly year by year. Therefore the need
26 to devise reliable and scalable techniques to explore the millions of transactions for the customer buying pattern
27 continues to be important. Above this, the increasing volume of data sets data demands for huge amounts of
28 resources in storage space and computation time. As it is not feasible to have huge storage spaces to store
29 the explosively growing data in a single location they are stored in distributed database and data warehouse
30 located in different geographical location. Inherently data distributed over a network with limited bandwidth
31 and computational resources motivated the development of distributed data mining (DDM).

32 Though mining process in DDM is carried out in distributed locations parallel and generates required
33 results in the local areas it is necessary to analyze these Author ? : Research Scholar,Bharathiar Univer-
34 sity,Coimbatore,India. Telephone: 09751149851 E-mail : anbarasi2@gmail.com Author ? : Prof,Bharatiary
35 University,School of Management and studies,Coimbatore India. E-mail : vivekbsmed@gmail.com local patterns
36 to obtain the global data model. Hence the knowledge derived from local distributed location is moved to the
37 central site and the local results are combined there to obtain the final result. This approach is less expensive
38 but may produce ambiguous and incorrect global results. Even though communication is a bottleneck problem
39 in a central data repository it guarantees accurate results of data analysis. To address the bottleneck problem
40 in central learning strategy, this work proposes a dimension reduction method which uses the concept of sum of
41 subset and scalable to very large databases. In this work the site which request data from different geographical
42 locations is treated as central site.

2 II.

3 EXISTING WORK

Principal component analysis (PCA) is a mathematical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of uncorrelated variables called principal components. The number of principal components is less than or equal to the number of original variables. This transformation is defined in such a way that the first principal component has as high a variance as possible), and each succeeding component in turn has the highest variance possible under the constraint that it be orthogonal to (uncorrelated with) the preceding components. Principal components are guaranteed to be independent only if the data set is jointly normally distributed.

Linear Discriminant Analysis (LDA) attempts to maximize the linear separability between data points belonging to different classes. In dissimilarity to most other dimensionality reduction techniques, LDA is a supervised technique. LDA finds a linear mapping M that maximizes the linear class separability in the lowdimensional representation of the data.

4 III.

5 PROBLEM DESCRIPTION

Data storage conversion algorithm transforms a transaction into a single dimension transaction with all attributes that appears in its original form. The encoded transactions are represented by a sequence of numbers, which is sum of subset approach. Any kind of combination of $2^1, 2^2, 2^3, \dots, 2^n$, the sum of different values gives as unique, This is way motivated to do the my research work. By this way, the new transaction is smaller than the original form and hence the cost of Storage is reduced. To offer highly specialized solutions for small parts of the general problem.

IV.

6 ENCODING AND DECODING TECHNIQUES FOR DATA STORAGE

A matrix is constructed with the given set of data items as shown in the Table-1. The order and dimensions of the matrix are user defined. The only constraint is that the number of columns should not exceed 14, as the value of 2^{14} will exceed the range of an 'int'.

7 Table 1: Display of item set

Of the entire set of data items, a transaction is always a subset of the data items. This subset of data items is encoded into a reduced database. If it reduced database minimized the memory area. The Table-2 explains the encoding process. For each data item in the transaction, the row 'i' and column 'j' is noted. The value 2^j is calculated and added to the i'th value in the transaction value E. The process is repeated for each data item one by one and final 'n' digits from the Table ?? The Transaction T1 = Chicken, Wheat bread, Dry fruits, Jam, Soft drink, Sugar, Pizza

In Table 4a, the items chosen in the row 1 are found by decoding the number '34'. Since the given matrix has 5 columns -the values of $2^5, 2^4, 2^3, 2^2$ and 2^1 are all subtracted from the value '34' one by one, cumulatively. Each time, the subtraction gives a positive value, the corresponding column's data item is chosen. The final list of chosen data items in this table indicates the original items from the transaction. The process is repeated on the other values of the reduced transaction form ie, 50 n 10 in Tables 4b and 4c respectively. Thus the remaining data items from the transaction are also decoded.

V.

8 CONCLUSION

The above stated technology appears to be the most fitting and forceful method adaptable in the distributed data as well as in the distributed data mining process in terms of speed and competence when we measure it up to the old methods. Another useful characteristic that is covered under this new technology is that it could be updated constantly when it is essential Encoding and Decoding Techniques for Distributed Data Storage Systems since the data is maintained at remote sites. The huge quantity of data is not needed to be stored to the much location for this purpose hence the storage spaces are used most favorably. To make complete use of the novel technology, the customary client server distributed data mining scheme must be entirely replaced with it. Methodology expansions for merging the accumulated information from different spots are in advancement. The purpose with this paper was to provide an overview of the specific of approach that can be employed for dimension reduction when processing high dimension data.

¹T © 2011 Global Journals Inc. (US)

²© 2011 Global Journals Inc. (US)



Figure 1:

Data Item	Matches with ith Row of	jth Col- umn of	E after adding 2j to the existing value in ith column		
Jam	2	4	0	0+16	0
Wheat bread	1	1	0+2	16	0
Chicken	1	5	2+32	16	0
Soft drink	2	1	34	16+2	0
Dry fruits	2	5	34	18+32	0
Sugar	3	3	34	50	0+8
Pizza	3	1	34	50	8+2

Figure 2:

2

Figure 3: Table 2 :

3

contains all the transactions in the encoded form. It will be these encoded values that will be transferred across the network between the client and the server.

34	50	10
58	16	02
18	16	18

Figure 4: Table 3

3

	Encoding and Decoding Techniques for Distributed Data Storage Systems				
pizza	sauce		sugar	sweet	
				bun	
soft	fruit wheat		honey	jam	Dry
					fruits
drink	bread				
wheat	bun		burger	butter	chickn
bread					

Figure 5: Table 3 :

4a

:
 $i=2; d1 = 50; n=5$

Figure 6: Table 4a

4b

Figure 7: Table 4b :

4c

Figure 8: Table 4c :

-
- 95 [Jackson ()] *A User's Guide to Principal Components*, J E Jackson . 1991. New York: John Wiley and Sons.
- 96 [Wu-Shan et al. (2005)] 'Distributed Data Mining on the Grid'. Wu-Shan , Ji-Hui Jiang , Yu . *Proceedings of the*
97 *Fourth International Conference on Machine Learning and Cybernetics*, (the Fourth International Conference
98 on Machine Learning and CyberneticsGuangzhou) August 2005. p. .
- 99 [Kulkarni et al. ()] 'Exploring the capabilities of Mobile Agents in Distributed Data Mining'. U P Kulkarni ,
100 K K Tangod , S R Mangalwede , A R Yardi . *10th International Database Engineering and Applications*
101 *Symposium (IDEAS'06)*, 2006. IEEE.
- 102 [Fukunaga ()] *Introduction to Statistical Pattern Recognition*, K Fukunaga . 1990. San Diego,CA, USA: Academic
103 Press Professional, Inc.
- 104 [Jimenez and Landgrebe ()] 'Supervised classification in high-dimensional space: geometrical, statistical,and
105 asymptotical properties of multivariate data'. L O Jimenez , D A Landgrebe . *IEEE Transactions on Systems,*
106 *Man and Cybernetics* 1997. 28 (1) p. .