Artificial Intelligence formulated this projection for compatibility purposes from the original article published at Global Journals. However, this technology is currently in beta. *Therefore, kindly ignore odd layouts, missed formulae, text, tables, or figures.*

A Survey on Clustering Techniques for Multi-Valued Data Sets

2	Lnc. Prakash K^1
3	1 Annamacharya institute of technology and sciences
4	Received: 10 December 2015 Accepted: 5 January 2016 Published: 15 January 2016

6 Abstract

 $_{7}$ $\,$ The complexity of the attributes in some particular domain is high when compare to the

 $_{\rm 8}$ $\,$ standard domain, the reason for this is its internal variation and the structure .their $\,$

⁹ representation needs more complex data called multi-valued data which is introduced in this

¹⁰ paper. Because of this reason it is needed to extend the data examination techniques (for

11 example characterization, discrimination, association analysis, classification, clustering, outlier

¹² analysis, evaluation analysis) to multi-valued data so that we get more exact and consolidated

¹³ multi-valued data sets. We say that multi-valued data analysis is an expansion of the standard

¹⁴ data analysis techniques. The objects of multi-valued data sets are represented by

¹⁵ multi-valued attributes and they contain more than one value for one entry in the data base.

¹⁶ An example for this type of attribute is ?languages known? .this attribute may contain more

than one value for the corresponding objects because one person may be known more than one
 language.

19

20 Index terms—

²¹ 1 I. Introduction

22 Analysis to such objects, called multi-valued Data analysis.

²³ 2 II. Input to Multi-Valued Data Analysis Algorithms

In general Columns of the source data table are variables (attributes) and rows (objects) are multi-valued descriptions. Each cell of this multi-valued data table may contain data of different types as given bellow.

? The sell may contain a quantitative value for example. If "height" is an attribute and a is a unit then height (a)=2.4.? The sell may contain discrete value, for instance, consider "village" is a variable and a is a unit then village (a) = Raja pet.? The sell may contain Multi-valued, for example, in the quantitative case: height (a)

29 = (3.2, 4.1, 5)

) which means that the weight of a maybe 3 .2 or 4 .1 or 5 . In the discrete case, language (a) = (Telugu, Hindi, Tamil) means that the language of a may be Telugu or Hindi or Tamil. Here it is clear that the above two are special cases of this case. ? The sell may contain Interval valued data: for instance height (a) = [12,15],

which tells that the height of a varies in the interval [12,15]. Some other types of variables may also exist with

³⁴ Multi valued weights, Taxonomic, Hierarchically dependent, logically dependent.

³⁵ 3 III. Result of Multi-Valued Data Analysis Algorithms

The output of any data analysis depending on the type of the data mining task. If the task is characterization it represents the nature of data and its character. When the task is association we can find the frequent items and

then form the association rules. If the task is classification we can find the different classes that are associated with the data. If the task is clustering we can partition the data into different similar groups depending on the

40 similarity measure.

⁴¹ 4 IV. Source of Multi-Valued Data

42 Multi-valued data is formed from the tables that are integrated from different sources so this type of data is 43 summarized and concise so the requirement of consolidation of different huge data sets causes multivalued data 44 sets. The result from several probability 43 Year 2016

45 **5** () C

Abstract-The complexity of the attributes in some particular domain is high when compare to the standard 46 domain, the reason for this is its internal variation and the structure their representation needs more complex 47 data called multi-valued data which is introduced in this paper. Because of this reason it is needed to extend the 48 data examination techniques (for example characterization, discrimination, association analysis, classification, 49 clustering, outlier analysis, evaluation analysis) to multi-valued data so that we get more exact and consolidated 50 multi-valued data sets. We say that multi-valued data analysis is an expansion of the standard data analysis 51 techniques. The objects of multi-valued data sets are represented by multi-valued attributes and they contain 52 more than one value for one entry in the data base. An example for this type of attribute is "languages known" 53 .this attribute may contain more than one value for the corresponding objects because one person may be known 54 more than one language. n the process of making more general surveys the persons are permitted to give more 55 than one answer for a particular question. The answer may be a set of categorical values or a set of numerical 56 values or sometimes it may be the combination of both. Now a days the data base is going on increasing and it 57 is needed the integration of different data bases as a result a very lengthy databases are formed. The data bases 58 that are formed in this way may contain the same objects that are repeated many times with different values .the 59 increasing importance is to reduce the data base size by summarizing the data without loss of the information 60 that is used for analysis. This can be achieved by introducing the concept of multi-valued attributes in the data 61 bases. Because the cells of such data may contain not only single numerical or categorical values, but much 62 more multifaceted information, such as subsets of categorical variable values, intervals of ordinal variable values, 63 dependencies and need rules to be specified. These new statistical units are called multi-valued objects and there 64 is a need for an extension of standard Data in order to answer to a query that is applied on different relational data 65 bases form the multi-valued data sets. Multi-valued data can also be generated from Data Analysis (functional 66 analysis, clustering, associations, neural networks . . .) of standard databases, from expert knowledge (scenario 67 of traffic accidents, type of employed . . .), from time series (in the intervals of time), from private data (in the 68 view of hiding the initial data by less precision), etc. 69

70 6 Table c

There are some approaches like ribeiro95 [1] and Thompson 91 [5] identify the knowledge directly from domains that are structured. The most straight forward way to form a linear file is joining of related tables. the above table-c is the result of both of the tables table-a and table-b. Table-c is obtained by applying the natural join to both of the tables table-a and table-b which has the problem that a single object cannot be identified unique, that means a single object can be represented by more than one row. The problem as a unique object even

though an object is identified by more than one tuple.

So, Therefore this construction between a structure database and linear file is understood by many analysis 77 frame works in efficiently This paper introduces a data mining framework and a information discovery to deal 78 with this drawback of the previous linear file representations and we finally concentrate on the issues of structured 79 80 database and discovery of a set of queries that describe the features of objects in the database. To overcome 81 from this problem mentioned above it is needed to develop a better illustration procedure which with this type of representation is most of the data mining and machine learning algorithms identify each row represents 82 information for objects that are co-related with other objects. One simple solution to solve the problem is 83 to combine the associated objects into a single object by applying some aggregate operations. The issue in this 84 concept is the selection of aggregate function that can be applied and minimizing the Loss of technical information 85 because of the aggregate function so finally there forms a cluster of related objects that represent a set of objects 86 as a single object. To solve the same problem another technique is introduction of multivalued attributes for 87 the objects. The generalization of this type of linear file arrangement is as multi-valued data set. in this types 88 of data bases the values of multivalued attribute are either a single value or a bin of values these bin of values 89 represent the data that associate with the object with respect to the multi-valued attribute. Table-d consists the 90 91 multi-valued data that is generated from the two tables table-a and table-b, in this table product type, purchase 92 place and purchase amount are multi-valued attributes, the objects have a set of values in between braces in the 93 corresponding cells of these multi-valued attributes, if there is no any value we represent it by null or if there is 94 a single value we put that value as it is. Now table-d contains the unique tuples for representing the data of one person which is different from the data that is denoted by table-c. So a useful set of queries can be applied to 95 this type of data for implementing data mining tasks, because of multivalued attributes are existed in the data 96 base the following problems are encountered during this process. 97

? It is needed to identify the set of queries that access the features of objects in the database. ? Applying the
 data mining techniques to the database that consists of multi-valued attributes.

Most of the data mining techniques that are existed cannot work properly on the multi-value datasets which consists of a set of values because they are developed for applying on single valued attributes. So the need is to develop the techniques that can be applied on multi valued attributes.

¹⁰³ Multi-valued Datasets that are available in UCI machine learning repository are given in the bellow table.

¹⁰⁴ 7 Name of the Dataset

- **105 8 Number of classes**
- ¹⁰⁶ 9 Number of normal attributes

107 10 Number of multivalued attributes

108 11 Number of binary attributes

¹⁰⁹ "Promoter" "Hayes-Roth" "Breast cancer" "Monks-1" "Monks-2" "Monks-3" "Balance" "Soya large" "Tic-tac-¹¹⁰ toe" "Car" "DNA" "Mushroom" "Nursery" In general the categorization of two types of attributes in the data ¹¹¹ bases are quantitative attributes and qualitative attributes to apply data mining tasks on these databases we are ¹¹² needed to find the similarity measures among these type of attributes .

VI. Qualitative Type Tversky (1977) proposed a contrast model and ratio model that is implemented by generalizing a several set of theoretical similarity models. Tversky described the objects as sets of features as a replacement of geometric points in a geometric space. To demonstrate his models, let a and b are two objects, and X and Y indicate the sets of features associated with the objects x and y respectively. Tversky proposed the following family of proximity measures which is called the contrast model:S(x, y) = ? n (X?Y) -? n(X - Y) -?n(X - Y)

For some ?, ?, ? ? 0; n is usually the cardinality of the set. In the previous models, the proximity between objects was determined only by their common features, or only by their distinctive features. In the contrast model, the proximity of a pair of objects is represented as a linear combination of the measures of the common and the distinctive features. The contrast model denotes proximity between objects as a weighted difference of the measures for their common and distinctive features. The given bellow family of similarity measures denotes

the ratio model: S(x, y) = n(X?Y) / [n(X?Y) + ? n(X-Y) + ? n(Y-X)], ?, ? ? 0In the ratio model, the proximity value is normalized to a value range of 0 and 1. In Tversky's set theoretic

similarity models, a feature usually denotes a value of a binary attribute or a nominal attribute but it

¹²⁷ 12 VII. Attributes That are Quantitative Type

To find the proximity within the group when the attributes are multi-valued type we use group mean for the particular attribute with respect to the object by using Euclidean distance like measures but the problem with this method is it should not consider the cardinality of the elements in a group. Another approach towards this is group average which can be used to calculate inter-group proximity. In this approach group similarity is calculated by taking the average of all the inter-object measures for those pairs of objects from which each object of a pair is in diverse groups.

For example, the average dissimilarity between group P and Q can be defined as given bellowD (P, Q) = ? ??(??, ??)/?? ?? ??=1

), where n is the cardinality of object-pairs, d (a, b) is the variation metric for the i th pair of objects p
and q where p? P And q ? Q . In calculating group similarity using on group average, decision on whether we
compute the average for every probable pair of similarity or the average for a subset of possible pairs of similarity
is required.

¹⁴⁰ 13 VIII. Algorithms for Clustering

141 IX.

¹⁴² 14 Partitioning Algorithms for Clustering

The data objects in these types of methods are initially considered as a single cluster and this can be partitioned into different clusters depending on the number of clusters required. The objects are assigned into partitions by iteratively placing the points between the partitions. There are different re arrangement schemes that iteratively reassign the points among the pre defined number of clusters not like hierarchical methods clusters not allowed revisiting the previously formed clusters. These algorithms Progressively develop the quality of clusters.

In all of the partitioning algorithms the number of clusters is decided initially, if we want to find the number of clusters automatically we can think data comes from the mixture of various probability distributions .the major advantage of probabilistic approach is interpretability of the clusters that are formed. The summarized cluster representation can be done by allowing the measures of intra cluster similarity, an another approach to measure

the quality of clustering is the definition of objective function . in general partitioning algorithms constructs the

153 clusters in non-convex shape.

Some partitioning algorithms for clustering are K_means, K_modes, K_medoids (PAM, CLARA, CLARANS, and FCM). Partition based algorithms can found clusters of Non convex shapes. The pros and cons of partitioning

clustering methods are described bellow.

157 15 Advantages:

158 ? Comparatively scalable and simple.

? These algorithms are suitable for the datasets which forms spherical shaped clusters that are wellseparatedDisadvantages:

? These algorithms causes efficiency degradation in the high dimensional spaces because almost all pairs of points are at a distance of its average the reason for this is the distance between points in high dimensional spaces

¹⁶³ is not defined properly. ? The description of cluster is less explained by these algorithms.

164 ? User needs to specify the number of clusters in advance.

¹⁶⁵ ? These are in efficient to deal with non-convex shaped clusters of unreliable size and density.

¹⁶⁶ 16 X. Hierarchical Algorithms for Clustering

Hierarchical clustering algorithms are of two types they are Agglomerative (top-bottom) and Divisive (bottom-167 top). In Agglomerative Hierarchical clustering initially one object is selected and the other objects are embedded 168 consecutively to form bigger clusters the merging of objects depends on the distance like maximum, minimum and 169 170 average .this procedure is repeated number of times until to form desired number of clusters. some Hierarchical 171 algorithms for clustering are BIRCH, CURE, ROCK, Chameleon, Wards, SNN, GRIDCLUST, CACTUS in which 172 clusters of Non convex shaped and, Arbitrary Hyper rectangular are produced. Some times to find the proximity between sub sets it is needed to generalize the proximity between individual points. Such proximity measure 173 is called linkage metric. This is also the case when the data base consists of multi-valued data. This type of 174 metric impacts the Hierarchical clustering algorithms because it impacts the closeness and relation of connectivity 175 among the attributes without considering its structure. 176

¹⁷⁷ Mostly used inter-cluster linkage metrics are single link, average_link, and complete link. The principal ¹⁷⁸ proximity of dissimilarity measure (usually distance) is computed for every pair of points with one point in ¹⁷⁹ the first set and another point in the second set. A specific functional operation such as minimum Average or ¹⁸⁰ maximum is applied to pair-wise proximity of dissimilarity measures an example for such functional operation is ¹⁸¹ given bellow:D (C 1 , C 2) = functional operation {d (x,y)/ x? C 1, y? C 1 }

182 The pros and cons of partitioning clustering methods are described bellow Advantages:

These algorithms are flexible with respect to the level of granularity. similarity or distance proximities are easily managed. These algorithms can be applicable to any types of attributes.

185 17 Disadvantages:

186 It is difficult to identify the termination point. Most of the hierarchical algorithms cannot be repeated if once 187 constructed. Clusters are formed depending on the reason of their improvement.

188 47 Year 2016

189 **18** () C

In density based methods the set of points which are in the space of Euclidean are partitioned into a set of ?
More sensitive to initialization phase that means final results depends on the clusters that are formed at initial,
also sensitive to noise and outliers.

points depending on the density among the points, connectivity between each pair of points and the boundary 193 of the points that are grouped. In density based clusters each point in the cluster closely related to its adjacent 194 neighbor, the clusters are formed in any direction that the density of the points in one cluster should be maximum. 195 Because of this reason the cluster that depends on density are of arbitrary shaped as a result it provides a good 196 security against outliers. These algorithms classify the data objects into core_points, border_points and noise 197 points to form dense regions. The clusters are formed by connecting the core points together. So there may 198 be a chance of forming non spherical or arbitrary shaped clusters. Some well known density based algorithms 199 for clustering are "DBSCAN", "OPTICS", "DBCLASD", "GDBSCAN", "DENCLUE" and "SUBCLUE". When 200 the data is of numeric type then only density based and partition based methods give results efficiently but 201 grid based techniques perform well for the attributes with different data types along with multivalued data. 202 For better performance the density based method "DENCLUE" initially uses grid structure. To structure real 203 clusters Grid based algorithms uses subspace and hierarchical techniques for clustering. When compare to all 204 Clustering algorithms Grid algorithms are very speedy and efficient processing algorithms. Arbitrary shaped 205 clusters are formed by Adaptive grid algorithms such as MAFIA and AMR by using the grid cells. 206

207 19 XV. Model Clustering Algorithms

A model is hypothesized for each of the clusters and tries to find Set of data points are related together based on different strategies like conceptual methods and statistical methods. For model based algorithms the well known approaches are one is neural network approach and another one is statistical approach. Algorithms such
as "EM", "CLASSIT", "COBWEB SOM", and "SLINK" are some Model based clustering algorithms.

212 20 XVI. Research Challenges

We have previously discussed that the problem of clustering becomes very demanding, attractive and at the same 213 time challenging, when the data is of type categorical attributes and multi-valued attributed. The number of 214 algorithms for the discovery of groups in such data is restricted, confined and limited, compared to the research 215 devoted on data sets with numerical data. Further, few algorithms deal with mixtures of values, i.e., attributes 216 of numerical and categorical The methods that partition the space into different cells or segments are called grid-217 based methods. Here each segment is related to a single value or a group of values (if the data is multi-valued 218 data) is considered as a unit, as a result the concentration is turned on to the partition space. The partition is 219 based on the characteristics of the grid that is formed from the data. The advantage of this type of arrangement 220 of data in the form of grids is it does not depend on the order of the points that are arranged in the form of grids. 221

222 **21** Global

223 22 XVII. Summary and Conclusion

In this paper, we analyzed the difficulty of generating single flat file arrangement to represent Data sets that have 224 been generated from structured databases, and pointed out its inappropriate representation to represent related 225 information, a fact that has been frequently overlooked by recent data mining investigation. To overcome these 226 difficulties, we used a better representation scheme, called multi-valued data set, which allows attributes of an 227 object to have a set of values, and studied how existing similarity measures for single-valued attributes could be 228 applied to measure group similarity for multi-valued data sets in clustering. We also proposed a unified framework 229 for proximity Once the target database is grouped into clusters with similar properties, the discriminate query 230 231 detection system, MASSON can find out useful characteristic information for a set of objects that fit in to a cluster. We claim that the planned representation scheme is suitable to cope with related information and that 232 it is more communicative than the traditional single flat file format. 233

More prominently, the relationship information in a structured database is actually considered in clustering procedure.



Figure 1: 1 C



Figure 2: Figure 1

ab

pan num- ber	City	Mode of payment	Date of pay- ment	Amount in Bs
11231	Visakhapatnam(V)	Physical cash	1/06/13	500.00
11231	Hyderabad(H)	Credit to bank	10/07/13	80.00
11231	Nellore(N)	Check	11/07/13	100.00
14561	Visakhapatnam(V)	Physical cash	19/08/13	400.00
14561	Anakapalli(A)	Credit to bank	20/08/13	200.00
17891	Anakapalli(A)	Credit to bank	5/09/13	50.00
11011	Vizianagaram(Vz)	Physical cash	8/09/13	25.00

Figure 3: Table a Table b

	1	I	
(J	L	
		_	

Name	Age	Product type category	Purchase place	Purchase
				amount
A.Ajay	21.0	{physical cash,credit to	Male {Visakhapatnam(V),Hyden	ra{15000(00),80.00,100.00}
kumar		bank, check	Nellore(N)	
A.Vijay	22.0	bank $\}$ { cash, credit to	Male {Visakhapatnam(V),anaka	1p{a100(.00),200.00}
kumar		physical		
A.Arun	23.0	{Credit credit to bank}	Male Anakapally(A)	50.00
kumar				
A.Varun	24.0	bank $\}$ { cash, credit to	Male {vizianagaram(Vz),Visakha	a \$25n90nf(W)\$ }
kumar		physical		

Figure 4: Table d :

 \mathbf{e}

[Note: . Similarity Measures for Multi-Valued Data Sets]

Figure 5: Table e V

²³⁶ .1 This page is intentionally left blank

- 237 [Gower ()] 'A general coefficient of similarity and some of its properties'. J C Gower . Biometrics 1971. 27 p. .
- 238 [Ryu et al. (1998)] 'Automated Discovery of Discriminate Rules for a Group of Objects in Databases'. T W Ryu
- 239 , C F Eick , O ; R , P E Duda , D G Hart , Stork . Conference on Automated Learning and Discovery,
- (Pittsburgh, PA; New York, NY, USA) 1998. June 11-13. 2012. Wiley. Carnegie Mellon University (Pattern Classification)
- [Everitt ()] Cluster Analysis, Edward Arnold, Copublished by Halsted Press and Imprint, B S Everitt . 1993. John
 Wiley & Sons Inc. (3rd edition)
- ²⁴⁴ [Thompson et al. ()] 'Concept formation in structured domains'. K Thompson , P Langley , D Fisher , M Pazzani
- P Langley , Morgan Kaufmann . Concept Formation: Knowledge and Experience in Unsupervised Learning,
 1991.
- 247 [Ryu et al. ()] 'Deriving Queries from Results using Genetic Programming'. T Ryu, C F Eick, Oregon J Portland
- A B Pei , X Jiang , Y Lin , Yuan . Procee dings of the 2nd Int'l Conf. on Knowledge Discovery and Data Mining, (ee dings of the 2nd Int'l Conf. on Knowledge Discovery and Data Mining) 1996a. 2007. (Proc. 33rd
- Int'l Conf. Very Large Data Bases (VLDB))
- [Li and Pang (2010)] 'Deterministic column-based matrix decomposition'. X Li , Y Pang . *IEEE Trans. Knowl. Data Eng* Jan. 2010. 22 (1) p. .
- [Nie et al. ()] 'Efficient and robust feature selection via joint L2, 1-norms minimization'. F Nie , H Huang , X
 Cai , C Ding . *Proc*, (null) 2010. 23 p. .
- 255 [Cheng et al. ()] 'Evaluating Probabilistic Queries over Imprecise Data'. R Cheng , D V Kalashnikov , S
- Prabhakar . Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD), (ACM SIGMOD Int'l Conf.
 Management of Data (SIGMOD)) 2003.
- 258 [Tversky (1977)] Features of similarity, Psychological review, A Tversky . 1977. July. 84 p. .
- [Koza ()] Genetic Programming: On the Programming of Computers by Means of Natural Selection, John R Koza
 . 1990. Cambridge, MA: The MIT Press.
- [Tao et al. ()] 'Indexing Multi-Dimensional Uncertain Data with Arbitrary Probability Density Functions'. Y
 Tao , R Cheng , X Xiao , W K Ngai , B Kao , S Prabhakar . Proc. Int'l Conf. Very Large Data Bases
 (VLDB), (Int'l Conf. Very Large Data Bases (VLDB)) 2005.
- [Ribeiro et al. ()] 'Knowledge Discovery from Multiple Data bases'. J S Ribeiro , K Kaufmann , L Kerschberg .
 Proc. Of the 1st Int'l Conf. On Knowledge Discovery and Data Mining, (Of the 1st Int'l Conf. On Knowledge Discovery and Data MiningQuebec, Montreal) 1995.
- [Yang et al. ()] 'L2, 1-norm regularized discriminative feature selection for unsupervised learning'. Y Yang , H T
 Shen , Z Ma , Z Huang , X Zhou . Proc. Int. Joint Conf, (Int. Joint Conf) 2011. 22 p. 1589.
- [He et al. ()] 'Laplacian score for feature selection'. X He, D Cai, P Niyogi. Proc, (null) 2005. 186 p. 189.
- 270 [Ryu et al. ()] 'MASSON: Discovering Commonalities in Collection of Objects using Genetic Programming'. T
- 271 Ryu, C F Eick, D V Cheng, S Kalashnikov, X Prabhakar, D He, P Cai, Niyogi. Proceedings of the
- 272 Genetic Programming 1996 Conference, (the Genetic Programming 1996 ConferenceSan FranciscoR) 1996b.
 2005. 186 p. 189. Stanford University (Proc.)
- 274 [Duda et al. ()] Pattern Classification, R O Duda, P E Hart, D G Stork. 2012. New York, NY, USA: Wiley.
- [Pei et al. ()] 'Probabilistic Skylines on Uncertain Data'. J Pei , B Jiang , X Lin , Y Yuan . Proc. 33rd Int'l Conf.
 Very Large Data Bases (VLDB), (33rd Int'l Conf. Very Large Data Bases (VLDB)) 2007.
- [Kotsiantis and Pintelas ()] 'Recent Advances in Clustering: A Brief Survey'. S B Kotsiantis , P E Pintelas .
 WSEAS Trans. Information Science and Applications 2004. 11 (1) p. .
- [Liu et al. ()] 'Transductive component analysis'. W Liu , D Tao , J Liu . Proc. 8th IEEE Int. Conf. Data Mining,
 (8th IEEE Int. Conf. Data Mining) 2008. p. .
- 281 [Cai et al. ()] 'Unsupervised feature selection for multi-cluster data'. D Cai , C Zhang , X He . Proc. 16th ACM
- 282 SIGKDD Int. Conf. Knowl. Discovery Data Mining, (16th ACM SIGKDD Int. Conf. Knowl. Discovery Data 283 Mining) 2010. p. .