

Clustering Data from Heterogeneous Dissimilarities

Prof. Sheetal More, Prof. Bajirao Shirole, Roshani Balkrushna Derle, Rinita Dilip Jadhav, Shraddha Chandrabhan Jadhav, Jyoti Yadav Kadale

Computer Engineering Department, Sanghavi College of Engineering, Pune University, Nashik, Maharashtra, India

ABSTRACT

The clustering model which, to handle the heterogeneity, uses all available dissimilarity matrices and identifies for groups of individuals clustering objects in a similar way. The model is a non-convex problem and difficult to solve exactly, and we thus introduce a Variable Neighborhood Search heuristic to provide solutions efficiently. In our Proposed System we use clustering mechanism to create groups of given heterogeneous datasets, in this system we process heterogeneous data like html and xml as well as numeric data and convert them to single vector by using correlation of values and then this single vector will be clustered by corr-k mean algorithm. output will be number of unlabelled clusters and these clusters will be more precise than what produced by Existing system. In the scope Find centroids and create clusters by automatic clustering method not by iterative method. It can be used to cluster different types of data.

Keywords: Clustering, Heterogeneous, Centroids, Heuristic, K Mean Algorithm.

I. INTRODUCTION

The Clustering algorithms determine groups of objects in such a way that objects in the same group, called clusters, only one dissimilarity matrix is available. For instance the Iris dataset (Fisher, 1936), The use of classical clustering algorithms (e.g. k -means, single-linkage, complete- linkage) on this dataset can provide excellent results. It is however possible that more than one dissimilarity matrix is available. In the context of the Iris dataset, one could envision asking a sample of multiple experts to measure the same owners, in case there has been significant measurement error. If there is heterogeneity in the data reported, clustering mechanism to create groups of given dataset, previous clustering mechanisms work on homogenous data like set of html pages but in this system we are processing heterogeneous data like set html pages and set of xml pages, we process these datasets separately using corresponding preprocessing mechanism and then we convert them to numeric values by identifying tf, tfidf and then we apply clustering algorithms like K-means clustering to define clusters. In our Proposed System we use clustering mechanism to create groups of given heterogeneous datasets, in this system we process heterogeneous data like html and xml as well as numeric

data and convert them to single vector by using correlation of values and then this single vector will be clustered by corr-k mean algorithm.

Motivation of the Project

Clustering of heterogeneous data provides better clusters containing mixed type of data, this is improved grouping than homogeneous clustering. Improved Clustering technique. In this system we process heterogeneous data like html and xml as well as numeric data and convert them to single vector by using correlation of values and then this single vector will be clustered by corr-k mean algorithm.

The organization of this document is as follows. In Section 2 gives literature survey, Section 3 gives details of system architecture. In Section 4 presents research findings and your analysis of those findings. Section 5 concludes the paper.

II. LITERATURE SURVEY

In our Existing System we use clustering mechanism to create groups of given dataset, previous clustering mechanisms work on homogenous data like set of html pages but in this system we are processing

heterogeneous data like set html pages and set of xml pages, we process these datasets separately using corresponding pre-processing mechanism and then we convert them to numeric values by identifying tf, tfidf and then we apply clustering algorithms like K-means clustering to define clusters. Each Dataset will have separate clusters created as an output, after all datasets converted to clusters, then we merge these clusters as a final clusters. In our attempts to solve the HCP exactly, we used Couenne. The first two are different implementations of the spatial Branch-and- Bound (sBB) algorithm (Liberti, 2006) for nonconvex mixed-integer nonlinear problems (MINLP). Much like a Branch-and-Bound (BB) algorithm for MIPs, sBB explores the feasible space exhaustively but implicitly, finding a guaranteed ϵ - approximate solutions for any given $\epsilon > 0$ in finite (potentially exponential) time. Unlike MIPs, whose continuous relaxation is a linear program, and unlike convex MINLPs, whose continuous relaxation of a nonconvex MINLP is usually difficult to solve. To find this issue, sBB algorithms used to solve convex relaxations of the given MINLP. The convexity gap between the original MINLP and its convex relaxation therefore stems from two factors: the relaxation of the integrality constraints and the relaxation of the nonconvex terms appearing in the MINLP. The third generic solver GloMIQO is a branch-and-cut algorithm based on generating tight convex relaxations from detecting special structures such as convexity and edge-concavity. The algorithm is specialized to address MIQCQPs to ϵ -optimality.

In the VNS framework, the neighborhoods are defined around types of moves, or perturbations, of the best current solution x –the center of the search. When looking for a better one in a minimization problem, a solution x_{new} is drawn at random in an increasingly wider neighborhood, and a local descent is performed from x_{new} leading to another local optimum x_{new} . If x_{new} is worse than x , then x_{new} is ignored and one chooses a new neighbor solution x_{new} in a more distant neighborhood of x . If instead x_{new} is better than x , the search is re-centered around x_{new} and the local search restarts in the closest neighborhood of the newly found best current solution. Once all neighborhoods of x have been explored without success, one begins again with the closest one to x , until a stopping condition (e.g. maximum CPU time) is met. As the size of neighborhoods tends to increase

with their distance from the current best solution x , close-by neighborhoods are explored more thoroughly than far away ones. This strategy takes advantage of the three observations 1–3 mentioned above, and yet can ensure with sufficient computational time that the algorithm is not stuck in a poor local optimum. We now turn to our implementation of VNS for the HCP.

III. PROPOSED SYSTEM

In our System we use clustering mechanism to create groups of given heterogeneous datasets, in this system we process heterogeneous data like html and xml as well as numeric data and convert them to single vector by using correlation of values and then this single vector will be clustered by corr-k mean algorithm. output will be number of unlabelled clusters and these clusters will be more precise than what produced by Existing system.

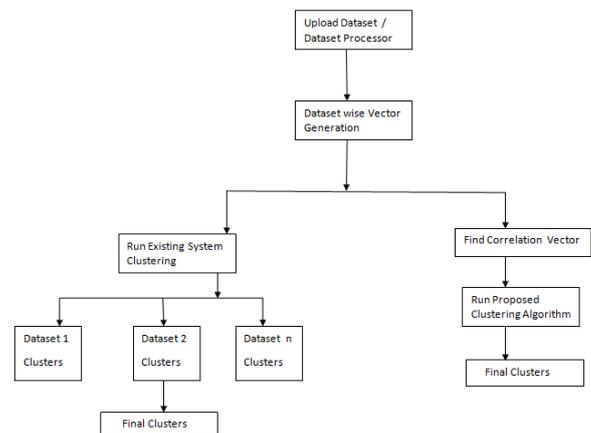


Figure 1: System Architecture

Heterogeneous Data consists of data with many dissimilarities, mostly data consists of objects with many properties and objects can be differentiated on different parameters. This data can be represented in a matrix form where n objects and total m dissimilarities form rows and columns where we can create or define clusters which have most common properties or most common 1's in rows, now again cluster can have separate matrix of objects and their properties we can create sub clusters and same process applied until no more or minimum dissimilarities are there. we can send query to level one cluster the to next level and to another level this reduces search time comparisons.

A. Dataset Processor:

Allows to Upload Dataset, Preprocessing of dataset.

B. Html Data Clustering:

Retrieve text and design term vector and find idf vector , output is clustered html data.

C. Xml-Tags Data Clustering:

Retrieve XML tag values and find tags vector output is Clustered xml documents.

D. Correlation based Clustering:

Proposed System finds correlation vector and find clusters on vector, output is clusters of mixed data

E. Result Analysis:

Clustering time and clustering levels comparison

IV. RESULTS

TABLE 1: System running time in ms

System	Running Time(in ms)
Html clustering	13558
Xml clustering	6158
Html + xml correlation	24781
Html + xml	16994
Html clustering R-K mean	5969
xml clustering R-K mean	5736
Html + xml correlation R-K mean	11935
Html + xml R-K mean	13968

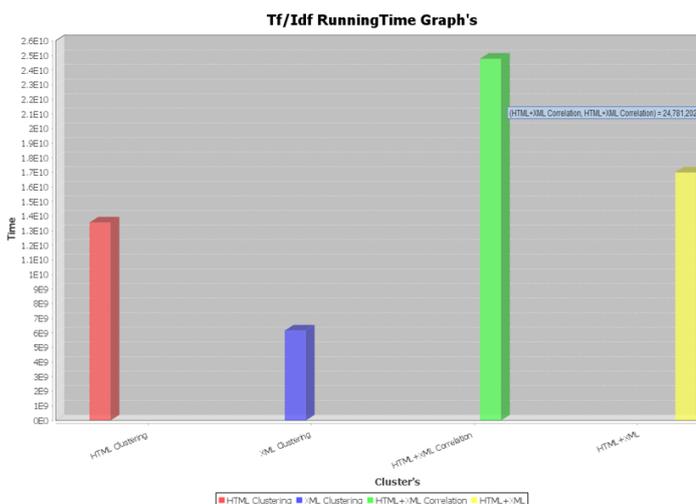


Figure 2: Tf/Idf running time graph

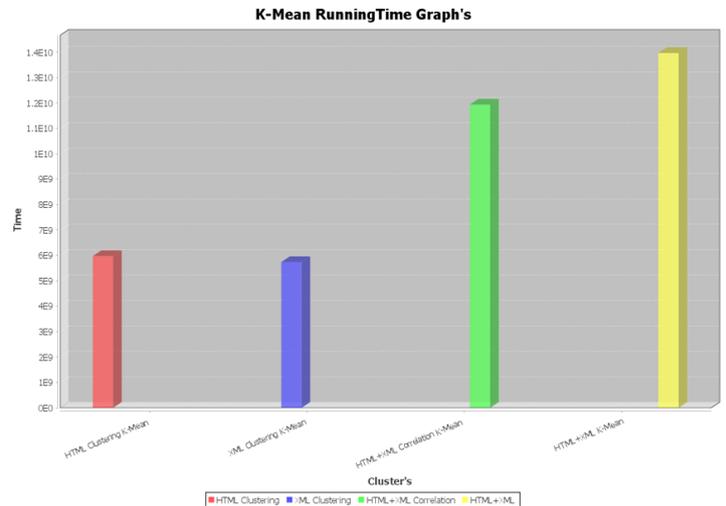


Figure 3: K-mean Running Time graph

V. CONCLUSIONS

Clustering of heterogeneous data provides better clusters containing mixed type of data, this is improved grouping than homogeneous clustering. clustering mechanism to create groups of given heterogeneous datasets, in this system we process heterogeneous data like html and xml as well as numeric data and convert them to single vector by using correlation of values and then this single vector will be clustered by corr-k mean algorithm. output will be number of unlabelled clusters and these clusters will be more precise than what produced by Existing system. Reduced Clustering time.

VI. REFERENCES

- [1]. Anstreicher, K. M. (2012). On convex relaxations for quadratically constrained quadratic programming. *Mathematical Programming* , 136 (2), 233–251 . Arora, R. (1982).
- [2]. Consumer involvement in retail store positioning. *Journal of the Academy of Marketing Science* , 10 (1-2), 109–124 . Audet, C. , Hansen, P. , Jaumard, B. , & Savard, G. (2000). A branch and cut algorithm for nonconvex quadratically constrained quadratic programming. *Mathematical Programming* , 87 (1), 131–152 . Avella, P. , Boccia, M. , Salerno, S. , & Vasilyev, I. (2012).
- [3]. An aggregation heuristic for large scale p-median problem. *Computers & Operations Research* , 39 (7), 1625–1632 .
- [4]. Belotti, P. , Lee, J. , Liberti, L. , Margot, F. , & Wächter, A. (2009). Branching and bounds

- tightening techniques for non-convex Minlp. *Optimization Methods and Software* , 24 (4–5), 597–634 . Bettman, J. R. , Luce, M. F. , & Payne, J. W. (1998).
- [5]. Constructive consumer choice processes. *Journal of Consumer Research* , 25 (3), 187–217 . Bettman, J. R. , & Park, C. W. (1980). Effects of prior knowledge and experience and phase of the choice process on consumer decision processes: A protocol analysis. *Journal of Consumer Research* , 7 , 234–248 . Bijmolt, T. H. , & Wedel, M. (1995).
- [6]. The effects of alternative methods of collecting similarity data for multidimensional scaling. *International Journal of Research in Marketing* , 12 (4), 363–371 . Bijmolt, T. H. , Wedel, M. , Pieters, R. G. , & DeSarbo, W. S. (1998).
- [7]. Judgments of brand similarity. *International Journal of Research in Marketing* , 15 (3), 249–268 . Billionnet, A. , Elloumi, S. , & Lambert, A. (2016). Exact quadratic convex reformulations of mixed-integer quadratically constrained problems. *Mathematical Programming* . in press.
- [8]. Blanchard, S. , & Banerji, I. (2016). Evidence-based recommendations for designing free-sorting experiments. *Behavior Research Methods* . in press. Blanchard, S. J. , Aloise, D. , & DeSarbo, W. S. (2012a). The heterogeneous p-median problem for categorization based clustering. *Psychometrika* , 77 (4), 741–762 . Blanchard, S. J. , & DeSarbo, W. S. (2013). A new zero-inflated negative binomial methodology for latent category identification. *Psychometrika* , 78 (2), 322–340 .
- [9]. Blanchard, S. J. , DeSarbo, W. S. , Atalay, A. S. , & Harmancioglu, N. (2012b). Identifying consumer heterogeneity in unobserved categories. *Marketing Letters* , 23 (1), 177–194 . Bomze, I. M. , & Locatelli, M. (2004). Undominated dc decompositions of quadratic functions and applications to branch-and-bound approaches. *Computational Optimization and Applications* , 28 (2), 227–245 .
- [10]. Brusco, M. J. , & Cradit, J. D. (2001). A variable-selection heuristic for k-means clustering. *Psychometrika* , 66 (2), 249–270 . Brusco, M. J. , & Cradit, J. D. (2005). Conpar: a method for identifying groups of concordant subject proximity matrices for subsequent multidimensional scaling analyses.
- [11]. *Journal of Mathematical Psychology* , 49 (2), 142–154 . Brusco, M. J. , Steinley, D. , Cradit, J. D. , & Singh, R. (2012). Emergent clustering methods for empirical OM research. *Journal of Operations Management* , 30 (6), 454–466 . Carpenter, G. S. , & Nakamoto, K. (1994).
- [12]. Reflections on “consumer preference formation and pioneering advantage”. *Journal of Marketing Research* , 31 (4), 570–573 . Coxon, A. P. M. (1999). *Sorting data: Collection and analysis*: 127. Sage . DeSarbo, 300.