

SUPERVISED CLUSTERING FOR CLASSIFICATION WITH FEATURE SPACE HETEROGENEITY

R.Pavani¹, K.Mounasree², B.Kiranmai³, J.V.Kirankumari⁴

Abstract: *Feature space heterogeneity often exists in many real world data sets so that some features are different importance for classification over different subsets. Moreover, the pattern of feature space heterogeneity might dynamically change over time as more and more data are accumulated. In this paper, we develop an incremental classification algorithm, Supervised Clustering for Classification with Feature Space Heterogeneity (SCCFSH), to address this problem. In our approach, supervised clustering is implemented to obtain a number of clusters such that samples in each cluster are from the same class. After the removal of outliers, relevance of features in each cluster is calculated based on their variations in this cluster. The feature relevance is incorporated into distance calculation for classification. The main advantage of SCCFSH lies in the fact that it is capable of solving a classification problem with feature space heterogeneity in an incremental way, which is favourable for online classification tasks with continuously changing data. Experimental results on a series of data sets and application to a database marketing problem show the efficiency and effectiveness of the proposed approach.*

Key words: *Classification, Clustering, Heterogeneity, Feature Relevance.*

I. INTRODUCTION

In classification problems, feature space heterogeneity is the phenomenon that a data set consists of some heterogeneous subsets, and the optimal features for classification are distinct over different subsets. The challenge of this problem is that we do not know how many heterogeneous subsets exist in the data set or which subset each sample belongs to. In the last decade, the problem of feature space heterogeneity in data has been addressed under different names, such as local feature relevance, case-specific feature weights, relevance in context, feature space and class heterogeneity, and attribute instability.

Feature space heterogeneity exists widely in various application fields of classification techniques, such as marketing, customs inspection decision, credit scoring, and medical diagnosis. For example, in marketing, a major concern of the market managers is to develop and implement efficient marketing programs by fully utilizing the customer databases and identifying the households that are most likely to be interested in the marketing programs. The above process can be formulated as a classification problem, in which the Features (attributes) are characteristics of the households such as demographic, psychographic, and behavioural information, and the target variable is whether a household responds to the marketing messages.

II. LITERATURE REVIEW

The phenomena that relevant features for classification vary across the data set have been observed by many researchers and practitioners. Until recently, a number of classification methods have been developed, which can be divided into two categories. In the first category, one of the best known methods is “bagging” [18]. In this approach, k subsets are generated by randomly sampling from the original set of samples. Consequently, relevant features might be different in the obtained subsets. Based on this approach, Puuronen. Proposed a Meta-Level Classification (MLC) method, which can be used to deal with the problem of feature space heterogeneity. MLC first divides the training sample set into some subsets and obtains the component classifiers based on these subsets. In the application phase, testing samples are put into the training sample set, and MLC dynamically selects the optimal component classifier for a testing sample by comparing the performance of different classifiers in its neighbourhood. Different from the method of sample partitioning, the Random Subspace Method (RSM) [20] divides the whole feature set into a number of feature subsets and constructs different classifiers based on the whole training samples with different feature subsets obtained. Feature space heterogeneity in testing samples is considered through synthesizing (usually by voting) the application results of all classifiers. These methods deal with the problem of feature space heterogeneity by firstly dividing sample set or feature set into different subsets in a random way and then training component classifiers in the subsets. These component classifiers are then combined for classification, mostly by major voting or selecting the optimal one. A major problem of these methods lies in the random set (sample set or feature set) partitioning, which may result in seriously biased component classifiers due to the feature redundancy and irrelevance in some subsets, especially for high dimensional data sets. In the second category, modified lazy learning methods are applied to classification problems with feature space Mathematical Problems in Engineering 3 heterogeneity. Friedman addresses the problem of feature space heterogeneity by investigating the variability of feature relevance in different data subsets. In his method, the local relevance of features in each subset is measured by the estimated reduction in classification error. Hastie and Tibshirani developed an adaptive form of nearest neighbour classification method for dealing with feature space heterogeneity. In their approach, distance metric for each sample is adaptively calculated in an iterative process using local discriminative information of features. Therefore, different relevant features are taken into account for classification in different subsets. Although both works report favourable results on their local approaches

compared to global ones, both of them are computationally expensive. Aredes and Vidal propose a locally weighted lazy learning approach for better classification accuracy. In their method, different samples would have different feature weights obtained by approximately minimizing the Leaving-One-Out (LOO) classification error of the given training set. However, the computational complexity of this method is high because of the gradient descent algorithm employed to search for the optimal weights. In spite of the fact that many researchers have been carried out for dealing with feature space heterogeneity in classification, we have not found any for incremental learning among them.

Researches on incremental classification are mainly focused on statistical methods, neural networks, and evolutionary algorithm. Instance-based learning, especially K-nearest neighbour (K-NN) learning, is a widely used nonparametric incremental classification approach where training or learning does not take place until a query is made. In contrast to complex learning algorithms such as neural networks or support vector machines, K- NN learning does not require a complex function fitting process or model training procedure. Thus, it is easy to do incremental learning. Nevertheless, once a query point with unknown class label is presented, conventional K-NN learning traverses the whole data set to find the K nearest neighbours of the query point. Therefore, the computational time and requirement of computer storage space of K-NN are not scalable to large amounts of data.

To solve this problem, Li and Ye propose a data mining algorithm based on supervised clustering to learn data patterns and use these patterns for classification. This algorithm enables a scalable and incremental learning of patterns from data with both numeric and nominal variables. However, it calculates the feature relevance by using squared correlation coefficient between predictor variables and target variable over the entire data set, regardless of the possible heterogeneity that exists in the feature space.

The Proposed Approach the Supervised Clustering for Classification with Feature Space Heterogeneity (SCCFSH) proposed in this paper is based on the ECCAS. However, SCCFSH differs significantly from ECCAS in that it takes feature space heterogeneity into consideration. SCCFSH first divides the data set into a number of subsets in a supervised way and then explores the feature relevance in each subset obtained. The main steps of SCCFSH include grid-based supervised clustering, supervised grouping of clusters, removal of outliers, calculation of feature relevance in each cluster, and distance based classification.

III. PROPOSED WORK

Proposed Approach is to develop a Supervised Clustering for Classification with Feature Space Heterogeneity (SCCFSH) to address this problem and consists of four main steps: grid-based supervised clustering, supervised hierarchical grouping of clusters, feature relevance evaluation in each cluster, and weighted distance calculation for classification.

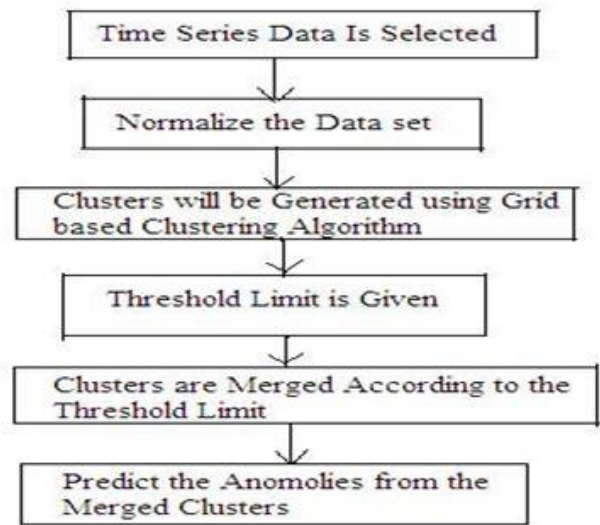


Fig 1: Flowchart of proposed work

- Grid-based supervised clustering: The grid-based supervised clustering procedure first divides the - dimensional space of samples into grid cells and then generates clusters within the grid cells
- Supervised hierarchical grouping of clusters: Supervised grouping of clusters plays an important role in the proposed SCCFSH. If some underlying clusters cover the area of several grid cells, grid-based supervised clustering would divide these large clusters into small clusters. Therefore, refinement of the clustering results is needed.
- Feature relevance evaluation in each cluster: After the above three steps, we could obtain a number of clusters, and samples in each cluster are with the same class label. Vucetic and Obradovic [29] argue that some features with the same values may result in quite different outputs (class labels) in different regions. Therefore, spatial characteristics of samples should be explored for better classification performance. Motivated by this statement, we investigate the feature relevance in the obtained clusters that represent different spatial distributions of samples.
- Weighted distance calculation for classification

Algorithm 1: Main procedure of supervised grouping of clusters

(INPUT: Clusters $C_1, C_2 \dots C_k$ obtained by grid-based supervised clustering.

OUTPUT: $R = \{New\ 1, New\ 2 \dots C_{New\ m}\}$. (1) Set $R =, R_0 = \{C_1, C_2 \dots C_3\}$ and $t=1$;

(2) Among all possible pairs of clusters (O_r, O_s) in R_{t-1} find the one, say (O_i, O_j) , such that $d(O_i, O_j) = \min_{r,s} \{d(O_r, O_s)\}$;

(3) IF Label $(O_i) =$ Label (O_j) THEN set $O_q = O_i \cup O_j$ and produce the new clustering

$R = (R_{t-1} - \{O_i, O_j\}) \cup O_q$; ELSE $R_t = R_{t-1} - \{O_i, O_j\}$, $R = R \cup O_i \cup O_j$; Set $t = t + 1$;

(4) Repeat steps 2-3 until $R_t = R_{t-1}$, return R.

Algorithm 2: Main steps of SCCFSH

INPUT: Training data set $S1 = \{X1, X2, \dots, XN\}$ with label set $L = \{y1, y2, \dots, yN\}$; Testing data set $S2 = \{XN+1, XN+2, \dots, XN+M\}$ with unknown labels.

OUTPUT: Predicted labels $\hat{y}^{N+1}, \hat{y}^{N+2}, \dots, \hat{y}^{N+M}$ for $S2$.

(1) Apply the grid-based supervised clustering to $S1$ and L ;

(2) Apply the algorithm shown in Algorithm 1 to Clusters $C1, C2, \dots, Ck$ obtained in step 1;

(3) Calculate relevance of each feature $R_k(d)$, ($d = 1, 2, \dots, p$) in $C_{New 1}, C_{New 2}, \dots, C_{New m}$ obtained in step 2;

(4) LOOP for each data point X_j , ($j = N + 1, N + 2, \dots, N + M$) in $S2$: (i). Calculate the distance between X_j , and $C_{New i}$, ($i = 1, 2, \dots, m$) obtained in step 2, by using the distance metric defined in (6); (ii). Set \hat{y}_j as the label of the nearest cluster ;

(5) Output the predicted labels $\hat{y}^{N+1}, \hat{y}^{N+2}, \dots, \hat{y}^{N+M}$ for $S2$.

Cluster	SENS.	SENS.	SENS.	SENS.	SENS.	SENS.	DEPSE	Prd	Dev%	Reco	isHead?	Anoma.
Cluster-0	164.14	30.3039	41.56	40.78	19.841	4.8016	75.16	75.16	0.0	0	true	false
Cluster-1	159.083	29.8316	41.56	40.78	19.9536	4.7184	77.56	75.16	3.0944	1977	false	true
Cluster-2	170.75	24.6072	41.56	40.78	19.9224	4.6895	73.5423	76.36	-3.8314	2239	false	true
Cluster-3	175.83	25.6072	41.56	40.78	19.9322	5.2092	73.5487	75.4206	-2.5453	2240	false	true
Cluster-4	161.532	25.0556	41.56	40.7822	19.9459	4.7268	72.7599	74.9527	-3.0138	2260	false	true
Cluster-5	165.10	24.8769	41.56	40.8342	19.9377	4.9082	72.7599	74.5142	-2.4111	2267	false	true
Cluster-6	167.64	26.2678	41.56	40.7664	19.9447	4.6884	72.7599	74.2218	-2.0002	2261	false	true
Cluster-7	165.83	24.0403	41.56	41.058	19.9364	4.7799	72.7599	74.0113	-1.7222	2269	false	false

Fig 5: Retrieving Cluster Report after Merging

IV. CONCLUSION

In this paper, we develop a Supervised Clustering for Classification with Feature Space Heterogeneity (SCCFSH) to address the problem in a scalable and incremental way. In this study, we considered Time series dataset which can be normalized and then by applying Grid based supervised clustering method clusters will be generated. Based on the threshold value clusters will be merged. In spite of the fact that we only consider classification problems. By this an employee can easily identify the anomalies and generate cluster reports.

REFERENCES

- [1] Y. Liu, Y. Liu, and K. C. C. Chan, "Dimensionality reduction for heterogeneous dataset in rushes editing," Pattern Recognition, vol. 42, no. 2, pp. 229–242, 2009.
- [2] J.-T. Wong and Y.-S. Chung, "Analyzing heterogeneous accident data from the perspective of accident occurrence," Accident Analysis and Prevention, vol. 40, no. 1, pp. 357–367, 2008. Mathematical Problems in Engineering 9
- [3] X. Li and N. Ye, "A supervised clustering and classification algorithm for mining data with mixed variables," IEEE Transactions on Systems, Man, and Cybernetics A, vol. 36, no. 2, pp. 396–406, 2006.
- [4] Scrypnik and T. K. Ho, "Feature selection and training set sampling for ensemble learning on heterogeneous data," Tech. Rep., DIMACS, 2003.
- [5] M. K. Lim and S. Y. Sohn, "Cluster-based dynamic scoring model," Expert Systems with Applications, vol. 32, no. 2, pp. 427–431, 2007.
- [6] R. Paredes and E. Vidal, "Learning weighted metrics to minimize nearest-neighbour classification error," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 28, no. 7, pp. 1100–1110, 2006.

supervised clustering for classification with Feature Space Heterogeneity

Cluster=FullData.csv
 Data Size: 5847
 CMM: L88T7
 Min Values=[18.7211, -58.9912, 38.14, 38.3832, 19.8486, 3.8488, 97.1798]
 Max Values=[24.1598, 62.0305, 48.38, 48.38, 25.4043, 1.1882, 95.9571]
 Mean Max Values=[0.2, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0]
 Mean Min Values=[0.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0]
 Mean Max Values=[0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0]
 Mean Min Values=[1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0]

RAW DATASET

F:0=[158.1428, 30.3039, 41.56, 40.78, 19.841, 4.8016, 75.16]
 F:1=[159.083, 29.8316, 41.56, 40.78, 19.9536, 4.7184, 77.56]
 F:2=[170.75, 24.6072, 41.56, 40.78, 19.9224, 4.6895, 73.5423]
 F:3=[175.83, 25.6072, 41.56, 40.78, 19.9322, 5.2092, 73.5487]
 F:4=[161.532, 25.0556, 41.56, 40.7822, 19.9459, 4.7268, 72.7599]
 F:5=[165.1, 24.8769, 41.56, 40.8342, 19.9377, 4.9082, 72.7599]
 F:6=[167.64, 26.2678, 41.56, 40.7664, 19.9447, 4.6884, 72.7599]
 F:7=[165.83, 24.0403, 41.56, 41.058, 19.9364, 4.7799, 72.7599]

Fig 2. Visualising Raw Data set

supervised clustering for classification with Feature Space Heterogeneity

Cluster=FullData.csv
 Data Size: 5847
 CMM: L88T7
 Min Values=[-0.5000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000]
 Max Values=[0.5000, 0.5000, 0.5000, 0.5000, 0.5000, 0.5000, 0.5000]
 Mean Max Values=[0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000]
 Mean Min Values=[0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000]
 Mean Max Values=[0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000]
 Mean Min Values=[0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000]

RAW DATASET

F:0=[0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000]
 F:1=[0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000]
 F:2=[0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000]
 F:3=[0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000]
 F:4=[0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000]
 F:5=[0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000]
 F:6=[0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000]
 F:7=[0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000, 0.0000]

Fig 3. Normalized Data set

Input

Enter Similarity Threshold

0.12

OK Cancel

Fig 4: Selecting Threshold Limit

- [7] S.-U. Guan and S. Li, "Incremental learning with respect to new incoming input attributes," *Neural Processing Letters*, vol. 14, no. 3, pp. 241–260, 2001.
- [8] S.-U. Guan and F. Zhu, "An incremental approach to genetic algorithms-based classification," *IEEE Transactions on Systems, Man, and Cybernetics B*, vol. 35, no. 2, pp. 227–239, 2005.