

GMADM-based attributes selection method in developing prediction model

Sue-Fen Huang · Ching-Hsue Cheng

Published online: 17 May 2012
© Springer Science+Business Media B.V. 2012

Abstract Attribute Selection is an important issue for developing a prediction model, however, how to determine an effective attribute selection algorithm is an important but difficult issue. Attribute selection can effectively delete the irrelevant and redundant attributes to increase the prediction accuracy, and evaluating attribute selection methods usually need to consider several criteria such as accuracy, type I error, and type II error. In this paper, the selected attribute process is modeled as a group multiple attributes decision making (GMADM) problem. In evaluating different GMADM methods, the most results usually are consistently, But there are some situations where the evaluated outcomes have different results. The GMADM method is useful tool for evaluating attribute selection algorithms, and the TOPSIS is capable of identifying a compromised solution when different GMADM method result in conflicting rankings. Therefore, this paper proposes an objective (persuasive) GMADM-based attributes selection method to solve this disagreement and help decision makers pick the most suitable method. After verification, the proposed model is more persuasive to evaluate the attributes selection methods for developing prediction model.

Keywords Attribute selection · Technique for order performance by similarity to ideal solution · Group multiple attributes decision making

S.-F. Huang
Department of Information Management, Central Taiwan University of Science and Technology,
666, Buzi Rd., Beitun Dist, Taichung 40601, Taiwan
e-mail: suefn2001@yahoo.com.tw

C.-H. Cheng (✉)
Department of Information Management, National Yunlin University of Science and Technology,
123 University Road, Section 3, Douliou, Yunlin 64002, Taiwan
e-mail: chcheng@yuntech.edu.tw

1 Introduction

Attribute selection is an important and widely studied topic in many disciplines, including statistics, artificial intelligent, operation research, computer science, data mining and knowledge discovery (Peng et al. 2010). Datasets usually include many redundant data, these redundant data not only increase the dimension of the dataset but also affect the subsequent analysis performance. Due to attribute selection can improve the quality of the dataset in a pre-processing stage, many attribute selection methods have been proposed in the last two decades, such as disease diagnosis, text categorization, credit analysis, software risk management, financial crisis, and network intrusion detection, a variety of methods (Chatfield and Janek 1972; Tan et al. 2009; Cornelis et al. 2010; Chuang et al. 2011) and algorithms (Maldonado and Weber 2009; Menjoge and Welsch 2010; Lee and Leu 2011; Li et al. 2011) have been developed for attribute selection in recent years.

The basic problem of attribute selection is optimized the measure function, with a performance measure for each subset of attributes is to measure its ability of classifying the samples. The problem is to search through the space of attribute subsets to identify the optimal or near-optimal subset, with respect to the performance measure. In dataset analysis, it is usual to have redundant attributes, and the selected attributes will affect the model performance. Therefore, attribute selection become an interesting research topic, and applied to extract a set of input attribute in predicting the target attributes.

Jong and Young (2010) derived 24 attributes which are related to credit card transaction period, number of transactions, and transaction amount by using credit card sales information. Among them, for the attributes of the month suspended, average of sales, maximum and minimum amount of sales, variance of sales, average of transaction and variance of transaction during 3 and 6 months period, t test statistics for the 6 months period are higher than those of 3 months period. Hence the 13 attributes including the derived attributes of 6 months period are selected for the input attributes.

Cho et al. (2010) use t statistic to test 56 financial ratio attributes turned out to be significantly different between the bankrupt and healthy groups. Then, they applied four attribute selection methods. Method #1: decision trees with the chi-square algorithm select five attributes. Method #2: decision trees with the entropy reduction algorithm select four attributes. Method #3: stepwise logistic regression obtain 26 significant attributes, then select six attributes based on the most frequently appearing attributes in the credit field expert analysis. Method #4: after obtaining 26 significant attributes by a stepwise logistic regression, select seven attributes based on the most frequently appearing attributes in the credit field expert analysis.

Ravisankar et al. (2010) using training data is fed to MLFF (multilayer feed forward neural network), PNN (probabilistic neural network), GP (genetic programming), t statistic and f statistic separately for feature selection. Top 10 features are selected in each fold for a given technique and it was observed that different folds yielded different features as the top features. In order to arrive at a unified and optimal feature subset, computed the frequency of occurrence all the features in the top 10 slots across all folds. Then, all the features are sorted in the descending order of the frequency of occurrence. This helps them in selecting feature subset for that particular technique. They repeated the same method for every other technique. In case of rough sets, we had taken the top 10 features from the results obtained by Bose (2006). In each case, top 10 features that contribute to high accuracy are selected from the 24 features. The feature subset so formed is fed separately to MLFF/PNN/GP for classification purpose in the second phase. In case of rough sets based approach, however,

the top 10 features are taken from Bose (2006), who followed hold-out method of 80–20 ratio.

Melek et al. (2009) apply t tests to 20 financial ratios which do not show a significant difference are taken out of the pure data set and a new data set is created. The 20 attributes there is not a significant difference regarding 11 of them considering the banks which were transferred to the SDIF (savings deposit insurance fund) and which were non-failed. Therefore the rest 9 ratios are thought to be more useful in making a difference between the failed and non-failed banks. Li et al. (2009) each year's features employed were selected separately from 35 original financial ratios account items through the stepwise discriminant analysis, final 16 feature sets employed.

According to that mentioned above, a good attribute selection method can remove unnecessary attributes which may affect both on the rule of comprehension and prediction performance, and good attribute selection method can influence on the accuracy of prediction. Thus, find a suitable attribute selection method to evaluate prediction model becomes a very important issue apparently. However, since there are so many attribute selection methods, how to select an effective one for a given task becomes an important yet difficult issue. The evaluation of attribute selection methods normally involves more than one criterion, such as accuracy, type I error, and type II error (Tsai 2009). Therefore attribute selection can be modeled as group multiple criteria decision making (GMCDM) problems.

A variety of GMCDM methods have been developed over recent years and it is a challenging task to decide which GMCDM methods are suitable for a problem. Because different GMCDM methods rank alternative using different approaches and may yield different results when applied to the same problem, one feasible way is to apply combinations of GMCDM methods is more trustful than one generated by a single GMCDM method. While many empirical studies show that the rankings of alternatives provided GMCDM methods may be in conflict result, there are situations where different GMCDM methods generate different rankings and how to reconcile these differences has not been fully investigated (Peng et al. 2011).

The decision making is complicated problem, as attribute selection becomes more challenging today. There is a need for simplicity, systematic, and logical methods or mathematical tools as guide for decision makers in considering among different selection attributes and their interrelations. The objective for attribute selection procedure is to identify appropriate selection attributes and acquire the most appropriate combination of attributes in cope with the real need. Thus, efforts need to be extended to identify those attributes that influence attribute selection for a given prediction to exclude unsuitable alternatives, and subsequently for the selection of the most appropriate alternative adopting simple and logical methods.

The purpose of this paper is to apply GMADM method to rank attribute selection methods to improve a compromised solution in conflicting rankings generated by different attribute selection methods. The proposed method helps the decision maker to arrive at a decision based on either the objective importance of the attributes or his/her subjective preferences, that is considering both the objective weights and the subjective preferences. The proposed approach is examined using examples Tsai (2009) six attribute selection methods are included to illustrate the proposed method. The results show that the proposed approach is able to generate an optimal ranking of attribute selection methods in different domains.

The rest of this paper is organized as follows. In Sect. 2 the attribute selection methods and TOPSIS (Technique for order preference by similarity to ideal solution) are described. The experimental design of study is provided in Sect. 3. Section 4 presents and analyzes the experimental results. The Sect. 5 summarizes the findings and discusses future research directions.

2 Related works

This section briefly reviews the related literatures which includes attribute selection methods, multiple criteria decision making, and TOPSIS.

2.1 Attribute selection methods

This section focuses on the statistical methods applied on quantitative data, which are correlation matrix, factor analysis, t test, stepwise regression, and principal component analysis.

(A) Correlation matrix

Correlation matrix (Sadanori 1979) is to confer the correlation of two quantitative groups, as well as to analyze whether one group affects the other one. A correlation coefficient is the result of a mathematical comparison of how closely related two attributes are. The relationship between two attributes is said to be highly correlated if a movement in one attribute results or takes place at the same time as a similar movement in another attribute. To select appropriate attributes affecting much more parts of the result by this technique could obtain related advantages (Atiya 2001).

(B) t test

The t test method (David et al. 1993) is often used to assess whether the means of two groups are statistically different from each other by calculating a ratio between the difference of two groups means and the variability of the two groups. It helps to answer the underlying question: do the two groups come from the same population, and only appear differently because of chance errors, or is there some significant difference between these two groups.

Three basic factors help determine whether an apparent difference between two groups is a true difference or just an error due to chance (Pagano 2001):

1. The larger sample, the less likely that the difference is due to sampling errors or chance.
2. The larger the difference between the two means, the less likely the difference is due to sampling errors.
3. The smaller variance among the participants, the less likely that the difference was created by sampling errors.

(C) Principle component analysis

PCA (Principle component analysis) can be used for reducing complexity of input attributes when there are large volumes of information and it is intended to have a better interpretation of attributes (Wang and Paliwal 2003; Noori et al. 2010d). In this method, the information of input attributes will present with minimum losses in PCs (Helena et al. 2000).

PCA model identification ends with the assessment of the number of PCs (principle components) as the result of a trade-off between dimension reduction and the relative cumulative variance (RCV). This choice is often made in a subjective way, especially in explorative studies. One approach is based on the eigenvalue scree plot. An alternative approach is based on cross-validation. Both will be used to come to a proper choice, next to a visual inspection of the candidate PCs.

(D) Factor analysis

Factor analysis (Schneeweiss and Mathes 1995; Alexander 1994) is a statistical method that is based on the correlation analysis of multi-attributes. The purpose is to reduce multiple attributes to a lesser number of underlying factors that are measured by the attributes. In other words, it uses fewer dimensions to present original structures of data and keeps the most information. Factors are formed by grouping the attributes that have a correlation with each other. Factor analysis is effective when the sample size is more than 300. Factor analysis has mainly four steps, which is described in the following.

- (1) Initial solution: Attributes are selected and an inter correlation matrix is generated for including all of the attributes. An inter-correlation matrix is a $k \times k$ (where k equals the number of attributes) array of the correlation coefficients of the attributes with each other. The correlation coefficients value should be significance value greater than 0.3. Kaiser–Meyer–Olkin (KMO) and Bartlett’s tests of sphericity (BTS) are then applied to the studied attributes in order to validate if the remaining attributes are factorable. The KMO value should be greater than 0.6 for a satisfactory factor analysis. The BTS should show that the correlation matrix is not an identity matrix by giving a significance value smaller than 0.05.
- (2) Extracting the factors: The PCA method is the most common form of factor analysis. The loadings values of PCA are easy to compute and explain the factors. An appropriate number of components (factors) are extracted from the correlation matrix based on the initial solution. In the initial solution, each attribute is standardized to have a mean of 0 and a standard deviation of ± 1 . Thus, the eigenvalue of the factor should be larger than one, if it is to be extracted.
- (3) Rotating the factors: In general, factors are rotated in order to clarify the relationship between the attributes and the factors. Factor rotation methods including orthogonal rotation (the varimax, quartimax, equimax method is the most commonly used) and oblique rotation (the oblimin, oblimax, quartimin method is the most commonly used).
- (4) Naming the factors: Results are then derived by analyzing the absolute factor loadings more than 0.5 or communality more than 0.5 of each attribute. Appropriate names are given to each factor by considering the factor loadings.

(E) Stepwise regression

Using regression to build models, one common technique to find the best combination of predictor attributes is stepwise regression. Although there are many variations, the most basic procedure is to find the single best predictor attribute and add attributes that meet some specified criterion. The result is a combination of predictor attributes, all of which have significant coefficients. A list of several potential explanatory attributes is available and this list is repeatedly searched for attributes which should be included in the model. The best explanatory attribute is used first, then the second best, and so on. This procedure is known as stepwise regression (Tsai 2009).

2.2 Group multiple criteria decision making

A decision making problem usually involves more than one criterion (factor), and criteria often conflict with each other. In GMCDM, it is usually assumed that the criteria are independent. The decision-making process involves identifying problems, constructing preferences,

evaluating alternatives, and determining the best alternative (Hwang and Yoon 1981). However, when decision-makers evaluate the alternatives with multiple criteria, many problems, such as the weights of the criteria, preference dependence, and conflicts among criteria, seem to complicate the decision-making process and should be resolved by more sophisticated methods.

Last decade, some of the popular solved GMCDM methods include: (1) Analytic hierarchy process (AHP) (Saaty 1980) can be enhanced with incremental analysis by a benefit–cost ratio. (2) Data envelopment analysis (DEA) (Charnes et al. 1978) evaluates the efficiency of decision making units through identifying the efficiency frontier and comparing each DMU with the frontier. (3) TOPSIS (Hwang and Yoon 1981) is chosen as the target for the analysis because of its stability and ease of use with cardinal information. (4) Decision-making trial and evaluation laboratory (DEMATEL) (Fontela and Gabus 1976) technique is used to detect complex relationships and build a network relation map among criteria for environment watershed measurement and evaluation.

Each GMCDM method reflects a different approach to solve GMCDM problems. How to obtain a final decision by integrating the different opinions from different experts is an important issue. In general, a GMCDM problem can be concisely expressed in matrix format.

$$\begin{matrix}
 & C_1 & \cdots & C_j & \cdots & C_n \\
 \\
 A_1 & \left[\begin{array}{cccc}
 x_{11} & \cdots & x_{1j} & \cdots & x_{1n} \\
 \vdots & & \vdots & & \vdots \\
 A_i & \left[\begin{array}{cccc}
 x_{i1} & \cdots & x_{ij} & \cdots & x_{in} \\
 \vdots & & \vdots & & \vdots \\
 A_m & \left[\begin{array}{cccc}
 x_{m1} & \cdots & x_{mj} & \cdots & x_{mn}
 \end{array} \right]
 \end{matrix} \right.
 \end{matrix} \tag{1}$$

where A_1, A_2, \dots, A_m are possible alternative, C_1, C_2, \dots, C_n are criteria with which performance of alternatives are measured, D is the rating of alternative A_i with respect to criteria C_j . The decision maker prefers to give his opinions. Hence, the rating x_{ij} of alternative A_i and the weights of the criteria are assessed in terms.

2.3 Technique for order preference by similarity to ideal solution

TOPSIS is one of the major techniques in dealing with GMADM problems. The TOPSIS was first developed by Hwang and Yoon (1981), and it concept of TOPSIS is that the most preferred alternative should not only have the shortest distance from the positive ideal solution, but also have the longest distance from the negative ideal solution. In short, the positive-ideal solution is composed of all the best values attainable of criteria, whereas the negative ideal solution is made up of all the worst values attainable of criteria. Despite many GMCDM techniques, the proposed use TOPSIS method is chosen as the target for the analysis because of its stability and ease of use with cardinal information.

3 The proposed method

Attribute selection is an important pre-process, the interesting patterns and useful relationships are extracted from the analysis of the input data. From Tsai’s paper (Tsai 2009), the

evaluated index (Accuracy, Type I Error Rate and Type II Error Rate) is utilized as the performance of attribute selection method under each dataset. For example, German credit dataset use "Accuracy" to rank the ordering, the result is {FA > t test > Stepwise regression > Correlation matrix > Baseline model > PCA}. And used the top two positions which have the highest scores are seen as the best attribute selection methods. i.e., FA and t test are selected as the best attribute selection methods. However, the ranking may be biased due to relative values or the ratio scale used in the evaluation. And most decision making problems have interactions in the real-world, these decision making problem have multiple attributes. Therefore, evaluation these decision making problems necessary to conduct an overall consideration of the relevant attributes and comprehensive assessment. For this reason, the past approaches are not good for the purpose of ranking and selection and cannot reflect the true dominance of attribute selection method.

For overcoming mentioned problem above, this paper proposes a GMADM-based attribute selection methods to deal with the disagreement of ranking and help decision makers pick the most suitable attribute selection method. These alternatives and criteria (factors) are formed GMADM framework. Then the decision matrix is calculated by TOPSIS method (Hwang and Yoon 1981) to find the optimized attribute selection method. Due to the result of TOPSIS is more robust fashion, the TOPSIS analysis is essential for solving GMADM ranking and selected attribute selection method. i.e., this paper utilizes six attribute selection methods: t test (A_1), Stepwise (A_2), correlation matrix (A_3), FA (A_4), PCA (A_5) and baseline model (A_6 , not reduce attribute) as GMADM alternatives, and the evaluated criteria (factors) are accuracy, Type I error and Type II error. Additionally, the ratings of GMADM are from the calculated five datasets as expert's ratings. The five datasets include Japanese credit, Australian credit, Bankruptcy dataset, German credit and UC Competition from Tsai (2009). Therefore, the process of building a GMADM model consists of alternatives and criteria, which forms the decision matrix and maximizes the decision problem by TOPSIS.

For easy understanding, the proposed method is shown in Fig. 1. the firstly five attribute selection methods belong to inferential statistics. With inferential statistics, we try to infer from the sample data about what the population might be like, or make judgments of the probability that an observed difference between groups is a dependable one. Applying those six methods to acquire the most representative attributes can affect the accuracy of prediction. For evaluated criteria, Type I error rate means that the error rate for the risk cannot categorize the normal company as a normal company, Type II error rate means that the error rate for the risk cannot categorize the failure company, and Accuracy means the obvious way to assess the quality of the learned model and to see on how long term the predictions given by the model are accurate.

Next, the proposed algorithm will be introduced in detailed, which includes 8 steps in the following:

Step 1: Construct the group decision matrix

Suppose that there are m attribute selection methods (A_1, A_2, \dots, A_m) will be evaluated in accordance with n decision criteria (C_1, C_2, \dots, C_n). Let $w = (w_1, w_2, \dots, w_n)$ be the weight vector of criteria. The group decision matrix D is constructed by aggregating the ratings of K decision maker as Eq. (2):

$$D = [x_{ij}]_{m \times n} \quad \text{where } x_{ij} = \frac{1}{K} \sum_{k=1}^K x_{ij}^k \quad k = 1, 2, \dots, K \quad (2)$$

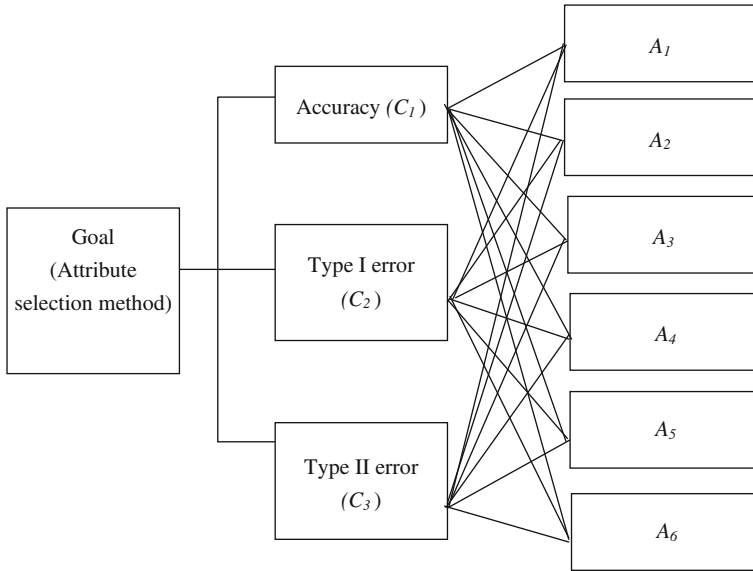


Fig. 1 The GMADM-based selection attribute method

Step 2: Normalize the group decision matrix $D = [x_{ij}]_{m \times n}$ using the Eq. (3):

$$r_{ij} = \frac{x_{ij}}{\sqrt{\sum_{i=1}^m x_{ij}^2}}, \quad i = 1, 2, \dots, m, \quad j = 1, 2, \dots, n \tag{3}$$

where r_{ij} is the normalized rating of A_i with respect to C_j . m and n denote the number of alternatives and the number of criteria, respectively. For alternative A_i , the performance measure of the i th criterion C_j is represented by x_{ij} .

Step 3: Calculate the weighted normalized group decision matrix

$$V = [v_{ij}]_{m \times n}, \quad v_{ij} = w_j^* r_{ij}, \quad i = 1, 2, \dots, j = 1, 2, \dots, n \tag{4}$$

where w_j is the relative weight of C_j , and $\sum_{j=1}^n w_j = 1$.

Step 4: Determine the positive ideal (A^*) and negative ideal (A^-) solutions as follows:

$$A^* = \{v_1^*, v_2^*, \dots, v_n^*\} \tag{5}$$

$$A^- = \{v_1^-, v_2^-, \dots, v_n^-\} \tag{6}$$

where $V_j^* = \max_i \{V_{ij}\}$ is associated with benefit criteria and $V_j^- = \min_i \{V_{ij}\}$ is associated with cost criteria.

Step 5: Calculate the Euclidean distances of each attribute selection method from the positive ideal solution (PIS) and the negative-ideal solution (NIS), respectively. It can be represented as

$$D_i^* = \sqrt{\sum_{j=1}^n (v_{ij} - v_j^*)^2}, \quad i = 1, 2, \dots, m \quad (7)$$

$$D_i^- = \sqrt{\sum_{j=1}^n (v_{ij} - v_j^-)^2}, \quad i = 1, 2, \dots, m \quad (8)$$

Step 6: Calculate the relative closeness of each attribute selection method. The relative closeness of the attribute selection method A_i is calculated as

$$CC_i = \frac{D_i^-}{D_i^* + D_i^-}, \quad i = 1, 2, \dots, n \quad (9)$$

Step 7: Rank the attribute selection method according to the relative closeness values. The higher the CC_i , the better the attribute selection method A_i . Thus, if $CC_i > CC_j$ then the method A_i is better A_j .

4 Numerical verification

In order to verify the proposed method, five datasets are calculated as expert's ratings to verify the proposed method. From Sect. 3, attribute selection has become a GMCDM problem in proposed method, six attribute selection methods is A_1 (t test), A_2 (Stepwise regression),

Table 1 Performance of attribute selection (unit %)

	t test	Stepwise	Correlation matrix	FA	PCA	Baseline models
Japanese credit						
Accuracy	63.53	82.64	60.16	74.22	74.00	85.88
Type I error	55.33	32.27	74.55	29.17	47.46	90.05
Type II error	17.29	6.77	3.49	23.75	10.37	22.40
Australian credit						
Accuracy	89.27	84.74	89.31	86.08	89.93	81.93
Type I error	9.38	12.80	13.33	14.58	7.93	21.89
Type II error	11.72	16.71	8.33	13.60	11.53	13.89
Bankruptcy dataset						
Accuracy	82.98	77	76.08	72.91	79.59	71.03
Type I error	7.69	37.27	22.76	22.50	16.55	12.85
Type II error	28.57	5.56	25.45	32.73	26	30.42
German credit						
Accuracy	75.87	75.51	74.84	78.76	67.03	74.28
Type I error	61.28	51.34	54.36	48.69	84.92	55.39
Type II error	8.62	12.25	12.04	10.66	6.27	9.63
UC competition						
Accuracy	97.25	96.33	96.70	97.30	96.47	96.92
Type I error	74.82	79.25	96.47	94.00	90.00	81.68
Type II error	0.16	0.35	0.04	0.08	0.13	4.05

*Data from Tsai (2009)

Table 2 Decision matrix

Five dataset	<i>t</i> test (A_1)	Stepwise regression (A_2)	Correlation matrix (A_3)	FA (A_4)	PCA (A_5)	Baseline model (A_6)
Accuracy (C_1)	81.7800	<i>83.2440</i>	79.4180	81.8540	81.4040	82.0080
Type I error (C_2)	<i>41.7000</i>	42.5860	52.2940	41.7880	49.3720	52.3720
Type II error (C_3)	13.2720	<i>8.3280</i>	9.8700	16.1640	10.8600	16.0780

Italic value denotes superior in same criteria

Table 3 Normalized decision matrix

Five dataset	<i>t</i> test (A_1)	Stepwise regression (A_2)	Correlation matrix (A_3)	FA (A_4)	PCA (A_5)	Baseline model (A_6)
Accuracy (C_1)	0.9824	<i>1.0000</i>	0.9540	0.9833	0.9779	0.9852
Type I error (C_2)	<i>0.5009</i>	0.5116	0.6282	0.5020	0.5931	0.6291
Type II error (C_3)	0.1594	<i>0.1000</i>	0.1186	0.1942	0.1305	0.1931

Italic value denotes superior in same criteria

A_3 (Correlation matrix), A_4 (FA), A_5 (PCA) and A_6 (baseline model), three criteria include: C_1 (Accuracy), C_2 (Type I error) and C_3 (Type II error), and five decision makers (datasets) are D_1 (Japanese credit), D_2 (Australian credit), D_3 (Bankruptcy dataset), D_4 (German credit) and D_5 (UC Competition). The performance of each dataset (decision maker) be calculated which is expressed as decision goal for each attribute selection method on each criterion. The results of *t* test, Stepwise, correlation matrix, FA, PCA and baseline model for the five datasets are listed in Table 1. The proposed method is currently applied to solve this decision making problem and the computational procedure is summarized in the following:

- (1) Construct the group decision matrix and normalize decision matrix.

From Step 1–2 of proposed algorithm, use Eq. (2) to construct the group decision matrix as Table 2. And normalize group decision matrix by Eq. (3), the normalized group decision matrix is shown in Table 3.

- (2) Calculate the weighted normalized group decision matrix.

In the process of decision making for focus different prediction result, it is evident that for each criterion as Accuracy, Type I error and Type II error, the perspective of the decision makers is not given the same importance. Therefore, a weighted vector W is introduced to denote the importance of weight for that criterion regarding the opinion of the decision maker. According prediction result focus of importance is assigned for each criterion. Assume that they express their preferences as Table 4. The weighted normalized decision matrix can be obtained from multiplying each assigned criterion weighting by normalized decision matrix, Table 5 shown the weighted normalized decision matrix under weighted vector $W_4 = [0.7 \ 0.1 \ 0.2]$.

- (3) Determine the positive ideal solution (A^*) and negative ideal solution (A^-)

The positive ideal solution and the negative ideal solution are then determined by Eqs. (5) and (6) shown in Table 6.

- (4) Calculate the separation measures using the Euclidean distance.

The separation of each alternative from PIS and NIS are then determined by Eqs. (7) and (8) shown in Table 7.

Table 4 The importance weight of the criteria

Weight	W_1	W_2	W_3	W_4	W_5	W_6	W_7	W_8	W_9
C_1	0.1	0.3	0	0.7	0.5	1	0.1	0.1	0
C_2	0.6	0.5	1	0.1	0.2	0	0.1	0.2	0
C_3	0.3	0.2	0	0.2	0.3	0	0.8	0.7	1

Table 5 The weighted normalized decision matrix

	A_1	A_2	A_3	A_4	A_5	A_6
$C_1(0.7)$	0.6877	0.7000	0.6678	0.6883	0.6845	0.6896
$C_2(0.1)$	0.0501	0.0512	0.0628	0.0502	0.0593	0.0629
$C_3(0.2)$	0.0319	0.0200	0.0237	0.0388	0.0261	0.0386

Table 6 The distance measurement

	A^*	A^-
C_1	0.7000	0.6678
C_2	0.0501	0.0629
C_3	0.0200	0.0388

Table 7 Each alternative from PIS and NIS

	A_1	A_2	A_3	A_4	A_5	A_6
D^*	0.3049	0.3135	0.2949	0.3035	0.3029	0.3014
D^-	0.2301	0.2338	0.2291	0.2279	0.2293	0.2242

(5) Calculate the relative closeness to the ideal solution and rank the preference according to CC_i .

The relative closeness to the ideal solution is determined by Eq. (9) shown in Table 8.

Rank the preference ordering in terms of the values of CC_i . In the numerical example the ranking of six attribute selection method is $A_3 > A_5 > A_1 > A_4 > A_2 > A_6$. Obviously, the best selection is candidate A_3 .

For sensitivity analysis, this paper repeatedly changes the criteria weight to rank attribute selection method by proposed algorithm under each different weight. The results of nine different weights are listed in Table 4, automatically one alters the preferences and the calculation of the impact of the attribute selection method, which, in turn, may modify the ranking of the attribute selection method. From Table 9, the correlation matrix method is relative important in overall consideration of the relevant attributes and comprehensive assessment (A_3 is first ordering with 4 times).

Table 8 The relative closeness to the ideal solution

	A_1	A_2	A_3	A_4	A_5	A_6
CC_i	0.4301	0.4272	0.4373	0.4288	0.4309	0.4265

Table 9 Ranking of the attribute selection method

Weight	Methods ranking
W_1	$A_4 > A_1 > A_2 > A_5 > A_6 > A_3$
W_2	$A_4 > A_1 > A_2 > A_5 > A_3 > A_6$
W_3	$A_1 > A_4 > A_2 > A_5 > A_3 > A_6$
W_4	$A_3 > A_5 > A_1 > A_4 > A_2 > A_6$
W_5	$A_3 > A_5 > A_1 > A_4 > A_2 > A_6$
W_6	$A_3 > A_5 > A_1 > A_4 > A_6 > A_2$
W_7	$A_3 > A_2 > A_5 > A_1 > A_6 > A_4$
W_8	$A_2 > A_1 > A_5 > A_3 > A_4 > A_6$
W_9	$A_2 > A_3 > A_5 > A_1 > A_6 > A_4$

Some findings can be figured out from Table 4 corresponding to Table 9:

- (1) The Accuracy (C_1) is higher weight (W_4, W_5, W_6), the correlation matrix method will be the best selection.
- (2) The Type I error (C_2) is given higher weight (W_1, W_2), the FA method (A_4) will be the best selection.
- (3) The Type II error (C_3) sets higher weight (W_8, W_9), the stepwise regression method (A_2) will be the best selection.

In simulation experiment, this paper used six attribute selection methods, three decision criteria and five decision makers only to verify the proposed method. We can see that the proposed method has found the best attribute selection method.

5 Conclusion

A variety of algorithms have been developed for attribute selection problems. How to select an effective and appropriate algorithm for a given task is an important yet difficult task. Since algorithms evaluation normally involves more than one criterion, it can be modeled as GMCMD problems. The proposed method has verified that attribute selection methods (t test, Stepwise, correlation matrix, FA, PCA and baseline model) are considered as GMADM to evaluate the optimal selection problem, the proposed method could resolve the problem of ranking conflicts.

From experimental results, we can obtain that applied attribute selection is better results than using all available attributes. And the proposed method can provide a compatible ranking when TOPSIS techniques yield conflicting results. It's a challenging task to find a general approach which can correctly rank attribute selection methods in any circumstances. This study has proposed GMADM-based method to rank attribute selection methods and has verified its feasible solution. As a future research direction, more datasets from various domains could be tested, and the other GMCMD methods also could be applied, such as VIKOR, DEA, and AHP.

References

- Alexander, T.B.: Statistical Factor Analysis and Related Methods: Theory and Applications. Wiley-Blackwell, (1994)

- Atiya, A.F.: Bankruptcy prediction for credit risk using neural networks: a survey and new results. *IEEE T. Neural Netw.* **12**, 929–935 (2001)
- Bose, I.: Deciding the financial health of dot-coms using rough sets. *Inf. Manage.* **43**, 835–846 (2006)
- Charnes, A., Cooper, W.W., Rhodes, E.: Measuring the efficiency of decision making units. *Eur. J. Oper. Res.* **2**(6), 429–444 (1978)
- Chatfield, D.C., Janek, E.J.: Attribute selection in concept identification. *J. Exp. Psychol.* **95**, 97–101 (1972)
- Cho, S., Hyojung, H., Byoung, C.H.: A hybrid approach based on the combination of variable selection using decision trees and case-based reasoning using the Mahalanobis distance: For bankruptcy prediction. *Expert Syst. Appl.* **37**, 3482–3488 (2010)
- Chuang, L.Y., Yang, C.H., Li, J.C.: Chaotic maps based on binary particle swarm optimization for feature selection. *Appl. Soft Comput.* **11**, 239–248 (2011)
- Cornelis, C., Jensen, R., Hurtado, G., Slezak, D.: Attribute selection with fuzzy decision reducts. *Inf. Sci.* **180**, 209–224 (2010)
- David, E.A., Giles, V., Srivastava, K.: The exact distribution of a least squares regression coefficient estimator after a preliminary *t*-test. *Stat. Probabil. Lett.* **16**, 59–64 (1993)
- Fontela, E., Gabus, A.: The DEMATEL Observer, DEMATEL Report, Battelle Geneva Research Center, Geneva (1976)
- Helena, B., Pardo, R., Vega, M., Barrado, E., Fernandez, J.M., Fernandez, L.: Temporal evolution of groundwater composition in an alluvial aquifer (Pisuerga river, Spain) by principal component analysis. *Water Res.* **34**, 807–816 (2000)
- Hwang, C.H., Yoon, K.: Multiple attribute decision making methods and applications. pp. 1–7. Springer, Berlin (1981)
- Jong, S.Y., Young, S.K.: A practical approach to bankruptcy prediction for small businesses: Substituting the unavailable financial data for credit card sales information. *Expert Syst. Appl.* **37**, 3624–3629 (2010)
- Lee, C.P., Leu, Y.: A novel hybrid feature selection method for microarray data analysis. *Appl. Soft Comput.* **11**, 208–213 (2011)
- Li, H., Sun, J.: Forecasting business failure in china using case-based reasoning with hybrid case representation. *J. Forecasting* **29**, 486–501 (2009)
- Li, S., Wu, H., Wan, D., Zhu, J.: An effective feature selection method for hyperspectral image classification based on genetic algorithm and support vector machine. *Knowl.-Based Syst.* **24**, 40–48 (2011)
- Maldonado, S., Weber, R.: A wrapper method for feature selection using support vector machines. *Inf. Sci.* **179**, 2208–2217 (2009)
- Melek, A.B., Yakup, K., Baykan, O.K.: Predicting bank financial failures using neural networks, support vector machines and multivariate statistical methods: A comparative analysis in the sample of savings deposit insurance fund (SDIF) transferred banks in Turkey. *Expert Syst. Appl.* **36**, 3355–3366 (2009)
- Menjoge, R.S., Welsch, R.E.: A diagnostic method for simultaneous feature selection and outlier identification in linear regression. *Comput. Stat. Data Anal.* **54**, 3181–3193 (2010)
- Noori, R., Sabahi, M.S., Karbassi, A.R., Baghvand, A., aati Zadeh, H.: T multivariate statistical analysis of surface water quality based on correlations and variations in the data set. *Desalination* **260**, 129–136 (2010)
- Pagano, R.R.: *Understanding Statistics in the Behavioral Sciences*, Sixth ed. Wadsworth/Thomson Learning, California (2001)
- Peng, Y., Kou, G., Wang, G., Shi, Y.: FAMCDM: a fusion approach of MCDM methods to rank multiclass classification algorithms. *Omega* **39**, 677–689 (2011)
- Peng, Y., Zhiqing, W., Jianmin, J.: A novel feature selection approach for biomedical data classification. *J. Biomed. Inform.* **43**, 15–23 (2010)
- Ravisankar, P., Ravi, V., Bose, I.: Failure prediction of dotcom companies using neural network–genetic programming hybrids. *Inf. Sci.* **180**, 1257–1267 (2010)
- Saaty, T.L.: *The Analytic Hierarchical Process*. McGraw-Hill, New York (1980)
- Sadanori, K.: Asymptotic expansions for the distributions of functions of a correlation matrix. *J. Multivar. Anal.* **9**, 259–266 (1979)
- Schneeweiss, H., Mathes, H.: Factor analysis and principal components. *J. Multivar. Anal.* **55**, 105–124 (1995)
- Shin, K., Lee, S.Y.J.: A genetic algorithm application in bankruptcy prediction modeling. *Expert Syst. Appl.* **23**, 321–328 (2002)
- Tan, K.C., Teoh, E.J., Yu, Q., Goh, K.C.: A hybrid evolutionary algorithm for attribute selection in data mining. *Expert Syst. Appl.* **36**, 8616–8630 (2009)
- Tsai, C.F.: Feature selection in bankruptcy prediction. *Knowl.-Based Syst.* **22**, 120–127 (2009)
- Wang, X., Paliwal, K.K.: Feature extraction and dimensionality reduction algorithms and their applications in vowel recognition. *J. Pattern Recogn. Soc.* **36**, 2429–2439 (2003)