# An Alternative Model to Fisher and Linear Programming Approaches in Two-Group Classification Problem: Minimizing Deviations from the Group Median

## (Review)

Hasan BAL, H.Hasan ÖRKCÜ[*], Salih ÇELEBİOĞLU

*Gazi University, Arts and Science Faculty, Department of Statistics Teknikokullar, 06500, Ankara, TURKEY,*

**ABSTRACT**

In this study, new classification models were developed which can be used in the solution to the problems of Discriminant Analysis having two groups. For the solution of these type of problems, Lam, Choo and Moy (1996) proposed a model regarding the minimization of deviations from the group means. The model examined by these authors loses its efficiency in respect of the hit ratio as the distributions of populations of samples considered go away from the normal distribution. For the samples drawn from non normal or skewed distributions, the median is a much more suitable descriptive statistic than the mean. The aim of the study is to consider the models of two-group classification problems by minimizing the deviations from the group medians. When these proposed approaches are applied to the data of real life or of simulation drawn from different distributions, it is observed that the attained performance of classification is better than both some important classification approaches in the literature and especially the classification performance minimizing the deviations from group means proposed by Lam, Choo and Moy.

*Key Words*: Statistical Discriminant Analysis, Goal Programming, median scores

## 1. INTRODUCTION

The discriminant analysis is a technique interested in determining the groups of objects based on their observed scores. Fisher's linear discriminant function is especially the most popular technique which is frequently used for the discriminant problem. As an alternative for the examination of classification problems using the statistical methods, recently a number of new efficient mathematical programming approaches have been developed, Fred and Glover (1), Glover (2), Lam and Moy (3), Lam and Moy (4), among others. In two-group and multigroup classification problems, Lam, Choo and Moy (5), and Lam and Moy (6) developed a satisfactory model of classification based on cluster analysis. In two-group classification problems, their approach minimizes the sum

of deviations of all object's classification scores from the mean group classification scores.

The approach presented here examines the two-group classification problem by minimizing the sum of deviations between the classification scores of all objects and the group median scores of objects. The reason that the median is used in this situation is that it represents results more accurately than the arithmetic mean in cases of skewed distributions like in Exponential, Gamma, Chi Square and F distributions and that most of the real life data sets are of that nature. The layout of the article is as follows; in the second part, linear programming approaches proposed previously in literature are explained, in the third part, the linear models based on median scores and in the fourth part goal programming approaches based on mean and median scores are proposed. In the fifth part, a real data set of economical variables from 91 countries is

[*] Corresponding Author, e-mail: hhorkcu@gazi.edu.tr

illustrated and in the sixth part, a simulation study is performed to compare samples from different distributions based on their hit-ratios. In the last part, we discuss the results and findings of the article and further research about this topic.

## 2. LINEAR PROGRAMMING APPROACHES

Consider the two-group classification problem with $k$ attributes. Let $x$ be the $k \times n$ matrix of attribute scores of the sample with size $n$ drawn from the groups $G_1$ and $G_2$. If $w_1, w_2, \ldots, w_k$ are the attribute weights, then the classification score is defined as $S_i = \sum_{j=1}^{k} x_{ij} w_j$, ($i = 1, 2, \ldots, n$). The assignment of an object into groups depends on the value of its classification score. The simple MSD (Minimum Sum of Deviations) classification score can be formulated as follows:

MSD

$$\min \sum_{i=1}^{n} d_i \qquad\qquad [1]$$

$$\sum_{j=1}^{k} w_j x_{ij} - d_i \leq c \qquad i \in G_1$$

$$\sum_{j=1}^{k} w_j x_{ij} + d_i \geq c \qquad i \in G_2$$

where $d_i \geq 0, \ (i = 1, 2, \ldots, n)$, $w_j \ (j = 1, 2, \ldots, k)$ and $c$ are unrestricted (positive or negative) variables (2). Solving this model gives us the $w_j$ and $c$ values, with which we can obtain the classification score of any object. An object will be classified into $G_1$ if its classification score is greater than or equal to c, otherwise into $G_2$.

Like this model, many of the other existing linear programming models determine the attribute weights and cut-off values taking place here at the same time. Lam, Choo and Moy (6) divide the process made by their model into two steps: the first constitutes the determination of attribute weights, and the second determines the cut-off values for the classification. Their model makes use of an objective function minimizing the sum of deviations from the group mean classification scores. The Modified model of Lam, Choo and Moy (6) can be formulated as follows (MLCM):

MLCM 1

$$\min \sum_{i=1}^{n} d_i \qquad\qquad [2]$$

$$\sum_{j=1}^{k} w_j \left( x_{ij} - \mu_{1j} \right) + d_i \geq 0 \qquad i \in G_1$$

$$\sum_{j=1}^{k} w_j \left( x_{ij} - \mu_{2j} \right) - d_i \leq 0 \qquad i \in G_2$$

$$\sum_{j=1}^{k} w_j \left( \mu_{1j} - \mu_{2j} \right) \geq 1$$

where $d_i \geq 0, \ (i = 1, 2, \ldots, n)$ and $w_j \ (j = 1, 2, \ldots, k)$ is the unrestricted variable and $\mu_{1j}$ is the median of the $j^{\text{th}}$ variable in group $G_1$ and $\mu_{2j}$ is the median of $j^{\text{th}}$ variable in group $G_2$.

In this model, to the weights are reached by making object scores close to the mean score of the group in which they take place. Then the object scores are used in the following model [3] and the classification is made:

MLCM 2

$$\min \sum_{i=1}^{n} h_i \qquad\qquad [3]$$

$$S_i + h_i \geq c \qquad i \in G_1$$

$$S_i - h_i \leq c \qquad i \in G_2$$

where $h_i \geq 0, \ (i = 1, 2, \ldots, n)$, and $c$ is an unrestricted variable. As seen classification is made in two steps. These models are easily solved by any simplex algorithm.

## 3. NEW LINEAR PROGRAMMING FORMULATIONS BASED ON MEDIAN MINIMIZATON

As mentioned, the above modified model of Lam, Choo and Moy (6) minimizes the sum of deviations of the classification scores for each unit. The model we present minimizes the sum of deviations of the classification scores from the median classification scores for the units in the two-group classification problem. Similar to the model of Lam, Choo and Moy (6) our model can be formulated as follows:

LMED 1

$$\min \sum_{i=1}^{n} d_i \qquad\qquad [4]$$

$$\sum_{j=1}^{k} w_j \left( x_{ij} - med_{1j} \right) + d_i \geq 0 \qquad i \in G_1$$

$$\sum_{j=1}^{k} w_j \left( x_{ij} - med_{2j} \right) - d_i \leq 0 \qquad i \in G_2$$

$$\sum_{j=1}^{k} w_j \left( med_{1j} - med_{2j} \right) \geq 1$$

where $d_i \geq 0, \ (i = 1, 2, \ldots, n)$ and $w_j \ (j = 1, 2, \ldots, k)$ is the unrestricted variable and $med_{1j}$ is the median of the $j^{\text{th}}$ variable in group $G_1$ and $med_{2j}$ is the median of $j^{\text{th}}$ variable in group $G_2$. In this model, in the first step the weights $w_j$ are found after the solution to LMED1. Here the weights are found by making the object scores close to the median score of the group in which they take place. Using these weights the classification scores ($S_i$) for each object are evaluated and then the designation of objects to groups are made by the LMED2 model.

LMED 2

$$\min \sum_{i=1}^{n} h_i \qquad [5]$$

$$S_i + h_i \geq c \qquad i \in G_1$$

$$S_i - h_i \leq c \qquad i \in G_2$$

where $h_i \geq 0$, $(i = 1, 2, \ldots, n)$, and $c$ is an unrestricted variable. These models are also solved by simplex algorithm.

## 4. GOAL PROGRAMMING APPROACHES

### 4.1. Goal programming approach based on the minimization of deviations from the mean

Since the goal programming is an extension of linear programming, we can consider the models MLCM1 and MLMC2 together as the priority linear programming models. In this model the first priority is to minimize the sum of deviations of objects' classification scores from the group mean scores. The second priority is to minimize the sum of deviations between classification scores and cut-off values. In this case the goal programming approach based on the minimization of deviations from the mean (GPM) model is given as follows:

GPM

$$Min\ a = \left\{ \sum_{i=1}^{n} d_i , \sum_{i=1}^{n} h_i \right\}$$

$$\sum_{j=1}^{k} w_j \left( x_{ij} - \mu_{1j} \right) + d_i \geq 0 \qquad i \in G_1 \qquad [6]$$

$$\sum_{j=1}^{k} w_j \left( x_{ij} - \mu_{2j} \right) - d_i \leq 0 \qquad i \in G_2$$

$$\sum_{j=1}^{k} w_j x_{ij} + h_i \geq c \qquad i \in G_1$$

$$\sum_{j=1}^{k} w_j x_{ij} - h_i \leq c \qquad i \in G_2$$

$$\sum_{j=1}^{k} w_j \left( \mu_{1j} - \mu_{2j} \right) \geq 1$$

where $d_i, h_i \geq 0$, $(i = 1, 2, \ldots, n)$, $w_j$ $(j = 1, 2, \ldots, k)$ and $c$ are unrestricted variables. An object will be classified into $G_1$ if its classification score is greater than or equal to c, otherwise into $G_2$. This model is solved by lexicographic goal programming algorithm.

### 4.2. Goal programming approach based on the minimization of deviations from the median

Similar to the foregoing model, we can solve the problems LMED1 and LMED2 by Pre-emptive goal programming model. Here the first priority is to minimize the sum of deviations of objects' classification scores from the group median scores. The second priority is to

minimize the sum of deviations between classification scores and cut-off values. The related GPM model is given as follows:

GMED

$$Min\ a = \left\{ \sum_{i=1}^{n} d_i , \sum_{i=1}^{n} h_i \right\}$$

$$\sum_{j=1}^{k} w_j \left( x_{ij} - med_{1j} \right) + d_i \geq 0 \qquad i \in G_1 \qquad [7]$$

$$\sum_{j=1}^{k} w_j \left( x_{ij} - med_{2j} \right) - d_i \leq 0 \qquad i \in G_2$$

$$\sum_{j=1}^{k} w_j x_{ij} + h_i \geq c \qquad i \in G_1$$

$$\sum_{j=1}^{k} w_j x_{ij} - h_i \leq c \qquad i \in G_2$$

$$\sum_{j=1}^{k} w_j \left( med_{1j} - med_{2j} \right) \geq 1$$

where $d_i, h_i \geq 0$, $(i = 1, 2, \ldots, n)$, $w_j$ $(j = 1, 2, \ldots, k)$ and $c$ are unrestricted variables. This model is also solved by lexicographic (pre-emptive) goal programming algorithm.

The purpose of this paper is to bring up the classification performances of the new LMED, GPM and GMED approaches we proposed above by using both the data of real life or of simulation.

## 5. AN APPLICATION

In order to evaluate the performances of the LMED and GMED models, we considered the real data set of 91 countries with six variables (7). The data of the year 2002 are provided from T.R. Undersecreteriat of The Prime Ministry for Foreign Trade. The data used are as follows:

$x_1$ : the portion of country in World export

$x_2$ : the portion of country in World import

$x_3$ : national income per capita ( in thousand $ )

$x_4$ : export / gross national product and national income per capita ( billion $ )

$x_5$ : import / gross national product and national income per capita ( billion $ )

$x_6$ : the annual inflation rate

All of the 91 countries are divided into two groups, those of export portion 1 and over to group 1 ($G_1$; 23 countries), and those of export portion less than 1 to group 2 ($G_2$; 68 countries).

The attribute weights for six methods are given in Table 1.

**Table 1:** Attribute weights of six approaches

| Method | Sample[a] | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | $x_6$ |
|--------|-----------|-------|-------|-------|-------|-------|-------|
| **FLDF** | W | 2.583 | 0.824 | 0.0032 | 0.112 | 0.006 | -0.024 |
|          | D | 2.231 | 0.926 | 0.0005 | 0.019 | 0.007 | -0.001 |
| **MSD**  | W | 0.625 | 0.215 | -0.378 | -0.012 | 0.175 | -0.245 |
|          | D | 0.752 | 0.117 | -0.214 | -0.003 | 0.254 | -0.179 |
| **MLCM** | W | 0.542 | 0.752 | 0.036 | 0.176 | 0.009 | -0.063 |
|          | D | 0.687 | 0.621 | 0.078 | 0.162 | 0.002 | -0.027 |
| **LMED** | W | 0.851 | 0.487 | 0.111 | 0.066 | 0.089 | -0.108 |
|          | D | 0.692 | 0.520 | 0.093 | 0.101 | 0.051 | -0.068 |
| **GPM**  | W | 0.941 | 0.252 | 0.178 | 0.063 | 0.287 | -0.117 |
|          | D | 0.827 | 0.178 | 0.134 | 0.057 | 0.195 | -0.085 |
| **GMED** | W | 0.789 | 0.428 | 0.013 | 0.178 | 0.082 | -0.024 |
|          | D | 0.782 | 0.303 | 0.092 | 0.090 | 0.110 | -0.005 |

[a] W: Whole sample ($n = 91$)    D: Development sample ($n = 50$)

An examination of Table 1 gives the weight values of the inflation rate ($x_6$) for both the whole sample and the development sample as negative in all methods. It is an expected result for the weight value of this variable to become negative and the results obtained from all models also match this finding. While the weight values of variables $x_3$ and $x_4$ are expected to be positive, they are evaluated negatively in MSD model. This state of affairs does not answer the expectations.

The hit-ratios of classification obtained from six approaches are given in Table 2.

**Table 2:** Hit-ratios of six approaches

| Method | Sample[a] | Correctly accepted | Erroneously Accepted | Correctly rejected | Erroneously rejected | Hit-ratio |
|--------|-----------|--------------------|--------------------|--------------------|--------------------|-----------|
| **FLDF** | W | 17 | 7 | 61 | 6 | 0.857 |
|          | D | 8  | 5 | 33 | 4 | 0.820 |
|          | H | 8  | 5 | 25 | 3 | 0.804 |
| **MSD**  | W | 18 | 4 | 64 | 5 | 0.901 |
|          | D | 10 | 4 | 34 | 2 | 0.880 |
|          | H | 8  | 4 | 26 | 3 | 0.829 |
| **MLCM** | W | 19 | 2 | 66 | 4 | 0.934 |
|          | D | 10 | 3 | 35 | 2 | 0.900 |
|          | H | 9  | 3 | 27 | 2 | 0.878 |
| **LMED** | W | 18 | 1 | 67 | 5 | 0.934 |
|          | D | 10 | 3 | 35 | 2 | 0.900 |
|          | H | 9  | 2 | 28 | 2 | 0.902 |
| **GPM**  | W | 19 | 1 | 67 | 4 | 0.945 |
|          | D | 10 | 2 | 36 | 2 | 0.920 |
|          | H | 9  | 2 | 28 | 2 | 0.902 |
| **GMED** | W | 21 | 2 | 66 | 2 | 0.956 |
|          | D | 11 | 2 | 36 | 1 | 0.940 |
|          | H | 10 | 2 | 28 | 1 | 0.926 |

[a] W: Whole sample ($n = 91$),  D: Development sample ($n = 50$), H: Holdout sample ($n = 41$)

Development sample: 12 countries from group 1 and 38 countries from group 2, totally 50 countries. Holdout sample: 11 countries from group 1 and 30 countries from group 2, totally 41 countries.

When Table 2 is examined it can be observed that the model GMED model has the best hit-ratio of classification for both cases, the whole sample and the holdout sample. The GPM and LMED approaches give better results than the MLCM, MSD and FLDF approaches.

## 6. SIMULATION EXPERIMENT

We carried out a simulation work to compare the performances of the FLDF, MSD, MLCM, LMED, GPM and GMED methods for the two-group classification problem. All of the results for these methods are obtained by using the MATLAB 7 program.

The distributions considered in this study and their parameters are as follows:

Mult. Normal   $G_1 \sim (6, 6, 6)$,  $G_2 \sim (5, 5, 5)$ , $\Sigma = I$

Uniform         $G_1 \sim (1, 5)$, $G_2 \sim (1, 4)$

Exponential    $G_1 \sim (2)$, $G_2 \sim (1)$

Gamma          $G_1 \sim (1, 5)$, $G_2 \sim (1, 4)$

Chi-Square     $G_1 \sim (2)$, $G_2 \sim (1)$

Fisher (F)       $G_1 \sim (1, 4)$, $G_2 \sim (1, 5)$

For each group a random sample of 100 objects having three attributes are drawn randomly from the above distributions. From totally 200 observations of objects, 100 of them are used as the development sample and the remaining 100 observations are used as the holdout sample. For example, in Table 3, for the case $n_1 = n_2 = 50$, 100 observations for group 1 and 100 observations for group 2, making a total of 200 observations is drawn from the related distribution and then 50 observations of the random sample of size 100 for group 1 and so far group 2 are used as the development sample and the holdout sample, respectively. In Table 4, for the case $n_1 = 20, n_2 = 80$, again 100 observations for group 1 and 100 observations for group 2, making a total of 200 observations is drawn, but 80 observations of the random sample of size 100 are used as the development sample and the remaining 20 observations as the holdout sample for group 1. The same process is repeated for group 2. In Table 5, for the case $n_1 = 80, n_2 = 20$, the process is applied in the reverse order.

To run the simulation experiments, 50 data sets are generated for each of the distributions. That is, 50 data of sample are used.

The hit-ratio numbers of 6 procedures based on different distributions are given in Table 3, 4 and 5.

**Table 3:** Average hit-ratios of six approaches in the holdout samples ( $n_1 = n_2 = 50$ )

| Distribution | M E T H O D | | | | | |
| | MSD | FLDF | MLCM | LMED | GPM | GMED |
|---|---|---|---|---|---|---|
| **Mult.Normal** | 77.85 (6.31)* | 79.12 (4.31) | 80.25 (4.52) | 80.12 (5.06) | 80.35 (4.62) | 81.38 (4.15) |
| **Exponential** | 70.12 (6.55) | 73.23 (4.37) | 74.81 (5.41) | 75.27 (4.92) | 76.36 (3.68) | 77.61 (3.96) |
| **Uniform** | 73.25 (6.92) | 81.46 (3.25) | 81.51 (3.45) | 80.31 (4.96) | 80.35 (4.78) | 83.24 (3.39) |
| **Gamma** | 71.09 (5.61) | 69.12 (5.42) | 74.01 (3.25) | 74.12 (4.93) | 74.57 (4.36) | 76.91 (4.80) |
| **Chi Square** | 75.13 (5.22) | 68.32 (4.38) | 73.76 (3.83) | 76.02 (5.11) | 76.86 (4.43) | 78.12 (3.62) |
| **Fisher (F)** | 75.01 (4.61) | 68.82 (3.86) | 75.24 (2.36) | 75.11 (5.06) | 77.13 (4.12) | 78.49 (4.12) |

*The values in parentheses are standard deviations.

**Table 4**: Average hit-ratios of six approaches in the holdout samples  ( $n_1 = 20, n_2 = 80$ )

| Distribution | M E T H O D | | | | | |
| | MSD | FLDF | MLCM | LMED | GPM | GMED |
|---|---|---|---|---|---|---|
| **Mult.Normal** | 78.21 (4.02) | 81.88 (4.85) | 82.75 (4.52) | 78.45 (4.23) | 81.65 (3.72) | 83.89 (3.59) |
| **Exponential** | 74.19 (3.87) | 77.47 (4.74) | 76.15 (4.41) | 80.07 (4.19) | 80.23 (3.09) | 82.79 (3.28) |
| **Uniform** | 76.96 (4.78) | 82.63 (4.12) | 82.46 (3.91) | 83.31 (3.17) | 83.43 (3.78) | 85.12 (3.14) |
| **Gamma** | 74.52 (4.64) | 71.65 (4.49) | 77.58 (4.65) | 77.23 (4.20) | 78.95 (4.41) | 81.32 (4.07) |
| **Chi Square** | 75.98 (4.25) | 73.85 (4.19) | 77.08 (4.16) | 79.14 (3.45) | 80.02 (3.78) | 82.85 (3.67) |
| **Fisher (F)** | 76.88 (4.04) | 73.56 (5.08) | 77.13 (3.53) | 79.41 (4.44) | 79.52 (3.96) | 81.28 (4.19) |

**Table 5**: Average hit-ratios of six approaches in the holdout samples ( $n_1 = 80, n_2 = 20$ )

| Distribution | M E T H O D | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | *MSD* | *FLDF* | *MLCM* | *LMED* | *GPM* | *GMED* |
| *Mult.Normal* | 79.12 (5.23) | 80.59 (5.27) | 82.14 (4.21) | 79.43 (4.78) | 80.78 (4.62) | 83.12 (3.34) |
| *Exponential* | 73.59 (4.25) | 76.32 (4.89) | 77.25 (5.11) | 79.27 (5.24) | 79.23 (3.68) | 81.75 (3.09) |
| *Uniform* | 77.25 (5.12) | 81.92 (3.65) | 82.44 (3.98) | 81.31 (4.23) | 82.93 (4.12) | 84.67 (3.62) |
| *Gamma* | 73.19 (5.02) | 71.32 (4.41) | 76.87 (3.87) | 78.66 (4.93) | 79.98 (5.23) | 80.08 (4.33) |
| *Chi Square* | 76.23 (4.59) | 71.23 (3.88) | 76.13 (4.34) | 79.70 (4.23) | 79.83 (3.25) | 81.45 (2.94) |
| *Fisher (F)* | 76.71 (4.28) | 71.56 (4.01) | 76.37 (3.25) | 78.11 (5.12) | 80.13 (3.72) | 82.78 (4.36) |

From tables 3, 4 and 5, the GMED model seems to have a larger hit-ratio than other 5 procedures. In skewed distributions like Gamma, Chi Square and F, the superiority of GMED model is very obvious.

There are three main hypotheses to be tested and these hypotheses are used for 6 distributions and selected for development and holdout samples case $n_1 = n_2 = 50$, $n_1 = 20, n_2 = 80$ and $n_1 = 80, n_2 = 20$.

The hypotheses are set up as follows. $H_0$: There is no difference between the mean hit-ratio of GMED and the mean hit-ratio of model $i$; $H_1$: The mean hit-ratio of GMED is greater than the mean hit-ratio of model $i$ ( $i$ = FLDF, MSD, MLCM ).

In Table 3, for example, 79.12 and 4.31 values pertaining to FLDF model for the multinormal distribution are the mean hit-ratio and the standard deviation of the FLDF model in 50 random samples. Similarly the values 81.38 and 4.15 of the GMED model are also the mean hit-ratio and the standard deviation in 50 samples, respectively. Using this knowledge, we perform the following hypothesis test.

$H_0$ : There is no difference between the mean hit-ratio of GMED and the mean hit- ratio of FLDF; $H_1$ : The mean hit-ratio of GMED is greater than the mean hit-ratio of FLDF.

Table 6 lists the results of hypotheses tests which claim the mean hit-ratio the GMED model is greater than the mean hit-ratios of FLDF, MSD and MLCM models. The value 2.67[a] in this table shows the value of test statistic t of the hypothesis that the mean hit-ratio of the GMED model is greater than the mean hit-ratio of

the FLDF method for the multivariate normal distribution. [a] shows that the calculated p-value corresponding to 2.67 and is less than 0.01. We can infer the results of the hypothesis according to the value 2.67 of the test statistic or the p-value. For the significance value 0.01, a $p-$value less than 0.01 results in the hypothesis $H_0$ being rejected. This means that statistically, the mean hit-ratio of GMED is found to be greater than the mean hit-ratio of FLDF.

When the results of the hypotheses tests are examined, it is observed that in the case of nonnormal distributions, the GMED model is superior to all models in respect of the average hit-ratio.

**7. CONCLUSION**

In this study, for the without outlier cases, three new mathematical approaches have been developed in solving the two-group classification problem. When these new proposed approaches are applied to a real life problem or to a simulation data, it is seen that the approaches are practicable and efficient. Moreover, it is seen that the GMED approach is superior to all other models in respect of classification performance (the true hit-ratio obtained from holdout samples) for the case of nonnormal distributions. Meanwhile it is seen that the LMED, GPM approaches are also useful to the extent of the other models. For a furher study, the performances of the LMED and GMED approaches may be investigated for multi-group classification problems and for outlying cases in their group.

**Table 6**: Hypothesis test results (*t*-values of paired *t*-tests) of GMED and other (FLDF, MSD, MLCM) using holdout samples

| | $n_1 = 50, n_2 = 50$ | $n_1 = 20, n_2 = 80$ | $n_1 = 80, n_2 = 20$ |
|---|---|---|---|
| *Test 1 (FLDF)* | | | |
| **Mult.Normal** | 2.67[a] | 1.74 | 2.66[a] |
| **Exponential** | 5.25[c] | 4.76[c] | 6.18[c] |
| **Uniform** | 2.68[a] | 2.53[a] | 3.03[a] |
| **Gamma** | 7.60[c] | 8.82[c] | 8.06[c] |
| **Chi-Square** | 12.19[c] | 8.79[c] | 13.08[c] |
| **Fisher(F)** | 11.83[c] | 6.31[c] | 10.45[c] |
| *Test 2 (MSD)* | | | |
| Mult. Normal | 3.30[b] | 5.76[c] | 4.23[c] |
| **Exponential** | 6.92[c] | 9.14[c] | 9.75[c] |
| **Uniform** | 9.16[c] | 7.23[c] | 7.50[c] |
| **Gamma** | 5.57[c] | 5.99[c] | 6.16[c] |
| **Chi-Square** | 3.34[b] | 6.63[c] | 6.28[c] |
| **Fisher(F)** | 3.98[c] | 4.35[c] | 5.58[c] |
| *Test 3 (MLCM)* | | | |
| **Mult. Normal** | 1.30 | 1.08 | 1.11 |
| **Exponential** | 2.95[a] | 6.30[c] | 5.03[c] |
| **Uniform** | 2.52[a] | 2.83[a] | 2.41[a] |
| **Gamma** | 3.53[b] | 3.29[b] | 3.02[a] |
| **Chi-Square** | 5.85[c] | 5.67[c] | 6.53[c] |
| **Fisher(F)** | 4.83[c] | 4.52[c] | 6.15[c] |

[a] Reject $H_0$ ($sig. < 0.01$), [b] Reject $H_0$ ($sig. < 0.001$), [c] Reject $H_0$ ($sig. < 0.0001$)

**REFERENCES**

1. Fred, N., Glover, F., "A linear programming approach to the discriminant problem", *Decision Sciences*, 12: 68-74, (1981).

2. Glover, F., "Improved Linear Programming Models for Discriminant Analysis", *Decision Sciences*, 21: 771-785, (1990).

3. Lam, K.F., Moy, J.W., "An experimental comparison of some recently developed linear programming approaches to the discriminant problem", *Computers and Operations Research,* Vol: 24, No 7, 593-599, (1997).

4. Lam, K.F., Moy, J.W., "Combining discriminant methods in solving classification problemsin two-group discriminant analysis", *European Journal of Operational Research*, 138: 294-301, (2002).

5. Lam, K.F., Choo, E.U., Moy, J.W., "Minimizing deviations from the group mean: A new linear programming approach for the two-group classification problem", *European Journal of Operational Research*, 88, 358-367, (1996).

6. Lam, K.F., Moy, J.W., "Improved Linear Programming Formulations for the Multi-group Discriminant Problem", *Journal of Operational Research Society*, 47, 1526-1529, (1996).

7. *www.foreigntrade.gov.tr*