



## **International Journal of Contemporary Educational Research (IJCER)**

[www.ijcer.net](http://www.ijcer.net)

### **A Comparison of Classification Performances between the Methods of Logistics Regression and CHAID Analysis in accordance with Sample Size**

**Mehmet Şata<sup>1</sup>, Fuat Elkonca<sup>2</sup>**  
<sup>1</sup>Agri Ibrahim Cecen University  
<sup>2</sup>Mus Alparslan University

#### **To cite this article:**

Şata, M. & Elkonca, F. (2020). A comparison of classification performances between the methods of logistics regression and CHAID analysis in accordance with sample size. *International Journal of Contemporary Educational Research*, 7(2), 15-26. DOI: <https://doi.org/10.33200/ijcer.733720>

This article may be used for research, teaching, and private study purposes.

Any substantial or systematic reproduction, redistribution, reselling, loan, sub-licensing, systematic supply, or distribution in any form to anyone is expressly forbidden.

Authors alone are responsible for the contents of their articles. The journal owns the copyright of the articles.

The publisher shall not be liable for any loss, actions, claims, proceedings, demand, or costs or damages whatsoever or howsoever caused arising directly or indirectly in connection with or arising out of the use of the research material.

## **A Comparison of Classification Performances between the Methods of Logistics Regression and CHAID Analysis in accordance with Sample Size \***

**Mehmet Şata<sup>1†</sup>, Fuat Elkonca<sup>2</sup>**

<sup>1</sup> Agri Ibrahim Cecen University

<sup>2</sup> Mus Alparslan University

### **Abstract**

The aim of the study is to analyze how classification performances change in accordance with sample size in logistic regression and CHAID analyses. The dataset used in this study was obtained by means of “Attentional Control Scale.” The scale was applied to 1824 students and the analyses were done by randomly choosing the samples from the dataset. Nine classification criteria were determined in order to evaluate classification performances of logistic regression and CHAID analyses, and the results were interpreted in consideration of these criteria. As a result of the analyses, it was found that classification performance in logistic regression showed no change as sample size increased, and performed a better classification in small sample size (N= between 25 and 900) than CHAID analysis. On the other hand, in the method of CHAID analysis it was seen that classification performance improved as sample size increased, and provided stronger findings in large sample size (N= 1000 and above). Moreover, in classification studies logistic regression analysis yielded more reliable results, and CHAID analysis provided stronger classifications. The results of this study are considered to suggest researchers to select the methods in classification studies based on sample size.

**Key words:** Logistic regression, CHAID analysis, Classification, Sample size.

### **Introduction**

Classification is a method which is commonly used in scientific studies as well as daily life due to its benefit to problem solution (Köktürk, 2012). The decisions for placing students to a certain academic program, determining the individuals with psychopathology and the customers with a credit risk in a bank are examples for the studies indicating the importance of group membership. Many practitioners in various disciplines use different statistical methods to estimate the group membership belonging to any property (Finch & Schneider, 2007). In the classification studies, different results were obtained because of the reasons such as the existence of a great number of classification algorithm, each algorithm's having different parameters within itself, each algorithm's having more than one type, different purposes of different algorithms, the use of different data source, algorithms' support for different data types and the dependence of pre-treatments of data on the practitioner (Akpınar, 2000; Berry & Linoff, 1997).

In classification studies, the methods based on data mining are mostly used. The data mining techniques are influential in estimating determining crucial data classifications and data tendency by classifying especially large datasets. Among these methods, Logistic Regression, Decision Trees and Artificial Neural Networks are often used in classification studies (Kıran, 2010). Logistic regression and CHAID analyses of decision trees will be mentioned below as they were used within the scope of the study.

Logistic regression assumes that there is a logit relation between dependent and independent variables; therefore, logistic regression can provide non-linear models. The reasons for common preference of logistic regression in social sciences are that they impose no restriction on variable's being continuous or non-continuous; there is no condition for the likelihood function distribution of independent variables; there is no

---

\* This study was presented as an abstract proceeding at the 26th International Conference on Educational Sciences held between 20-23 April 2017.

† Corresponding Author: *Mehmet Şata*, [mehmetsata@gmail.com](mailto:mehmetsata@gmail.com); [msata@agri.edu.tr](mailto:msata@agri.edu.tr)

obligation for a linear relation between independent and dependent variable; there are many statistical package programs (Tabachnick & Fidell, 2013).

In logistic regression analysis, there is no assumption regarding the distribution of independent variables. However, when logistic regression analysis is used, there are some assumptions. First, the sample rates of variables in the analysis is an important necessity. Second, logistic regression analysis uses goodness of fit tests to evaluate model-data fit. The goodness of fits tests includes the values expected for each cell in the dataset of the combination of discrete variables. If any expected frequency in cells is very low (usually  $ef < 5$ ), the strength of the analysis is very low. Third, there should not be a problem with multicollinearity. Last, the extreme values of independent variables should be carefully examined (Tabachnick & Fidell, 2013).

CHAID analysis is one of the methods in decision trees which are commonly used in data mining. CHAID analysis uses chi-square ( $\chi^2$ ) if a dependent variable is discontinuous and F statistics if it is a continuous as it is concerned especially with the relation between independent variables and their interactions. As it is known, chi-square test statistics deals with the dependence between the variables. Due to this property, it is consistent that CHAID analysis establishes mathematical models on a chi-square ground. In addition, CHAID analysis helps to provide objective and robust results in the evaluation of the effects of sociodemographic variables in the sample and the sub factors of assessment instrument on dependent variables (Kayri & Boysan, 2007).

Classification studies are often seen in the disciplines of education, medicine and banking. In classification studies, there are a great number of models and different algorithms belonging to these models. The answers to the questions which one of these algorithms provide more accurate results, which algorithms are more successful in certain disciplines will increase the practices' success and enhance proficiency of the work. Therefore, an evaluation of algorithms by comparison is of great importance. However, classification studies are usually conducted with one sample, and also as a result of the difference in sample size classification studies provide different results. Therefore, in classification studies it is not determined which method is more effective in accordance with sample size. For example, classification performances of logistic regression, artificial neural networks and decision trees methods were compared based on small samples by Sabzevari, Soleymani and Noorbakhsh (2007), and logistic regression was indicated the most successful method. On the other hand, in the study conducted by Kiran (2010) decision trees were found to make a better classification than logistic regression. In the study by Neuilly, Zgoba, Tita, and Lee (2011) decision trees were found to have a lower classification rate than logistic regression. In another study by Heckert and Gondolf (2005), binary and multinomial logistic regression among logistic regression methods were compared with CHAID, Exhaustive CHAID, CART and QUEST analyses among decision trees methods. As a result of the analyses, logistic regression analysis has a better classification rate than decision trees. In addition, in other several classification studies, these methods were found to have different classification methods (Ekici, 2012; Karakış, 2009; King, Feng & Sutherland, 1995; Kurt & Türe, 2005; Zurada & Lonial, 2005). These differing results in the literature cause a confusion regarding which method is better and complicate the method preference of researchers in classification studies.

Therefore, in this study the classification performances of these methods from small sample size ( $n=25$ ) to large sample size ( $n= 1800$ ) were analyzed, and it is aimed to lead researchers which method to use in accordance with their sample size. In this sense, the study is expected to contribute to the literature. Considering that there is no empirical study comparing classification performances in accordance with sample size, and there are only studies with simulative data based on simulation methods (Dolgun, 2014), the current study is considered to have great importance for the relevant literature. In this study, it is hypothesized that logistic regression provides better results in small sample size and classifies within the acceptable error limit ( $\alpha=0.05$ ), CHAID analysis yields better results in a large sample size and classifies within the acceptable error limit ( $\alpha=0.05$ ).

## Method

### Research Type

In this study, how classification performances of logistic regression and CHAID analysis changed in accordance with sample size. As the relevant study determines the existent situation, it is a type of descriptive study.

### Study Group

The study group of the study is composed of total 1824 students of 700 females and 1124 males in the high schools of Batman Provincial Directorate of National Education through convenience sampling which is one of

the non-random sampling methods. The samples were randomly created by the study group. While the samples were created, the condition for being five people minimum per each cell in the dataset for independent variable which is an assumption of logistic regression is taken into consideration. As there are five different variables in the study, the smallest sample size was determined as 25 and selected randomly from the total dataset. After this procedure, the other sample sizes were determined as multiples of 25 and it continued to reach 1800 participants. In this way, 72 samples in total were created and both analysis methods were applied to these samples.

### **Data Tools**

In this study, the “Attentional Control Scale” developed by Fajkowska and Derryberry (2010) and translated into Turkish by Akın et al. (2013) was used as an instrument for data collection. Before the scale was applied, the necessary permission was obtained from the researchers translating into Turkish. Moreover, “Personal Information Form” was developed by the researcher to collect demographic information in accordance with the aim of the study. The items were scaled from negative to positive in a 4-point Likert scale. The points to be taken from the scale change between 20 and 80. The high score of total points from the scale indicate students’ high level of attentional control while the low score of total points from the scale indicate students’ low level of attentional control. The reliability and validity analyses were done and the reliability of the scale was found McDonald  $\omega$  0.894 (%CI 0.887 – 0.900). The confirmatory factor analysis was done for validity and it was found RMSEA = 0.077, SRMR = 0.057, IFI = 0.95, RFI = 0.94, NNFI = 0.94 and CFI = 0.95. In other words, the data collection instruments were found to be reliable and valid.

### **Data Analysis**

In the study, logistic regression and CHAID analyses were conducted in order to determine classification performances in accordance with the sample size. For the analysis, SPSS (Version 25.0) and Mplus (Version 7) and Origin (Version 8) package programs were utilized. Polynomial fit was drawn to the graphics by means of origin package program in order to make these obtained graphics more understandable for interpretation. This polynomial fit was gained by creating a polynomial regression model between dependent and independent variables.

Two-step cluster analysis was used to categorize the scale scores as the total scores from the Attentional Control Scale is continuous. In this way, the total score of the scale was formed into categorical two clusters (high level and low level). The main reason for transferring total continuous score of the scale into categorical score is that CHAID analysis provided better results in categorical data (Pehlivan, 2006), and it is aimed to compare classification performances with logistic regression.

Before the analyses were conducted, multivariate normal distribution hypothesis, outlier, missing value and the multicollinearity were analyzed. As a result of the analyses, it is found that there are no outliers and multicollinearity, but there is a missing value and multivariate normal distribution could not be ensured.

For missing value, a new value was attained by using the EM (Expectation Maximization) algorithm. As there is no obligation to counter to normal distribution hypothesis of both CHAID and logistic regression analyses, analysis were done without any procedure for multivariate normality. In the current study, the same model was used to provide consistent results in the considered analysis methods.

### **Classification Criteria**

In classification studies, the performances of analysis methods were determined by means of certain criteria. When the literature is examined, it is seen that there are various criteria. The most used criteria are as follows: positive likelihood rate (PLR), negative likelihood rate (NLR), Type I error rate, Type II error rate, confidence level, power of test, sensitivity, specificity and total accurate classification percentage. The classification table is used to calculate these parameters as shown in Table 1 (Koyuncu, 2015).

Table 1. Classification Table

		Expected (Actual) State				Total
		Positiive	N	Negative	N	
Observed (Test Result) State	Positive	Correct Positive (Confidence Level)	a	Incorrect Positive (Type I Error)	c	a+c
	Negative	Incorrect Negative (Type II Error)	b	Correct Negative (Power of Test)	d	b+d
	Total		a+b		c+d	a+b+c+d

Positive likelihood rate (PLR) is a classification performance measurement by combining sensitivity and specifying classification studies. It is obtained by dividing correct positive rate into incorrect positive rate. This value indicates how many incorrect positive results were provided for each correct positive result in classification study. It is usually demanded for PLR's being as high as possible (Deeks & Altman, 2004; Grimes & Schulz, 2005; Medcalc, 2018, s.222). PLR values are calculated as given below:

$$\text{PLR} = \text{Sensitivity}/(1-\text{Specificity}) \quad (1)$$

Negative likelihood rate (NLR) is obtained by dividing incorrect negative rate into correct negative rate. This value indicates how many correct negative results were provided for each incorrect negative result in classification study. It is usually demanded for NLR's being as low as possible (Deeks & Altman, 2004). NLR values are calculated as given below:

$$\text{NLR} = (1-\text{Sensitivity})/\text{Specificity} \quad (2)$$

Type I error rate ( $\alpha$ ) is an error rate when the test result is accepted as positive though it is negative in reality. Type I error is also known as  $\alpha$  error (Cohen, 1988, p.4). It is accepted to have low values in practice. Considering that Type II error rate will be high as the Type I error rate is lower, it is necessary to use large samples in order to lower both error types (Tan, 2016, p.265). Type I error rate is calculated in 2x2 tables as follows:

$$\alpha = c/(a+b+c+d) \quad (3)$$

Type II error rate ( $\beta$ ) is an error rate when the test result is accepted as negative though it is positive in reality (Tan, 2016, p.265). Type II error is also known as  $\beta$  error (Cohen, 1988, p.5). It is accepted to have low values in practice similar to Type I error rate. In the studies, while Type II error rate is not stated while Type I error rate is calculated. Type II error rate is calculated in 2x2 tables as follows:

$$\beta = b/(a+b+c+d) \quad (4)$$

The confidence level of the test is the likelihood of not having Type I error, and the power of the test is the likelihood of not having Type II error (Tan, 2016, p.265). The confidence level of the test is obtained by subtracting Type I error rate from 1, while the power of the test is obtained by subtracting Type II error rate from 1.

Sensitivity is described as the positive likelihood though it is positive in reality (Deeks & Altman, 2004; Medcalc, 2018, p.222). At the same time, it is known as an accurate positive rate. Sensitivity is calculated in 2x2 tables as follows:

$$\text{Sensitivity} = a/(a+b) \quad (5)$$

Specificity is described as the negative likelihood as it is negative in reality (Deeks & Altman, 2004; Medcalc, 2018, p.222). At the same time, it is known as an accurate negative rate. Specificity is calculated in 2x2 tables as follows:

$$\text{Specificity} = d/(c+d) \quad (6)$$

In the current study, classification performances of logistic regression and CHAID analyses were compared in accordance with the nine classification criteria mentioned above. The graphics were used for comparison and polynomial fit were drawn to the graphics to make it more understandable.

## Findings

In the created regression model for both CHAID and logistic regression analyses, the grouped state (high control/low control) of total score obtained from the scale as dependent variable were used. As independent variable, the students' sex, age, state of living with their families, monthly income of their families and use of smart phone were included in the model. This created regression model was randomly selected and was tested in total 72 samples. Then, the results were provided by graphics for each classification criteria. First, the positive likelihood in accordance with sample size in CHAID and logistic regression analyses is shown in Figure 1.

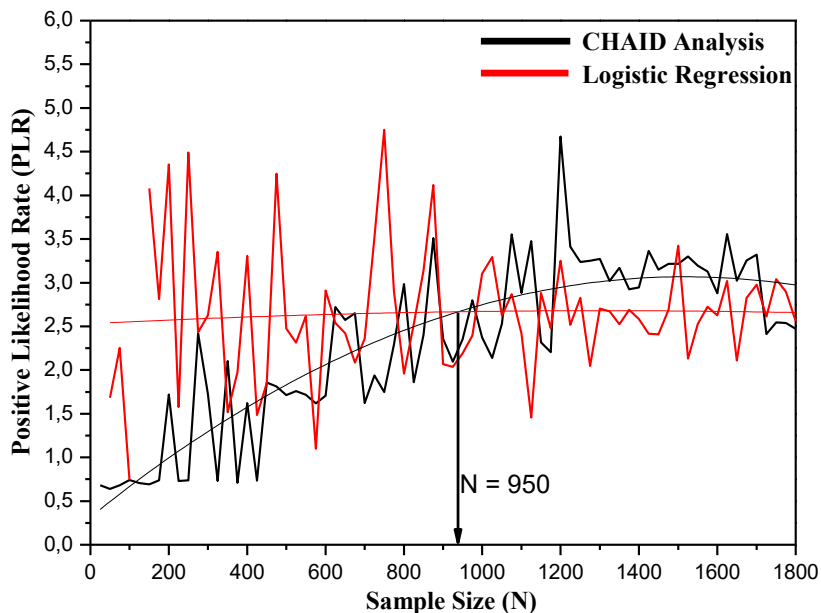


Figure 1. Positive likelihood rate for both analyses

When Figure 1 was examined, logistic regression analysis showed no change in accordance with sample size; however, in CHAID analysis PLR value is seen to increase as sample size increases. It is clear that logistic regression provides more desired results in classification studies in sample size up to 1000. Therefore, the sample size can be argued an important factor in selecting the method of analysis used in classification studies. The change in negative likelihood rate (NLR) in accordance with sample size, which is another classification criterion, is shown in Figure 2.

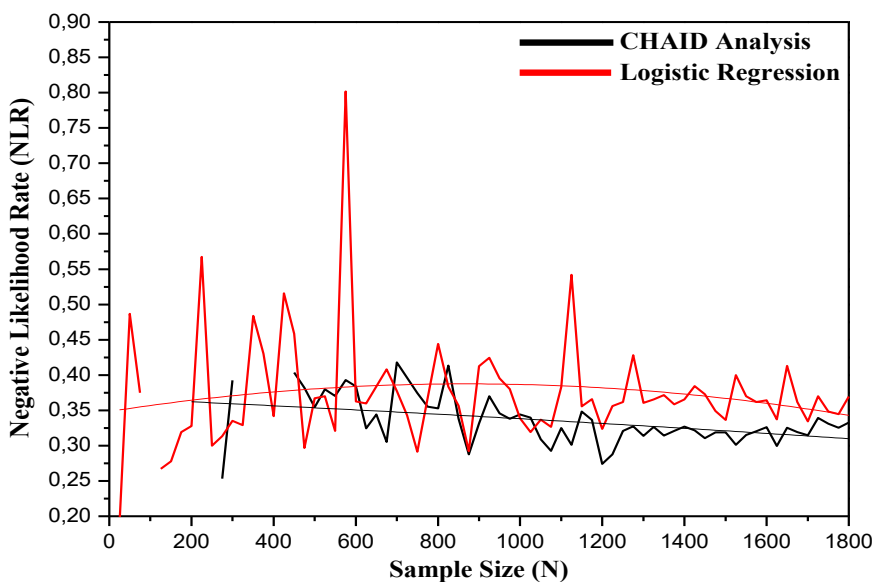


Figure 2. Negative likelihood rate for both analyses

When Figure 2 was examined, both analysis methods did not show much change as sample size increased. CHAID analysis is seen to have discrete values when sample size is approximately at 600. The main reason is derived from the fact that specificity percentage of in all the samples except two samples in CHAID analysis is zero. After the analyses of PLR and NLR values, the change of Type I error rate, which is a criterion often used in classification studies, in accordance with sample size was analyzed and the obtained values are shown in Figure 3.

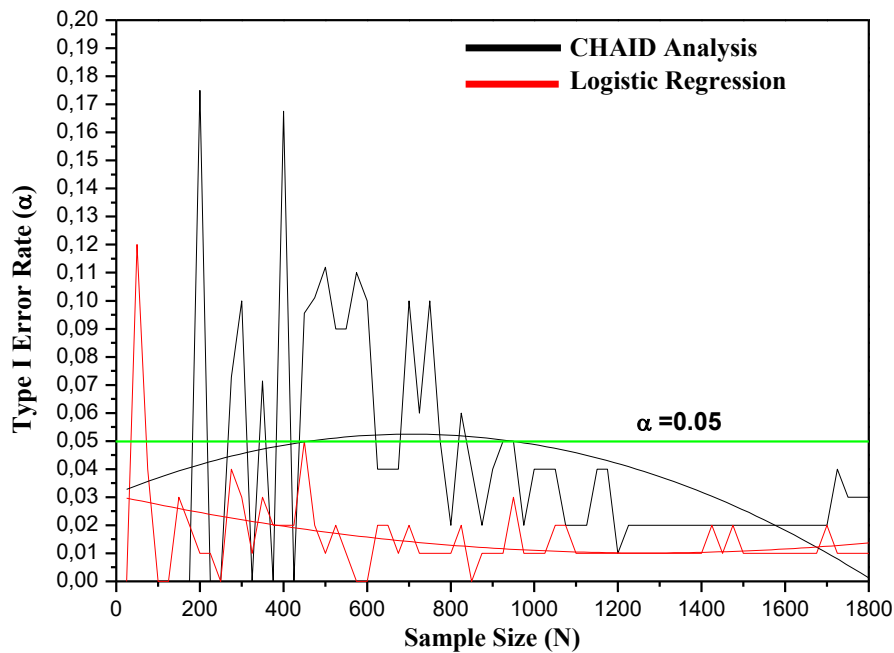


Figure 3. Type I error rate for both analyses

In Figure 3, it is seen Type I error rate decreases in logistic regression as the sample size increases while in CHAID analysis Type I error rate increases to a certain value and then decreases. Moreover, when the sample size is approximately between 440 and 920, it is seen that the error rate obtained from CHAID analysis exceed the level  $\alpha = 0.05$  which is usually accepted as a range in social sciences. In the study, the change of Type II error rate in accordance with sample size was analyzed, and the graphic regarding the results is shown in Figure 4.

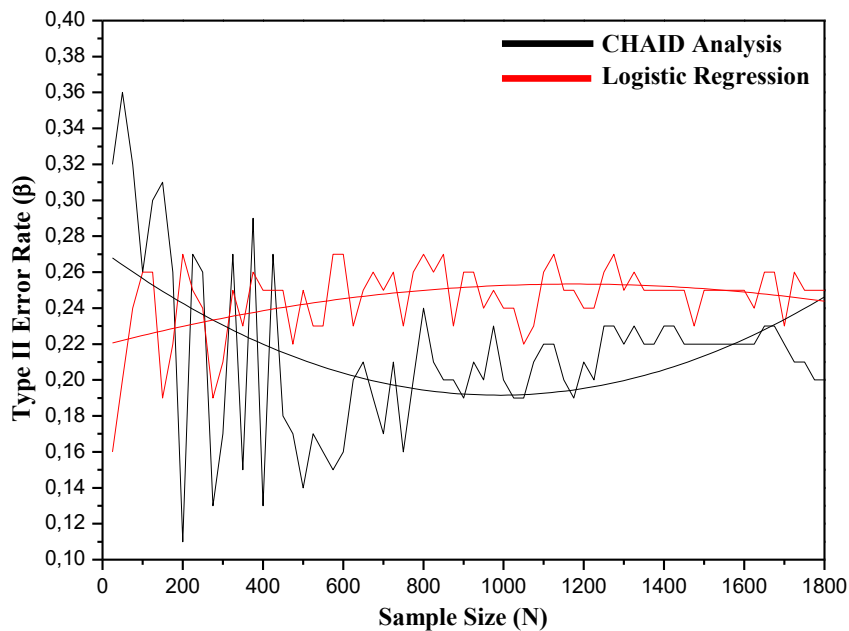


Figure 4. Type II error rate for both analyses

When Figure 4 was analyzed, it is seen Type II error rate did not show much change in logistic regression as the sample size increased while in CHAID analysis Type II error rate decreased and then increased once again. The reason is considered to derived from the considerably high level of specificity percentage while CHAID analysis classifies to a certain sample size (approximately 600 samples). After analyzing the change of Type I and Type II error rates in accordance with sample size, how the confidence level of the test changed was analyzed, and the results are shown in Figure 5.

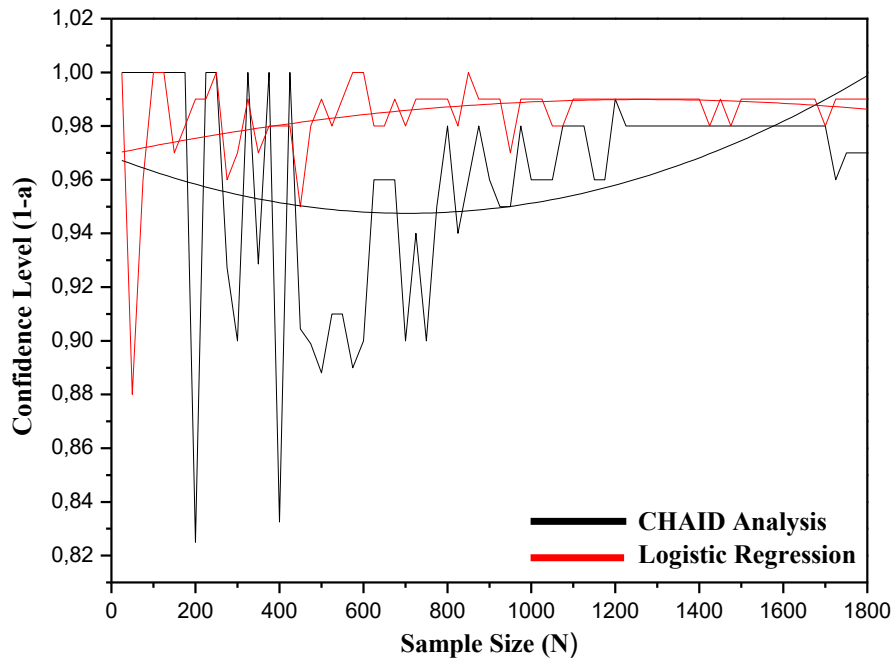


Figure 5. The confidence level of the test for both analyses

When Figure5 was analyzed, the confidence level of the classification through logistic regression analysis was found to be higher. However, as the sample size increased, the confidence level of the classifications through CHAID analysis was found to increase and provide more reliable classifications after n= 1600 samples than logistic regression. In the study, the change of power of classification tests in accordance with sample size was analyzed and the results are shown in Figure 6.

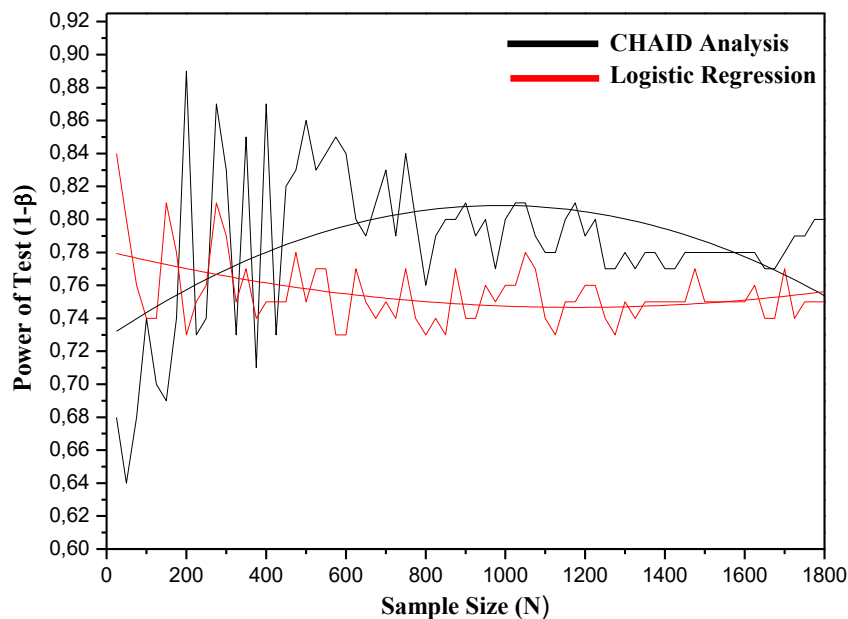


Figure 6. The power of test for both analyses

When Figure 6 was analyzed, the power of classification in logistic regression did not show much change as the sample size increased, while it first increased and then decreased in CHAID analysis. It is seen that CHAID



analysis has a weak power at approximately  $n=400$ . In the study, the change of specificity percentages in accordance with sample size was analyzed, and the results are shown in Figure 7.

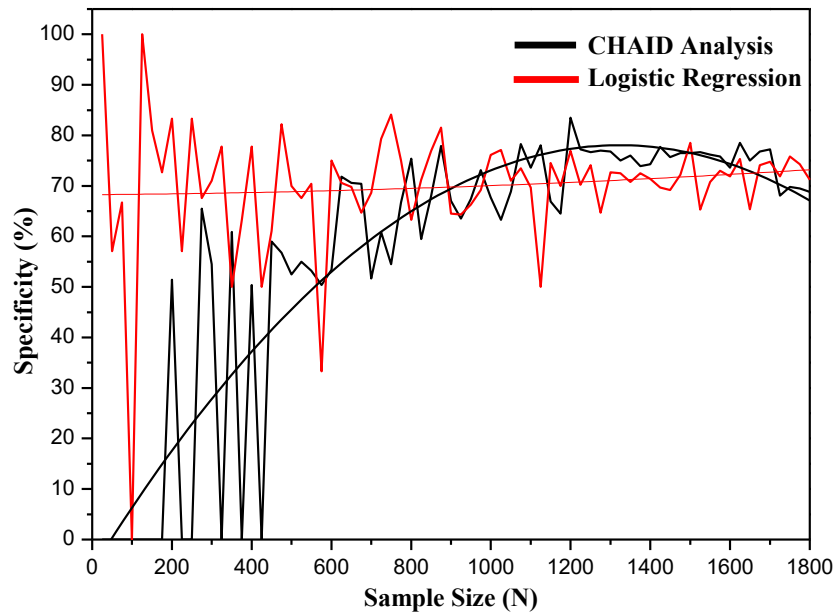


Figure 7. Specificity percentages for both analyses

When Figure 7 was analyzed, it is seen that specificity was very low in CHAID analysis with small samples and logistic regression was almost independent from sample size. In classification studies when CHAID analysis was used, it should be remembered that specificity would be low at size between 0 and 700 samples. If a specificity is an important criterion in a classification study and the sample size is small, the preference of logistic regression might be more reliable and valid. In the study, the change of classification sensitivity of both analysis methods was analyzed, and the results are shown in Figure 8.

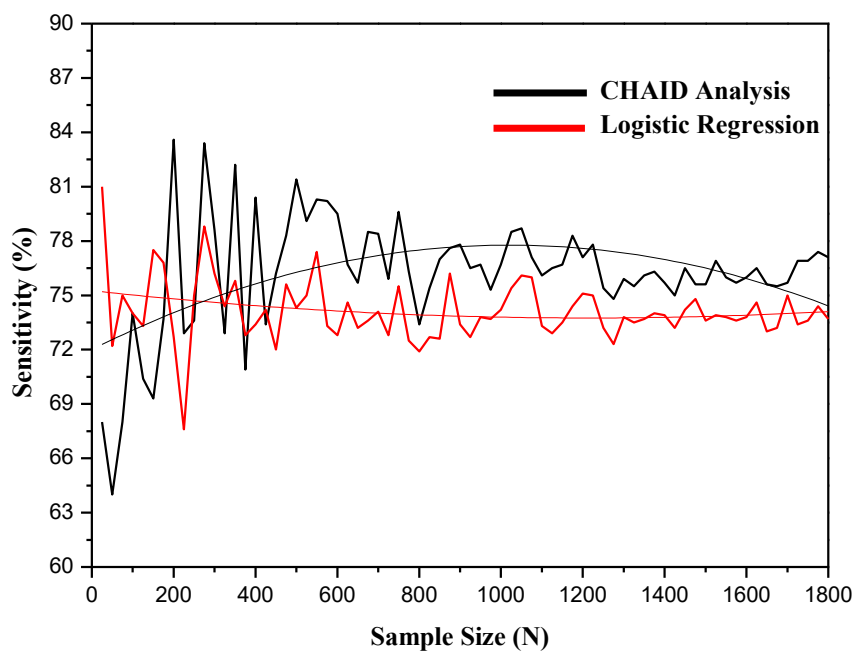


Figure 8. Sensitivity percentages for both analyses

As seen in Figure 8, as the sample size increased, the sensitivity percentage in CHAID analysis increased, logistic regression hardly changed and was independent from sample size. Last, the change of total classification percentages of both analysis methods was analyzed, and the results are shown in Figure 9.

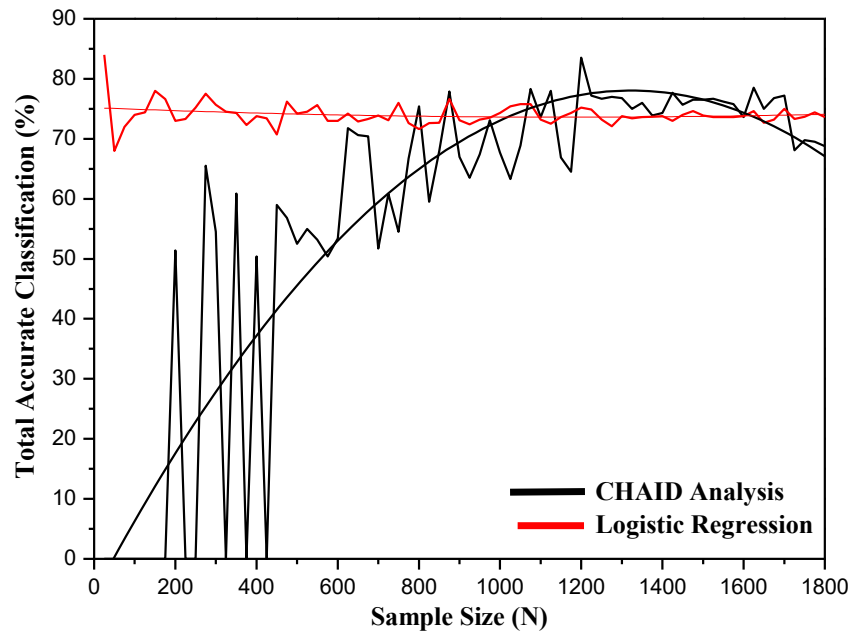


Figure 9. Total accurate classification percentages for both analyses

When Figure 9 was analyzed, it is seen that accurate classification percentage of CHAID analysis increased as the sample size increased and there was no considerable change in logistic regression. Moreover, it is found that logistic regression provided more accurate classification at small samples.

## Discussion, Conclusion and Suggestions

The aim of this study is to compare classification performances of CHAID and logistic regression analysis in accordance with sample size. Therefore, nine classification criteria were determined to compare classification performances of both analyses. In addition, at the beginning of the study based on the literature the hypothesis was developed that logistic regression analysis made better classification with small samples and CHAID analysis made better classification with large samples.

In a scientific study, “good sample” rather than large sample” should be aimed (Balci, 2015, p.105). This is related to be systematic and sensitive in the process of selection. The sample should be as small as possible, but should represent the population at a high level in order to accomplish its goal (Balci, 2015, p.105). Therefore, the sample size especially in classification studies in scientific studies is a crucial factor. However, inappropriate sample size causes inaccurate or subjective results. This leads us to take incorrect decisions and use resources inefficiently (Stafford et al., 2006).

Based on the findings of the study, the sample size in accordance with the analysis methods used in classification studies can be considered as an important factor. It is observed that logistic regression analysis should be preferred in classification studies rather than decision trees methods particularly in small samples. In a similar way, in the study of Nemes et al., (2009) which the classification performance of the logistic regression in accordance with sample size was analyzed, logistic regression has more accurate classification with small samples. In the study conducted by Demidenko (2007), the likelihood rate in logistic regression was seen to decrease as sample size increased.

When the literature was analyzed, there are various studies which compare classification performances of logistic regression and decision trees. However, it was observed that classification performances of the methods used in these studies showed variance among studies (Ekici, 2012; Brewer, 2012; Neuilly et al., 2011; Kıran, 2010; Çakır, 2008; Heckert & Gondolf, 2005). Therefore, it has become difficult to choose the method to be used in the classification studies. In this sense, the current study provides evidence for which method researchers would choose in accordance with the sample size they can reach. According to the analyses, it is seen more appropriate to choose logistic regression analysis in small sample sizes and one of both analyses in large sample sizes.

In classification studies, as Type I error rate is an important criterion it is required to determine the threshold value of Type I error rate. In the current study, Type I error rate was analyzed and the threshold value was identified as  $\alpha = 0.05$ . As a result of the analyses, logistic regression classified below the identified error rate, but CHAID analysis was seen to classify above the error rate at small samples. Considering this finding, the use of logistic regression is seen to be more appropriate with small samples and in the studies whose error rate is important. When the international literature was examined, in classification studies it is seen that Type I error rate was taken into consideration while the studies at national level do not consider Type I error rate in classification studies. In the current study, CHAID analysis exceeded the  $\alpha = 0.05$  level for some sample groups. Therefore, the accuracy of the classification is suspected. It is understood that Type I and Type II error rates should be consulted in classification studies.

The power of classification of the used analysis is also an important criterion in classification studies. Particularly at large datasets, the speed of classification methods becomes an important factor. Therefore, the obtained findings suggested that CHAID analysis had stronger classification performance. Similarly, decision trees were concluded to classify more quickly in the study by Çakır (2008).

When the graphics, in which the nine criteria within the scope of the study, were analyzed, it was found that there was much fluctuation in both analysis methods in small sample sizes. This fluctuation causes doubt in the accuracy of the classification studies with small samples (between 0 and 400). Therefore, it is considered necessary to determine the sample size before conducting classification study.

There are some limitations of the study. First, although the dependent variable is continuous by nature, it was transformed into an artificial categorical form by means of two-step cluster analysis. Therefore, the analysis was conducted by the classification accuracy of two-step clustering analysis. The second limitation is that this study is limited to merely two of the classification methods. Finally, the distribution of the data used in this study is not normal. Therefore, it is useful for further studies to consider these limitations and conduct their research.

To sum up the changes of both analysis methods in accordance with sample size for the considered nine criteria within the scope of the study;

- It was observed that positive likelihood rate in logistic regression showed no change and positive likelihood rate in CHAID analysis increased as the sample size increased. In this case, at the beginning of the study, it can be accepted that the CHAID analysis classifies better in large samples.
- Considering the negative likelihood rate, there was no significant change in both analysis methods.
- Considering the change of Type I error rate, which is an important criterion in classification studies, in accordance with sample size, it was observed that Type I error rate decreased in both analysis methods despite increasing sample size. However, it should not be forgotten the reliability and validity of the classification done with the related samples were weak because CHAID analysis exceeded error limits in some sample sizes.
- Considering another criterion Type I error rate, it was seen that logistic regression was independent from sample size and CHAID analysis decreased as sample size increased. Bulut (2015) stated that Type I and Type II error rates decreased as sample size increased. It can be argued that the current study was supported with the literature.
- Considering the confidence level of the test, it was seen that it showed no change as sample size increased while it increased in CHAID analysis. Moreover, when the confidence level was examined, logistic regression was found to have a higher level.
- Considering the power of the test, logistic regression showed no change in accordance with sample size, while it increased in CHAID analysis. Moreover, it was found that CHAID analysis had stronger classification performance than logistic regression. Therefore, if the power of the test is important in classification studies, CHAID analysis might be useful; on the other hand, if the confidence level of the test is important, logistic regression might be useful.
- Last, considering specificity, sensitivity and total accurate classification, logistic regression showed no change in accordance with sample size while CHAID analysis increases. In addition, it was observed that logistic regression made a better classification with small samples.

To conclude, it was observed that logistic regression showed no change, but CHAID analysis showed change as sample size increased. The hypothesis which logistic regression made a better classification with small samples within acceptable error limits ( $\alpha = 0.05$ ) can be accepted. In a similar vein, the hypothesis which CHAID analysis made a better classification with large samples within acceptable error limits ( $\alpha = 0.05$ ) can be accepted.

Based on the results obtained from the study, these suggestions can be made:

- In this study, two analysis methods used in the classification studies were used. There are many methods used in classification studies ; therefore, similar studies can be conducted for different analysis methods.
- It is expected that logistic regression would be useful in a study which the confidence level is important (for example, clinical studies), and CHAID analysis would be useful in a study which the power of the test is important (for example, large datasets).
- Considering the limitations of the current study, further studies might contribute to the literature.

## Acknowledgements or Notes

A part of this study is presented as an oral presentation at 26th International Conference on Educational Sciences in Antalya.

## References

- Akın, A., Kaya, Ç., Uysal, R., Çardak, M., Çitemel, N., Özdemir, E., & Gülşen, M. (2013). *Dikkat Kontrol Ölçeği Türkçe Formu: Geçerlik ve Güvenirlik Çalışması [The Turkish version of the attentional control scale: the validity and reliability study]*. Paper presented at VI. National Graduate Education Symposium. Retrieved from [http://www.academia.edu/download/43723223/Eitim\\_Modelinin\\_renci\\_zerindeki\\_Etkilili20160314-25744-1i99q7c.pdf#page=19](http://www.academia.edu/download/43723223/Eitim_Modelinin_renci_zerindeki_Etkilili20160314-25744-1i99q7c.pdf#page=19)
- Akpınar, H. (2000). Veri tabanlarında bilgi keşfi ve veri madenciliği [Knowledge discovery and data mining in databases]. *Istanbul Business Research*, 29(1), 1-22. Retrieved from <https://dergipark.org.tr/tr/pub/ibr/archive>
- Balcı, A. (2015). *Sosyal bilimlerde araştırma yöntem, teknik ve ilkeler [Research methods, techniques and principles in social sciences]*. Ankara: Pegem Akademi.
- Berry M., & Linoff G., (1997). *Data Mining Techniques for Marketing Sales and Customer Support*. John Wiley & Sons.
- Brewer S. L. (2012). *An empirical comparison of logistic regression to decision tree induction in the prediction of intimate partner violence reassault*. (Doctoral dissertation). Retrieved from <https://www.proquest.com/>
- Bulut, N. (2015). *İzleme amaçlı klinik araştırmalarda öngörülen ölçütlere göre örneklem büyüklüğünün belirlenmesi [Determination of sample size by criterias proposed on monitoring in clinical research]*. (Master thesis). Retrieved from <https://tez.yok.gov.tr/UlusalTezMerkezi/>
- Çakır, Ö. (2008). *Veri madenciliğinde sınıflandırma yöntemlerinin karşılaştırılması "bankacılık müşteri veri tabanı üzerinde bir uygulama" [ Comparison of classification methods in data mining "an application on banking customer database"]*. (Doctoral dissertation). Retrieved from <https://tez.yok.gov.tr/UlusalTezMerkezi/>
- Cohen, J. (1988). *Statistical power analysis for the behavioral sciences*. NJ: Erlbaum Hillsdale.
- Deeks, J. J., & Altman, D. G. (2004). Diagnostic tests 4: likelihood ratios. *Bmj*, 329(7458), 168-169. <https://doi.org/10.1136/bmj.329.7458.168>
- Demidenko, E. (2007). Sample size determination for logistic regression revisited. *Statist. Med.*, 26, 3385–3397. <https://doi.org/10.1002/sim.2771>
- Ekici, E. (2012). *Farklı sınıflandırma yöntemlerinin karşılaştırılması ve bir uygulama [An application on the comparison of various classification methods]*. (Master thesis). Retrieved from <https://tez.yok.gov.tr/UlusalTezMerkezi/>
- Fajkowska, M. & Derryberry, D. (2010) . Psychometric properties of Attentional Control Scale: The preliminary study on a Polish sample. *Polish Psychological Bulletin*, 41(1), 1-7. <https://doi.org/10.2478/s10059-010-0001-7>
- Finch, H., & Schneider, M. K. (2007). Classification accuracy of neural networks vs. discriminant analysis, logistic regression, and classification and regression trees. *Methodology*, 3(2), 47-57. <https://doi.org/10.1027/1614-2241.3.2.47>
- Grimes, D. A., & Schulz, K. F. (2005). Refining clinical diagnosis with likelihood ratios. *The Lancet*, 365(9469), 1500-1505. [https://doi.org/10.1016/S0140-6736\(05\)66422-7](https://doi.org/10.1016/S0140-6736(05)66422-7)
- Heckert, D.A., & Gondolf, E.W. (2005). Do multiple outcomes and conditional factors improve prediction of batterer reassault? *Violence and Victims*, 20 (1), 3-24. <https://doi.org/10.1891/vivi.2005.20.1.3>
- Karakış, R., (2009). *Yapay sinir ağları ve lojistik regresyon yöntemleri ile meme kanseri koltuk altı lenf nodu*

- durumunun belirlenmesi[Prediction of the axillary lymph node status in breast cancer using artificial neural network and logistic regression analysis methods]. (Master thesis). Retrieved from <https://tez.yok.gov.tr/UlusalTezMerkezi/>
- Kayri, M., & Boysan, M. (2007). Araştırmalarda CHAID analizinin kullanımı ve baş etme stratejileri ile ilgili bir uygulama[Using Chaid analysis in researches and an application pertaining to coping strategies]. *Ankara University Journal of Faculty of Educational Sciences*. 40(2), 133-149. [https://doi.org/10.1501/Egifak\\_0000000146](https://doi.org/10.1501/Egifak_0000000146)
- King, R. D., Feng, C., & Sutherland, A. (1995). Statlog: comparison of classification algorithms on large real-world problems. *Applied Artificial Intelligence an International Journal*, 9(3), 289-333. <https://doi.org/10.1080/08839519508945477>
- Kıran, Z. B. (2010). Lojistik regresyon ve CART analizi teknikleriyle sosyal güvenlik kurumu ilaç provizyon sistemi verileri üzerinde bir uygulama[An application on pharmacy provision system data of social security institution by logistic regression and CART analysis technics]. (Master thesis). Retrieved from <https://tez.yok.gov.tr/UlusalTezMerkezi/>
- Köktürk, F. (2012). K-en yakın komşuluk, yapay sinir ağları ve karar ağaçları yöntemlerinin sınıflandırma başarılarının karşılaştırılması[comparing classification success of k-nearest neighbor, artificial neural network and decision trees]. (Doctoral dissertation). Retrieved from <https://tez.yok.gov.tr/UlusalTezMerkezi/>
- Koyuncu, M. S., (2015). Psikolojik ölçeklerde ROC analizi yöntemiyle standart belirleme[Standard determination in psychological scales using ROC analysis]. (Master thesis). Retrieved from <https://tez.yok.gov.tr/UlusalTezMerkezi/>
- Kurt, İ. & Türe, M.(2005). Tıp öğrencilerinde alkol kullanımını etkileyen faktörlerin belirlenmesinde yapay sinir ağları ile lojistik regresyon analizi'nin karşılaştırılması[Comparison of artificial neural networks and logistic regression analysis in determining factors affecting alcohol consumption among medicine students]. *The Balkan Medical Journal*. 22(3), 142-153. Retrieved from <https://dergipark.org.tr/en/pub/bmj/issue/3749/49838>
- Medcalc. (2018). *Software manual*. Retrieved from <https://www.medcalc.org/download/medcalcmanual.pdf>
- Nemes, S., Jonasson, J.M., Genell, A., & Steineck, G. (2009). Bias in odds ratios by logistic regression modelling and sample size. *BMC Medical Research Methodology*, 56(9), 1-5. <https://doi.org/10.1186/1471-2288-9-56>
- Neuilly, M. A., Zgoba, K. M., Tita, G. E., & Lee, S. S. (2011). Predicting recidivism in homicide offenders using classification tree analysis. *Homicide Studies*, 15(2), 154-176. <https://doi.org/10.1177/1088767911406867>
- Pehlivan, G. (2006). CHAID analizi ve bir uygulama[CHAID analysis and an application]. (Master thesis). Retrieved from <https://tez.yok.gov.tr/UlusalTezMerkezi/>
- Sabzevari, H., Soleymani, M., & Noorbakhsh, E. (2007). A comparison between statistical and data mining methods for credit scoring in case of limited available data. In *Proceedings of the 3rd CRC Credit Scoring Conference* (pp. 1-5).
- Stafford, J.D., Kaminski, R.M., Reinecke K.J., & Gerard, P.D., (2006). Multi-stage sampling for large scale natural resources surveys: a case study of rice and waterfowl. *Journal of Environmental Management*, 78, 353-361. <https://doi.org/10.1016/j.jenvman.2005.04.029>
- Tabachnick, B.G. & Fidell, L.S. (2013). *Multivariate statistics*. New Jersey: Pearson Education Inc.
- Tan, Ş. (2016). SPSS ve excel uygulamalı temel istatistik-I[Basic statistics-I with SPSS and excel application]. Ankara: Pegem Akademi. <https://doi.org/10.14527/9786053183877>
- Zurada, J., & Lonial, S. (2005). Comparison of the performance of several data mining methods for bad debt recovery in the healthcare industry. *Journal of Applied Business Research*, 21(2), 37-54. <https://doi.org/10.19030/jabr.v21i2.1488>