

Logistic regression and classification tree methods as elements of diagnosis in cardiology

Metody regresji logistycznej i drzewa klasyfikacyjnego jako elementy procesu diagnostycznego w dziedzinie kardiologii

Anna Spychała, Michał Skrzypek, Ewa Niewiadomska

Department of Environmental Health, School of Public Health in Bytom, Medical University of Silesia in Katowice, Poland

ABSTRACT

INTRODUCTION: The purpose of statistical analysis in research is to identify accurate and reliable conclusions where the researcher has a great deal of sources and information. Usually, one can point to a few different methods that allow the task to be fulfilled, but each time the question arises: which one to choose?

MATERIAL AND METHODS: The study was conducted using a database that included 3246 patients in the Second Department of Cardiology, Silesian Medical Centre in Katowice-Ochojec in 2003–2008. We were A model in which the STROKE dependent variable was considered was subjected to statistical analysis, and the results of the analysis suggested selecting the following variables: gender, transfusion, PTCA, IVA, IVM, SVA, aneurysm and hematocrit.

RESULTS: The essential factors affecting the occurrence of stroke, according to logistic regression are: aneurysm, transfusion of blood components, prior treatment with PTCA and according to the classification tree: aneurysm and level of hematocrit.

CONCLUSIONS: The results achieved by both the two statistical models complemented each other, and by combining them one is able to obtain reliable information to use as a base for the decision-making process.

KEY WORDS

stroke, logistic regression, statistical analysis, classification tree, ROC curve

STRESZCZENIE

WSTĘP: Zadaniem analiz statystycznych w badaniach naukowych jest wskazanie trafnych i maksymalnie wiarygodnych wniosków w sytuacji, gdy badacz dysponuje wieloma informacjami. Zwykle można wskazać kilka różnych metod, które pozwalają to zadanie spełnić, jednak za każdym razem nasuwa się pytanie, którą z nich wybrać?

Received: 21.02.2016

Revised: 24.03.2016

Accepted: 11.04.2016

Published online: 12.08.2016

Address for correspondence: Mgr Anna Spychała, Department of Environmental Health, School of Public Health in Bytom, Medical University of Silesia in Katowice, Poland, ul. Piekarska 18, 41-902 Bytom, tel. 32 397 65 44, e-mail: aspychala@sum.edu.pl

Copyright © Śląski Uniwersytet Medyczny w Katowicach
www.annales.sum.edu.pl

MATERIAŁ I METODY: Badania zostały przeprowadzone na bazie danych, która obejmowała 3246 pacjentów przebywających na II Oddziale Kardiologii Górnośląskiego Centrum Medycznego w Katowicach-Ochojcu w latach 2003–2008. Analizie statystycznej poddano model, w którym za zmienną zależną uznano zmienną UDAR, natomiast wyniki przeprowadzonych analiz zasugerowały dobór następujących zmiennych objaśniających: płeć, przetoczenie, PTCA, IVA, IVM, SVA, tętniak i hematokryt.

WYNIKI: Czynniki istotnie wpływającymi na wystąpienie udaru, według regresji logistycznej, są: tętniak, przetoczenie składników krwi i przebyty zabieg PTCA, natomiast według drzewa klasyfikacyjnego – tętniak i poziom hematokrytu.

WNIOSKI: Wyniki uzyskane za pomocą obydwu modeli statystycznych dopełniały się, a ich łączenie pozwala na uzyskanie wiarygodnych informacji, stanowiących podstawę procesu decyzyjnego.

SŁOWA KLUCZOWE

udar, regresja logistyczna, analiza statystyczna, drzewo klasyfikacyjne, krzywa ROC

INTRODUCTION

In the process of decision making, it is extremely important to record and evaluate information, then apply the relevant selection criterion, which ultimately results in the issuance of an appropriate decision [1,2]. In the case of medical diagnostics, determining a clinical diagnosis is associated with the so-called diagnostic process, which is increasingly supported by tools for assessing the correctness of action, based on models and mathematical techniques [3,4,5,6]. These tools include: logistic regression, the classification tree and the ROC curve. The application of logistic regression as a mathematical model with a dichotomous variable explained by independent variables is particularly useful for medical research, whose aim is to estimate the risk factors for disease or death [7]. The basis for the logistic regression model is the logistic function, which takes the values from 0 to 1, the shape of which resembles a stretched letter S, and this shows minimal changes to the function value (between 0 and 1) [8]. In contrast, one of the most important logistic regression model parameters is the odds ratio, defined as the ratio of the 'chance' of occurrence A to the 'chance' of occurrence B [9].

The classification tree is a mathematical model presented using a graph, which is generated by dividing the total number of observations [10,11,12]. The basic premise for creating a classification tree model is to obtain subsets that are the most homogeneous from the point of view of the characteristic [13,14,15]. This method owes its popularity to the possibility of presenting data using graphical tools, making it easier to interpret the results. Moreover, it has two important features that make it so frequently used, namely its hierarchical nature and flexibility [16].

On the other hand, the ROC curve (Receiver Operating Characteristic Curve) is used for comparing classification models, based on the relation of the association measures of decision-making: the specificity and sensitivity [17,18]. To interpret the ROC curve, we

should look at the surface area created under the curve, which is determinant for exact demarcation of normal and abnormal results. The area under the ROC curve, designated as area under curve (AUC), ranges from 0 to 1 and the higher or closer the value is to 1, the better the model is [19].

AIM

The aim of the study was to compare two methods of statistical analysis: logistic regression and the classification tree using the ROC curve for suitability in analyzing medical data.

MATERIAL AND METHODS

As part of the comparative analysis of the mathematical models, the database of clinical trials conducted in 2003–2008 at the Second Department of Cardiology, Silesian Medical Centre in Katowice-Ochojec, the health status of 5781 adults aged 16 and over undergoing hospitalization after qualifying for surgery were used. It included the following information: vital records (birth date, gender), anthropometric (weight, height, BMI), environmental (smoking) and health status (morphological parameters, diabetes, coronary heart disease, hypertension, cardiovascular disease, stroke, heart attack, aneurysm, cardiac parameters, transfusion of blood components). The analyses included only 3246 comprehensive records presented in the non-personalized form, ensuring patient anonymity.

The analyzed models were developed using logistic regression and the classification tree with an assumed independent variable, which takes into account division of the testing population as a group with a history of stroke and its lack of a preoperative interview. In carrying out the construction of logistic regression, we use a backward stepwise approach to preserve

significance of the parameters with a large number of independent variables. In the process of determining the classification tree, the C&RT (Classification and Regression Trees) algorithm was used. To compare the results of the multivariate models, ROC curves were used. Simple univariate analysis was based on the U Mann-Whitney test and the χ^2 Pearson test. Compatibility with normal distribution was analyzed based on the Shapiro-Wilk test results.

The material was developed statistically using Statistica 10.0 StatSoft Poland, with the addition of Medical Kit 2.0. Statistical significance was set at $p < 0.05$.

RESULTS

Selection of variables to model

In order to identify the factors affecting the occurrence of stroke in the patients we studied, a simple analysis was performed. Among the persons diagnosed with cerebrovascular disease (I60–I67 according to ICD-10), statistically significant disturbances in hematocrit levels were observed (Tab. I).

Moreover, equally important was the impact of the following variables: gender, blood component transfusions carried out, diagnosed cardiac aneurysm and factors such as percutaneous transluminal coronary angioplasty (PTCA), aortic stenosis (SVA), aortic regurgitation (IVA), mitral regurgitation (IVM) (Tab. II).

Table I. Level of morphological indicators in groups of hospitalized people
Tabela I. Poziom wskaźników morfologicznych w grupach osób hospitalizowanych

Morphological parameters	Total N = 3246 (100%)		STROKE 0	STROKE 1	Mann-Whitney p-value
	X ± SD	M; IQR	M; IQR	M; IQR	
Hemoglobin [g/dl]	13.80 ± 1.52	14; 1.80	14; 1.80	13.80; 2.5	0.07 NS
Hematocrit [%]	40.66 ± 4.42	41; 5.30	41; 5.30	40.15; 7	$p < 0.05$
Platelets [thous/mm ³]	198.86 ± 66.27	189; 68	189; 68	180; 83.50	0.10 NS
Glucose [mg/dl]	116.31 ± 43.71	103; 35	103; 35	106; 53.50	0.37 NS
Creatinine [mg/dl]	1.12 ± 0.24	1.02; 0.30	1.02; 0.30	1; 0.27	0.24 NS

Table II. Clinical characteristics of studied groups of hospitalized people
Tabela II. Charakterystyki kliniczne w badanych grupach osób hospitalizowanych

Clinical characteristics	Total N = 3246 (100%)	STROKE 0 N = 3154	STROKE 1 N = 92	χ^2 Pearson test p-value	
	1	2	3		4
Gender					
Woman	989 (30.5%)	948 (30.1%)	41 (44.6%)		$p < 0.0001$
Man	2257 (69.5%)	2206 (69.9%)	51 (55.4%)		
Age					
< 30	19 (0.6%)	19 (0.6%)	0 (0%)		0.29 NS
31–44	89 (2.7%)	89 (2.8%)	0 (0%)		
45–59	1081 (33.3%)	1055 (33.5%)	26 (28.3%)		
60–74	1711 (52.7%)	1657 (52.5%)	54 (58.7%)		
>=75	346 (10.7%)	334 (10.6%)	12 (13.0%)		
BMI					
underweight	24 (0.7%)	23 (0.7%)	1 (1.1%)		0.30 NS
standard	761 (23.4%)	736 (23.3%)	25 (27.2%)		
overweight	1590 (48.9%)	1541 (48.9%)	59 (53.3%)		
obese	871 (26.8%)	854 (27.1%)	17 (18.4%)		

	1	2	3	4	5
Smoking					0.80
No		2148 (66.2%)	2086 (66.1%)	62 (67.4%)	NS
Yes		1098 (33.8%)	1068 (33.9%)	30 (32.6%)	
Transfusion of blood components					p < 0.01
No		2361 (72.8%)	2329 (73.8%)	32 (34.8%)	
Yes		885 (27.2%)	825 (26.2%)	60 (65.2%)	
Coronary heart disease					0.58
No		490 (15.1%)	478 (15.2%)	12 (13.0%)	NS
Yes		2756 (84.9%)	2676 (84.8%)	80 (87.0%)	
Aneurysm					p < 0.01
No		3215 (99.0%)	3127 (99.0%)	88 (95.7%)	
Yes		31 (1.0%)	27 (1.0%)	4 (4.3%)	
Heart attack					0.82
No		1761 (54.3%)	1710 (54.2%)	51 (55.4%)	NS
Yes		1485 (45.7%)	1444 (45.8%)	41 (44.6%)	
Reoperation					0.35
No		3113 (96.0%)	3023 (96.0%)	90 (98.0%)	NS
Yes		133 (4.0%)	131 (4.0%)	2 (2.0%)	
Hypertension					0.80
No		810 (25.0%)	786 (24.9%)	24 (26.0%)	NS
Yes		2436 (75.0%)	2368 (75.1%)	68 (74.0%)	
Cardiovascular disease					0.83
No		2708 (83.4%)	2632 (83.5%)	76 (82.6%)	NS
Yes		538 (16.6%)	522 (16.5%)	16 (17.4%)	
Diabetes					0.35
No		2404 (74.0%)	2332 (74.0%)	72 (78.0%)	NS
Yes		842 (26.0%)	822 (26.0%)	20 (22.0%)	
PTCA					p < 0.01
No		3185 (98.0%)	3102 (98.0%)	83 (90.0%)	
Yes		61 (2.0%)	52 (2.0%)	9 (10.0%)	
SVA					p = 0.05
No		2914 (90.0%)	2837 (90.0%)	77 (84.0%)	
Yes		332 (10.0%)	317 (10.0%)	15 (16.0%)	
IVA					p < 0.05
No		2873 (88.5%)	2800 (88.8%)	73 (80.0%)	
Yes		373 (11.5%)	354 (11.2%)	19 (20.0%)	
SVM					0.07
No		3116 (96.0%)	3031 (96.0%)	85 (92.4%)	NS
Yes		130 (4.0%)	123 (4.0%)	7 (7.6%)	
IVM					p < 0.01
No		2568 (80.0%)	2508 (80.0%)	60 (65.0%)	
Yes		678 (20.0%)	646 (20.0%)	32 (35.0%)	
IVT					0.28
No		3050 (94.0%)	2966 (94.0%)	84 (91.0%)	NS
Yes		196 (6.0%)	188 (6.0%)	8 (9.0%)	

Table III. Results of estimating best logit model for probability of having stroke
Tabela III. Wyniki estymacji najlepszego modelu logitowego dla prawdopodobieństwa wystąpienia udaru

Logit model $\chi^2 = 84.75$; p-value < 0.0001 <i>Bayesian inference</i>	Parameter estimation	p-value	Odds ratio	95% CI for odds ratio	
Free term	-4.49	< 0.0001	0.01	0.008	0.016
Transfusion of blood components	1.63	< 0.0001	5.0	3.2	7.8
PTCA	1.73	< 0.001	4.9	2.2	10.6
IVA	0.71	0.01	2.0	1.2	3.1
Aneurysm	1.44	0.014	5.9	1.9	18.3

Logistic regression

Table III shows the results of estimating the best logit model for the probability of having a stroke. The risk of stroke was:

- nearly 5 times higher in patients who received blood component transfusions and on whom a percutaneous transluminal coronary angioplasty (PTCA) procedure was performed
- 2 times higher in patients with aortic regurgitation
- 5.9 times higher in patients diagnosed with an aneurysm.

Other variables (hematocrit, gender, aortic stenosis-SVA, mitral regurgitation-IVM) have been removed from the model of the reasons for the lack of statistical significance for the dependent variable model STROKE.

Classification tree

In the process of determining the classification tree, the C&RT algorithm was used. It enabled the appointment of 8 divisions giving 9 end nodes (Fig. 2). In the first division, transfusion plays a significant role (no/others) – Node 1. The next division was made based on the hematocrit level (Hct) – Node 2, 3, 13 and 14. In further divisions, the variable indicator of coronary angioplasty of cardiac vessels – PTCA – (yes/other) began to play a significant role as well as the occurrence of aortic stenosis in patients (SVA).

The meaning of the aneurysm variable (yes/other) in assessing the risk of stroke is shown by node 16, while the gender variable (female/male) was rejected due to lack of significance.

The importance of ranking the predictors – independent variables, is presented in Figure 2. The ranking scale estimates the values from 0 to 1, where values closer to 1 represent the greatest influence of each variable on a subsidiary variable. In this case, the risk of having a stroke for patients is related to aneurysm diagnosis and the hematocrit level (Hct).

ROC curve

To find a reference point when comparing logistic regression and a classification tree, an ROC curve is generated. The top-matching classifier in terms of a parameter space of the ROC curve (AUC) is the transfusion variable (0.695), which is also characterized by relatively high sensitivity and specificity in comparison with other variables. The variable with the lowest values for AUC (only 0.427) is gender, which still has high sensitivity in comparison with other classifiers. The other variables have very similar values of the parameter characterizing the area under the ROC curve (AUC) and the specificity of the test, but the variable relating to patient diagnosis, aneurysm has the highest specificity (0.991) (Tab. IV). Figure 3 is the graphic confirmation of these observations.

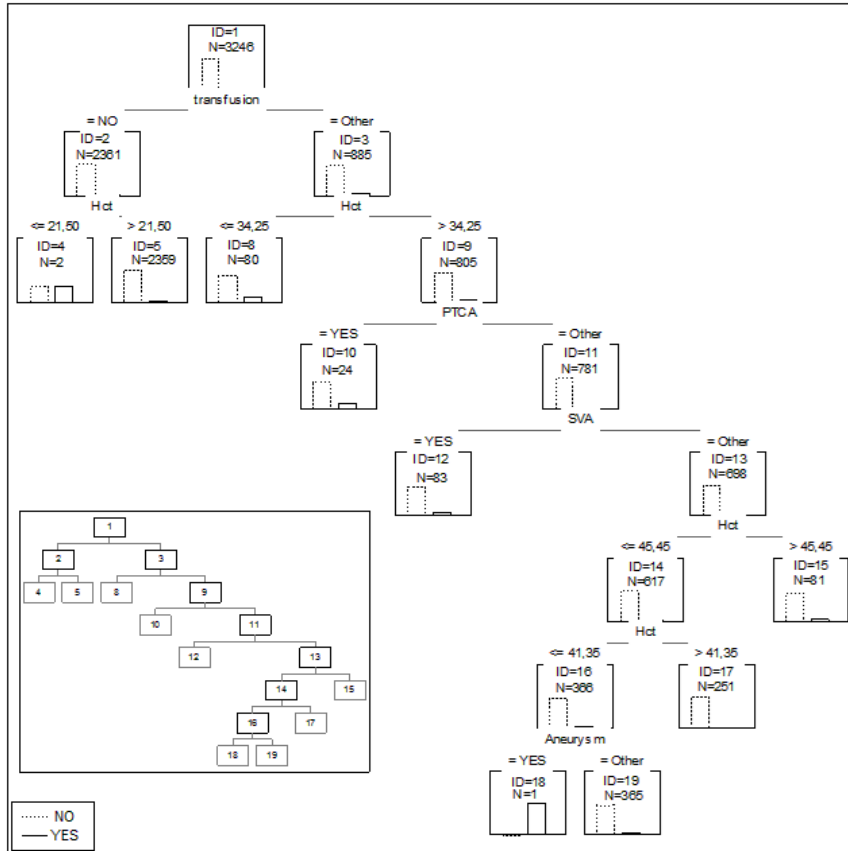


Fig. 1. Decision tree for STROKE dependent variable model.
Ryc. 1. Drzewo decyzyjne dla modelu zmiennej zależnej UDAR.

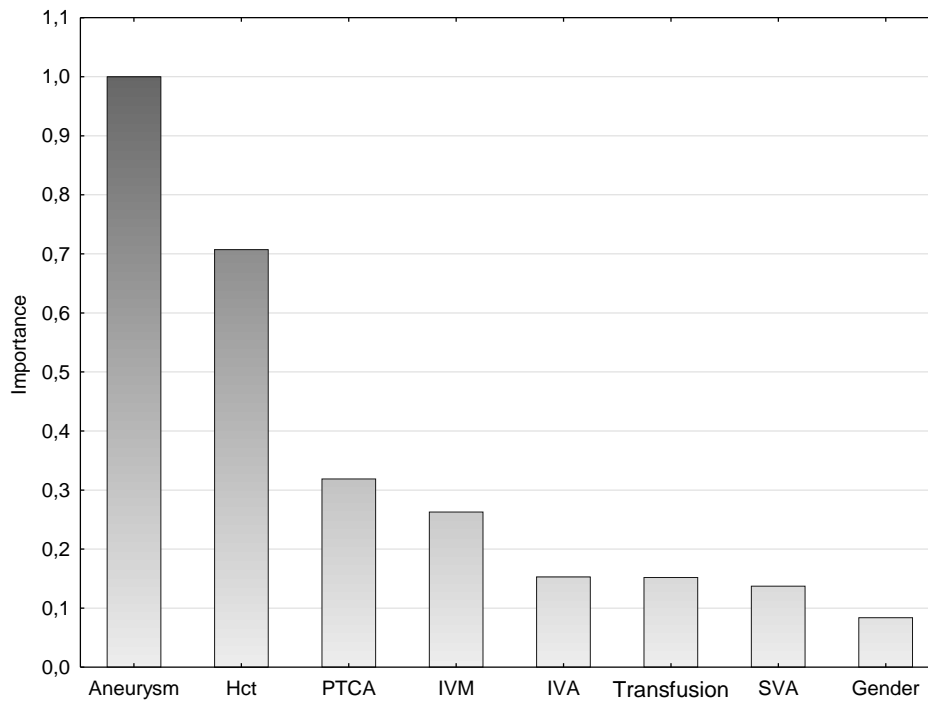


Fig. 2. Importance of predictors for STROKE dependent variable model.
Ryc. 2. Ważność predyktorów dla modelu zmiennej zależnej UDAR.

Table IV. ROC curve indicators for STROKE dependent variable model
Tabela IV. Wskaźniki krzywej ROC dla modelu zmiennej zależnej UDAR

STROKE					
Explanatory variables	Sensitivity	Specificity	Positive predictive value	Negative predictive value	Area under ROC curve AUC
Transfusion of blood components	0.652	0.738	0.068	0.986	0.695
IVM	0.348	0.795	0.047	0.977	0.572
IVA	0.207	0.888	0.051	0.975	0.547
PTCA	0.098	0.984	0.148	0.974	0.541
SVA	0.163	0.899	0.045	0.974	0.531
Aneurysm	0.043	0.991	0.129	0.973	0.517
Hematocrit	0.269	0.713	0.098	0.911	0.430
Gender	0.554	0.301	0.026	0.959	0.427

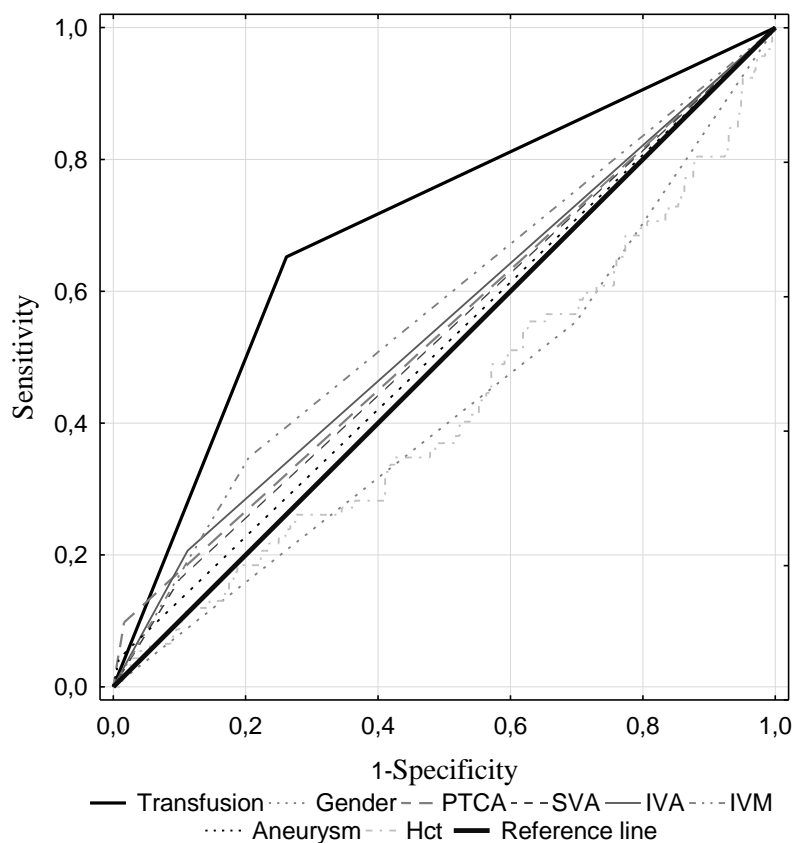


Fig. 3. ROC curve graph for STROKE dependent variable model.
Ryc. 3. Wykres krzywej ROC dla modelu zmiennej zależnej UDAR.

DISCUSSION

Statistics as a science allows for collecting, summarizing, presenting and interpreting data, and thus verifying hypotheses. The priorities in medical statistics include: hospital statistics, epidemiology and evaluation of measurement accuracy and quality of diagnostic tests [1,5,20]. In the modern world, a biostatistician should be a partner for the doctor and provide him with reliable statistical analyzes, testing hypotheses and interpretation of results because improving the quality of health care can be indirectly linked to progress in the field of statistical analysis [21].

The presented work was designed to assess the two methods of data analysis: logistic regression and the classification tree for their usefulness in medical research. The selection of appropriate statistical methods for the analysis of medical data can affect the planning and implementation processes of individual patient therapy [22]. The classification tree can be considered a method that does not require the researcher to have broad statistical knowledge to carry it out and interpret it because the results are presented using a decision tree graph and the conclusions that the researcher can draw are visible at first glance. In addition, the classification tree used as a method of data analysis due to its simpler form, as compared to logistic regression, requires less time to carry it out. An example of research which uses decision trees described in the literature was to predict adverse cardiac events in patients with chronic heart failure. During this study the data on patients with chronic heart failure diagnosed at least 6 months before the beginning of the study was analyzed and the main aim of this study was to evaluate the usefulness of different classification tree models in predicting adverse cardiac events and death among patients [23]. An important advantage of classification trees is also the possibility of introducing a wide range of predictors, especially when they are not familiar with the appropriate criteria for their selection [24].

Logistic regression analysis in medical sciences has grown to the status of the ideal method. However, it is challenging for both the researcher with specialized analytical skills, and also for the reader, who

must have adequate knowledge of statistics to properly analyze the presented results. Nevertheless, it is recognized as a more precise and more certain method due to the large number of objectives for which the variables must be checked for the analysis to be performed [2]. Logistic regression was used in the analysis of the results of studies describing the use of three-dimensional ultrasound in the diagnosis of thyroid nodules among 65 patients. This method was used to assess the independent risk factors for thyroid cancer [25].

An example of studies that have used both methods of data analysis discussed in this paper is a study on the identification of insulin resistance among young people living in Asia. Scientists tried to prove that the simultaneous use of classification trees and logistic regression as methods that complement each other, will develop technology that is simple, not too expensive, and at the same time gives very satisfactory results in terms of its accuracy. The results of this study showed that both methods can be used to predict insulin resistance and may be useful in the development and implementation of prevention programs, which will cover the youth population living in Asia [26].

CONCLUSIONS

Comparing the model determined using logistic regression and the classification tree graph, a group of factors can be identified which significantly affect the occurrence of cerebrovascular disease. They are: previous transfusion of blood components, prior treatment with PTCA and the presence of an aneurysm. It should be noted that several variables that logistic regression considered statistically insignificant for the model and consequently removed, the classification tree treated those variables as its major hubs.

Analysis of the results obtained using ROC curves confirms the appropriate choice of independent variables in the case of the methods used, however, at the same time the analysis highlights the complementarity between them. Combining this type of analysis allows one to obtain reliable information which acts as a basis for the decision-making process.

Author's contribution

Study design – A. Spychała, M. Skrzypek

Data collection – M. Skrzypek

Data interpretation – A. Spychała, M. Skrzypek, E. Niewiadomska

Statistical analysis – A. Spychała

Manuscript preparation – A. Spychała, M. Skrzypek, E. Niewiadomska

Literature research – A. Spychała, E. Niewiadomska

REFERENCES

1. Koźmiński A.K., Piotrowski W. Zarządzanie. Teoria i praktyka. Wydawnictwo Naukowe PWN. Warszawa 2002.
2. Szydło R. Komu jest potrzebny statystyk medyczny? *Onkol. Prakt. Klin.* 2005; 1(3): 129–131.
3. Michalski T. *Statystyka*. Wydawnictwo Szkolne i Pedagogiczne. Warszawa 2004, s. 5–7.
4. Sobczyk M. *Statystyka*. Wydawnictwo Uniwersytetu Marii Curie-Skłodowskiej. Lublin 2001, s. 7–9.
5. Dobbson A.J. The role of the statistician. *Int. J. Epidemiol.* 1983; 12(3): 274–275.
6. Lenz R., Reuchert M. IT suport for healthcare processes-premises, challenges, perspectives. *Data & Knowledge Engineering* 2007; 61: 39–58.
7. Stanisław A. Przystępny kurs statystyki z zastosowaniem STATISTICA PL na przykładach z medycyny. Tom II. StatSoft Polska. Kraków 2006, s. 217–254.
8. Stanisław A. *Biostatystyka*. Wydawnictwo Uniwersytetu Jagiellońskiego. Kraków 2005, s. 325–335.
9. Stanisław A. Regresja logistyczna. *Medycyna Praktyczna* 200110, http://www.mp.plartykulyindex.php?aid=11813&_tc=32F4510F44D54B93654578ACE3F57E99 [dostęp z dnia: 02.01.2013].
10. Łapczyński M. Drzewa klasyfikacyjne i regresyjne w badaniach marketingowych. Uniwersytet Ekonomiczny. Kraków 2011.
11. Stanisław A. Przystępny kurs statystyki z zastosowaniem STATISTICA PL na przykładach z medycyny. Tom III. StatSoft Polska. Kraków 2006, s. 113–153.
12. Łapczyński M. Drzewa klasyfikacyjne w badaniach satysfakcji i lojalności klientów. StatSoft Polska. Kraków 2003.
13. Gatnar E., Walesiak M. *Statystyczna analiza danych z wykorzystaniem programu R*. Wydawnictwo Naukowe PWN. Warszawa 2009, s. 238–260.
14. Stanisław A. *Statistica w badaniach naukowych i nauczaniu statystyki*. StatSoft Polska. Kraków 2010, s. 61–80.
15. StatSoft Electronic Statistic Textbook. Drzewa klasyfikacyjne, http://www.statsoft.pl/textbookstathome_stat.html?http%3A%2F%2Fwww.statsoft.pl%2Ftextbook%2Fstclatre.html [dostęp z dnia: 02.01.2013].
16. Koronacki J., Ćwik J. *Statystyczne systemy uczące się*. 2 wyd. EXIT. Warszawa 2008, s. 129–164.
17. Harańczyk G. Krzywa ROC, czyli ocena jakości klasyfikatora i poszukiwanie optymalnego punktu odcięcia. StatSoft Polska. Kraków 2010, s. 79–89.
18. Swets J.A., Dawes S.M., Monahan J. Better decisions through science. *Sci. Am.* 2000; 283(4): 82–87.
19. Stanisław A. Przystępny kurs statystyki z zastosowaniem STATISTICA PL na przykładach z medycyny. Tom I. StatSoft Polska. Kraków 2006, s. 340–370.
20. Harańczyk G., Stępień M. Ilustrowana sztuka podejmowania decyzji. *Matematyka Społeczeństwo Nauczanie* 2008; 41: 12–15.
21. Lemon S.C., Roy J., Clark M.A., Friedmann P.D., Rakowski W. Classification and regression tree analysis in public health: methodological review and comparison with logistic regression. *Ann. Behav. Med.* 2003; 26(3): 172–181.
22. Long W.J., Griffith J.L., Selker H.P., D'Agostino R.B. A comparison of logistic regression to decision-tree induction in a medical domain. *Comput. Biomed. Res.* 1993; 26: 74–97.
23. Owczarek A. Drzewa decyzyjne oraz analiza hazardu proporcjonalnego Coxa w przewidywaniu niekorzystnych zdarzeń sercowych u chorych z przewlekłą niewydolnością serca. Rozprawa habilitacyjna nr 12/2011. Śląski Uniwersytet Medyczny w Katowicach. Katowice 2011.
24. Classification and Regression Trees. <http://documents.software.dell.com/StatisticsTextbookClassification-and-Regression-Trees> [dostęp z dnia: 15.06.2015].
25. Słapa R.Z. *Zastosowanie ultrasonografii trójwymiarowej w diagnostyce zmian ogniskowych tarczycy*. Akademia Medyczna w Warszawie. Warszawa 2007.
26. Goel R., Misra A., Kondal D., Pandey R.M., Vikram N.K., Wasir J.S., Dhingra V., Luthra K. Identification of insulin resistance in Asian Indian adolescents: classification and regression tree (CART) and logistic regression based classification rules. *Clin. Endocrinol. (Oxf)* 2009; 70(5): 717–724.