

Original Article

Application of Survival Tree Model in Determining Affecting Factors in Breastfeeding Duration

Ameneh Sadat Sheykholeslami¹ Nasser Behnampour² Reza Ali Mohammadpour^{3*} Fatemeh Abdollahi⁴

1. Department of Biostatistics, Faculty of Health, Mazandaran University of Medical Sciences, Sari, Iran
2. Assistant professor, Department of Biostatistics, Faculty of Health, Golestan University of Medical Sciences, Gorgan, Iran
3. Associate professor, Department of Biostatistics, Faculty of Health, Mazandaran University of Medical Sciences, Sari, Iran
4. Associate professor, Department of Public Health, Faculty of Health, Health Sciences Research Center, Addiction Institute, Mazandaran University of Medical Sciences, Sari, Iran

*Correspondence to: Reza Ali Mohammadpour
Mohammadpour2002@yahoo.com

(Received: 4 Jan. 2021; Revised: 11 Apr. 2021; Accepted: 30 Apr. 2021)

Abstract

Background and Purpose: Survival tree model is a nonparametric method which can be used to identify the affecting factors from a specific time to the onset of an event. In this method, the categories are selected according to the most important factors. The purpose of this study was to determine the factors affecting the duration of breastfeeding in mothers and introduce the homogeneous subgroups using a survival tree model.

Materials and Methods: It was a historical cohort study analyzing the survival data of mothers with healthy single childbirths referring to the rural and urban health centers of Agh-Ghala County since 2011 until 2014. Data analyses and groupings of breastfeeding survival were performed using survival tree model with conditional inference algorithm in R Software. A separation criterion (SEP) confirmed the relevance of the model.

Results: Survival tree model results revealed that the type of consumed milk with the complementary nutrition, ethnicity and the time interval between current childbirth and the previous delivery were the most important factors affecting the duration of breastfeeding. The SEP's criterion was 2.082. Thus, due to the significant difference between the subgroups and the value of more than 1 for SEP criterion, the efficiency of the model was confirmed.

Conclusion: Survival tree model could be introduced as a suitable and powerful method for ranking the duration of breastfeeding rate which presents four homogeneous subgroups for analysis in addition to identifying the predictive variables.

Keywords: Survival Tree; Conditional Inference Algorithm; Homogeneous Subgroups; Breastfeeding

Citation: Sheykholeslami AS, Behnampour N, Mohammadpour RA*, Abdollahi F. Application of Survival Tree Model in Determining Affecting Factors in Breastfeeding Duration. Iran J Health Sci. 2021; 9(2): 9-17.

1. Introduction

Regression models are commonly used in conventional statistical methods in which the response variables are quantitative and affected by numerous factors (1). Survival models could be performed when the response variable is time, such as the duration of breastfeeding (BF), and the data are censored (2). In the last few years, in addition to classic survival models such as Cox proportional hazards analysis survival tree model has been introduced and extended for censored data (3-5).

Tree-based models as nonparametric regression methods are of the most flexible, intuitive and powerful tools in data analysis to discover the complex data structure. These methods provide suitable solutions to provide diagnostic tests in medicine and biomedical sciences for detecting the exact causes of diseases (6-8). Survival tree model is a nonparametric method which can be used to identify the affecting factors from a specific time to the onset of an event, such as BF duration (9).

The most suitable food for the infant is breast milk. The world health organization (WHO) recommendation since 2003 is exclusively BF for the first 6 months and BF up to 2 years of life for babies (10, 11). For the first 6 months of life, breast milk alone is the ideal nourishment for infants (10, 12-13). The findings of various literature reviews showed that BF duration depends on a large number of factors. Maternal age, age of infants, multiparity, education level, religion, and occupation are associated with BF duration (10, 14-15). Determinants of BF duration and complementary feeding practices of Iranian mothers were studied by using the Cox proportional hazard model in the previous published article (16).

In various studies about BF, researchers employed classical Cox proportional hazards model for the BF but these models, in addition to the need to establish the necessary preconditions, are not able to determine a proper ranking of the studied subjects. The

acceptance of methods based on tree models in medical sciences is due to the need of clinical researchers to define categorical criteria for diagnosis and prediction of disease (7). In fact, the basic idea of these models is the formation of subgroups of individuals, in which these subgroups are homogeneous according to the desired event (17). Various recursive partitioning methods were introduced in the literature (18-19).

In the present study, we used the conditional inference tree algorithm (18) for subgroup identification in mothers to investigate the factors affecting the BF duration. According to the huge number of effective variables and their diversity, the survival tree method introduces a classification model for existing observations with a simple and understandable structure for decision making, and this model may facilitate the research process by identifying the affecting factors and homogeneous subgroups to provide solutions in removing the barriers in this field. Therefore, the findings of this study can play a significant role in promoting BF in the area of Iran with socio-economic, demographic and multi-cultural variables.

2. Materials and Methods

The data were collected in the range of 2011 to 2014 from a historical cohort study on 501 registered healthy and single birth deliveries since 20th march 2011 until 20th September 2012 who were follow at least for 24 complete months. The study area was located in rural and/or urban districts of Agh-Ghala County, Golestan Province, north of Iran with various ethnicities and cultures. The BF duration was determined from the time of birth until the cases were 24 months (optimum BF duration). The BF was discontinued or not (right censorship) at this time.

Infants older than two years of age who were still BF were considered as right censorship. A checklist was used to collect the data of all children on the basis of available information in the infant's file and the mother's maternity booklet as registered in health centers. Other

information was collected by phone call or instant interview. The checklist was comprised of 68 questions including demographic data of mother, father and child, socio-economic factors, obstetric factors and delivery conditions, maternal and/or childhood health complications, smoking, supplementary nutrition, and BF information.

We used the conditional inference algorithm for survival tree model and a measure of separation criterion which confirmed the relevance of the model (20).

The conditional inference tree algorithm was used to analyze censored data. The method of this algorithm was based on the partition method to increase the differences between nodes. R Software and *rpart* or *party* packages which are freely available at website <https://cran.r-project.org/> were also used to implement these methods (21-22).

Statistical analysis

The conditional distribution $D(Y|X)$ of the response variable Y depends on the function f of the covariates variables on the condition of the covariates X .

$$D(Y|X) = D(Y|X_1, \dots, X_m) = D(Y|f(X_1, \dots, X_m))$$

A regression relation model is fitted based on the n training sample, i.e. a random sample and co-distribution of observations

$$\mathcal{L}_n = \{(Y_i, X_{1i}, \dots, X_{mi}); i = 1, 2, \dots, n\}$$

In this type of tree structure, the generic algorithm for recursive binary partitioning can be formulated using integer non-negative values of observation weights $W = (w_1, \dots, w_n)$ (19).

The following algorithm performs recursive partitioning based on conditional inference;

1. Test the hypothesis of overall independence between each covariate and the response. If this assumption is not rejected, the separation will stop. Otherwise, the variable X_j^* is selected with the strongest relation to Y .

2. Choose a set A^* in covariate space to divide into two disjoint sets with left and right weights for $i = 1, \dots, n$.

$$w_{\text{left},i} = w_i I(X_{j^*i} \in A^*)$$

$$w_{\text{right},i} = w_i I(X_{j^*i} \notin A^*)$$

3. Recursively repeat step 1 and 2 with modified case weights.

Variable selection and stopping criteria

The first step is the issue of independence. At this stage, it must be decided whether the covariate variable contains information about the response variable or not. For this test, the relationship between Y and X_j where $j = 1, \dots, m$ is measured by linear statistics in the following form which is known as the permutation test method (19, 21-22).

$$T_j(\mathcal{L}_n, w) = \text{vec} \left(\sum_{i=1}^n w_i g_j(X_{ij}) h(Y_i, (Y_1, \dots, Y_n))^T \right) \in \mathbb{R}^{p \times q}$$

This function is a rank logarithm function for the survival tree for survival data, which includes censored data (23).

Splitting criteria

To find the optimal separation on the selected variable X_j^* , the goodness of the separation is evaluated by linear statistics of two samples for all possible subsets A of the sample space X_j^* .

$$T_{j^*}^A(\mathcal{L}_n, w) = \text{vec} \left(\sum_{i=1}^n w_i I(X_{j^*i} \in A) h(Y_i, (Y_1, \dots, Y_n))^T \right) \in \mathbb{R}^q$$

The important point in this statistical method is to ensure that a tree of the right size is built, so there is no need for the pruning and validation process of the model or any other method to select the right size of the tree. Selecting the input variable for p-value-based separation is a process to prevent bias in selecting relative to the input variables with different number of possible cut-off points (21, 22). Further details in CTree algorithm evaluation were described in Schlosser et al. (23).

The study was approved by the committee of ethics, Mazandaran University of Medical Sciences, Mazandaran, Iran (Code of Ethics: IR. MUZUMS.1394.1456).

3. Results

The duration of BF in 501 children was analyzed using the stated model. Of all the cases, 253 newborns (50.6%) were male and 248 were female (49.4%). A total of 425 (84.8%) cases were from Sunnites mothers and 76 (15.2%) cases were Shiites. 363 mothers (72.6%) were the residents of rural areas, while 138 mothers (27.4%) were the residents of city regions. The mothers' ages ranged between 15

to 42 years. The average age of auxiliary nutrition initiation was 6 months in 389 (78.1%) children. The median duration of BF was 22 months with a standard error of 0.388, while the average of BF duration was 20.44 months with a standard error of 0.206. The minimum duration of BF was reported to be one day. A univariate tree model was initially implemented for all covariate variables to find the variables which were related to the response variable of BF duration. Significant variables in a univariate tree model using log-rank test are demonstrated in Table 1.

Table 1. Significant variables in a univariate tree model by Log-Rank test

Variable	Log-rank statistics	p-value
Mother's age at marriage	3.92	0.021
Mother's age at child birth	3.85	0.015
Current weight of child	5.2	0.004
Current age of child	8.49	0.049
Child birth rank	12.72	0.041
Mother's religion	5.6	<0.001
Father's religion	5.8	<0.001
Ethnicity	5.1	0.002
The period between last delivery and previous one	7.91	0.009
Type of consumed milk with supplementary nutrition	127.39	<0.001

The purpose of univariate analysis of the survival tree was to examine the independence assumption for the desired variables and the response variable. Cutting points may be different from multivariate analyses due to the existence of different variables and the interaction between them.

The most important variable to be selected for data partitioning was the type of milk used with supplementary nutrition to fit the model in the first step. This variable divided individuals into

two subgroups according to the smallest possible p-value of all possible cutting points, in which 96.6% of the subjects were graded in the left node and 3.4% in the right. Ethnicity, with the strongest relation with the response variable, entered into the model. On the basis of the log-rank statistics, the most convenient distribution for this variable was Sistani ethnicity.

This phase of partitioning demonstrated that the most important factor which can affect the prognosis of BF advancement or discontinuing was ethnicity among the infants whose type of milk was not supplemented with exclusive nutritional supplementations. The cases with ethnicities other than Sistani, with a p-value = 0.005 derived from ordinal logarithm statistics, were categorized in the left node. The interval between last delivery and previous one was divided into two categories of more than 60 months and less than or equal to 60 months based on the distance variables. At this stage, due to the rejection of independence assumption, no other variable for the partition of these two subgroups was included in the model, and the two final nodes were obtained. Therefore, the variable of the interval between last delivery and previous one was introduced into the model as the most important factor in the discontinuation of BF after the ethnicity. It resulted in two different categories with

p-value = 0.005 derived from the logarithm of the rank.

The individuals with Sistani ethnicity were categorized in the right node. Because there was no rejection of the independence assumption, there were no other input variables for this subgroup, and the other final node was obtained. At the end of the node 7, no fraction which may have produced homogeneous subgroups was made for the infants who had been fed merely with milk powder in addition to the supplementary nutrition. Accordingly, a model was created which constituted of four distinguished subgroups of BF duration and prognostic factors, and four subgroups of homogeneity which were obtained. Figure 1 demonstrates the model of survival tree and subgroups of patients with the most important variables as factors affecting BF duration.

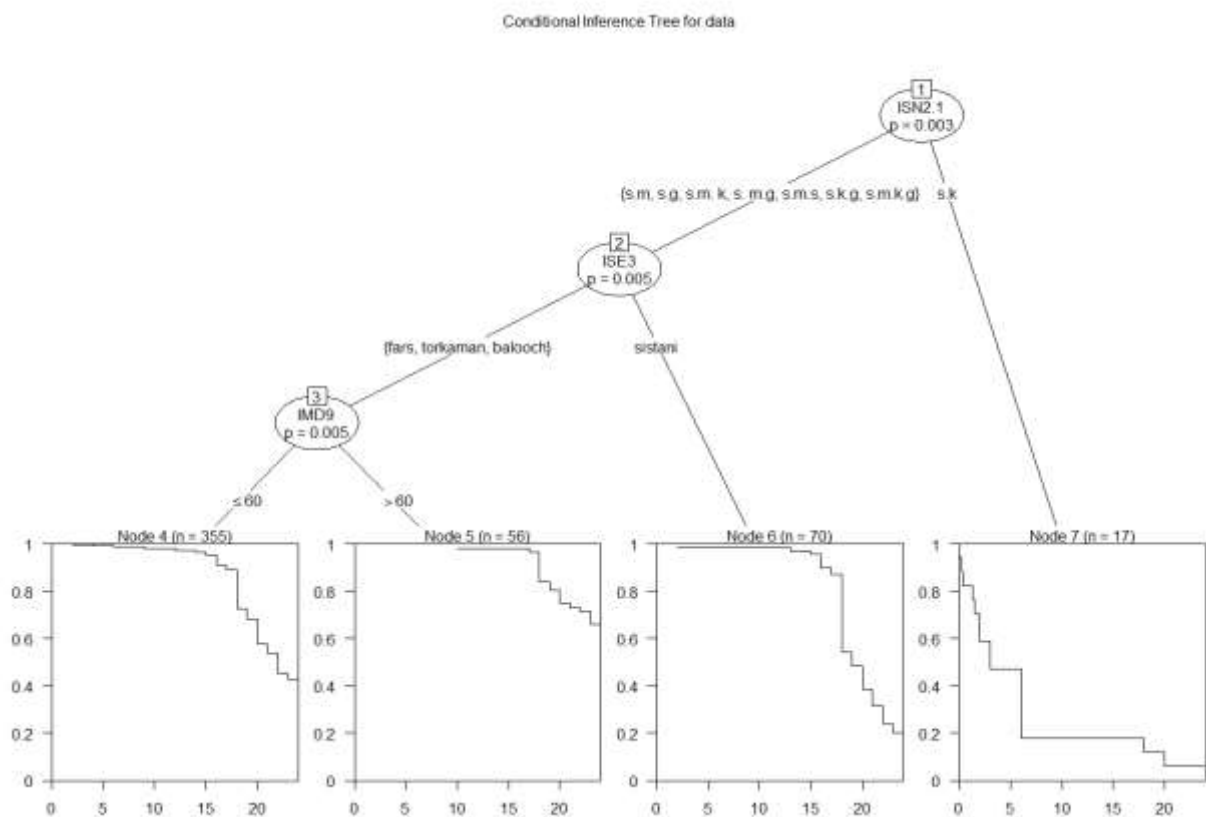


Figure 1. Multivariate survival tree, Kaplan-Meier curve inside each terminal node for factors associated with breastfeeding in the area of north of Iran (ISN2.1: Type of consumed milk with supplementary nutrition; ISE3: Ethnicity; IMD9: The period between last delivery and previous one)

Implementing the survival tree model, in addition to determining the factors influencing the BF duration can also introduce homogeneous subgroups. In the present study, the implementation of this model resulted in the creation of four subgroups: the first subgroup included infants who did not exclusively use milk powder in addition to the nutritional supplementary and had Baluch, Turkmen or Fars ethnicity, and the interval between last delivery and previous one was 60 months or less. This group had the highest average survival rate in comparison to other subgroups. The second subgroup included the infants who did not use milk powder in addition to the supplementary nutrition and

had non-Sistani ethnicity, and the interval between last delivery and previous one was more than 60 months. This subgroup had the highest average survival rate. The third subgroup included infants who did not use milk powder in addition to the supplementary nutrition and had Sistani ethnicity. The fourth subgroup included infants who were fed merely with milk powder in addition to the supplementary nutrition which had the lowest survival rate compared to other subgroups. The KM survival functions of these subgroups are presented in Figure 1. The median and mean survival rates in the subgroups derived from the survival tree model are also reported in Table 2.

Table 2. The median and mean survival rates in the subgroups derived from the survival tree model

Group	Number	Mean	SE*	Median	SE	log-rank test (P-value)	Tarone-Waire test (P-value)
1	355	20.65	0.23	21.0	0.4	117.55	136.66
2	56	22.28	0.39	-	-	(p<0.001)	(p<0.001)
3	70	19.74	0.42	19.0	0.75		
4	17	6.2	1.72	3.0	1.18		

*SE: Standard Error

Tarone-Waire (24) and log-rank statistics were used to address the difference between these subgroups. Assuming no difference in survival function in four at risk sub-groups as the zero assumption (p-value <0.001) indicated a significant difference in survival between subgroups. The SEP's criterion was found to be 2.082. Thus, due to the significant difference between the subgroups and the value of more than 1 for SEP criterion, the efficiency of the model was confirmed. Therefore, the obtained tree model in the previous section reached the main goal of this study which was the formation of homogeneous subgroups with maximum difference.

4. Discussion

In the present study, the type of consumed milk in addition to the supplementary nutrition, ethnicity and the interval between last delivery and previous one were accounted as the affecting factors on early discontinuation of BF using a survival tree model with a conditional inference algorithm on the basis of increasing the difference between the nodes. Tree models are implemented to create subgroups with the greatest possible differences. These subgroups could be used for subsequent clinical analyzes only if there are significant differences. The findings of this model by survival tree analysis revealed

that the type of consumed milk in addition to the supplementary nutrition, ethnicity, the interval between last delivery and previous one were the most important factors affecting the duration of BF. This finding was in accordance with some other studies about BF in other populations which were performed through Cox Model (14-16).

Various studies have been conducted using survival tree model to investigate the predictors of different diseases (9). However, to the best of our knowledge, no record was found on BF factors in the literature. Therefore, the studies which have used survival tree models for other applications are discussed in this section. In the majority of these studies, parametric and semi-parametric classical models were used. These models were unable to determine a proper ranking of the studied subjects in addition to the necessity of establishing pre-assumptions (7). The acceptance of tree-based models in the field of medical sciences is due to the need of clinical researchers to define the categorization criteria for the diagnosis and prediction of the disease. The principal idea of these models is the formation of subgroups of individuals which should be homogeneous according to the desired event (17).

Saki et al. (25) used tree models with a conditional inference algorithm on the basis of increasing differences between nodes to determine at risk homogeneous groups in patients with colorectal cancer and the important prognostic factors for predicting survival. Four important factors including the tumor stage at the time of cancer diagnosis, the age of the diagnosis, the morphology type, and the degree of tumor were determined as the major factors of disease prognosis. A survival tree model

was also used by Parizadeh et al. (26) to identify risk factors for ischemic stroke as a complementary method to the conventional Cox-analysis. DBP was shown to be the most important predictor of ischemic stroke in middle-aged and old subjects upon identification of risk patterns in ischemic stroke and exploration of the interactions between risk factors. Another study was also designed by Ramezankhani et al. (27) regarding the application of survival tree analysis in exploring potential interactions between predictors of chronic kidney disease.

LeBlanc and Crowley (5) designed a study to evaluate the survival in Leukemic patients on the basis of survival tree-centered algorithm and minimum internode deviances. As a result, four factors of serum albumin, creatinine, calcium, and age were determined as important prognostic factors. Valera et al. (28) conducted a study for patients with colorectal cancer using tree models on the basis of the maximum differences between nodes. Metastasis variables, involved lymph nodes, tumor size, and ki-67 marker, which were identified as prognostic factors. Schumacher et al. (17) presented a diagnostic model using breast cancer data from the German Breast Cancer Study Group since 1984 to 1989 in which the variables of the number of involved lymph nodes and the progesterone receptor were set as prognosis factors. The algorithm used for data partitioning was similar to our study and based on increasing the differences between groups using log-rank statistics. Hothorn et al. (19) reanalyzed the data from Schumacher et al. (17) using a conditional inference tree model. In addition to previously mentioned variables, hormone therapy was also mentioned, and patients were classified into four

homogeneous subgroups (20). In fact, this model showed the relative efficiency of this algorithm in comparison to the other methods. An advantage of this algorithm is the ability to obtain a tree with a proper size without editing and using methods for estimating cross-validation. It should be noted that other tree-structured algorithms use an editing process to determine the tree with an appropriate size after constructing a large tree with the highest number of nodes. In another study, Zhou and McArdle (8) showed that the conditional inference survival tree by Hothorn et al. (19) seems to be reliable among the survival tree algorithms. However, it is agreed that a single survival tree has some restrictions in situations with complex data. According to the literature review, it seems that the conditional inference algorithm using log-rank test is a suitable model for subgroup identification and variable prediction in BF (29).

Based on the findings of the present study, through using the survival tree analysis, we concluded that the type of powder milk used with supplementary nutrition, ethnicity, and the last delivery interval were the most important factors affecting the duration of BF.

In conclusion, survival tree-based models were found not to need any kind of presumptions to establish a model. Therefore, they were documented not to have the limitations of parametric and semi-parametric models. In this study, the survival tree model provided homogeneous subgroups of newborns in addition to introducing important variables in the prognosis of BF duration. Analysis of these subgroups is very important for clinical researchers and healthcare management. However, one of the limitations of survival tree-based models is the poor performance

of the model if there is coherence between variables. In this case, the best variable and cut-off point for the partition will not be well-defined, and dichotomizing the score values may decrease the power of CTree algorithm (23).

Acknowledgements

This article was derived from a Master's thesis submitted by Ameneh Sadat Sheykholeslami in biostatistics with the project code of 1456. The study was approved by the Research Council of the Deputy of Research and Technology of Mazandaran University of Medical Sciences, Mazandaran, Iran. Hereby, we extend our gratitude to all personnel of healthcare centers of Agh-Ghala for their cooperation.

Conflicts of Interest

The authors confirm that there are no known conflicts of interest associated with this publication.

Reference

1. Neter J, Kutner MH, Nachtsheim CJ, Wasserman W. Applied linear statistical models. Vol 4: Irwin Chicago; 1996.
2. Kleinbaum DG, Klein M. Survival analysis. Vol 3: Springer; 2010.
3. Breiman L, Friedman J, Stone CJ, Olshen RA. Classification and regression trees. CRC press; 1984.
4. Segal MR. Regression trees for censored data. *Biometrics* 1988; 44(1):35–47.
5. LeBlanc M, Crowley J. Relative risk trees for censored survival data. *Biometrics*. 1992;411-25.
6. Banerjee M, Noone AM. Advances in the Biomedical Sciences. New Jersey: John Wiley& Sons; 2008.
7. Noon AM, Banerjee M. Computational Methods in Biomedical Research. 2008:77-101.
8. Zhou Y, McArdle JJ. Rationale and applications of survival tree and survival ensemble methods. *Psychometrika*. 2015;80(3):811-833. doi: 10.1007/s11336-014-9413-1. [PubMed] [Cross Ref].

9. Bou-Hamad I, Larocque D, Ben-Ameur H. A review of survival trees. *Stat Surv.* 2011; 5:44–71. doi: 10.1214/09-SS047. [Cross Ref]
10. World Health Organization, Global Strategy on Infant and Young Child Feeding, World Health Organization, Geneva, Switzerland, 2003.
11. Robert E, Coppieters Y, Swennen B and Dramaix M. Breastfeeding Duration: A Survival Analysis—Data from a Regional Immunization Survey. *BioMed Research International.* Volume 2014, Article ID 529790, 8 pages, <http://dx.doi.org/10.1155/2014/529790>
12. World Health Organization. Guiding principles for complementary feeding of the breastfed child. Division of Health Promotion and Protection. Geneva. 2003. available at <http://apps.who.int/iris/bitstream/10665/42590/1/9241562218.pdf>
13. Kasahun A W, Wako W G, Gebere M V and Neima G H. Predictors of exclusive breastfeeding duration among 6–12 month aged children in gurage zone, South Ethiopia: a survival analysis. *International Breastfeeding Journal* (2017) 12:20 DOI 10.1186/s13006-017-0107-z.
14. Teresa S, Abada J, Trovato F, Lalu N. Determinants of breast feeding in the Philippines: a survival analysis. *Soc Sci Med.* 2001; 52:71–81. doi: 10.1016/S0277-9536(00)00123-4. [PubMed] [Cross Ref]
15. Chaves RG, Lamounier JA, César CC. Factors associated with duration of breastfeeding. *J Pediatr (Rio J).* 2007;83(3):241-246. doi 10.2223/JPED.1610
16. Mohamadpour R, Behnampour B, abdollahi F, Sheykholeslami A, Mehrbakhsh Z, Barzanuni S. Determination of effective factors in breastfeeding duration using survival analysis. *J Res Dev Nurs Midwifery.* 2017;14(2):45-50.
17. Schumacher M, Hollander N, Schwarzer G, Sauerbrei W. Handbook of Statistics in Clinical Oncology. In J Crowley (ed.), *Prognostic Factor Studies.* New York: Marcel Dekker Basel; 2001
18. Zhang H, Singer B. *Recursive partitioning in the health sciences: Springer Science & Business Media;* 2013.
19. Hothorn T, Hornik K, Zeileis A. Unbiased recursive partitioning: A conditional inference framework. *Journal of Computational and Graphical statistics.* 2006;15(3):651-674.
20. Holfprnd TR. *Multivariate Methods in Epidemiology.* New York: Oxford University Press; 2002.
21. Hothorn T. "Package party. 2009; Available from: <http://cran.R-projec.org>
22. Hothorn T, Hornik K, Strobl C, Zeileis A. Package 'party' Version 1.3-7: A Laboratory for Recursive Partitioning. URL <http://party.R-forge.R-project.org> 2021-03-03 17:10:13 UTC
23. Schlosser L, Hothorn T, Zeileis A (2019). "The Power of Unbiased Recursive Partitioning: A Unifying View of CTree, MOB and GUIDE", arXiv:1906.10179, arXiv.org E-Print Archive. <https://arXiv.org/abs/1906.10179>
24. Tarone RE, Ware J. On distribution-free tests for equality of survival distributions. *Biometrika* 1977; 64: 156-60. Google Scholar, Crossref
25. Saki Malehi A, Hajizadeh E, Fatemi R. Evaluation of prognostic variables for classifying the survival in colorectal patients using the decision tree. *Iranian Journal of Epidemiology.* 2012;8(2):13-19.
26. Parizadeh D, Ramezankhani A, Momenan AA, Azizi F, Hadaegh F. Exploring risk patterns for incident ischemic stroke during more than a decade of follow-up: a survival tree analysis. *Computer methods and programs in biomedicine.* 2017; 147:29-36.
27. Ramezankhani A, Tohidi M, Azizi F, Hadaegh F. Application of survival tree analysis for exploration of potential interactions between predictors of incident chronic kidney disease: a 15-year follow-up study. *Journal of translational medicine.* 2017;15(1):240.
28. Valera VA, Walter BA, Yokoyama N, et al. Prognostic groups in colorectal carcinoma patients based on tumor cell proliferation and classification and regression tree (CART) survival analysis. *Annals of surgical oncology.* 2007;14(1):34-40.
29. Shimokawa A, Kawasaki Y, Miyaoka E. Comparison of splitting methods on survival tree. *Int J Biostat.* 2015 May;11(1):175-88. doi: 10.1515/ijb-2014-0029.