# COMPUTER PROGRAMMING BASED: AN OPTIMIZED PROGNOSTIC MODEL OF THE FUZZY SURVIVAL ORAL CANCER USING SAS

WAN MUHAMAD AMIR W AHMAD[1], MUHAMMAD AZEEM YAQOOB[1], NURFADHLINA ABDUL HALIM[2], NOR AZLIDA ALENG[2], ZALILA ALI[3], RUHAYA HASAN[1]

_**Abstract.** There are various techniques to determine the associated factors for oral cancer data. One of the most common techniques (in medical research) that always used is the prognostic model of the survival. Through the cox proportional hazards regression, we can investigate the association between the survival times of patients with the potential predictor variables. Current research nowadays is being focused on the most efficient model for the right decision making. To optimize the gained results, a combination method is being considered which consist of the fuzzy approach to the prognostic model of the survival. This promising technique had given a full attention because of its ability to give the best results. The model which gained from the proposed method is employed to estimate the potential factors that contribute to the oral cancer case. This provides useful information for the determination of potential factors that lead to oral cancer._

_**Keywords:** Prognostic Cancer Model, Fuzzy linear model, SAS._

## 1. INTRODUCTION TO THE PROGNOSTIC SURVIVAL CANCER MODEL

Survival analysis is a method for analyzing data where the outcome variable is the time until the occurrence of an event of interest. In survival analysis, the event will not necessarily have occurred in all individuals (patients) by the time the study ends, and for these patients, their full survival times are unknown [1]. Survival analysis is useful in clinical research because it focuses on comparing the survival distributions and the identification of risk factors [2]. In survival analysis, it is common for a proportion of patient to remain alive and disease-free at the end of the follow-up period [1]. The prognostic model of the survival was introduced in 1972 by Cox in order to estimate the effects of different covariates influencing the times to the failures of a system. According to Cox in 1972, the prognostic model of survival has been used rather extensively in biomedicine and engineering [3]. Survival model is a powerful tool that is used frequently in studies of clinical outcomes. Survival model can involve a mixture data which consist of categorical and continuous variables and they can handle partially observed (censored) responses. However, uncritical application of modeling techniques can result in models that poorly fit the dataset at

---

[1] Universiti Sains Malaysia, School of Dental Sciences, 11800 Gelugor, Penang, Malaysia.
E-mail: wmamir@usm.my.
[2] Universiti Malaysia Terengganu, School of informatics and Applied Mathematics, 11800 Gelugor, Penang, Malaysia.
[3] Universiti Sains Malaysia, School of Mathematical Sciences, 21030 Kuala Terengganu, Malaysia.

hand, or, even more likely, inaccurately predict outcomes in new subjects. Measurement of predictive accuracy can be hard for survival time data in the presence of censoring [4]. The model of the survival was used to inference a prognostic model of metastatic hormone-refractory prostate cancer patients (HRPC) from 1991 to 2001 which is consisted of 1,101 patients. Calibration of the survival model predictions was assessed by comparing the predicted probability with the actual survival probability [5].

In medical analysis for example, Seker, Odetayo et al. in 2003 [6] had investigated the used of fuzzy k-nearest neighbor (FK-NN) classifier as a fuzzy logic method that provides a certainty degree for prognostic decision and assessment of the markers. They also combined their proposed techniques with logistic regression as a statistical method and multilayer feed forward backpropagation neural networks an artificial neural-network tool, the latter two techniques having been widely used for oncological prognosis because the gained result leads to an accurate and right reliable decision making [6]. The rapid increased in the methodology development of prognostic model has dictated the need for developing reliable methods for extracting clinically decision making. To proof that, a result based on breast and prostate cancer data has indicated that the FK-NN-based method yields the highest predictive accuracy and also has produced a more reliable prognostic marker model by combining both statistical and artificial neural-network-based methods [6].

In 1995, a survival model was developed using the following predictor variables: diagnosis, age, number of days in the hospital before study entry, presence of cancer, neurologic function, and 11 physiologic measures recorded on day 3 after study entry. Physicians were interviewed on day 3. Patients were followed for survival for 180 days after study entry [7]. In 1985, Chen and George [8] investigated the stability of a stepwise selection procedure in the framework of the Cox proportional hazard regression model based on bootstrap resampling procedure. They develop a bootstrap-model selection procedure, combining with with existing selection techniques for the best variable selection and illustrate the proposed strategy using data from two cancer clinical trials featuring two different situations [9]. Chen and George in 1985 [8], describes the use of the bootstrap in prognostic survival model for acute lymphocytic leukemia patients using computer-based statistical methodology. To validate the accuracy of the prognostic survival cancer model, they use a bootstrap resampling technique (100 bootstrap samples) to select the important prognostic factors via a stepwise regression. At the second stage, it is involved 400 bootstrap samples for the estimate the corresponding regression parameters. The bootstrap result suggests that the model constructed from the training set is reasonable [8].

## 2. DATA AND ALGORITHM FOR PROGNOSTIC SURVIVAL CANCER MODEL

Data from medical record unit were review and related information was extracted. The sampling frame was the list of patient diagnosed with oral cancer admitted to HUSM. The details of the studied variable as shown as follows:

**Table 1. Description of Data**

| Num. | Variables | Explanation of user variables |
|------|-----------|-------------------------------|
| 1 | Gender | 1= Male , 2= Female |
| 2 | Smoking | Smoking Status [0= Never, 1= Yes] |
| 3 | Alcohol | Alcohol used [1= Never, 2= Stop, 3=Current] |
| 4 | Betel | Betel Quid [0= Never, 1= Yes] |
| 5 | Nerve | Nerve Invasion [0= Never, 1= Yes] |
| 6 | Time | Time in months |

## 2.1 *FLOW CHART AND ALGORITHM FOR PROGNOSTIC ORAL CANCER MODEL*

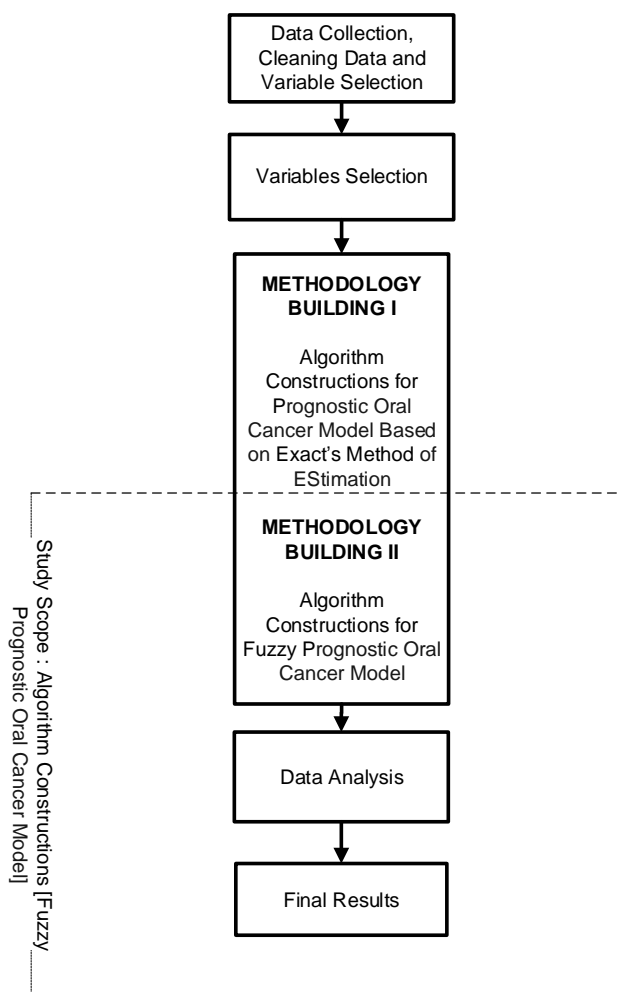Below is the flow chart of prognostic fuzzy model for oral cancer model using SAS algorithm.



**Figure 1. Study Scope for Prognostic Oral Cancer Model.**

## Methodology of Algorithm for Prognostic Oral Cancer Model

**Part I:** `/*PROGRAMMING FOR ORAL CANCER MODEL BASED ON EXACT METHODS */`

```
%MACRO bootstrap(data=_last_, booted=booted, boots=10, seed=1234);
DATA &booted;
** randomly picks an integer from 1 to n;
pickobs = INT(RANUNI(&seed)*n)+1;
** POINT tells SAS to read value pickobs
** NOBS sets n to number of obs in &Data;
** when the point option is used SAS will loop through the data step forever;
SET &data POINT = pickobs NOBS = n;
** saves number of current bootstrap;
REPLICATE=int(i/n)+1;
        i+1;
** stop will leave data set when n*&boots obs have been created;
IF i > n*&boots THEN STOP;
RUN;
%MEND bootstrap;

/* INPUT DATA */
data Cancer;
input  Gender Smoking Alcohol Betel Nerve Time;
cards;
2.00     1.00     2.00     2.00     0.00     87.00
1.00     1.00     3.00     1.00     0.00     18.00
1.00     2.00     1.00     1.00     0.00     65.00
1.00     3.00     2.00     1.00     0.00     69.00
1.00     3.00     2.00     2.00     0.00     42.00
1.00     3.00     2.00     1.00     0.00     44.00
2.00     1.00     1.00     3.00     1.00     13.00
2.00     1.00     1.00     3.00     1.00     15.00
2.00     1.00     1.00     3.00     0.00     19.00
2.00     1.00     1.00     3.00     0.00     18.00
2.00     1.00     1.00     3.00     0.00     77.00
1.00     2.00     2.00     2.00     0.00     11.00
1.00     3.00     3.00     1.00     0.00     37.00
2.00     1.00     1.00     3.00     1.00     16.00
1.00     2.00     2.00     1.00     1.00      9.00
2.00     1.00     1.00     3.00     1.00      7.00
;
run;
ods rtf file='abc.rtf' style=journal;

/**GENERATE BOOTSTRAP SAMPLE**/
%bootstrap(data=Cancer, boots=10);
run;

/**PRINT DATA **/
proc print data=booted;
run;

/**SURVIVAL ANALYSIS**/
Proc lifetest data=booted plots= (s);
Title 'Survival by Treatment';
Time Time*Nerve(1);
Strata Gender;
run;
```

```
proc lifetest data= booted plots=(s,ls,lls) censoredsymbol=none;
time Time*Nerve(1);
strata Gender;
run;



/******** EXACT PROCEDURE ********/
Proc phreg data= booted;
model Time*Nerve(1) =  Gender Smoking Alcohol Betel /ties=exact;
BASELINE OUT=set1 SURVIVAL=st LOGSURV=lst LOGLOGS=llst;
OUTPUT OUT=resid1 DFBETA=dfgred RESSCH=scgred RESDEV=deres
                  RESMART=mares XBETA=linpred STDXBETA=cipred;
RUN;
PROC PRINT DATA=set1;
RUN;
PROC PRINT DATA=resid1;
RUN;
PROC GPLOT DATA=resid1;
PLOT dfgred*Time;
RUN;
ods rtf
close;
```

**Part II:** /*PROGRAMMING FOR PROGNOSTIC OF ORAL CANCER MODEL BASED NONLINEAR MODEL */

```
/* For Gender Smoking Alcohol Betel_ Quid */
proc nlp;
min Y;
decvar a0c a0w a1c a1w a2c a2w a3c a3w a4c a4w;
bounds a0w>=0, a1w>=0, a2w>=0, a3w>=0, a4w>=0;
lincona0c+2*a1c+1*a2c+2*a3c+2*a4c-a0w-2*a1w-1*a2w-2*a3w-2*a4w<=-8.26;
lincona0c+1*a1c+1*a2c+3*a3c+1*a4c-a0w-1*a1w-1*a2w-3*a3w-1*a4w<=-3.19;
lincona0c+1*a1c+2*a2c+1*a3c+1*a4c-a0w-1*a1w-2*a2w-1*a3w-1*a4w<=-6.40;
lincona0c+1*a1c+3*a2c+2*a3c+1*a4c-a0w-1*a1w-3*a2w-2*a3w-1*a4w<=-7.07;
lincona0c+1*a1c+3*a2c+2*a3c+2*a4c-a0w-1*a1w-3*a2w-2*a3w-2*a4w<=-4.37;
lincona0c+1*a1c+3*a2c+2*a3c+1*a4c-a0w-1*a1w-3*a2w-2*a3w-1*a4w<=-7.07;
lincona0c+2*a1c+1*a2c+1*a3c+3*a4c-a0w-2*a1w-1*a2w-1*a3w-3*a4w<=-6.41;
lincona0c+2*a1c+1*a2c+1*a3c+3*a4c-a0w-2*a1w-1*a2w-1*a3w-3*a4w<=-6.41;
lincona0c+2*a1c+1*a2c+1*a3c+3*a4c-a0w-2*a1w-1*a2w-1*a3w-3*a4w<=-6.41;
lincona0c+2*a1c+1*a2c+1*a3c+3*a4c-a0w-2*a1w-1*a2w-1*a3w-3*a4w<=-6.41;
lincona0c+2*a1c+1*a2c+1*a3c+3*a4c-a0w-2*a1w-1*a2w-1*a3w-3*a4w<=-6.41;
lincona0c+1*a1c+2*a2c+2*a3c+2*a4c-a0w-1*a1w-2*a2w-2*a3w-2*a4w<=-2.86;
lincona0c+1*a1c+3*a2c+3*a3c+1*a4c-a0w-1*a1w-3*a2w-3*a3w-1*a4w<=-6.22;
lincona0c+2*a1c+1*a2c+1*a3c+3*a4c-a0w-2*a1w-1*a2w-1*a3w-3*a4w<=-6.41;
lincona0c+1*a1c+2*a2c+2*a3c+1*a4c-a0w-1*a1w-2*a2w-2*a3w-1*a4w<=-5.56;
lincona0c+2*a1c+1*a2c+1*a3c+3*a4c-a0w-2*a1w-1*a2w-1*a3w-3*a4w<=-6.41;
lincona0c+2*a1c+1*a2c+2*a3c+2*a4c+a0w+2*a1w+1*a2w+2*a3w+2*a4w>=-8.26;
lincona0c+1*a1c+1*a2c+3*a3c+1*a4c+a0w+1*a1w+1*a2w+3*a3w+1*a4w>=-3.19;
lincona0c+1*a1c+2*a2c+1*a3c+1*a4c+a0w+1*a1w+2*a2w+1*a3w+1*a4w>=-6.40;
lincona0c+1*a1c+3*a2c+2*a3c+1*a4c+a0w+1*a1w+3*a2w+2*a3w+1*a4w>=-7.07;
lincona0c+1*a1c+3*a2c+2*a3c+2*a4c+a0w+1*a1w+3*a2w+2*a3w+2*a4w>=-4.37;
lincona0c+1*a1c+3*a2c+2*a3c+1*a4c+a0w+1*a1w+3*a2w+2*a3w+1*a4w>=-7.07;
lincona0c+2*a1c+1*a2c+1*a3c+3*a4c+a0w+2*a1w+1*a2w+1*a3w+3*a4w>=-6.41;
lincona0c+2*a1c+1*a2c+1*a3c+3*a4c+a0w+2*a1w+1*a2w+1*a3w+3*a4w>=-6.41;
lincona0c+2*a1c+1*a2c+1*a3c+3*a4c+a0w+2*a1w+1*a2w+1*a3w+3*a4w>=-6.41;
lincona0c+2*a1c+1*a2c+1*a3c+3*a4c+a0w+2*a1w+1*a2w+1*a3w+3*a4w>=-6.41;
lincona0c+2*a1c+1*a2c+1*a3c+3*a4c+a0w+2*a1w+1*a2w+1*a3w+3*a4w>=-6.41;
lincona0c+1*a1c+2*a2c+2*a3c+2*a4c+a0w+1*a1w+2*a2w+2*a3w+2*a4w>=-2.86;
lincona0c+1*a1c+3*a2c+3*a3c+1*a4c+a0w+1*a1w+3*a2w+3*a3w+1*a4w>=-6.22;
```

```
lincona0c+2*a1c+1*a2c+1*a3c+3*a4c+a0w+2*a1w+1*a2w+1*a3w+3*a4w>=-6.41;
lincona0c+1*a1c+2*a2c+2*a3c+1*a4c+a0w+1*a1w+2*a2w+2*a3w+1*a4w>=-5.56;
lincona0c+2*a1c+1*a2c+1*a3c+3*a4c+a0w+2*a1w+1*a2w+1*a3w+3*a4w>=-6.41;
Y=a0w*16+24.00*a1w+27.00*a2w+26.00*a3w+33*a4w;
run;
```

## 3. RESULTS AND DISCUSSION

PART I: PROGNOSTICS CANCER MODEL WITH SURVIVAL CURVE ESTIMATION TIME

Fig. 2 shows the survival probabilities for nerve invasion scenario according to gender. The plot shows that the survival probability is about lower for females compared to male at all times point to develop nerve invasion among oral cancer patient which registered in Hospital University Sains Malaysia (HUSM).
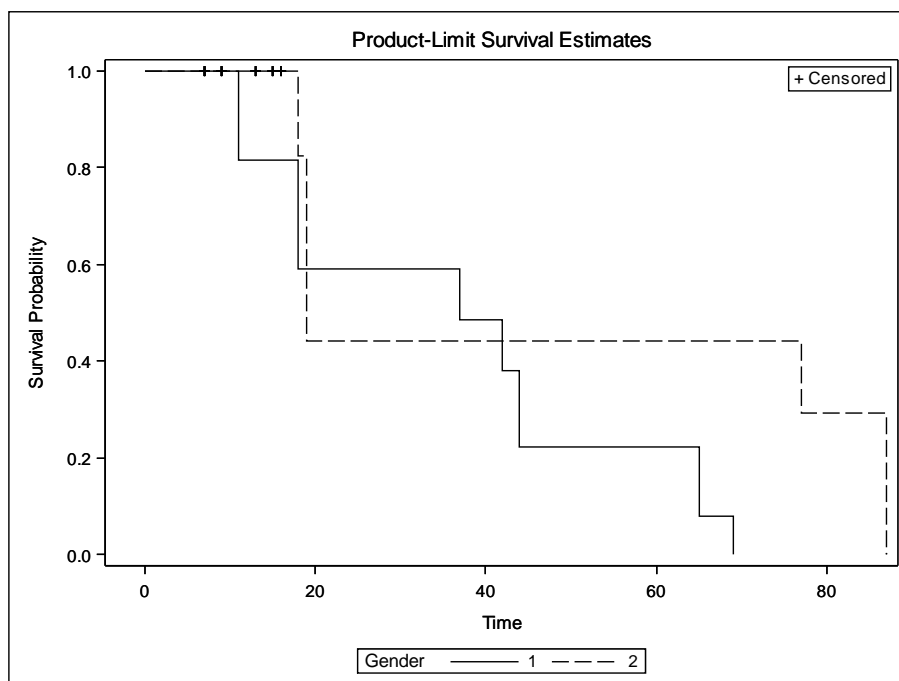


**Figure 2. Survival Probabilities for Nerve Invasion Scenario According to Gender.**

**Table 2. Discrete's Method for Prognostics Cancer Estimation.**

| Analysis of Maximum Likelihood Estimates | | | | | | |
|---|---|---|---|---|---|---|
| **Parameter** | **DF** | **Parameter Estimate** | **Standard Error** | **Chi-Square** | **p-value**[*] | **Hazard Ratio** |
| Gender | 1 | -6.92020 | 0.74079 | 87.2662 | <.0001 | 0.001 |
| Smoking | 1 | -1.51533 | 0.22273 | 46.2881 | <.0001 | 0.220 |
| Alcohol | 1 | 0.84705 | 0.21668 | 15.2819 | <.0001 | 2.333 |
| Betel | 1 | 2.70080 | 0.30624 | 77.7806 | <.0001 | 14.892 |

[*]*significant at $p < 0.05$.*

Table 2 shows the results of Discrete's Method estimation for prognostics cancer. The finding shows that there are five factors were associated to the survival of oral cancer towards nerve invasion. Four factors Gender ($\beta_1 = $ -6.92020, p < 0.0001), Smoking ($\beta_2 = $ -1.51533, p < 0.0001), Alcohol ($\beta_3 = 0.84705$, p < 0.0001) and Betel quid ($\beta_4 = 2.70080$, p < 0.0001) were significant at $\alpha = 0.05$. The prognostics cancer model using Discrete's Method is given by

$$\text{HR = Exp [-6.92020 (Gender) -1.51533 (Smoking) + 0.84705 (Alcohol)} + 2.70080 \text{ (Betel)]}$$

Taking natural logarithm to the equation above, we can achieve the equation as shown as

$$\text{HR= Exp [-6.92020 (Gender) -1.51533 (Smoking) + 0.84705 (Alcohol)} + 2.70080 \text{ (Betel)]}$$

Then we obtained the model as follow:

$$\text{Ln [HR] = [-6.92020 (Gender) -1.51533 (Smoking) + 0.84705 (Alcohol)} + 2.70080 \text{(Betel)]} \tag{1}$$

## PART II: FUZZY REGRESSION

**Table 3. Procedures Nonlinear Programming: Nonlinear Minimization**
**(Optimization Results for Parameter Estimates).**

| N | Parameter | Estimate | Gradient Objective Function | Variables |
|---|---|---|---|---|
| 1 | a0c | 0.001250 | 0 | Coefficient |
| 2 | a0w | 0 | 16.000000 | |
| 3 | a1c | -6.916250 | 0 | Gender |
| 4 | a1w | $-2.1684 \times 10^{-19}$ | 24.000000 | |
| 5 | a2c | -1.515000 | 0 | Smoking |
| 6 | a2w | 0 | 27.000000 | |
| 7 | a3c | 0.845625 | 0 | Alcohol |
| 8 | a3w | 0.001875 | 26.000000 | |
| 9 | a4c | 2.697500 | 0 | Betel quid |
| 10 | a4w | 0 | 33.000000 | |

*Value of Objective Function = 0.04875*

## Prognostic Fuzzy Model for Survival

In fuzzy regression methodology, parameters are estimated by minimizing total vagueness in the model. The prediction equations for computing upper and lower limits, lower and upper widths of prediction interval for fitted fuzzy linear regression models are computed respectably as:

## Upper limit of prediction interval for fuzzy model for survival:

$$\text{Hazard} = (0.001250 +0) + (-6.916250 + (-2.1684 \times 10^{-19})) \text{ Gender} + (-1.515000+0)$$
$$\text{Smoking} + (0.845625+ 0.001875) \text{ Alcohol} + (2.697500+0) \text{ Betel quid}$$

$$\text{Hazard} = 0.001250 - 6.916250\text{Gender} -1.515000 \text{ Smoking}+ 0.8475 \text{ Alcohol}$$
$$+ 2.697500 \text{ Betel quid} \tag{2}$$

## Lower limit of prediction interval for fuzzy model for survival:

$$\text{Hazard} = (0.001250 -0) + (-6.916250-(-2.1684 \times 10^{-19})) \text{ Gender} + (-1.515000-0)$$
$$\text{Smoking} + (0.845625-0.001875) \text{ Alcohol} + (2.697500-0) \text{ Betel quid}$$

$$\text{Hazard} = (0.001250) -6.916250 \text{ Gender} -1.515000\text{Smoking}+ 0.84375\text{Alcohol}$$
$$+ 2.697500\text{Betel quid} \tag{3}$$

## Prognostic Model for Survival

From equation (1) from the Part I of an analysis, the prognostic oral cancer model:

Ln [HR] = [-6.92020 (Gender) -1.51533 (Smoking) + 0.84705 (Alcohol) + 2.70080 (Betel)]
Std. Error: [0.74079]          [0.22273]          [0.21668]          [0.30624]

Further, the prediction equations for the upper prognostic oral cancer model regression model, is given by:

$$\text{Ln [HR]} = [-6.92020+ 0.74079] \times (\text{Gender}) + [-1.51533+ 0.22273] \times (\text{Smoking}) +$$
$$[0.84705+ 0.21668] \times (\text{Alcohol}) + [2.70080+ 0.30624] \times (\text{Betel})]$$

$$\text{Ln [HR]} =-6.17941 (\text{Gender}) -1.2926 (\text{Smoking}) + 1.06373(\text{Alcohol})$$
$$+ 3.00704(\text{Betel})$$

The prediction equations for the lower prognostic oral cancer model regression model, is given by:

$$\text{Ln [HR]} = [-6.92020-0.74079] \times (\text{Gender}) + [-1.51533-0.22273] \times (\text{Smoking})$$
$$+ [0.84705-0.21668] \times (\text{Alcohol}) + [2.70080-0.30624] \times (\text{Betel})]$$

$$\text{Ln [HR]} = -7.66099 (\text{Gender})-1.73806 (\text{Smoking}) + 0.63037(\text{Alcohol})$$
$$+2.39456(\text{Betel})$$

**Table 4. Fitting of prognostic oral cancer model regression model and fuzzy oral cancer model regression.**

| Prognostic Oral Cancer Model | | | Fuzzy Regression (FR) Model | | |
|---|---|---|---|---|---|
| **Lower limit** | **Upper limit** | **Width** | **Lower limit** | **Upper limit** | **Width** |
| -11.01 | -5.51 | 5.50 | -8.26 | -8.26 | 0.008 |
| -5.11 | -1.27 | 3.84 | -3.19 | -3.20 | 0.011 |
| -8.11 | -4.69 | 3.42 | -6.40 | -6.40 | 0.004 |
| -9.22 | -4.92 | 4.30 | -7.07 | -7.08 | 0.008 |
| -6.83 | -1.92 | 4.91 | -4.37 | -4.38 | 0.008 |
| -9.22 | -4.92 | 4.30 | -7.07 | -7.08 | 0.008 |
| -9.25 | -3.57 | 5.68 | -6.41 | -6.41 | 0.004 |
| -9.25 | -3.57 | 5.68 | -6.41 | -6.41 | 0.004 |
| -9.25 | -3.57 | 5.68 | -6.41 | -6.41 | 0.004 |
| -9.25 | -3.57 | 5.68 | -6.41 | -6.41 | 0.004 |
| -9.25 | -3.57 | 5.68 | -6.41 | -6.41 | 0.004 |
| -5.09 | -0.62 | 4.46 | -2.86 | -2.86 | 0.008 |
| -8.59 | -3.86 | 4.73 | -6.22 | -6.23 | 0.011 |
| -9.25 | -3.57 | 5.68 | -6.41 | -6.41 | 0.004 |
| -7.48 | -3.63 | 3.85 | -5.55 | -5.56 | 0.008 |
| -9.25 | -3.57 | 5.68 | -6.41 | -6.41 | 0.004 |
| Average width | | 4.942 | Average width | | 0.006 |

The width of prediction intervals in respect of prognostic oral cancer model and fuzzy regression corresponding to each set of observed explanatory variables is computed in SPSS and the results are reported in the following Table 4. From this table, average width for former was found to be 4.942 while that for latter was 0, indicating thereby the superiority of fuzzy regression methodology.

## 4. CONCLUSIONS

This paper demonstrates an algorithm that has a high potential to determine factors that contribute to prognostic oral cancer model regression. Combining the exact method of estimation, bootstrap and nonlinear fuzzy regression leads to the high-performance result.

For the better, we provided the dataset in the SAS programming language. For the information, the result will be slightly different when we running the SAS for the second time, this is because the method will choose the sample randomly and differently. But all the result will lead to the most accuracy of parameter estimation.

Our second aim is to share the algorithm and also provide the researcher with an alternative programming that suit for the case of time event data. This proposed method can be applied to as a second option on the hazard estimation for the oral cancer patient. With more detailed of programming language and step by step explanation, the estimation of

prognostic oral cancer model will more precise and obvious and this will indirectly open a new idea for the current stage of research.

## REFERENCES

[1]     Bradburn, M.J., et al. *British Journal of Cancer*, **89**(3), 431, 2003.
[2]     Lee, S.J., et al. *American Journal of Orthodontics and Dentofacial Orthopedics*, **137**(2), 194, 2010.
[3]     Cox, D.R., *J. Royal Statist. Soc.*, **134**, 187, 1972.
[4]     Harrell, F.E., et al., *Statistics in Medicine*, **15**(4), 361, 1996.
[5]     Halabi, S., et al., *Journal of Clinical Oncology*, **21**(7), 1232, 2003.
[6]     Seker, H., et al., *IEEE Transactions on Information Technology in Biomedicine*, **7**(2), 114, 2003.
[7]     Knaus, W.A., et al., *Annals of Internal Medicine*, **122**(3), 191, 1995.
[8]     Chen, C.H., George, S.L., *Statistics in Medicine*, **4**(1), 39, 1985.
[9]     Sauerbrei, W., Schumacher, M., *Statistics in Medicine*, **11**(16), 2093, 1992.