

Applying Logistic Regression Model to The Second Primary Cancer Data.

Amr I. Abdelrahman

**Department of Statistics, Mathematics, and Insurance.
Faculty of Commerce, Ain Shams University, Egypt**

Abstract

The logistic regression model is used to determine the social-demographic risk factors which affect the second cancer occurrence for 200 patients who were initially treated for first primary cancer stage I and were cancer free for at least 1 year after first primary cancer treatment. The 200 patients were classified as "having a second cancer", and "not having a second cancer". The social-demographic risk factors used are age at first cancer, gender, area the patient lives in, marital status, family history, smoking, education and obesity in addition to treatment by radiation. The binary Logistic regression model is used in this study to estimate the probability of the occurrence and to determine the effective risk factors that cause the second cancer occurrence. The odds ratio analysis compare whether the probability of having a second primary cancer is the same for each covariate groups. Significance testing for the logistic coefficients using Wald test and likelihood ratio show that five risk factors were significant. To assess the fitness of the model the Hosmer and Lemeshow test is used. The logistic regression model proved to have a lower sensitivity level due to the clinical risk factors not considered in this study.

Keywords: Logistic regression model; Wald test, Odds Ratio, Cross-Validation; Roc curve; Second primary cancer.

1. Introduction

Early detection and evaluation of the risk factors which might cause the occurrence of second cancer is very important. The prediction of risk factors is an important pivot of the war against cancer. The use of statistical methods to identify risk factors would help to identify the probability of second cancer occurrence.

We distinguish between two medical cases: a) Recurrence case: Cancer that has recurred (come back), usually after a period of time during which the cancer could not be detected. The cancer may come back to the same place as the original (primary) tumor or to another place in the body. Also called recurrent cancer, and

b) Second cancer: a new primary cancer in a person with a history of another cancer.

According to DeVita et al. (2008) Second cancers can reflect the late sequel of treatment, as well as the influence of lifestyle factors, environment exposures, host determinants and gene-gene interactions. The main life style factors are tobacco and alcohol; the environmental factors are: contaminants and viruses; and the host factors are gender, age, genetics, immune function and hormonal factors.

A statistical model is proposed to explain the association between the studied covariates and its effect on the probability of the second cancer occurrence. Data included 200 patients were have a first primary cancer stage I, and have at least one year free cancer after first cancer treatment. Covariates used in the analysis were Age at first Cancer, Gender, Marital status, Area the patient lives in, Treatment by Radiation, Family History, Smoking, Obesity, and Education status. This study proposes to:

- a. Determine the effective risk factors that cause the second cancer occurrence and propose a statistical model to explain the association between the studied covariates and second cancer occurrence.
- b. Explain the relative risk for each studied covariate and its effect on the probability of the second cancer occurrence.

In Section 2, we present the logistic regression model to estimate the probability of occurrence of second cancer; the Wald test, likelihood ratio test, Hosmer-Lemeshow test, cross validation methods and ROC curve are also introduced in section 2. In Section 3, we apply the binary regression model to the data; SPSS is used for the analysis. Summary and conclusions are given in Section 4.

2. The Binary Logistic Regression Model

The logistic regression model has been used in many disciplines including medical studies. It has been used in the social research (Ingles et al., 2009; King and Zeng, 2002; Saijo et al., 2008; and Garcia-Ramirez et al., 2005), in market research (Neagu and Hoerl, 2005; Kleijnen et al., 2004; Barone et al., 2007; Sallis and Sharma, 2009; and Kirkos, 2009), also become an important tool at the commercial applications (Erhart et al., 2009; O'Leary

, 2009; and Weber et al., 2008); and in medical studies(Sanchez et al., 2008; Kaufman et al., 2000; Rubino et al., 2003).

The dependent variable of the logistic model is classified into two basic types (Afifi et al., 2004);

- a- Continuous Variable: can assume any value within a specified range.
- b- Discrete Variable: can only assume certain values and there are usually “gaps” between values(categorical response has two main categories: success (occurrence) and fail (no occurrence)).

Everitt (1998) gave the following definition for logistic distribution:" the limiting probability distribution as n tends to infinity, of the average of the largest to smallest sample values, of random samples of size n from an exponential distribution".

The logistic distribution is given by

$$f(x) = \frac{\exp[(x - \alpha) / \beta]}{\beta \{1 + \exp - [(x - \alpha) / \beta]\}^2} \quad -\infty < x < \infty, \beta > 0$$

The location parameter α is the mean. The variance of the distribution is $\pi^2 \beta^2 / 3$, its skewness is zero and its kurtosis is 4.2 The standard logistic distribution with $\alpha = 0, \beta = 1$, with cumulative probability function. F(x), and probability distribution, f(x), has the property

$$f(x) = F(x) [1 - F(x)]$$

see also, (Evans et al.,1993).

The logistic regression is a form of regression analysis used when the response variable is a binary variable (Altman,1991and Everitt, 1998). The method is based on the logistic transformation or logit proportion, namely;

$$\text{Logit}(p) = \frac{p}{1-p}$$

Where ;

$$\begin{aligned} p &= \text{Pr}(y = 1) \\ (1-p) &= \text{Pr}(y=0) \end{aligned}$$

As p tends to 0, $\text{Logit}(p)$ tends to $-\infty$ and as p tends to 1, $\text{Logit}(p)$ tends to ∞ . The function $\text{Logit}(p)$ is a sigmoid curve that is symmetric about $p = 0.5$

The logistic regression makes no assumption about the distribution of the independent variables. They do not have to be normally distributed, linearly related or of equal variance within each group. The relationship between the predictor and response variables is not a linear function in logistic regression.

The logistic regression function is the logit transformation of P , where;

$$\text{Logit}(P) = \ln \frac{p}{1-p} = \beta_0 + \beta_1 x_1 + \dots + \beta_q x_q$$

Where β_0 = the constant of the equation and, β_i = the coefficient of the predictor variables i . Using the logistic transformation in this way overcomes problems that might arise if p was modeled directly as a linear function of the explanatory variables; in particular it avoids fitted probabilities outside the range (0, 1). The parameters in the model can be estimated by maximum likelihood estimation.

The slope coefficient β_j associated with an explanatory variable x_j represents the change in log odds for an increase of one unit in x_j .

To assess the significance of the logistic regression coefficients the Wald statistic and likelihood ratio test are used (Afifi et al., 2004). The Wald statistic takes the form:

$$\left(\frac{\hat{\beta}}{s.e(\hat{\beta})} \right)^2$$

Where $\hat{\beta}$ represents the estimated coefficient β and $s.e(\hat{\beta})$ is its standard error. Under the null hypothesis of zero slope and based on asymptotic theory, this quantity follows a chi-square distribution with one degree of freedom. If the estimated value of the slope is small and its estimated variability is large, then we can not conclude that the slope is significantly different from zero and vice versa (Afifi et al., 2004).

The likelihood ratio test for overall significance of the beta's coefficients for the independent variables in the model is used

(Hosmer and Lemeshow, 2000; Fienberg,1998). The test based on the statistic " G" under the null hypothesis that the beta's coefficients for the covariates in the model are equal to zero. G statistic takes the form:

$$G = -2\ln\left[\frac{\text{Likelihood} \cdot \text{without the variable}}{\text{Likelihood} \cdot \text{with the variable}}\right]$$

The distribution of "G" is a chi-square with q degree-of-freedom, where q is the number of covariates in the logistic regression equation. Hauck and Donner (1977) and Jennings (1986) examined the performance of the Wald test and found that the test often failed to reject the null hypothesis when the coefficient was significant. They recommended that the likelihood ratio test to be used.

The likelihood statistic L is used to assess the fitness of the model. The sampling distribution of the $-2 \log L$ has a chi-square distribution with q degrees of freedom under the null hypothesis that all regression coefficients of the model are zero (Fienberg, 1998). A significant p-value provides evidence that at least one of the regression coefficients for an explanatory variable is non zero.

Hosmer and Lemeshow (2000) developed a goodness-of-fit test for logistic regression models with binary responses. They proposed grouping based on the value of the estimated probabilities. This test is obtained by calculating the Pearson chi-square statistic from the $2 \times g$ table of observed and expected frequencies, where g is the number of groups. The statistic is written

$$x_{HW}^2 = \sum_{i=1}^g \frac{(o_i - N_i \bar{\pi}_i)^2}{N_i \bar{\pi}_i (1 - \bar{\pi}_i)}$$

Where;

N_i Is the number of observation in the i^{th} group.

o_i Is the number of event outcomes in the i^{th} group.

$\bar{\pi}_i$ Is the average estimated probability of an event outcome for the i^{th} group.

The Hosmer and Lemeshow statistic is then compared to a chi-square distribution with $(g - 2)$ degree of freedom. However, Christensen (1997) gave the following warnings about the Hosmer and Lemeshow goodness-of-fit test;

1. If too few groups are used to calculate the statistic (<5) it will always indicate that the model fits the data. That is why Hosmer and Lemeshow (2000) advocated that, before finally

accepting that a model fits; an analysis of the individual residuals and relevant diagnostic statistics be performed (pp.151-156).

2. It is highly dependent on how the observations are grouped.
3. It is a conservative test.
4. It has low power to detect specific types of lack of fit (such as nonlinearity in an explanatory variable).

The odds ratio

The odds ratio is a measure of association for 2×2 contingency table (Agresti, 2007). In 2×2 tables the probability of "success" is π_1 in row 1 and π_2 in row 2. Within row 1, the odds of success are defined to be:

$$odds_1 = \frac{p_1}{1-p_1} \quad \text{and} \quad odds_2 = \frac{p_2}{1-p_2}$$

Evaritt (1998) and Agresti (2002) define the odds ratio in two groups of subjects as "the ratio of odds". Thus;

$$\theta = \frac{odds_1}{odds_2} = \frac{p_1/(1-p_1)}{p_2/(1-p_2)}$$

For the binary regression model, the odd ratio is the exponent (e^{β_j}) is the ratio of odds for a one-unit change in x_j (Hosmer and Lemeshow, 2000). The change in Log odds, and the corresponding change in the odds ratio, for a c units is estimated $\exp[c \hat{\beta}_j]$ (Fleiss, 1981). When the two groups of odds are identical then the odds ratio is equal to one.

The corresponding lower and upper confidence limits for odds ratio for a c units change are $\exp[c L_j]$ and $\exp[c U_j]$, respectively, for ($c>0$), or $\exp[c U_j]$ and $\exp[c L_j]$ respectively, for ($c<0$), where (L_j, U_j) ; can be either the likelihood ratio-based confidence interval or the Wald confidence interval for β_j (Agresti, 2002 and The SAS system ,1995).

Cross Validation Techniques

Cross - validation is a general procedure used in statistical model building. It can be used to decide on the order of a statistical model including time series models, regression models, mixture distribution models, and discrimination models (Chernick, 2008).

Cross validation is performed in different ways, some of them are:

1. Take two random subsets of the data. Models are fit or various statistical procedures are applied to the first subset and then are tested on the second subset.
2. Leave - one – out technique is performed by fitting to all but one observation and then testing on the remaining one and has also been called "cross - validation by Efron" (1983), but it does not provide an adequate test.
3. Fit the model n times, each time leaving out a different observation and testing the model on estimating or predicting the observation left out each time. This provides a fair test by always testing an observations not used in the fit. It also is efficient in the use of the data for fitting the model since n – 1 observation are always used in the fit.

Hit ratio is the percentage of objects (individuals, respondents, firms, etc.) correctly classified by the logistic regression model. It is calculated as the number of objects in the diagonal of the classification matrix (H_o) divided by the total number of objects (N).

Also known as the Percentage correctly classified (Hair et al. 2009).

This can be compared with the maximum chance and proportional chance criterion to determine the discriminating power of the function. Maximum chance criterion is the percentage of the total sample represented by the larger of the two groups (H_c). The proportional chance criterion is obtained from the actual occurrence of second cancer by the equation $p^2 + (1 - p)^2$, where p = proportion of individuals in group (having a second cancer) and 1-p = proportion of individuals in group (not having a second cancer).

According to Marcoulides (1997) the difference between H_o and H_c may be tested by the following statistic

$$z = \frac{H_o - H_c}{\sqrt{H_c(N - H_c) / N}}$$

Where the significance of z is found by comparison with a critical value from a standard normal distribution.

Classification Accuracy: The ROC curve

ROC (Receiver Operating Characteristic) analysis is being used as a method for evaluation and comparison of classifiers (Ferri et. Al. 2002). The ROC gives complete description of classification

accuracy as given by the area under the ROC curve. The ROC curve originates from signal detection theory (Hosmer and Lemeshow, 2000); the curve shows how the receiver operates the existence of signal in the presence of noise.

The ROC curve plots the probability of detecting true signal (sensitivity) and false signal ($1 - \text{specificity}$) for an entire range of possible cut points.

The sensitivity and specificity of a classifier also depend on the definition of the cut-off point for the probability of predicted classes. In many situations, not all misclassifications have the same consequences, and misclassification costs have to be taken into account. A ROC curve demonstrates the trade-off between true positive rate and false positive rate in binary classification problems.

To draw a ROC curve, the true positive rate (TPR) and the false positive rate (FPR) are needed.

- TPR determines the performance of a classifier or a diagnostic test in classifying positive cases correctly among all positive samples available during the test.
- FPR, on the other hand, defines how many incorrect positive results, which are actually negative, there are among all negative samples available during the test.
- Because TPR is equivalent to sensitivity and FPR is equal to $(1 - \text{specificity})$, the ROC graph is sometimes called the sensitivity vs. $(1 - \text{specificity})$ plot.

The area under the ROC curve has become a particularly important measure for evaluating classifiers' performance because it is the average sensitivity over all possible specificities (Bradley 1997). The larger the area, the better the classifier performs. If the area is 1.0, the classifier achieves both 100% sensitivity and 100% specificity. If the area is 0.5, then we have 50% sensitivity and 50% specificity, which is no better than flipping a coin. This single criterion can be compared for measuring the performance of different classifiers analyzing a dataset. (Hanley, 1982; Bamber, 1975)

After a classifier has been made, it is also useful to measure its calibration. Calibration evaluates the degree of correspondence between the estimated probabilities of a specific outcome resulting from a classifier and the outcomes predicted by domain experts. This can then be tested using goodness-of-fit statistics. This test examines the difference between the observed frequency and the expected frequency for groups of patients and can be used to

determine if the classifier provides a good fit for the data. If the p-value is large, then the classifier is well calibrated and fits the data well. If the p-value is small, then the classifier is not well calibrated.

3. Statistical Analysis and Results

Data used for the analysis comprised of 1500 registered patients in Ain shams university hospitals, Cairo, Egypt, by different stages of cancer in 2006; 200 patients met the study assumptions were classified as:

- 1- Has a first primary cancer stage I.
- 2- Has at least one year free cancer after first cancer treatment.

The dependent variable used in the study was the classification variable (0 for those not has a second primary cancer, 1 for those who has a second primary cancer), explanatory variables used in this study were: age at first cancer occurrence, gender(male-female), marital status(married –single), area (urban or rural), radiation treatment of first cancer(yes- no) ,family history of cancer (yes, no), smoking (yes-no), Obesity before first cancer (yes-no), and education (Yes-no)for patients above 18 or parents for patients less than 18 .

SPSS software package is used for the analysis. The maximum likelihood method is used to estimate the coefficient and its standard error in addition the Newton-Raphson method to solve the nonlinear equations for the logistic model maximum likelihood estimations, table 1 shows the SPSS output.

Table 1 : The estimated coefficient , its S.E and Wald test

<i>Covariate</i>	<i>Beta estimate</i>	<i>S.E</i>	<i>Wald</i>	<i>P-value</i>
<i>Age (x₁)</i>	<i>0.007</i>	<i>0.019</i>	<i>0.154</i>	<i>0.695</i>
<i>Gender (x₂)</i>	<i>- 0.518</i>	<i>0.598</i>	<i>0.751</i>	<i>0.386</i>
<i>Marital Status (x₃)</i>	<i>1.274</i>	<i>0.562</i>	<i>5.146</i>	<i>0.023</i>
<i>Living Area(x₄)</i>	<i>0.299</i>	<i>0.426</i>	<i>0.492</i>	<i>0.483</i>
<i>Treatment.by.Radiation (x₅)</i>	<i>- 1.311</i>	<i>0.411</i>	<i>10.192</i>	<i>0.001</i>
<i>Family.History (x₆)</i>	<i>1.187</i>	<i>0.393</i>	<i>9.129</i>	<i>0.003</i>
<i>Smoking (x₇)</i>	<i>2.720</i>	<i>0.743</i>	<i>13.398</i>	<i>0.000</i>
<i>Obesity (x₈)</i>	<i>0.083</i>	<i>0.486</i>	<i>0.029</i>	<i>0.864</i>
<i>Education (x₉)</i>	<i>- 1.472</i>	<i>0.394</i>	<i>13.955</i>	<i>0.000</i>
<i>Constant</i>	<i>- 0.651</i>	<i>0.755</i>	<i>0.745</i>	<i>0.388</i>

At the .05 level of significant, Table 1 shows that " Education" ," Smoking", " Treatment by radiation", " family history", and " marital status" were highly significant.

The coefficients estimate are used to estimate the probability of the second cancer occurrence (Ashour and Abo Elfotouh 2005) as follows:

$$P (y=1 \mid x) = \frac{e^z}{1+e^z} \quad \text{or} \quad \frac{1}{1+\exp^{-z}}$$

Where ;

$$Z = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

Hence:

$$Z = -0.651 + 0.007 x_1 - 0.518 x_2 + 1.274 x_3 + 0.299 x_4 - 1.311 x_5 + 1.187 x_6 + 2.720 x_7 + 0.083 x_8 - 1.742 x_9$$

The sign of the coefficients of the estimated logistic function in Table 1 above gives an explanation of the explanatory variables used, as given in Table 2 .

Table 2 : The sign analysis

Covariate	Codes	Sign	Explanation
Age at first Cancer	Quantitative	Positive	Older age increases the probability of second cancer
Gender	1 Male 0 Female	Negative	Male decreases the probability of second cancer
Marital status	1 Married 0 Single	Positive	Married increases the probability of second cancer
Area	1 Urban 0 Rural	Positive	Living in urban increases the probability of second cancer
Treatment by Radiation	1 Yes 0 No	Negative	Radiation decreases the probability of second cancer
Have a Family History	1 Yes 0 No	Positive	family history increases the probability of second cancer
Smoking	1 Yes 0 No	Positive	Smoking increases the probability of second cancer
Obesity	1 Yes 0 Not	Positive	obesity increases the probability of second cancer
Education	1 Educated 0 Illiterate	Negative	Education decreases the probability of second cancer

The odds Ratio Results

The following odds ratios were calculated using the formula;

$$\theta = \frac{odds_1}{odds_2} = \frac{p_1 / (1 - p_1)}{p_2 / (1 - p_2)}$$

For every covariate used in the study, results are given in Table 3.

Table 3: Odds Ratios and 95% Confidence Intervals for Covariates

Variable	Odds ratio	95% Confidence interval
Gender	0.596	{0.184 to 1.923}
Marital Status	3.577	{ 1.189 to 10.756}
Area	1.348	{0 .585 to 3.104}
Radiation	0.270	{0.121 to 0 .603}
Family History	3.278	{ 1.518 to 7.079}
Smoking	15.181	{ 3.538 to 65.141}
Obesity	1.087	{0.419 to 2.819}.
Education	0.229	{0.106 to 0 .497}

From Table 3, it is evident that patients who smoke, patients with family history and married persons are highly susceptible for a second cancer occurrence.

Table 4 gives the classification table. Using the obtained Z function observations are classified as follows, using a prior probability of 0.50

Table 4: Classification Table

Observed	Predicted		
	Have.a.Second.Cancer		Percentage Correct
	Not have	Have	
Have.a.Second.Cancer Not have	80	20	80.0
Have	32	68	68.0
Overall Percentage			74.0

From Table 4, we conclude that;

- a- 80% of all patients not have a second cancer are correctly classified, and 20% are incorrectly classified.
- b- 68% from all patients who have a second cancer are correctly classified, 32% are incorrectly classified.

- c- The overall correct percentage was 74% , which reflects the model's overall explanatory strength

3.2 Model Assessment

The -2 log likelihood for the constant only model obtain by fitting the constant only model was 277.259;and the -2 log likelihood for the overall model was 194.585.

Thus the value of the likelihood ratio test is;

$$G = 277.259 - 194.585 = 82.674$$

And the p-value for the test is $p [\chi^2(9) > 82.674] = 0.00000002$ which is highly significant at the $\alpha < 0.001$ level. The null hypothesis is rejected and we conclude that at least one and perhaps all beta's coefficient are different from zero.

The likelihood ratio tests for all covariates and for each covariate are given in Table 5.

Table 5: likelihood ratio test

Model	-2loglikelihood	G	P-value
Model with constant only	277.259		
Model with all covariates (full model)	194.585	82.674	0.000
Model without family history	204.284	9.699	0.002
Model without smoking	209.602	15.017	0.000
Model without education	209.815	15.230	0.000
Model without Age at first cancer	194.739	0.154	0.695
Model without Treatment by radiation	205.701	11.116	0.001
Model without Gender	195.354	0.769	0.381
Model without Marital status	200.023	5.438	0.020
Model without area	195.080	0.495	0.482
Model without obesity	194.614	0.029	0.865

From table 5 we note that the covariates (family history, smoking, education, treatment by radiation and marital status) are statistically significant; while the covariates (gender, age at first cancer, area and obesity) are statistically non-significant.

The Wald test is obtained by comparing the maximum likelihood estimate of the beta's, $\hat{\beta}_i$, to its standard error. The resulting ratio, under the hypothesis that $\beta_i = 0$ are given in Table 5.

It is evident that the covariates (family history, smoking, education, treatment by radiation and marital status) are statistically significant; while the covariates (gender, age at first cancer, area and obesity) are statistically not-significant.

Stepwise logistic regression analysis is used to reduce number of covariates. results are summarized the results as in table 6.

Table 6: Step-wise Binary Logistic Regression Results

	B	S.E.	Wald	d.f	p-value	Exp(B) Odds ratio
Marital status	1.540	0.432	12.690	1	0.000	4.667
Radiation	-1.250	0.402	9.658	1	0.002	0.286
Family History	1.268	0.378	11.265	1	0.001	3.555
Smoking	2.279	0.503	20.528	1	0.000	9.767
Education	-1.426	0.375	14.447	1	0.000	0.240
Constant	-0.515	0.517	0.992	1	0.319	0.598

And the logit is:

$$Z = -0.515 + 1.540 (\text{Marital status}) - 1.250 (\text{Radiation}) + 1.268 (\text{Family History}) + 2.279 (\text{Smoking}) - 1.426 (\text{Education})$$

The Logit (Z) above indicates that: married patients are more susceptible to develop a second cancer ; treatment by radiation decreases the susceptibility; a patient with family history is more susceptible to develop second cancer; smokers are more susceptible than non-smokers, and educated patients are less susceptible to develop a second cancer.

The exponent (Exp (B)) in Table 6 is the odds ratio, thus:

1. The odds for married patients to single patients to develop second cancer are 4.667.
2. The odds for patients with family history to patients with no family history to develop second cancer is 3.55.
3. The odds for smokers to nonsmokers to develop second cancer is 9.76.

Table 7 gives the classification table. Using the obtained Z function observations are classified as follows, using a prior probability of 0.50.

Table 7 : Classification Table

Observed		Predicted		
		Have.a.Second.Cancer		Percentage Correct
		Not have	Have	
Have.a.Second.Cancer	Not have	82	18	82.0
	Have	33	67	67.0
Overall Percentage				74.5

- a- 82% of all patients not have a second cancer are correctly classified, and 20% are incorrectly classified.
- b- 67% from all patients who have a second cancer are correctly classified, 32% are incorrectly classified.
- c- The overall correct percentage was 74.5%, which reflects the model's overall explanatory strength.

The value of the Hosmer – Lemeshow goodness-of-fit statistic computed for the full model was $C = 4.060$ and the corresponding p-value computed from the chi-square distribution with 8 degree of freedom is 0.852 this indicates that the model seems to fit quite well.

Cross Validation Results

Using Efron (1983) leave-one-out Cross Validation goodness-of-fit statistic the results for the full model was (using prior probability of .50) summarized in table 8 .

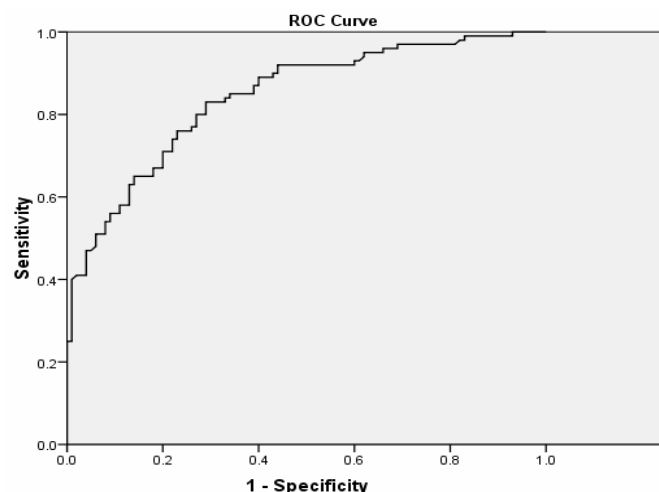
Table 8: Cross Validation Results;

Table 3: Cross Validation Results,				
Group	Actual no. of cases	Predicting group		%Correctly classified
		membership		
		Having a second cancer	Not having a second cancer	
Having a second cancer	100	67	33	72.5 %
Not having a second cancer	100	22	78	

The classification matrix shows the accuracy of second cancer occurrence prediction in the cross validation leave-one-out sample as presented in table 8 above. In this sample of 200 patients, actual occurrence of second cancer was 50%. Of the 100 patients that 67 or 67% were correctly classified into group having a second cancer. Of the 100 patient that not having a second cancer, 78 or 78% were correctly classified into group not having a second cancer. The total correctly classified was 145 of 200 or 72.5%. The maximum chance Criterion is 50% and the proportional chance criterion is 50% also. Because the percentage correctly classified was 72.5% (22.5% greater than proportional chance), Z test evident that difference are statistically significant (p-value <0.001).

Using ROC curve for the classification accuracy, it is found that the area under the ROC curve, which ranges from zero to one, provides a measure of the model's ability to discriminate between those subjects who experience the response of interest versus those who do not.

Plotting sensitivity versus (1 – specificity) over all possible cut-points is shown in the Figure below .The area under the ROC curve for the full model was is 0.843 this is considered excellent discrimination



4. Summary and Conclusions

In this study, social-demographic risk factors of developing a second primary cancer using logistic regression model were studied. The social-demographic risk factors used are age at first

cancer, gender, area the patient live in, marital status, family history, smoking, education and obesity in addition to treatment by radiation. The binary logistic regression model is used to estimate the probability of having second primary cancer. Significance testing for the logistic coefficients using Wald test and likelihood ratio show that smoking, family history, marital status, and education are the significant factors. The odds ratio for each covariate compare whether the probability of having a second primary cancer is the same for each covariate groups. The odds ratio for smokers to non-smokers ranges between 3 times to 65 times with confidence 95%. To assess the fitness of the model the maximum likelihood test and Hosmer and Lemeshow test are used. The logistic regression model proved to have a lower sensitivity level due to some other clinical risk factors not considered in this study.

The study concludes that: *married patients are more susceptible to develop a second cancer; treatment by radiation decreases the susceptibility; a patient with family history is more susceptible to develop second cancer; smokers are more susceptible than non-smokers, and educated patients are less susceptible to develop a second cancer.*

The researcher recommends the following:

- 1- Replicate the same study with an increased sample size.
- 2- Develop a logistic regression model that contains repeated measures.
- 3- Replicate the same study to include repeated measures on the same patients, especially when some demographic factors change, and age develops.
- 4- Use the reached significant factors and add more clinical risk factors which was not available at the hospitals records when the research was conducted.
- 5- Apply Classification and Regression Tree (CART) and compare the results with the binary logistic regression model.

References

1. Afifi, A., Clark, V. A., and May, S. (2004). Computer- Aided Multivariate Analysis. Fourth Edition, Champman and Hall/CRC.

2. Agresti, A. (2007). An Introduction to Categorical Data Analysis. Second Edition, Wiley, Inc., New York.
3. Agresti, A. (2002). Categorical Data Analysis. Second Edition, Wiley, Inc., New York.
4. Altman, D. G. (1991). Practical statistics for medical research. Champman and Hall, London.
5. Ashour, S., and Abo Elfotouh, S. (2005). Presentation and statistical analysis using SPSSWIN. Second Part, Advanced Applied Statistics, Institute of Statistical Studies and Research. Cairo university, Egypt (in Arabic).
6. Bamber, D. (1975) .**The area above the ordinal dominance graph and the area below the receiver operating characteristic graph.** Journal of mathematical psychology, 12, pp. 387-415.
7. Barone S., Lombardo A. and Tarantino P.. (2007). A weighted logistic regression for conjoint analysis and Kansei engineering. Quality and Reliability Engineering International, Vol. 23 Issue 6, pp. 689 – 706, John Wiley & Sons, Ltd.
8. Bradley, A. P. (1997). The use of the area under the roc curve in the evaluation of machine learning algorithms. Pattern Recognition. Jul; 30(7): pp.1145-59.
9. Chernick ,M., R. (2008). Bootstrap Methods: A Guide for Practitioners and Researchers. Second Edition, Wiley, Inc., New York.
10. Christensen, R. (1997) Log-Linear models and logistic regression. Second Edition. Springer-verlag. New York.
11. DeVlta, Hellman, and Rosenberg's. (2008). Cancer Principles and practice of oncology. Eighth edition, vol. 6, wolters kluwer, lippincott Williams&wilkins.
12. Efron, B. (1983). Estimating the error rate of a prediction rule: improvements on cross validation. Journal of the American Statistical Association, 78, pp. 316-331.

13. Erhart, M., Hagquist C., Auquier P., Rajmil L., Power M., Ravens-Sieberer U. and the European KIDSCREEN Group. (2009). A comparison of Rasch item-fit and Cronbach's alpha item reduction analysis for the development of a Quality of Life scale for children and adolescents. *Child: Care, Health and Development*, Blackwell Publishing Ltd.
14. Evans, M. Hastings, N. and Peacock, B. (1993). *Statistical Distributions*. Second Edition, Wiley, New York.
15. Everitt, B. S. (1998). *The Cambridge Dictionary of Statistics*. Cambridge University Press.
16. Ferri C., Flach P., Hernandez-Orallo J. (2002). Learning Decision Trees Using the Area under the ROC Curve. Nineteenth International Conference on Machine Learning (ICML 2002); Morgan Kaufmann; pp. 46-139.
17. Fienberg, S. E (1980). *The analysis of cross-classified categorical data*. Second Edition, The MIT Press, Cambridge, Massachusetts.
18. Fleiss, J., L. (1981). *Statistical Methods For Rates And Proportions*. Second Edition. John Wiley & Sons, Inc.
19. Garcia-Ramirez M., Martinez, M., F. M., Balcazar F., E., Suarez-Balcazar Y., Albar M., Domínguez E. and Santolaya, F.,J. (2005). Psychosocial empowerment and social support factors associated with the employment status of immigrant welfare recipients. *Journal of Community Psychology*, Volume 33 Issue 6, Pages 673 – 690, Wiley Periodicals, Inc., A Wiley Company.
20. Hair, J. F., Anderson, R. E., Babin, B. J., and Black, W. C.(2009) *Multivariate Data Analysis*. Seventh Edition. Maxwell Macmillan International, New York.
21. Hanley, J.A. and McNeil, B., J. (1982). The meaning and the use of the Area under a receiver operating characteristic curve (Roc).' *Radiology*, 143, pp. 29-36.
22. Hauck, W.W., and Donner, A. (1977). Wlad's test as applied to hypotheses in logit analysis. *Journal of the American Statistical Association*, 72, pp.851-853.

23. Hosmer, D. W., Lemeshow, S. (2000). *Applied Logistic Regression*, Second Edition, Wiley, Inc., New York.
24. Ingles, C., J.; Garcia-Fernandez, j., M.; Castejon, J., L. ; Valle Antonio, B., D., and Marzo, J., C. (2009). Reliability and validity evidence of scores on the Achievement Goal Tendencies Questionnaire in a sample of Spanish students of compulsory secondary education. *Psychology in the Schools*, Vol. 46 Issue 10, pp. 1048 – 1060, Wiley Periodicals, Inc., A Wiley Company
25. Jennings, D.E. (1986a). Judging inference adequacy in logistic regression. *Journal of the American Statistical Association*, 81, pp. 471-476.
26. Kaufman, E., L., Jacobson, J.; S. Hershman, D., L.; Desai M., and Neugut, A., I. (2008). Effect of Breast Cancer Radiotherapy and Cigarette Smoking on Risk of Second Primary Lung Cancer. *Journal of clinical oncology*, 26(3): pp. 392-398.
27. King G. and Zeng L. (2002). Estimating risk and rate levels, ratios and differences in case-control studies. *Statistics in Medicine*, Vol. 21 Issue 10, pp. 1409 – 1427, John Wiley & Sons, Ltd.
28. Kirkos E., Spathis C., and Manolopoulos Y. (2009). Audit-firm group appointment: an artificial intelligence approach. Article on line in advance of print, *Intelligent Systems in Accounting, Finance & Management*, John Wiley & Sons, Ltd.
29. Kleijnen M., De Ruyter K., Wetzels M. (2004). Consumer adoption of wireless services: Discovering the rules, while playing the game. *Consumer adoption of wireless services: Discovering the rules, while playing the game. Journal of Interactive Marketing*, Vol. 18 Issue 2, pp. 51 – 61, Wiley Periodicals, Inc., A Wiley Company, and Direct Marketing Educational Foundation, Inc.
30. Marcoulides, A. George., and Hershberger, L. Scott. (1997). *Multivariate statistical methods: A first course*. Lawrence Erlbaum Associates, Mahwah, New Jersey.
31. Neagu R. and Hoerl R. (2005). A Six Sigma Approach to Predicting Corporate Defaults. *Quality and Reliability Engineering International*. Vol. 21 Issue 3, pp. 293-309, John Wiley & Sons, Ltd.

32. O'Leary, D., E. (2009). Downloads and citations in Intelligent Systems in Accounting. Finance and Management. Intelligent Systems in Accounting, Finance & Management, Vol. 16 Issue 1-2, pp. 21 – 31, John Wiley & Sons, Ltd.
33. Rubino, C., De Vathaire, F., Shamsaldin, A., Labbe, M., and le M. G. (2003). Radiation dose, chemotherapy, hormonal treatment and risk of second cancer after breast cancer treatment. British Journal of Cancer; 89(5): pp. 840–846.
34. Saijo, Y., Ueno T., Hashimoto, Y. (2008). Twenty-four-hour shift work, depressive symptoms, and job dissatisfaction among Japanese firefighters. American Journal of Industrial Medicine, Vol. 51 Issue 5, pp. 380 – 391. Wiley-Liss, Inc., A Wiley Company.
35. Sallis, J., E. and Deo Sharma, D. (2009). Knowledge seeking in going abroad. Thunderbird International Business Review, Vol. 51 Issue 5, pp. 441 – 456, Wiley Periodicals, Inc., A Wiley Company.
36. Sanchez, L., A.; Lana, A. B. ; Hidalgo, A. A.; Rodriguez, M., J. C.; Del Valle, M., D.; Cueto, A., b; Folgueras, M., C.; Belyakova, E., C.; Comendador, M. D.; Lopez, M., L. (2008). Risk factors for second primary tumours in breast cancer survivors. European Journal of Cancer Prevention. 17(5): pp. 406-413.
37. The SAS System (1995) .Logistic regression examples using the SAS system. Version 6, First Edition. SAS institute Inc., Cary. NC, USA.
38. Weber, S., O.; and Michalik G., W., (2008). Incorporating sustainability criteria into credit risk management. Business Strategy and the Environment, Vol. 19 Issue 1, pp. 39 – 50, John Wiley & Sons, Ltd. and ERP Environment.

Acknowledgments

The researcher would like to thank Mr. Mahmoud R. Nooh and the officials of the radiotherapy department; Faculty of the Medicine, Ain shams university hospitals for providing data for this research paper.