

RESEARCH ARTICLE

Application of Data Mining Techniques to Explore Predictors of HCC in Egyptian Patients with HCV-related Chronic Liver Disease

Dalia Abd El Hamid Omran^{1*}, AbuBakr Hussein Awad², Mahasen Abd El Rahman Mabrouk¹, Ahmad Fouad Soliman¹, Ashraf Omar Abdel Aziz¹

Abstract

Background: Hepatocellular carcinoma (HCC) is the second most common malignancy in Egypt. Data mining is a method of predictive analysis which can explore tremendous volumes of information to discover hidden patterns and relationships. Our aim here was to develop a non-invasive algorithm for prediction of HCC. Such an algorithm should be economical, reliable, easy to apply and acceptable by domain experts. **Methods:** This cross-sectional study enrolled 315 patients with hepatitis C virus (HCV) related chronic liver disease (CLD); 135 HCC, 116 cirrhotic patients without HCC and 64 patients with chronic hepatitis C. Using data mining analysis, we constructed a decision tree learning algorithm to predict HCC. **Results:** The decision tree algorithm was able to predict HCC with recall (sensitivity) of 83.5% and precision (specificity) of 83.3% using only routine data. The correctly classified instances were 259 (82.2%), and the incorrectly classified instances were 56 (17.8%). Out of 29 attributes, serum alpha fetoprotein (AFP), with an optimal cutoff value of ≥ 50.3 ng/ml was selected as the best predictor of HCC. To a lesser extent, male sex, presence of cirrhosis, $AST > 64$ U/L, and ascites were variables associated with HCC. **Conclusion:** Data mining analysis allows discovery of hidden patterns and enables the development of models to predict HCC, utilizing routine data as an alternative to CT and liver biopsy. This study has highlighted a new cutoff for AFP (≥ 50.3 ng/ml). Presence of a score of > 2 risk variables (out of 5) can successfully predict HCC with a sensitivity of 96% and specificity of 82%.

Keywords: HCC - HCV-related chronic liver disease - data mining - decision tree - prediction - AFP

Asian Pac J Cancer Prev, 16 (1), 381-385

Introduction

Hepatocellular carcinoma (HCC) is the most common primary malignant tumor of the liver and the fifth most frequent malignant tumor in the world (third in terms of mortality) (Parkin et al., 2001). An association was observed between the occurrence of HCC and viral hepatitis (either HBV or HCV) (Saeed et al., 2012). Studies had demonstrated that unlike other countries of the Middle East, the attributable fraction of HCC due to hepatitis C virus (HCV) is quite high in Egypt (Ezzat et al., 2005). Egypt has the highest prevalence of hepatitis C virus (HCV) in the world, ranging from 15% to 25% in rural communities (Strickland et al., 2002; Kurosaki M et al., 2010). This high prevalence in a country with a population census of about ninety million results in a massive number of HCV-infected patients (12-22 million). This, in turn, results in massive numbers of HCV related HCC. All HCV related liver cirrhosis patients are recommended for HCC surveillance according to

European Association for the Study of Liver Diseases (EASL), the American Association for the Study of Liver Diseases (AASLD) and Asian Pacific Association for the Study of Liver Diseases (APASL) (Bruix et al., 2001; Omata et al., 2010; Bruix et al., 2011).

Regular HCC surveillance by ultrasound examination (US) of the liver and/or alpha fetoprotein (AFP) may identify early HCC and reduce mortality (Wong et al., 2008). Given the constrained economy and lack of resources in Egypt, such large-scale screening is almost impossible. Moreover, repeated examinations are costly for the patients and reduce their compliance. In clinical practice, it had been reported that over 60% of HCCs were diagnosed at late stages, suggesting failures in the surveillance. This failure is attributed mainly to poor patients' compliance and to failure to detect small lesions (Altekruse et al., 2005).

These findings emphasize the need for a simple, cheap and noninvasive method that can predict HCC. Identification of the risk factors associated with HCC

¹Endemic Medicine and Hepatology Department, Faculty of Medicine, ²Computer Science Department, Faculty of Computers and Information, Cairo University, Cairo, Egypt *For correspondence: daliaomran2007@yahoo.com; daliaomran@kasralainy.edu.eg

development in HCV related chronic liver disease (CLD) is essential for formulating personalized surveillance programs. Decision-tree analysis is a core component of data mining analysis that can be used to build predictive models (Breiman et al., 1980). The major advantage of decision-tree analysis over logistic regression analysis is that the results of analysis are easy to understand. The simple allocation of patients into subgroups by following the flowchart form could define the predicted possibility of outcome (LeBlanc and Crowley, 1995)

This method had been used to define prognostic factors in various diseases such as prostate cancer (Garzotto et al., 2005), diabetes (Miyaki et al., 2002), melanoma (Averbook et al., 2002), colorectal carcinoma (Valera et al., 2007), liver failure (Baquerizo et al., 2003) and for the prediction of virological response in HCV patients (Salim, 2009). The objective of the current study is to develop an economic, reliable mathematical model to predict HCC in HCV related CLD patients using routine workup. In areas with limited resources like Egypt, it is wise to restrict the semiannual surveillance by ultrasound scan to risky patients with HCV-related liver cirrhosis. This restriction will eventually reduce unnecessary costs caused by screening all cirrhotic patients. We believe that the field of data mining can be used to solve real health problems that are currently facing Egypt with great success.

Materials and Methods

This cross sectional retrospective study enrolled 315 HCV related chronic liver disease (CLD) patients of both sexes, recruited from Endemic Medicine Department, Cairo University Hospital. Informed consent was obtained from each patient according to the 1975 Helsinki Declaration and the study was approved by Cairo University ethical committee.

One hundred thirty five (135) patients were diagnosed to have HCC according to the criteria of the European association for the study of the liver (EASL) (Bruix et al., 2001). One hundred and sixteen (116) patients were diagnosed to have liver cirrhosis on the basis of clinical, biochemical, and ultrasound findings. Sixty four (64) patients were diagnosed to have chronic hepatitis C. HCV infection was diagnosed by anti-HCV antibodies, HCV-RNA (Cobas Amplicor HCV Monitor v 2.0, Roche Diagnostic systems, CA).

Data Collection, Feature selection and reduction

A subset of 29 features including routine laboratory workup (categorical or numerical) was used for the model building process (Table 1). The dataset was created containing two demographic variables (age, gender), three hematological variables (hemoglobin, white blood cells, and platelets), eight biochemical variables (total bilirubin, albumin, AST, ALT, ALP, INR, creatinine and AFP), viral markers (for HBV and HCV) Child class status (A, B or C), in addition to clinical examination for the presence of splenomegaly and ascites.

A number of data transformation techniques have been used to format and prepare the patient records to be

Table 1. Summary of Features (Attributes) Included in The Study

Attribute	Represented As
Gender	Categorical (Male, Female)
Age	Numeric
Smoking (history/ongoing)	Categorical (Negative, Positive)
Alcohol intake (history/ongoing)	Categorical (Negative, Positive)
History of operations	Categorical (Negative, Positive)
History of tarter emetic intake	Categorical (Negative, Positive)
History of blood transfusion	Categorical (Negative, Positive)
History of dental procedures	Categorical (Negative, Positive)
Diabetes mellitus	Categorical (Negative, Positive)
Hypertension	Categorical (Negative, Positive)
Abdominal pain	Categorical (Negative, Positive)
Loss of weight and appetite	Categorical (Negative, Positive)
WBC	Numeric
Hb	Numeric
PLT	Numeric
Albumin	Numeric
ALP	Numeric
AST	Numeric
ALT	Numeric
Total Bilirubin	Numeric
AFP	Numeric
INR	Numeric
Creatinine	Numeric
Anti HCV Ab	Categorical (Negative, Positive)
HBsAg	Categorical (Negative, Positive)
HB core Ab	Categorical (Negative, Positive)
Child class	Categorical (A, B or C)
Spleen size	Categorical (normal, enlarged)
Ascites	Categorical (absent, present)

WBC, white blood cell count; Hb hemoglobin; PLT, platelet count; ALP, alkaline phosphatase; AST, aspartate aminotransferase; ALT, alanine aminotransferase; AFP, alpha fetoprotein; INR, international normalized ratio; Anti HCV Ab, anti hepatitis C antibodies; HBsAg, hepatitis B surface antigen; HB core Ab, hepatitis B core antibodies

processed by the learning algorithms.

Data mining

Using the data mining analysis, we constructed a decision tree learning algorithm C4.5 (weka J48). The C4.5 which was published by Ross Quinlan in 1993 is an example of commonly used decision trees of high accuracy in medical classification, which can handle both categorical and numerical data. The data set was evaluated to determine which variables can yield the most significant diagnosis of HCC. Internal validation was performed with test mode: 10-fold cross-validation which is a technique for assessing how the results of a statistical analysis will generalize to an independent data set.

Statistical Analysis

Patients were categorized into HCC and non HCC.

Qualitative variables were expressed by number, percent and compared by chi square or fisher's exact test.

Quantitative variables were expressed by mean and standard deviation (SD) and compared by t student.

Optimum cutoff values for serum AFP, serum AST, age were determined by data mining analysis. Sensitivity, specificity, PPV, NPV and accuracy were calculated subsequently.

Table 2. Baseline Characteristics of Patients

Variables	All patients	HCC	No HCC
Number of patients. N (%)	315	135 (42.9%)	180 (57.1%)
Gender (M/F)	222/93	112/23	110/70
Age (mean±SD)	51.5±11.4	55.96±7.69	47.6±12.6
AST (IU/ml) (mean±SD)	84.37±87.7	128.5±124.69	62.85±50.2
ALT (IU/ml) (mean±SD)	53.68±48.1	73.46±59.85	44.1±37.87
Total Bilirubin (mg/dl) (mean±SD)	2.67±3.4	3.78±4.23	1.94±2.47
Albumin (gm/dl) (mean±SD)	3.19±0.95	2.748±0.69	3.46±0.98
AFP (ng/dl) (mean±SD)	3134.66±1.38	5450.68±1.78	8.05±9.544
Spleen size (normal/Enlarged)	(119/196)	14/112	105/84
Ascites (Absent/present)	(119/150)	19/80	100/70

AST, aspartate aminotransferase; ALT, alanine aminotransferase; AFP, alpha fetoprotein

Table 3. Accuracy of The Studied Variables to Predict HCC

Variables	OR	95% CI of OR	Sensitivity	Specificity	PPV	NPV	Accuracy
AFP >50.3 ng/dl	252	34-1877	72%	99%	99%	72%	83%
AST >64 IU/ml	9.5	4.7-19.2	86%	61%	52%	90%	69%
Cirrhosis	2.2	1.9-2.5	100%	36%	54%	100%	86%
Ascites	6	3.3-10.8	81%	59%	53%	84%	60%
Male gender	2.9	1.7-5.1	83%	38%	53%	73%	50%
Number of risk variables >2	103.4	22.6-473	96%	82%	76%	97%	87%

OR, odds ratio; 95% CI, 95% confidence interval; PPV, positive predictive value; NPV, negative predictive value; AFP, alpha fetoprotein ; AST, aspartate aminotransferase

Results

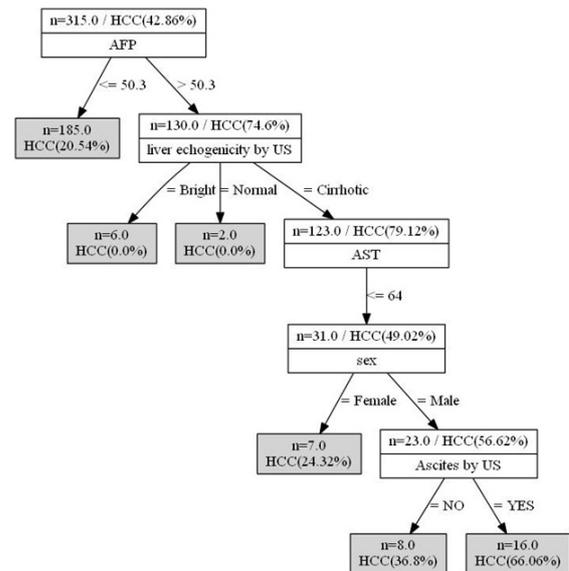
Our study included 315 patients suffering HCV related CLD (222 males, 93 females). Their mean age: 51.5±11.4 years. Baseline characteristics of patients are summarized in Table 2.

Decision tree algorithm was able to diagnose HCC with recall (sensitivity) 83.5% and Precision (specificity) 83.3% using only routine data. The correctly classified instances were 295 (82.2%) and the incorrectly classified instances were 56 (17.8%). Out of 29 attributes, serum AFP, with an optimal cutoff value of ≥ 50.3 ng/ml was selected as the best predictor (most decisive) of HCC according to the decision-tree models. To a lesser extent, male sex, presence of cirrhosis, AST >64U/L, and ascites were variables associated with HCC (Figure 1).

AFP was found to be the accurate single predictor of HCC (Table 3). Moreover, it was found that the presence of more than 2 of the studied five variables (i.e. having score >2) was associated with an increased risk for HCC development by 103.4 times and can successfully predict HCC with a sensitivity of 96% and specificity of 82%.

Discussion

HCC reduces quality of life and causes death within 6 months -1 year from the diagnosis (Bosch et al., 2005). In Egypt, there was a dramatic increase in the number of HCC cases in the last few years. The registry

**Figure 1. Decision Tree Algorithm to Predict HCC**

of the pathology department, National Cancer Institute (NCI), Cairo (2003-2004) showed that HCC was the 2nd malignancy in males after carcinoma of the urinary bladder and the 4th in females (Mokhtar et al., 2007) The stage of cancer dictates the therapeutic choice, making early detection a primary objective. Surveillance of HCC aims at detection of small tumors for curative treatment, which may be translated to improved patient survival. Regular screening of large number of cirrhotic patients have a high cost impact and may add burden to a country like Egypt having the highest prevalence of HCV worldwide (Frank et al., 2000; Arguedas et al., 2003).

Data mining analysis has been integrated into bioinformatics research in order to explore hidden patterns in large datasets and thus can be used to make prediction models or certain hypotheses (Han and Kamber, 2006). On one hand, conventional statistics can examine certain hypothesis while on the other hand; data mining analysis can set an algorithm by using a large amount of data. Decision tree model is considered to be rather superior over the traditional regression models as it can be readily interpreted by medical professionals simply by following the flowchart form without any specific knowledge of statistics (Witten and Frank, 2005).

In the current study, a decision tree model based on routinely available clinical and laboratory parameters was constructed for HCC prediction in patients with HCV related CLD. Serum AFP ≥ 50.3 ng/ml was the best predictor of HCC and to a lesser extent male sex, presence

of cirrhosis, AST >64 U/L, and ascites were variables associated with HCC.

Serum AFP, a well-known noninvasive marker for the development of HCC (Daniele et al., 2004), was the first split variable (most decisive) in the predictive model for HCC and was significantly associated with HCC development in the multivariate analysis as well. HCV related CLD patients with AFP serum level of 50.3 ng/ml or more are 252 times more liable to develop HCC. Our study proposed a cutoff of 50.3 ng/ml with a sensitivity of 72%, specificity of 99%, a positive predictive value of 99% and a negative predictive value of 72%. The AUROC was 0.833.

In previously published studies, AFP (with different cutoffs) had a sensitivity of 39%-65%, a specificity of 76%-94%, and a positive predictive value of 9%-50% for the presence of HCC (Franca et al., 2004.). The variation in sensitivity and specificity of AFP in the studies performed may be due to the diversity of patient populations examined, varying study designs and differing cut-off values for normality. There was a debate in defining the AFP cut-off level for the diagnosis of HCC. An AFP value >400 ng/mL had been considered to be diagnostic for HCC in cirrhotic patients. However, such a cut-off value is problematic in absolute diagnostic terms, since such high levels are not common in the presence of small tumors (<5 cm) and only 30% of HCC patients have levels higher than 100 ng/mL, furthermore, up to 20% patients with HCC do not produce AFP (Tao et al., 2010).

Recently, the American Association for the Study of Liver Diseases (AASLD) Practice Guidelines Committee recommended that ultrasound (US) examination alone (without AFP) should be used for HCC surveillance (Bruix et al., 2011). However, the interpretation of ultrasound is operator-dependent and can be difficult in obese persons; moreover, it cannot differentiate malignant from benign nodules in the small cirrhotic liver, and the detection of small HCC in a cirrhotic liver by US is much more difficult than the detection of metastases in a normal liver, owing to disturbed parenchymal architecture (Saar and Kellner-Weldon, 2008). An effectiveness study recently demonstrated that ultrasound only had a sensitivity of 32% for early stage tumors, which was significantly increased to 63% when used in combination with AFP (Singalat et al., 2012).

According to our study, male patients are 2.9 times more liable to develop HCC. Previous studies reported that primary liver cancer is more prevalent among men than women. The gender-specific age-adjusted incident rate ratio ranges from 1.3 to 3.6 worldwide (Ferlay et al., 2001). Male sex had been consistently shown to increase the risk for HCC development in HCV-infected persons (Degos et al., 2000). Androgens and androgen receptors had been suggested to induce and promote HCC (Strickland et al., 2002).

In the current study, liver cirrhosis increases the risk of HCC by 2.2 times. It is well known that HCC development is restricted largely to patients with cirrhosis or advanced fibrosis (Shiratori et al., 1995; El-Serag, 2000).

Presence of ascites in our study was associated with an increased risk of HCC development. It is well known

that ascites signifies advanced liver disease. Shiratori et al (1995) analyzed the characteristics of 205 HCC cases from Japan and noted that HCV-associated HCC occurred in the presence of more severe liver disease than in hepatitis B virus (HBV)-associated HCC.

According to our study, having AST serum level >64 IU/L increase the risk of HCC by 9.5 times. High AST serum levels reflect severe hepatic fibrosis, portal inflammation and piecemeal necrosis that may eventually progress to HCC (Saar and Kellner-Weldon, 2008).

Identifying patients at risk of progressing to hepatocellular carcinoma may help in justifying health resource allocation by targeting high risk patients. In areas with constrained economy like Egypt, it is wise to restrict the HCC surveillance to patients having >2 risk variables (out of the five risk variables proposed by our decision-tree model). This restriction will eventually reduce unnecessary costs caused by screening all cirrhotic patients. This is also useful in the rural areas of Egypt where CT facilities are not readily present. Depending on these risk variables, physicians can easily identify high risk patients and refer them to specialized tertiary care centers. Limitation of our study included the small number of patients and the lack of evaluation of tumor markers other than AFP.

Conclusion: Data mining analysis explores data to discover hidden patterns, trends and enables the development of models to diagnose HCC utilizing simple laboratory data, without imposing extra costs for additional examinations. This study was the first one (up to our knowledge) to highlight a new cutoff value of AFP for diagnosis of HCC (≥ 50.3 ng/ml). This low cutoff together with other unfavorable (risk) variables when used in combination can help in early diagnosis of HCC. The field of data mining can be used to solve real health problems that Egypt is currently facing with great success. More studies are needed exploring more variables that may be associated with progression of HCV related CLD to HCC. Identification of risk factors associated with HCC development will result in better targeting of patient, thus having the utmost benefit from HCC surveillance programs with the least possible cost by restricting screening to risky patient only.

Acknowledgements

The authors would like to thank Dr Wafaa Al Akel for her assistance in statistical analysis of the data.

References

- Altekruse SF, McGlynn KA, Reichman ME (2009) Hepatocellular carcinoma incidence, mortality, and survival trends in the United States from 1975 to 2005. *J Clin Oncol*, **27**, 1485-91.
- Arguedas MR, Chen VK, Eloubeidi MA, Fallon MB (2003). Screening for hepatocellular carcinoma in patients with hepatitis C cirrhosis: a cost-utility analysis. *Am J Gastroenterol*, **98**, 679-90.
- Averbook BJ, Fu P, Rao JS, Mansour EG (2002) A long-term analysis of 1018 patients with melanoma by classic Cox regression and tree structured survival analysis at a major referral center: Implications on the future of cancer staging.

- Surgery*, **132**, 589-602.
- Baquerizo A1, Anselmo D, Shackleton C, et al (2003) Phosphorus as an early predictive factor in patients with acute liver failure. *Transplantation*, **75**, 2007-14.
- Bosch FX, Ribes J, Cleries R, Diaz M (2005). Epidemiology of hepatocellular carcinoma. *Clin Liver Dis*, **9**, 191-211.
- Breiman L JH, Friedman RA, Olshen CJ, Stone CM (1980). Classification and regression trees. CA: Wadsworth. ???
- Bruix J1, Sherman M, Llovet JM, et al (2001). Clinical management of hepatocellular carcinoma. Conclusions of the Barcelona-2000 EASL conference. European Association for the Study of the Liver. *J Hepatol*, **35**, 421-30.
- Bruix J, Sherman M. (2011). American Association for the Study of Liver Diseases. Management of hepatocellular carcinoma: an update. *Hepatology*, **53**, 1020-2.
- Daniele B, Bencivenga A, Megna AS, Tinessa V (2004). Alpha fetoprotein and ultrasonography screening for hepatocellular carcinoma. *Gastroenterology*, **127**, 108-12
- Degos F1, Christidis C, Ganne-Carrie N, et al (2000). Hepatitis C virus related cirrhosis: time to occurrence of hepatocellular carcinoma and death. *Gut*, **47**, 131-6.
- El-Serag HB (2002). Hepatocellular Carcinoma and Hepatitis C in the United States. *Hepatology*, **36**, 74-83
- Ezzat S, Abdel-Hamid M, Eissa SA, et al (2005). Associations of pesticides, HCV, HBV, and hepatocellular carcinoma in Egypt. *Int J Hyg Environ Health*, **208**, 329-39
- Ferlay J BF, Pisani P, Globocan (2000). Cancer incidence, mortality and prevalence worldwide. IARC Cancer Base No. 5, IARC Nonserial Publication, 2001.
- França AV, Elias Junior J, Lima BL, et al (2004). Diagnosis, staging and treatment of hepatocellular carcinoma. *Braz J Med Biol Res*, **37**, 1689-1705
- Frank C1, Mohamed MK, Strickland GT, et al (2000). The role of parenteral antischistosomal therapy in the spread of hepatitis C virus in Egypt. *Lancet*, **355**, 887-91.
- Garzotto M, Beer TM, Hudson RG, et al (2005). Improved detection of prostate cancer using classification and regression tree analysis. *J Clin Oncol*, **23**, 4322-9.
- Han J, Kamber M (2006). Data mining: concepts and techniques. San Francisco San, CA: Morgan Kaufmann Publishers. , 550.
- Kurosaki M, Matsunaga K, Hirayama I, et al (2010). predictive model of response to peginterferon ribavirin in chronic hepatitis C using classification and regression tree analysis. *Hepatol Res*, **40**, 251-60.
- LeBlanc M, Crowley J. (1995). A review of tree-based prognostic models. *Cancer Treat Res*, **75**, 113-24.
- Miyaki K, Takei I, Watanabe K, Nakashima H, Omae K (2002). Novel statistical classification model of type 2 diabetes mellitus patients for tailor-made prevention using data mining algorithm. *J Epidemiol*, **12**, 243-8.
- Mokhtar N, Gouda I, Adel I (2007). Cancer Pathology Registry and time trend analysis (2003-2004). PP.76
- Omata M, Lesmana LA, Tateishi R, et al (2010). Asian Pacific Association for the Study of the Liver consensus recommendations on hepatocellular carcinoma. *Hepatol Int*, **4**, 439-74
- Parkin DM, Bray F, Ferlay J, Pisani P (2001). Estimating the world cancer burden: Globocan 2000. *Int J Cancer*, **94**, 153-6.
- Salim EI1, Moore MA, Al-Lawati JA, et al (2009). Cancer epidemiology and control in the Arab world - past, present and future. *Asian Pac J Cancer Prev*, **10**, 3-16.
- Saar B, Kellner-Weldon F (2008). Radiological diagnosis of hepatocellular carcinoma. *Liver Int*, **28**, 189-99.
- Saeed NM, Bawazir AA, Al-Zuraiqi M, Al-Negri F, Yunus F (2012). Why is hepatocellular carcinoma less attributable to viral hepatitis in Yemen? *Asian Pac J Cancer Prev*, **13**, 3663-7.
- Shiratori Y, Shiina S, Imamura M, et al (1995). Characteristic difference of hepatocellular carcinoma between hepatitis B- and C- viral infection in Japan. *Hepatology*, **22**, 1027-33.
- Singal AG, Conjeevaram HS, Volk ML, et al (2012). Effectiveness of hepatocellular carcinoma surveillance in patients with cirrhosis. *Cancer Epidemiol Biomarkers Prev*, **21**, 793-9.
- Strickland GT, Elhefni H, Salman T, et al (2002). Role of hepatitis C infection in chronic liver disease in Egypt. *Am J Trop Med Hyg*, **67**, 436-42.
- Tao LY, Cai L, He XD, Liu W, Qu Q (2010). Comparison of serum tumor markers for intrahepatic cholangiocarcinoma and hepatocellular carcinoma. *Am Surg*, **76**, 1210-3.
- Valera VA1, Walter BA, Yokoyama N, et al (2007). Prognostic groups in colorectal carcinoma patients based on tumor cell proliferation and classification and regression tree (CART) survival analysis. *Ann Surg Oncol*, **14**, 34-40.
- Witten IH, Frank E (2005). Data mining: practical machine learning tools and techniques. San Francisco, CA: Morgan Kaufmann Publishers, p 525.
- Wong GL1, Wong VW, Tan GM, et al (2008). Surveillance programme for hepatocellular carcinoma improves the survival of patients with chronic viral hepatitis. *Liver Int*, **28**, 79-87