ARTICLE IN PR

Clinical Biochemistry xxx (2015) xxx-xxx



Contents lists available at ScienceDirect

Clinical Biochemistry





journal homepage: www.elsevier.com/locate/clinbiochem

Improve discrimination power of serum markers for diagnosis of cholangiocarcinoma using data mining-based approach 2

Sirorat Pattanapairoj^{a,1}, Atit Silsirivanit^{b,e,1}, Kanha Muisuk^{d,e}, Wunchana Seubwai^{d,e}, Ubon Cha'on^{c,e}, 02

Kulthida Vaeteewoottacharn ^{c,e}, Kanlayanee Sawanyawisuth ^{c,e}, Danaipong Chetchotsak ^{a,b,*}, Sopit Wongkham ^{c,e,**} 4

5

^a Department of Industrial Engineering, Faculty of Engineering, Khon Kaen University, Khon Kaen, Thailand

^b System Modeling for Industry Research Group, Faculty of Medicine, Khon Kaen University, Khon Kaen, Thailand

^c Department of Biochemistry, Faculty of Medicine, Khon Kaen University, Khon Kaen, Thailand

^d Department of Forensic Medicine, Faculty of Medicine, Khon Kaen University, Khon Kaen, Thailand 9

10 e Liver Fluke and Cholangiocarcinoma Research Center, Faculty of Medicine, Khon Kaen University, Khon Kaen, Thailand

ARTICLE INFO 1 1

12Article history:

- Received 9 January 2015 13
- Received in revised form 26 March 2015 14
- Accepted 30 March 2015 15
- 16Available online xxxx
- 17 Keywords:
- 18 Tumor markers
- Combined analysis 19
- 20C4.5 decision tree
- 21 Single neural network
- 22Hepatobiliary
- 23ALP 24Mucin

ABSTRACT

Objective: Cholangiocarcinoma (CCA) is usually fatal because of the absence of tests for early detection and 25 lack of effective therapy. Tumor markers with adequate diagnostic values are of clinical significance. This study is 26 aimed to improve the diagnostic power of serum markers using the computational data mining technique to de- 27 velop a combined diagnostic model that yielded the best diagnostic values for CCA.

Design and methods: Eight CCA-associated markers—carcinoembryonic antigen, carbohydrate antigen 19-9, 29 alkaline phosphatase (ALP), and gamma glutamyl transferase, biliary-ALP, mucin5AC, CCA-associated carbohy- 30 drate antigen (CCA-CA) and CA-S27-were used as the inputs for the C4.5 decision tree classification model 31 and the selected model was confirmed by ANN analyses. Eight serum markers for CCA were determined in the 32 training set of 85 histologically proven-CCA patients and 82 control subjects. The chosen set of combined markers 33 that gave the best diagnostic values for CCA was then validated in the testing set of 22 CCA patients and 60 34 controls. 35

Results: A decision tree diagram built by the C4.5 algorithm suggested the serial analysis of CCA-CA and ALP 36 for distinguishing CCA patients from non-CCA subjects with all diagnostic parameters ≥95%. The combined tests 37 showed a precise diagnosis in the testing set.

Conclusions: The C4.5 model indicates the combined markers of CCA-CA and ALP that produced the more 39 precise diagnosis for CCA.

© 2015 The Canadian Society of Clinical Chemists. Published by Elsevier Inc. All rights reserved. 41

49 44

Introduction 46

Cholangiocarcinoma (CCA), a malignancy of bile duct epithelium, is a 47 48relatively rare cancer worldwide; however, its incidence is extremely high in Southeast Asia, especially in the northeastern part of Thailand 49 [1]. CCA is normally difficult to diagnose until the disease develops to 5051the advanced or disseminated stage, which makes CCA a tumor with an extremely poor prognosis. As a result, approximately 35/100,000 or 525314,400 CCA patients die per year in Thailand [2]. Therefore, it is an urgent need to search for tumor markers that are sensitive and specific 54 enough to diagnose individuals with CCA.

At present, serum alkaline phosphatase (ALP), carbohydrate antigen 56 19-9 (CA19-9), carcinoembryonic antigen (CEA) and gamma-glutamyl 57 transferase (GGT) are the most routine markers for diagnosis of CCA 58 in general clinical practices. A number of new CCA-associated markers 59 have also been reported, such as biliary-alkaline phosphatase (BALP) 60 [3], mucin5AC (MUC5AC) [4], CCA-associated carbohydrate antigen 61 (CCA-CA) [5], CA-S27 [6], cytokeratin 19 fragment (CYFRA21-1) [7,8], 62 receptor-binding cancer antigen expressed in SiSo cells IRCAS1 [9], 63 and M2-pyruvate kinase [10]. However, none of these markers alone 64 has impressive clinical diagnostic value [11]. 65

Improvement in the diagnosis of cancer using combined measure- 66 ments of multiple tumor markers has been demonstrated in many stud- 67 ies [11-13]. This approach provides a better sensitivity, specificity, 68 and accuracy for detection of cancers. Computer-based diagnostic and 69 referral systems, such as artificial neural network and decision tree 70

http://dx.doi.org/10.1016/j.clinbiochem.2015.03.022

0009-9120/© 2015 The Canadian Society of Clinical Chemists. Published by Elsevier Inc. All rights reserved.

Please cite this article as: Pattanapairoj S, et al, Improve discrimination power of serum markers for diagnosis of cholangiocarcinoma using data mining-based approach, Clin Biochem (2015), http://dx.doi.org/10.1016/j.clinbiochem.2015.03.022

Correspondence to: D. Chetchotsak, Department of Industrial Engineering, Faculty of Engineering, Khon Kaen University, 40002, Thailand. Fax: +66 43 362 145.

^{*} Correspondence to: S. Wongkham, Department of Biochemistry, Faculty of Medicine, Khon Kaen University, Khon Kaen 40002, Thailand. Fax: +66 43 348 386.

E-mail addresses: cdanai@kku.ac.th (D. Chetchotsak), sopit@kku.ac.th (S. Wongkham). ¹ Co-first author.

2

71classifications, have been applied to simulate the combined-diagnostic 72models or equations from the measured tumor markers [12-14].

The computational data mining C4.5 model, and its most recent 73 74 version C5.0, have been used extensively as a classification model in many different areas such as dental therapy [15], mammogram analysis 7576[16], land-use planning [17], construction accidents [18], and individual 77 credit evaluation [19]. In this study, we aimed to use the computational 78data mining C4.5 model to develop a combined diagnostic model that 79yielded the best diagnostic values for CCA from four routinely used 80 serum markers for CCA-CEA, CA19-9, ALP, and GGT, and four recently reported markers-BALP, MUC5AC, CCA-CA and CA-S27. The C4.5 creat-81 ed several combined marker models and the model that had the highest 82 discrimination power was then compared with those of the reference 83 technique, single-hidden-layer feed-forward neural network, common-84 ly known as artificial neural network (ANN). Finally, the discrimination 85 86 power of the combined tests was validated in the second cohort subjects. The combination test may serve as a tool to provide more 87 precise diagnosis of CCA. 88

Materials and methods 89

Patients and serum samples 90

All serum samples were obtained from the specimen bank of The 91 92Liver Fluke and Cholangiocarcinoma Research Center, Khon Kaen University, Khon Kaen, Thailand. Informed consent was obtained from 93 all subjects, and the study protocol was approved by The Khon Kaen 94University Ethics Committee in Human Research. Serum from 85 histo-95logically proven-CCA patients (Table 1) and from 82 age and sex 96 97 matched control subjects, including 11 healthy persons (HE), 30 98 opisthorchiasis subjects (OV), 14 benign hepatobiliary disease patients 99 (BHD), and 27 gastrointestinal cancer patients (GICA) were used in the first cohort as a training set. A second cohort used as a validation 100 set was composed of 22 CCA patients and 60 controls (HE = 22, 101 OV = 7, BHD = 16, GICA = 15). Cancer and BHD patients were diag-102 nosed using histology. Opisthorchiasis in the OV group was confirmed 103 by a positive fecal result for Opisthorchis viverrini eggs, and healthy 104 controls were persons whose medical check-up, liver function tests 105 and complete blood count were normal. 106

107 Determination of serum markers

ALP (U/L) and GGT (U/L) were determined using enzymatic assay 108 109 (Roche Diagnostics GmbH, Mannheim, Germany). CEA (ng/mL) and CA19-9 (U/mL) were analyzed using enzyme linked immunosorbent 110 111 assay (ELISA, Roche Diagnostics GmbH) according to the International Federation of Clinical Chemistry and Laboratory Medicine (IFCC). 112 Serum BALP (U/L), MUC5AC (Absorbance Unit, AU), CCA-CA (AU) and 113 CA-S27 (AU) were determined using lectin capture ELISA as described 114 previously [3–6]. 115

1.1	Table	1

1.2 Demographical data of CCA patients of the training	g se
--	------

Characteristics		Number of patients	
Age (years)	≤56	46	54.1%
	>56	39	45.9%
Sex	Male	56	65.8%
	Female	29	34.2%
Histological types	Papillary	17	20.0%
	Non-papillary	68	80.0%
Staging	I–III	12	14.1%
	IVA	25	29.4%
	IVB	48	56.5%
Jaundice ^a	Jaundice	49	57.6%
	Non-jaundice	36	42.4%

Iaundice = total bilirubin > 2 mg/dL.

Construction of C4.5 classification model

116

In this study, C4.5 decision tree algorithm was used as a classification 117 model while an artificial neural network was used for benchmarking 118 purposes. C4.5 is an algorithm used to construct a decision tree classifi- 119 cation model in a logical form [20]. The model creates a decision tree 120 based on the information gain concept [21], where information gain re- 121 flects the ability to classify the data using a particular marker. Starting at 122 the top node of the tree, C4.5 selects the marker with the highest infor- 123 mation gain to classify the data in a "top-down recursive divide-and- 124 conquer manner" [22]. At the branch node, C4.5 selects another marker 125 from the remaining set in the same manner. These procedures take 126 place until no more markers are left. More details of C4.5 are provided 127 in [20]. 128

An artificial neural network known as (ANN) [23] is an information 129 processing and computational system which is modeled after a biologi- 130 cal nervous system. Like a human brain, ANN can process, learn, and re- 131 member information. The simplest and most common form of ANN is a 132 single-hidden layer feedforward neural network trained using the 133 backpropagation algorithm and was used in this study. ANN can map 134 the relationship between the inputs and outputs through adjusting 135 the connection weights and hence it has been used successfully in 136 both regression and classification tasks. 137

To construct a classification model using both C4.5 and ANN, all com- 138 bined markers were used as inputs and the diagnostic result as "CCA" or 139 "nonCCA" was the output. For C4.5, we conducted an experiment to 140 build a model with the best combination of markers. In this step, binary 141 splits of at least five instances per leaf and a confidence threshold for 142 pruning of 0.75 were used as settings. To construct ANN, the following 143 settings were used: 20 hidden units, sigmoid function, and 50,000 learn- 144 ing cycles. The results obtained were compared to those from the C4.5 145 model. Both C4.5 and ANN were built using an open source machine 146 learning software, WEKA [24] and were validated using the five-fold 147 cross validation method. 148

Performance evaluation and experimental trials

149

169

For consistency, the CCA samples were referred to as "positive" 150 while the non-CCA samples were labeled as "negative". The confusion 151 matrix and defined performance measures for a binary class problem, 152 used to measure performance of the classification models were: sensi- 153 tivity (SEN), specificity (SPEC), positive predictive value (PPV), negative 154 predictive value (NPV) and accuracy (ACC). These measures were 155 defined as follows: the terms TP and TN denoted correct classification 156 for positive and negative samples while FP and FN defined incorrect 157 classification for positive and negative samples, respectively. The diag- 158 nostic values were determined as SEN = TP / (TP + FN); SPEC = 159TN / (TN + FP); PPV = TP / (TP + FP); NPV = TN / (TN + FN); and 160 ACC = (TP + TN) / (TP + FN + FP + TN).161

To make sure that the training data was chosen randomly to avoid 162 bias, C4.5 model construction in this study was replicated ten times. 163 This was done by training and validating the classification model ac- 164 cording to the five-fold cross validation method. Performance of each 165 model was then evaluated. These procedures were subsequently re- 166 peated ten times. A 95% confidence interval for all performances was 167 then constructed and used for comparison and evaluation purposes. 168

Statistical analysis

Statistical analysis was performed using the SPSS 14.0 for Windows 170 Evaluation software (SPSS Inc., Chicago, USA). Mann-Whitney U 171 test was used to compare the serum level of marker of the CCA and 172 the control groups. P-values <0.05 were considered statistically 173 significant. 174

Please cite this article as: Pattanapairoj S, et al, Improve discrimination power of serum markers for diagnosis of cholangiocarcinoma using data mining-based approach, Clin Biochem (2015), http://dx.doi.org/10.1016/j.clinbiochem.2015.03.022

Download English Version:

https://daneshyari.com/en/article/8317227

Download Persian Version:

https://daneshyari.com/article/8317227

Daneshyari.com