CSBJ-00155; No of Pages 11

ARTICLE IN PR

Computational and Structural Biotechnology Journal xxx (2016) xxx-xxx





COMPUTATIONAL ANDSTRUCTURAL BIOTECHNOLOGY JOURNAL



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

A Hybrid Computer-aided-diagnosis System for Prediction of Breast Cancer Recurrence (HPBCR) Using Optimized Ensemble Learning 2

Mohammad R. Mohebian^a, Hamid R. Marateb^{a,b}, Marjan Mansourian^{c,*}, 04 Miguel Angel Mañanas^b, Fariborz Mokarian^{d,e} 4

^a Biomedical Engineering Department, Engineering Faculty, University of Isfahan, Hezar Jerib St., 81746-73441, Isfahan, Iran

^b Department of Automatic Control, Biomedical Engineering Research Center, Universitat Politècnica de Catalunya, BarcelonaTech (UPC), C. Pau Gargallo, 5, 08028 Barcelona, Spain 6

^c Department of Biostatistics and Epidemiology, School of Public Health, Isfahan University of Medical Sciences, Hezar Jerib St., 81745 Isfahan, Iran

^d Cancer Prevention Research Center, Isfahan University of Medical Sciences, Isfahan, Iran 8

^e Department of Internal Medicine, School of Medicine, Isfahan University of Medical Sciences, Isfahan, Iran 9

7

10

ARTICLE INFO 1 1

12Article history:

13Received 30 August 2016

Received in revised form 24 November 2016 14

15 Accepted 26 November 2016

16Available online xxxx

20

- 41 Keywords:
- 42 Breast cancer

43Cancer recurrence Computer-assisted diagnosis

- 44 45 Machine learning
- 46
- Prognosis

ABSTRACT

Cancer is a collection of diseases that involves growing abnormal cells with the potential to invade or spread to 22 the body. Breast cancer is the second leading cause of cancer death among women. A method for 5-year breast 23 cancer recurrence prediction is presented in this manuscript. Clinicopathologic characteristics of 579 breast 24 cancer patients (recurrence prevalence of 19.3%) were analyzed and discriminative features were selected 25 using statistical feature selection methods. They were further refined by Particle Swarm Optimization (PSO) as 26 the inputs of the classification system with ensemble learning (Bagged Decision Tree: BDT). The proper combi-27 nation of selected categorical features and also the weight (importance) of the selected interval-measurement- 28 scale features were identified by the PSO algorithm. The performance of HPBCR (hybrid predictor of breast cancer 29 recurrence) was assessed using the holdout and 4-fold cross-validation. Three other classifiers namely as sup- 30 ported vector machines, DT, and multilayer perceptron neural network were used for comparison. The selected 31 features were diagnosis age, tumor size, lymph node involvement ratio, number of involved axillary lymph 32 nodes, progesterone receptor expression, having hormone therapy and type of surgery. The minimum sensitivity, 33 specificity, precision and accuracy of HPBCR were 77%, 93%, 95% and 85%, respectively in the entire cross- 34 validation folds and the hold-out test fold. HPBCR outperformed the other tested classifiers. It showed excellent 35 agreement with the gold standard (i.e. the oncologist opinion after blood tumor marker and imaging tests, and 36 tissue biopsy). This algorithm is thus a promising online tool for the prediction of breast cancer recurrence. 37 © 2016 Mohebian et al.. Published by Elsevier B.V. on behalf of the Research Network of Computational and 38 Structural Biotechnology. This is an open access article under the CC BY license 39

(http://creativecommons.org/licenses/by/4.0/). 40

1. Introduction 51

40 49

Computer-aided diagnosis (CAD) is using computers and software to 5253interpret medical information. The purpose of CAD is to improve the diagnosis accuracy. In fact, CAD is used as a second opinion by the 54physicians to make the final diagnosis decision [1,2]. 55

05

Corresponding author.

E-mail addresses: van_mohebian@yahoo.com (M.R. Mohebian),

h.marateb@eng.ui.ac.ir, hamid.reza.marateb@upc.edu (H.R. Marateb),

j_mansourian@hlth.mui.ac.ir (M. Mansourian), miguel.angel.mananas@upc.edu (M.A. Mañanas), mokarian@med.mui.ac.ir (F. Mokarian).

Nowadays, CAD is used in many different fields in medicine includ- 56 ing, but not limited to, early detection of breast cancer [3], lung cancer 57 diagnosis [4], arrhythmia detection [5] and dental and maxillofacial 58 lesions diagnosis [6]. Several studies have been reported in the 59 literature focusing on the application of CAD for cancer diagnosis and 60 prognosis [7,8].

Cancer is a collection of diseases that involves growing abnormal 62 cells with the potential to spread to other parts of the body [9]. There 63 are over 200 types of cancer. The most common types of cancer in 64 women are breast, colorectal, lung, and cervical [10].

Cancer is the leading cause of death worldwide, accounting for 66 8.2 million deaths in 2012 [11]. The most common causes of cancer 67 death are: lung (1.59 million deaths), liver (745,000), stomach 68 (723,000), colorectal (694,000), breast (521,000), and esophageal 69 (400,000) [11]. More than 60% of world's total new annual cancer 70 cases occur in Africa, and Central and South America [12]. The cancer in-71 cidence in Asia is also high [13]. It is expected that annual cancer cases 72

2001-0370/© 2016 Mohebian et al.. Published by Elsevier B.V. on behalf of the Research Network of Computational and Structural Biotechnology. This is an open access article under the CC BY license (http://creativecommons.org/licenses/by/4.0/).

Please cite this article as: Mohebian MR, et al, A Hybrid Computer-aided-diagnosis System for Prediction of Breast Cancer Recurrence (HPBCR) Using Optimized Ensemble Learning, Comput Struct Biotechnol J (2016), http://dx.doi.org/10.1016/j.csbj.2016.11.004

Abbreviations: CAD, computer-aided diagnosis; DT, decision tree; FH, family history of cancer; HPBCR, the proposed hybrid predictor of breast cancer recurrence; HRT, hormone therapy; I. Node, number of involved axillary lymph nodes; NR, lymph node involvement ratio; T. Node, number of dissected axillary lymph nodes; TS, tumor size; XRT, radiotherapy.

http://dx.doi.org/10.1016/j.csbj.2016.11.004

2

ARTICLE IN PRESS

will rise from 14 million in 2012 to 22 million within the next 2 decades
[11,12]. Worthwhile, among all types of cancer, breast cancer is the
second leading cause of cancer death among women [14].

76 In Iran, cancer is the third cause of death after coronary heart disease, and accidents [15]. Breast cancer is the leading type of cancer 77 in Iranian females, accounting for 24.6% of all cancers. In Iran, the 78 79average women age with breast cancer is 49.6 years [16]. Among all of 80 the provinces in Iran, Isfahan Province is the biggest and most important 81 area located at a desert border in the center of Iran [17]. The rate of 82 cancer is increasing rapidly in Isfahan Province. Cancer control and 83 comprehensive prevention plan is thus necessary [18].

Breast cancer, a complex heterogeneous disease, occurs with a set of different clinical symptoms. It is usually diagnosed using blood tests, MRI, mammography, and CT scan and biopsy. The pathology results from biopsy samples indicate whether the suspicious area is cancerous. Cancer patients, then, undergo a systematic treatment procedure dependent on the cancer stage and additional lab tests such as hormone receptor status [19].

Cancer staging (assigning an ordinal numbers I-IV) is in fact the 91 process to determine the extent to which a cancer has been spread 92and is based on the following four characteristics: 1) the tumor size, 93 942) whether the cancer is invasive or non-invasive, 3) the spread of the 95tumor into the lymph nodes, and 4) the spread of the tumor into 96 other parts of the body (i.e. metastasis). Stage I is an isolated cancer (tumor size ≤20 mm in greatest dimension) while stage IV is a metasta-97 sis cancer. Most cancer deaths are due to cancer that has spread from its 98 primary site to other organs [12,20]. 99

100 Breast cancer is usually treated with surgery, which may be followed by chemotherapy, radiation, and hormone therapies [21]. When the 101 cancer patient is initially treated, the disease can recur at any time. 102However, most recurrences happen in the first 5 years after treatment 103 104 [22]. The recurrence could be local (near the mastectomy scar), regional 105(spread to nearby lymph nodes) or metastatic (spread to other parts of 106 the body not near the breast). Some of the most common sites of recurrence outside the breast are the lymph nodes, bones, liver, lungs, and 107 brain [23]. 108

An important issue is whether we can optimize the treatments to increase the therapeutic efficacy. In fact, 5-year recurrence-free survival is an important treatment quality measure. In principle, it is possible to predict 5-year cancer recurrence using clinicopathologic characteristics of cancer patients [24]. Such a prediction could be used by doctors to make proper treatment plan to considerably prolong patient life [25].

The prediction systems were used for cancer diagnosis in the literature [7]. However, there are few studies focusing on cancer prognosis (including recurrence or survival analysis). Since the focus of the current study is prediction of cancer recurrence, the literature review on cancer recurrence prediction models is provided. Meanwhile, the table of available methods was provided in the Supplementary material S1.

Zeng. suggested a mixture classification model containing a twolayer structure called mixture of rough set and support vector machine
(SVM) for breast cancer prognosis with the average accuracy of 91%
[26].

The Nottingham prognostic index (NPI) is a prognostic regression model proposed by Galea et al. [27] based on tumor size, histological grade, and lymph node status. The NPI calculation equation is as follows: tumor size (cm) \times 0.2 + histological grade + lymph node point (negative nodes = 1; 1–3 positive nodes = 2; \geq 4 positive nodes = 3). The patients were then classified into the low-risk (NPI point <3.4) and high-risk groups (NPI point \geq 3.4) [27,28].

Kim et al. [28] used normalized mutual information index for feature
selection and supported vector machines (SVM), Cox-proportional hazard regression model, and artificial neural network classifiers for classification in a sample size of 679 patients (the recurrence prevalence of
28.6%). The following features were used in their prognosis system:
local invasion of tumor, number of tumors, number of metastatic
lymph nodes, the histological grade, tumor size, estrogen receptor,

and lymphovascular invasion and reached the sensitivity, specificity 139 and area under the curve of 89%, 73% and 0.85, respectively for the 140 best classifier (SVM). Although the statistical power of their system 141 is acceptable (Power = 89% > 80%), the Type I error is beyond the 142 acceptable range ($\alpha = 0.17 > 0.05$). 143

Ahmad et al. [29] used the SVM, decision tree, and multilayer 144 perceptron artificial neural network classifiers with feature selection 145 in a sample size of 547 patients (the recurrence prevalence of 21.4%). 146 The predictors were age at diagnosis, menarche and menopause, 147 tumor size, number of involved and dissected axially lymph nodes, 148 grade and HER2. The best classifier (SVM) had sensitivity, specificity 149 and accuracy of 96%, 91% and 94%, respectively. Despite the high perfor- 150 mance of the algorithm proposed by the authors, cases with primary 151 metastasis, a metastasis diagnosed at the time of registry, were not 152 excluded from the analysis. In fact, having metastasis was one of the 153 attributes primarily used by their proposed system. This necessary 154 exclusion has been implemented in similar studies [28,30]. The reason 155 for such exclusion is that the survival rate of patients with stage IV 156 (metastasis) breast cancer within 5 years is about 10% to 15%. Thus, 157 most of patients with primary metastasis experience cancer recurrence. 158

The discriminative prognosis predictors vary regionally in different 159 studies [27-29]. Also, such prognosis tools have been recently used to 160 optimize the treatment protocol. An example is including the cancer 161 recurrence risk to breast cancer treatment guidelines by the American 162 Society of Clinical Oncology (ASCO) and the National Comprehensive 163 Cancer Network (NCCN). Thus, there is a need to design prognosis 164 tools in regions that are different based on cancer epidemiology. 165 Meanwhile, it is necessary to validate the cancer prognosis system 166 using extensive diagnosis validation criteria. Therefore, the aim of this 167 study was to reliably and accurately predict breast cancer recurrence 168 using pathological and demographics features of the patients. In the 169 proposed CAD system, HPBCR (hybrid predictor of breast cancer 170 recurrence), a hybrid technique including statistical features selection, 171 meta-heuristic population-based optimization and ensemble learning 172 were used to predict breast cancer recurrence in the first 5 years after 173 the diagnosis. 174

The rest of the paper is organized as follows: in the next section, information about the experimental protocol and the pattern recognition methods used in this study is presented. Section 3 provides the results of HPBCR and its comparison with the state-of-the-art. The discussion is provided in Section 4 and finally, the conclusions are summarized in Section 5. 180

181

182

2. Materials and Methods

2.1. The Cohort Dataset

In this study we used a 16-year registry cohort database (1998-183 2014) on 1085 women who diagnosed with breast cancer in Isfahan 184 Sayed-o-Shohada cancer research center. History of clinical conditions 185 and therapy of patients was continued until death or lost to follow up. 186 Characteristics of variables and tumor have been recorded by interview 187 and reported as pathology results. Time to recurrence was based on the 188 physicians' opinion. The following information was extracted from each 189 patient: age at diagnosis of breast cancer, lymph node involvement ratio 190 (NR) defined as ratio of involved to dissected lymph nodes [31], age of 191 menarche, number of pregnancy (No. Preg), primary tumor size (TS), 192 cellular marker for proliferation (Ki67), number of involved (I. Node) 193 and dissected (T. Node) axillary lymph nodes, number of chemother- 194 apies (No. Chemo) and categorical covariates of family history of cancer 195 (FH: negative/positive), having more than one tumor in the breast 196 (multifocal: negative/positive), estrogen receptor expression (ER: nega- 197 tive (also known as absent)/positive), progesterone receptor expression 198 (PR: negative (also known as absent)/positive), tumor protein 53 (p53: 199 negative/positive), type of surgery (MRM: modified radical mastectomy, 200 BCS: breast-conserving surgery, Mast: mastectomy), epidermal 201

Please cite this article as: Mohebian MR, et al, A Hybrid Computer-aided-diagnosis System for Prediction of Breast Cancer Recurrence (HPBCR) Using Optimized Ensemble Learning, Comput Struct Biotechnol J (2016), http://dx.doi.org/10.1016/j.csbj.2016.11.004 Download English Version:

https://daneshyari.com/en/article/8408337

Download Persian Version:

https://daneshyari.com/article/8408337

Daneshyari.com