

# Predicting breast screening attendance using machine learning techniques

Baskaran, V. , Guergachi, A. , Bali, R. and Naguib, R.

Author post-print (accepted) deposited in CURVE January 2012

## Original citation & hyperlink:

Baskaran, V. , Guergachi, A. , Bali, R. and Naguib, R. (2011) Predicting breast screening attendance using machine learning techniques. IEEE Transactions on Information Technology in Biomedicine, volume 15 (2): 251-259.

<http://dx.doi.org/10.1109/TITB.2010.2103954>

**Publisher statement:** © 2011 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

**Copyright © and Moral Rights are retained by the author(s) and/ or other copyright owners. A copy can be downloaded for personal non-commercial research or study, without prior permission or charge. This item cannot be reproduced or quoted extensively from without first obtaining permission in writing from the copyright holder(s). The content must not be changed in any way or sold commercially in any format or medium without the formal permission of the copyright holders.**

This document is the author's post-print version, incorporating any revisions agreed during the peer-review process. Some differences between the published version and this version may remain and you are advised to consult the published version if you wish to cite from it.

**CURVE is the Institutional Repository for Coventry University**

<http://curve.coventry.ac.uk/open>

# Predicting Breast Screening Attendance Using Machine Learning Techniques

Vikraman Baskaran, Aziz Guergachi, *Member, IEEE*, Rajeev K. Bali, *Senior Member, IEEE*,  
and Raouf N. G. Naguib, *Senior Member, IEEE*

**Abstract**—Machine learning-based prediction has been effectively applied for many healthcare applications. Predicting breast screening attendance using machine learning (prior to the actual mammogram) is a new field. This paper presents new predictor attributes for such an algorithm. It describes a new hybrid algorithm that relies on back-propagation and radial basis function-based neural networks for prediction. The algorithm has been developed in an open source-based environment. The algorithm was tested on a 13-year dataset (1995–2008). This paper compares the algorithm and validates its accuracy and efficiency with different platforms. Nearly 80% accuracy and 88% positive predictive value and sensitivity were recorded for the algorithm. The results were encouraging; 40–50% of negative predictive value and specificity warrant further work. Preliminary results were promising and provided ample amount of reasons for testing the algorithm on a larger scale.

**Index Terms**—Breast screening, cancer, machine learning, neural networks, prediction, screening attendance.

## I. INTRODUCTION

**B**REAST cancer is the most common cancer for women in North America [1]. In the U.K., over 40 000 women are being diagnosed with breast cancer each year [2], [3]. Mortality due to breast cancer is also one of the highest in the world [1], [4], and is the second highest of all cancers in the Canada [7]. Breast cancer should ideally be diagnosed at the earlier stages of its development to considerably reduce mortality. Possible treatments include removing or destroying the cancer cells to avoid the spread of the affected cells. Breast self-examination is an effective and noninvasive type of checking for any lumps in the breast tissue. Unfortunately, this greatly depends on the size of the lump, technique, and experience in carrying out a self-examination procedure by a woman [9]. An ultrasound test, examining breast tissue using sound waves, can be utilized to detect lumps but this is usually suited for women aged below 35

Manuscript received April 11, 2010; revised October 12, 2010; accepted December 23, 2010. Date of publication; date of current version. This work was supported in part by the NHS Cancer Screening Programs, U.K.

V. Baskaran and A. Guergachi are with Ryerson University, TRSM, Toronto, ON M5B 2K3, Canada (e-mail: vikraman@ryerson.ca; a2guerga@ryerson.ca).

R. K. Bali is with the KARMAH Group, Health Design and Technology Institute (HDTI), Coventry University, Coventry University Technology Park, Coventry CV1 2TT, U.K. (e-mail: r.bali@coventry.ac.uk).

R. N. G. Naguib is with the Biomedical Computing and BIOCORE, HDTI, Coventry University, Coventry University Technology Park, Coventry CV1 2TT, U.K. (e-mail: r.naguib@coventry.ac.uk).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TITB.2010.2103954

owing to the higher density of breast tissue [1]. Having a tissue biopsy via a fine needle aspiration or an excision is often used to examine the cells histopathologically and to diagnose if the growth, lump, is benign or cancerous. These investigations are mostly employed in treatments or post-treatment examination and as second rung diagnostic confirmation methods [10]. Performing a computed tomography or an MRI scan would result in a thorough examination of the breast tissue but these techniques are not favored due to reasons which include cost, needs preparation, noise, time, and images that may not be clear [10].

Mammography is a technique for detecting breast tissue lumps using a low dosage of X-ray. This technique can even detect a 3-mm-sized lump. The X-ray image of the breast tissue is captured and the image is thoroughly read by experienced radiologists and specialist mammogram readers [10]. Preliminary research suggests that women aged 50 and above are more susceptible to breast cancer; mammography is more suited to women in this age range due to the lower density of breast tissue [11]. Even though mammography has its critics—mainly due to its high rate of false positives and false negatives [13]—it has become the standard procedure for screening women by the NHS National Breast Screening Program in the U.K. [3], [15]. Mammography is the best and most viable tool for mass screening to detect cancer in the breast at an early stage [17]; however, the effectiveness of diagnosis through screening is directly dependent on the percentage of women attending the screening program [18]–[20]. The NHS Breast Screening Program, catering to the entire eligible women population, is funded by the Department of Health in the U.K. It covers 2.5 million women every year and detected nearly 16 500 cancers in the screened population for the year 2007–2008 [3]. Currently, the screening program routinely screens women between the ages 50 and 70.

Early breast cancer detection through screening is fundamental for increasing the efficacy of cancer treatment [11], [21]. Mammography has been accepted as the best and most economically viable tool for population screening [22]. Maximizing coverage for the target population is crucial for the success of such screening programs [11]. Currently, the breast cancer screening attendance rates are below expectations in many countries that have publicly funded healthcare programs [24]. This paper proposes a set of protocols to increase breast screening attendance for the U.K.'s NHS breast screening program. Based on this protocol, a new software prototype was created and tested. The prototype tests the prediction algorithm and shares the prediction results with multiple healthcare stakeholders for initiating opportunistic interventions on nonattendees. This prototype is a radical new idea that uses machine learning techniques for

86 predicting screening attendance and shares this knowledge by  
87 adopting the health informatics initiative of the NHS.

## 88 II. CHALLENGE

89 The NHS Breast Screening Program Annual Review (2008)  
90 states that, out of invited women, only 74% attend the screen-  
91 ing program [3]. This sizeable nonattendance could result in  
92 missed cancer detection for nearly 4 000 women (based on the  
93 cancer detection rate within screened women) [3]. This large  
94 percentage of nonattendance not only result in loss of life due  
95 to breast cancer but also result in loss of screening resources  
96 through costly imaging equipment laying idle, underutilization  
97 of specialist-imaging expertise, wasted screening slots, and so  
98 forth. Screening units are unable to arrange buffered attendees  
99 for the idle slots since the units do not know *a priori* which  
100 women will attend and which will not. In addition, there is a  
101 sizeable cost factor involved in sending repeat screening ap-  
102 pointments letters to nonattending women.

103 Reasons for nonattendance may be largely attributed to dis-  
104 interest in attending a mammography session, prior or current  
105 medical problems, and fear of X-rays [11], [24]. These rea-  
106 sons can be negated by proper education provided to women.  
107 Education has to be directed at explaining the advantages and  
108 importance of screening and assist in removing the sociocultural  
109 and personal barriers [25]. Other possible options include con-  
110 venience in terms of time, place, and dates provided to women  
111 for encouraging their attendance.

112 In spite of the expedient measures provided to the women,  
113 nonattendance has been a grave concern for the NHS—National  
114 Screening Program. This scenario can be properly addressed if  
115 those women who may probably not attend a screening appoint-  
116 ment can be identified in advance so that additional resources  
117 can be directed at interventions that can increase screening  
118 attendance.

119 A proposal enumerating the complete software solution is  
120 summarized at the end of Section IV. The National Screening  
121 Program has been constantly striving to provide better services  
122 to the public and one of the new enhancements offered by the  
123 screening services is to increase the screening age limit from  
124 64 to 70 [26]. This effectively increases the number of screen-  
125 ing episodes and results in augmenting the need for effective  
126 use of the already stretched NHS resources. All the aforemen-  
127 tioned factors underline the need to increase the breast screening  
128 attendance.

## 129 III. SOLUTION PROPOSED

130 To address these challenges, a set of protocols were devel-  
131 oped as part of the ongoing research. The protocols are based on  
132 two components: 1) machine learning algorithms for knowledge  
133 creation; and 2) health informatics for knowledge sharing. This  
134 paper elaborates on how the prediction-based knowledge was  
135 created through a machine learning algorithm. Machine learning  
136 [Artificial Intelligence (AI)-based algorithm] was implemented  
137 through the creation of a prototype software based on open

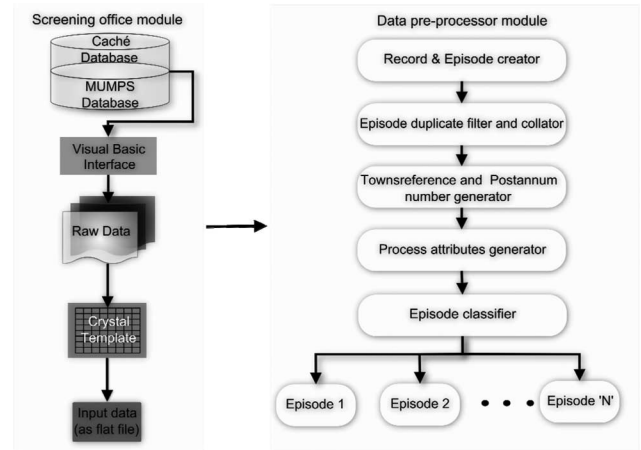


Fig. 1. Data filtering, preparation, and preprocessing.

138 source technologies. The prototype software was automated to  
139 produce the preprocessed data and eventually normalize the  
140 data for neural network (AI) assimilation. These activities were  
141 performed sequentially without human involvement for repeata-  
142 bility, reliability, and accuracy.

143 The AI-based neural network incorporates all additional  
144 transformations that occurred within the screening process (in-  
145 cluding the change in the screening upper age limit). The pro-  
146 totype framework was called JAABS—Java-based attendance  
147 prediction by AI for breast screening. The prototype combines  
148 the demographic data pertaining to the nonattending women  
149 and information related to their family physician as a package.  
150 This package then triggers the generation of an electronic mes-  
151 sage based on the Health Level 7 (HL7) standards and utilizes  
152 web services as the message delivering technology. This paper  
153 focuses on the machine learning techniques used within the pro-  
154 totype and subsequent testing of the algorithm for its prediction  
155 accuracy.

### 156 A. Data Preprocessing Module

157 The prototype was constructed using two main modules: 1)  
158 data preprocessing module; and 2) AI module. The data prepro-  
159 cessing module (see Fig. 1) consists of “Screening office mod-  
160 ule” that accomplishes data extraction from the screening unit’s  
161 database. The demography details for the three-year call/recall  
162 were downloaded (extraction date—Jan 2008) from the local  
163 health care authority’s database. The downloading is affected  
164 via the health link network onto a standalone system within  
165 the breast screening unit. The historical data related to screen-  
166 ing, appointments, and results pertaining to screening women  
167 are retained within the screening unit’s “Massachusetts Gen-  
168 eral Hospital Utility Multi-Programming System” (MUMPS)  
169 database. MUMPS, also known as the Oxford system, is one of  
170 the earliest programming languages used since the 1960s [27].  
171 This language was extensively employed to write database ap-  
172 plications explicitly for the healthcare domain.

*Generate input data as flat file from "Crystal Report" template*  
*For every record*  
     *Separate records for each woman*  
         *Remove duplicate episodes*  
         *Collate episodes into one record*  
     *Generate townsend reference and post annum numbers*  
     *Generate attributes*  
     *Classify and save record into their respective episode groups*  
     *End*

Pseudo-code 1. Pseudo-code for filtering raw data and preprocessing it to generate predictor attributes and classify them based on their episode details.

TABLE I  
THIRTEEN-YEAR DATASET DETAILS

Description	Number of records
Total valid women's record	159,412
Number of records deleted due to multiple entries	15,778
Records with missing values	9,799
CR template output records	540,539

173 The MUMPS database is based on the disk operating system  
 174 (DOS) and employs character-based user interface for database  
 175 interrogation [27]. The cumbersome DOS-based system is prone  
 176 to erroneous data entry and hence warranted a change in the  
 177 system. A new software package, the National Breast Screen-  
 178 ing Computer System (NBSS), was developed in 2002–2003  
 179 to address these issues [28]. This NBSS consists of a Visual  
 180 Basic (VB) front end connected to a "Caché" database which  
 181 is seamlessly integrated with the MUMPS database [29]. Due  
 182 to the aforementioned factors, an unstable environment, thus,  
 183 resulted in considerable complexities during data extraction for  
 184 the current research. The screening office module (see Fig. 1)  
 185 is executed with the existing software programs available in the  
 186 breast screening office.

187 The VB front end made data extraction straightforward from  
 188 the MUMPS database through Structured Query Language  
 189 (SQL) queries directed at the Caché database. Currently, the  
 190 breast screening office is employing "Crystal Report" (CR) as  
 191 part of the NBSS to generate reports for all the screening activ-  
 192 ities, including screening, administration, invitation, etc. Part of  
 193 the data preprocessing was implemented through the CR soft-  
 194 ware. The screening unit had earlier indicated that the routine  
 195 functioning of the screening office should not be affected during  
 196 the data extraction process.

197 Hence, prior to data extraction, a CR template was created to  
 198 reflect the format of the data to be exported (see pseudo-code  
 199 1). This template was used to export the data as a flat file to  
 200 negate any system instability. All the screening units around the  
 201 country were expected to have some form of minimum facility  
 202 for creating datasets in a flat file format. Coupled with this, a  
 203 need for a low overhead on the existing IT system and minimum  
 204 additional complexities was considered as fundamental for the  
 205 prototype. All the aforementioned rationale strengthened the  
 206 need for adopting a compromised strategy that exports data as  
 207 a flat file, so that the mode of data transfer can be standardized  
 208 across the country with minimum or no interrogation with the  
 209 screening database.

210 The SQL query generated details for all the women in as  
 211 many records, pertaining to the demography and episodes. The  
 212 demographic data were incomplete and only the first record of  
 213 a particular woman had the complete dataset and the remaining  
 214 records of the women corresponded to the historical episode  
 215 details (see Table I). The women's address and name were ex-  
 216 cluded from the study to address data protection and maintain

anonymity. In spite of its necessity for the messaging module, 217  
 the complete dataset was generated without the personal infor- 218  
 mation of the screening women. The post code of the women 219  
 is indispensable for the current study, as it generates the im- 220  
 portant predictor variable in the form of Townsend's reference 221  
 (Townsend deprivation score denotes the socioeconomic status 222  
 of a given postcode) and post annum number. 223

To address this without compromising the research work, 224  
 variables related to postcode, such as the Townsend score, post 225  
 annum (post annum is an arbitrary number associated with the 226  
 women's postcode) and screening distance, were all processed to 227  
 generate categorical variables within the screening unit and then 228  
 the data were ported to the AI module. The individual women 229  
 were identified by their SX number (pseudo-anonymised unique 230  
 identifier). The AI module generated the attendance prediction, 231  
 which formed the core of the knowledge transfer. The recipient 232  
 of the knowledge transfer is the woman's family physician; 233  
 hence, family physician information in the form of surname, 234  
 surgery address, and postcode was later collated for sending the 235  
 HL7-based message. 236

*For each episode group*  
     *Normalize data for AI module*  
     *Generate networks (BPNN and RBFN) and train*  
     *For each network*  
         *Validate data*  
     *Test data*  
     *Generate screening attendance prediction*  
     *Collate the best and save output with women's detail*  
     *End*

Pseudo-code 2. Pseudo-code for the AI module and results collation for the final output

One "Record" object was associated with one or more 237  
 "Episode" objects (see Fig. 2). The gaps in the demographic 238  
 record have to be filled and the episode details were associ- 239  
 ated with the women's demographic data. Exhaustive analyses 240  
 of the data indicated that the CR report had duplicate episode 241  
 details and are to be removed before further processing can be 242  
 implemented (see Table I). Each record read from the CR re- 243  
 port has to be first partitioned into episode details and stored 244  
 as "Episode" objects. They are finally collated and associated 245  
 with the women's demographic details (represented as "Record" 246  
 object). In addition to this, all the records have to be automat- 247  
 ically validated. The earlier work by Arochena had identified 248  
 all the contributing predictor attributes through comprehensive 249



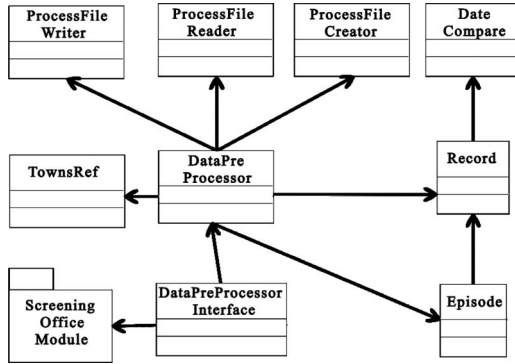


Fig. 2. UML class diagram for data preprocessing module (with I/O processing submodule).

TABLE II  
DATASET SPREAD ACROSS THE EPISODES AND ITS TRI-FURCATED DATA

Episode number	Total records	Train set	Valid set	Test set
Episode 1	23,277	4653	4708	13916
Episode 2	33,765	6838	6734	20193
Episode 3	29497	5868	5891	17738
Episode 4	43584	8792	8839	25953
Episode 5	26669	5340	5338	15991
Episode 6	2366	473	485	1408
Episode 7	238	36	39	163
Episode 8	16	3	3	10

250 statistical analyses [30]. After generating the required attributes,  
 251 the preprocessor module classifies the “Record” objects based  
 252 on the number of “Episode” objects it contains (see Fig. 2). This  
 253 dataset was then written as an in-process flat file for reference.  
 254 All errors generated during the execution of the preprocessing  
 255 module are written in a log (error) and is also saved as a flat file  
 256 for future reference.

257 The data preprocessing module identified episodes with miss-  
 258 ing data and removed them from the study. In total 2% (9 799)  
 259 were removed as records with missing data (see Table I). It fur-  
 260 ther deleted almost 3% (15 778) of the total records due to dupli-  
 261 cate entries. The valid records constituted 86% (159 412) of the  
 262 extracted dataset; on an average, each record had 3.2 episodes.  
 263 Table II depicts the spread of data for each episode. The highest  
 264 number of records was reached for the fourth episode. The first  
 265 to fifth episodes had an average of 31 000 records. For the re-  
 266 maining episodes (sixth, seventh, and eighth) the average is only  
 267 800 records. This might have a significant impact on the actual  
 268 prediction capacity of the JAABS algorithm for these episodes.

## 269 B. AI Module

270 JAABS is the new algorithm designed and developed in a  
 271 JAVA environment. As the design process was based on more  
 272 of an evolutionary type, a modular design strategy was selected.  
 273 This assists in parallel development of the implementation and  
 274 also enables testing as modules rather than as one single mono-  
 275 lithic program. The modular design also ensured that any addi-  
 276 tions or changes happening within the screening unit’s business

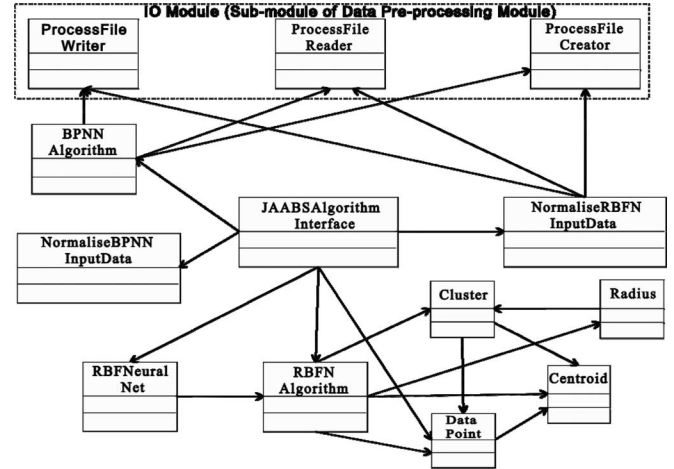


Fig. 3. UML class diagram of JAABS algorithm showing back propagation-based neural network and radial-basis function-based neural.

logic can be implemented without affecting the other modules  
 (see pseudo-code 2.). The “AI Module” encompasses the data  
 normalizer; the neural networks; and the results collator (see  
 Fig. 3). The Java-based algorithm implements two different  
 neural networks: feed-forward back-propagation neural network  
 (BPNN) and radial basis function neural network (RBFN).

The neural network algorithm requires the input data vector  
 classified as binary values; hence, the input data are normalized.  
 The input data in the RBFN are first passed through a radial basis  
 function algorithm, to identify the clusters and assign a radius  
 for cluster classification. These cluster centers are calculated  
 and the real-time data are checked against these established  
 cluster centers. Once the distance is calculated, the input dataset  
 is then associated with its nearest cluster. These data then trigger  
 a neural network for performing the prediction on attendance.  
 Each episode has a different set of predictor attributes; hence,  
 each episode is fed through separate neural networks that were  
 trained with their respective training dataset.

The results module collects the collated prediction for each  
 episode and submits it to a “Pooler” based classifier (see Fig. 4).  
 The “Pooler” finds the best prediction for the given episode  
 and generates the final prediction output based on the confi-  
 dence value of the prediction. This is fed into the prediction  
 result collator for all the input (women) based on each episode.  
 The consolidated result is used to generate the nonattendance  
 list and written as a flat file for processing by the “messaging  
 module” for message generation. The final output is associated  
 with the women’s SX number so that general physician details  
 can be added for knowledge sharing and to initiate physician  
 intervention.

## IV. ANALYSES

The predictor attributes (PA: post annum is an arbitrary num-  
 ber associated with the women’s postcode, TS: townsend depri-  
 vation score denotes the socioeconomic status of a given post-  
 code, AttBin: previous episode’s attendance, NumTest: number  
 of tests in the previous episodes, Cancer: denotes if cancer was  
 diagnosed in previous episodes, FP: false positive in previous

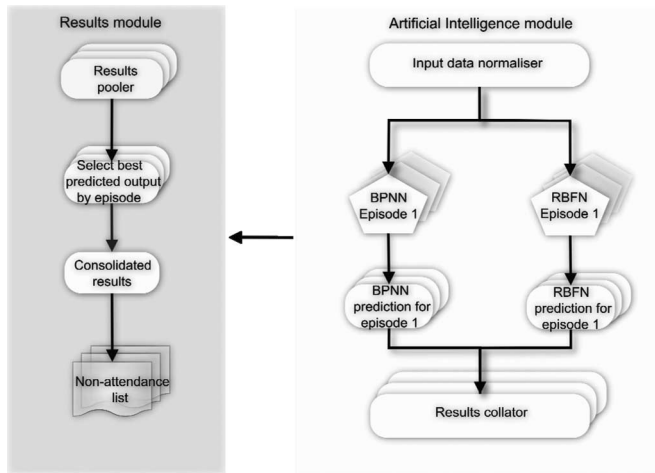


Fig. 4. Machine learning algorithm containing artificial intelligence and results module.

TABLE III  
PREDICTOR ATTRIBUTES AND THEIR ASSOCIATION TO THE SCREENING ATTENDANCE EPISODE WISE

Independent variables	Epi1	Epi2	Epi3	Epi4	Epi5	Epi6	Epi7	Epi8
PA	❖	❖	❖	❖	❖	✓	✓	✓
TS	❖	❖	❖	❖	❖	❖	❖	✓
AttBin		✓	✓	✓	✓	✓	✓	✓
NumTest		✓	✓	✓	✓	✓	✓	✓
Cancer		●	●	❖	●	●	●	
FP		●	●	●	●	●	●	
HFP			●	●	●	●	●	✓
HC			●	●	●	●	●	
AttTypeBin	✓	✓	✓	✓	✓	✓	✓	✓
AgeBand	❖	❖	●	●	●	●	●	✓
Slip	✓	✓	✓	✓	✓	✓	✓	✓
ScrDist	●	●	●	●	●	❖	✓	

✓ Association more than 0.2  
 ❖ Association more than 0.1 and less than 0.2  
 ● Association more than 0.0 and less than 0.1  
 No association is left blank

314 episodes, HFP: history of false positive, HC: history of cancer,  
 315 AttTypeBin: type of attendance like first or later episodes, Age-  
 316 Band: age categories, Slip: difference in days between screening  
 317 appointment and actual screening date, ScrDist: distance traveled  
 318 by the women for getting a mammogram) were initially  
 319 verified for their association with the screening attendance (see  
 320 Table III). The variables, being categorical, were analyzed  
 321 through parameters such as Lambda, Uncertainty, Phi (), Cram-  
 322 mer's V, and Contingency (confidence level at 95%).

323 These tests for association were conducted for establishing  
 324 some kind of linear relationship between the dependent and in-  
 325 dependent variables. Even though an association was not strong,  
 326 it was used only to establish some form of relationship between  
 327 the variables. This was used as an indication and as a first step  
 328 for resolving the real problem space which is multispatial. This  
 329 strategy assisted in filtering out the nonparticipating attributes  
 330 and to reduce the introduction of background noise.

331 Episode 1 lacked the historical variables and had to rely  
 332 only on demographic details. The rest of the episodes have

TABLE IV  
ROC FOR ALL EPISODES—AIATT AND JAABS (JAVA AND CLEMENTINE)

AI-ATT- Clementine (version 5)					
AIATT	ACC	NPV	PPV	SPC	SEN
Episode 1	67.01	20.45	87.48	41.81	71.43
Episode 2	87.76	56.1	92.85	58.91	93.14
Episode 3	86.49	50.54	92.91	55.99	91.32
Episode 4	81.65	41.26	92.51	64.59	85.42
Avg. for 4					
Episodes	80.73	42.09	91.44	55.33	85.33
JAABS- Java					
JAABS	ACC	NPV	PPV	SPC	SEN
Episode 1	67.29	42.07	76.71	40.22	78.05
Episode 2	69.38	47.65	77.87	45.66	79.22
Episode 3	69.95	39.45	76.46	26.29	85.59
Episode 4	79.17	39.25	87.06	37.37	87.93
Episode 5	76.23	51.61	83.84	49.64	84.89
Episode 6	57.79	46.51	64.77	44.92	66.21
Episode 7	51.39	30.02	76.53	60.05	48.18
Avg. for 4					
Episodes	71.45	42.11	79.53	37.39	82.7
Average	67.31	42.37	77.61	43.45	75.72
JAABS-Clementine (version 12)					
JAABS	ACC	NPV	PPV	SPC	SEN
Episode 1	68.16	52.58	69.35	11.57	95.04
Episode 2	79.61	74.59	81.33	57.93	90.28
Episode 3	81.24	72.56	83.86	57.63	90.99
Episode 4	85.73	74.91	88.45	62	93.34
Episode 5	80.81	74.43	82.56	53.88	92.18
Episode 6	67.88	63.8	70.36	56.7	76.16
Episode 7	78.99	86.49	77.61	41.56	96.89
Avg. for 4					
Episodes	78.68	68.66	80.75	47.28	92.41
Average	77.49	71.34	79.08	48.75	90.7

333 both the demographic and historical attributes as predictors; es-  
 334 pecially the new attribute in the form of screening distance  
 335 was found to increase the prediction efficiency for all the  
 336 episodes. The JAABS algorithm and its predictor attributes  
 337 were compared with its predecessor [AI-based attendance pre-  
 338 diction algorithm(AI-ATT)] for validation [30]. The AI-ATT  
 339 algorithm was developed in a visual modeling environment—  
 340 Clementine [30]. This off-the-shelf software assisted in design-  
 341 ing and implementing the algorithm rapidly, but created new  
 342 functional challenges such as the need for licensing the software  
 343 for all the screening units, specialist requirement for running the  
 344 algorithm, as it was not automated, and is based on outdated data  
 345 and semantics (1989–2001) to name just a few.

346 AI-ATT provided a base line for comparison and a reference  
 347 for validating the JAABS algorithm. To make the validation  
 348 more up-to-date, the same dataset that was applied to the JAABS  
 349 algorithm was also tested on Clementine (version 12.0). The  
 350 dataset was trifurcated into training, validating, and test sets (see  
 351 Table II). The training set contained equal numbers of women  
 352 categorized as attendees and nonattendees. The validating set  
 353 contained data that were never exposed during the training and  
 354 contained an equal number of attendees and nonattendees. The  
 355 test set contained skewed data, where nonattendees were only a  
 356 small proportion. This ensures that the test set reflects the real-  
 357 time dataset that would also be skewed (less nonattendees). The  
 358 JAABS algorithm was tested with the complete set of episodes  
 359 after appropriate training and validation.

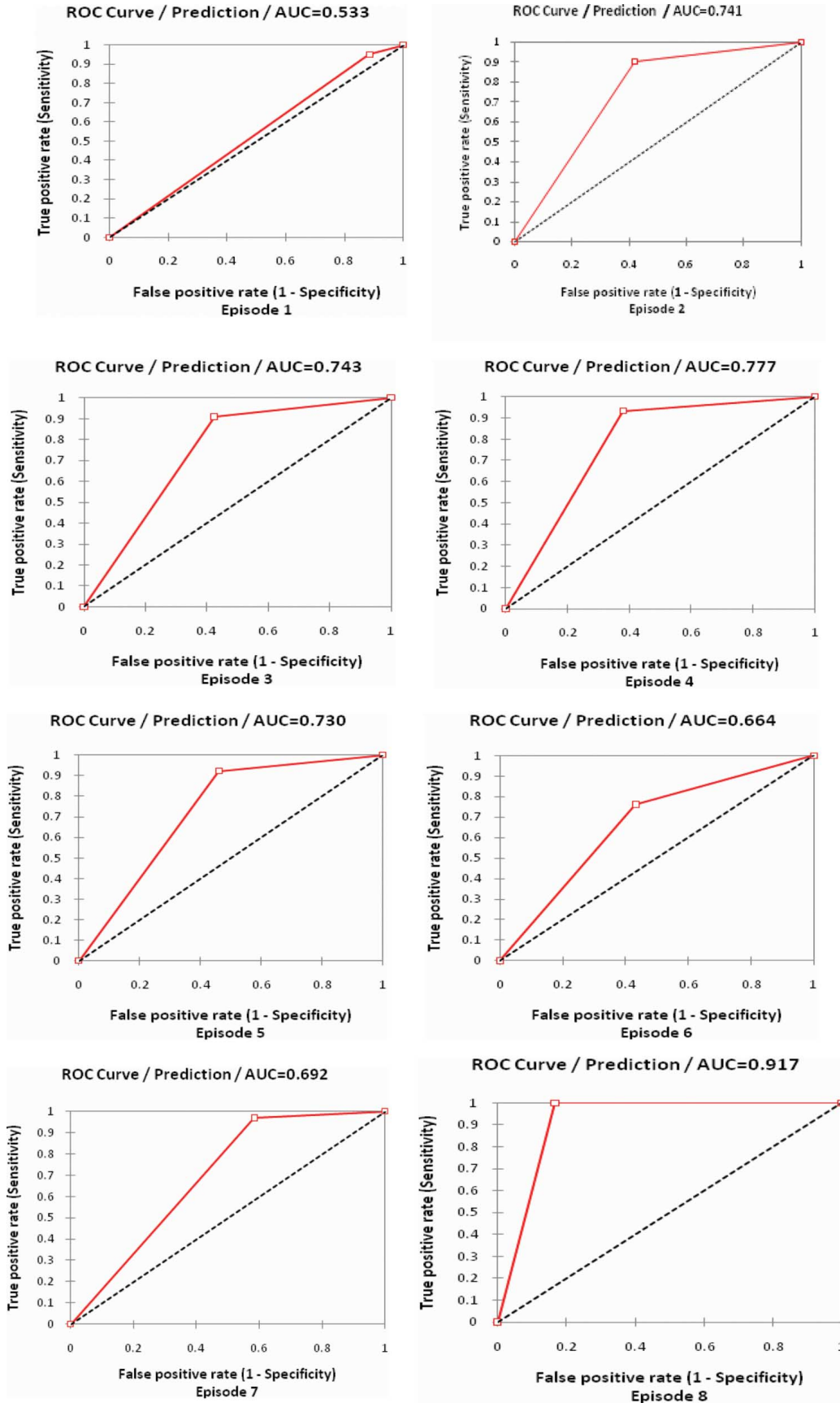


Fig. 5. ROC curve for Episodes one to eight for the machine learning algorithm.



The receiver operator characteristics (ROC) are summarized in Table IV (ACC: accuracy, NPV: negative predictive value, PPV: positive predictive value, SPC: specificity, SEN: sensitivity). The algorithm's final prediction of the screening attendance was based on a polling strategy that relies on the prediction confidence. The accuracy of the algorithm was around 68% for the first three episodes. Episode 4 had the maximum accuracy at 79%, closely followed by the fifth episode. The accuracies of the sixth and seventh episodes were lowest (57% and 51%, respectively). The NPV was the maximum at 51% for the fifth episode. The rest of the episodes had NPV values between 41% and 47%.

Episode 7 had the lowest NPV (30%). These lower NPVs were expected as the proportion of nonattendees was lesser in the test set (unbalanced). The PPVs for the fourth and fifth episodes were higher between 83% and 87%. The remaining episodes had values in the seventies range, except for the sixth episode where it was 64%. Specificity was highest for the seventh episode at 60%, but this may not be a true indicator as this episode had only 238 records in total. The next highest value was in the fifth episode at 49%. Episodes 1, 2, and 6 had values between 40% and 45%. Episodes 3 and 4 had lower values at 26% and 37%, respectively. The sensitivity was around 80% for the first four episodes, peaking at 85% for Episode 3. The higher the training set of records, the higher the sensitivity values. Since the previous algorithm (AI-ATT) had only four episodes, the averages for the first four episodes were used for comparing the JAABS and AI-ATT algorithms. The same set of attributes, when presented to commercial software (Clementine), generated improved results (see Table IV).

The first three episodes show an almost 10% increase in accuracy. Similarly, the later episodes (Episodes 4 and 5) when predicted by the JAABS–Clementine model, on average, do 6% better than the JAABS–Java algorithm, whereas Episodes 6 and 7 illustrated the maximum difference in accuracy (10–27%); this shows that the commercial software performed better even with a reduced training dataset. The NPV was lowest for the first episode, but was double when compared to AI-ATT and nearly 10% more than JAABS (Java). The NPV for the rest of the episodes (second to fifth) was around 73%. The remainder (sixth and seventh) were at 63% and 86%, respectively. The NPV is the metric that corresponds to the prediction of nonattendance and this was much better than that was achieved by the AI-ATT. Specificity is the next important measure and tests on Clementine showed promising results for all the episodes except for the first one.

The ROC curves for JAABS (Clementine) showed good prediction characteristics for all episodes except for Episode 1 (see Fig. 5). From the model's performance perspective, all these prediction characteristics were positive. The AI model proposed (JAABS—implemented in both Java and Clementine) was consistent and even outperformed the earlier model (AI-ATT) in many aspects. This could be attributed to the larger database and more complete attribute set and even the new predictor variable (screening distance) assisting in improving the algorithm's efficiency. The knowledge creation by applying AI (JAABS) is not only consistent, repeatable, and economical, but also ensures

minimal human intervention. This is ideal for automating the whole process.

The proposed AI network (JAABS) for predicting screening nonattendance would be incorporated in a new breast screening software model that connects to the screening database to generate the screening batch. Based on the prediction, an automated message would be sent to the women's healthcare stakeholders (GPs, nurses, and other clinical specialists). These messages would be assimilated by the clinical system used by the stakeholders and would eventually flag the women as a nonattender. When a woman's clinical record is opened, a flag/pop-up window would trigger opportunistic interventions that are aimed at educating the woman. This knowledge transfer would empower the woman to make an informed decision toward screening. This multistakeholder-based opportunistic intervention strategy would increase the overall breast screening attendance.

## V. CONCLUSION

This paper discussed the details of how a machine learning-based prediction tool can be effectively applied to increase the breast cancer screening attendance. The need for a high degree of automation was highlighted to simplify the algorithm's adoption; such automation would also reduce overheads and make integration as seamless as possible [31]. From the model's performance perspective, all the prediction characteristics were positive. The machine learning-based AI model (JAABS—implemented in both Java and Clementine) proposed was consistent and even outperformed the earlier model (AI-ATT) in many aspects. The performance improvement could be attributed to the larger database, more complete attribute set and even the new predictor variable (screening distance). The knowledge creation by applying AI (JAABS) is not only reliable, repeatable, and economical, but also ensures minimal human intervention. There is still scope for improving the prediction efficiency and this can be achieved through better predictor attributes and/or improved machine learning techniques. The former would be difficult to achieve as the data source itself may not be available but the latter would be possible as better AI models, such as support vector machines, fuzzy logic, and genetic algorithms or a combination of these, would enable further investigation for increasing the efficiency.

## ACKNOWLEDGMENT

The authors would like to thank J. Patnick CBE, Director, NHS Cancer Screening Programs (U.K.), for funding this research, Dr. M. Wallis, Consultant Radiologist, Cambridge Breast Unit team, and Margot Wheaton, Program Manager for the Warwickshire, Solihull and Coventry Breast Screening Service at Coventry and Warwickshire Hospital, for their excellent support and guidance throughout this research.

## REFERENCES

- [1] American Cancer Society. (2010, Feb. 10). *Breast Cancer Facts & Figures 2009–2010* [Online]. Available: [http://www.acsevents.org/downloads/STT/F861009\\_final%209-08-09.pdf](http://www.acsevents.org/downloads/STT/F861009_final%209-08-09.pdf).



- 469 [2] Cancer Research U.K. (2010, Feb. 10). *Breast Cancer—U.K.*  
470 *Mortality Statistics*. [Online]. Available: <http://info.cancerresearchuk.org/cancerstats/types/breast/mortality/index.htm>.  
471
- 472 [3] NHS Breast Screening Programme—Cancer Screening Programmes  
473 Annual Review 2009. (2010, Feb. 10). [Online]. Available: <http://www.cancerscreening.nhs.uk/breastscreen/publications/nhsbsp-annualreview2009.pdf>.  
474  
475
- 476 [4] K. Turner, J. Wilson, and J. Gilbert, “Improving breast screening uptake:  
477 Persuading initial non-attenders to attend,” *J. Med. Screening*, vol. 1,  
478 pp. 199–202, 1994.
- 479 [5] A. Majeed, R. Given-Wilson, and E. Smith, “Impact of follow up letters  
480 on non-attenders for breast screening: A general practice based study,” *J.*  
481 *Med. Screening*, vol. 4, pp. 19–20, 1997.
- 482 [6] J. P. Sin and A. S. Leger, “Interventions to increase breast screening  
483 uptake: Do they make any difference?,” *J. Med. Screening*, vol. 6, no. 1,  
484 pp. 170–181, 1999.
- 485 [7] Canadian Cancer Society. (2006). *Canadian Researchers Find Common*  
486 *Breast Cancer Chemotherapy Regime Inferior at Preventing Disease*  
487 *Recurrence* [Online]. Available: [http://www.cancer.ca/Canadawide/About%20us/Media%20centre/CW-Media%20releases/CW2006/Canadian%20Researchers%20Find%20Common%20Breast%20Cancer%20Chemo%20therapy%20Regime%20Inferior%20at%20Preventing%20Disease%20Recurrence.aspx?sc\\_lang=en](http://www.cancer.ca/Canadawide/About%20us/Media%20centre/CW-Media%20releases/CW2006/Canadian%20Researchers%20Find%20Common%20Breast%20Cancer%20Chemo%20therapy%20Regime%20Inferior%20at%20Preventing%20Disease%20Recurrence.aspx?sc_lang=en).  
488  
489  
490  
491
- 492 [8] Canadian Cancer Society. (2008, Mar. 22). *Canadian Cancer Statistics*  
493 *2008* [Online]. Available: [http://www.cancer.ca/Canada-wide/About%20cancer/Cancer%20statistics/~media/CSC/Canada%20wide/Files%20List/English%20files%20heading/pdf%20not%20in%20publications%20section/Canadian%20Cancer%20Society%20Statistics%20PDF%202008\\_614137951.ashx](http://www.cancer.ca/Canada-wide/About%20cancer/Cancer%20statistics/~media/CSC/Canada%20wide/Files%20List/English%20files%20heading/pdf%20not%20in%20publications%20section/Canadian%20Cancer%20Society%20Statistics%20PDF%202008_614137951.ashx).  
494  
495  
496  
497
- 498 [9] A. Oikonomou, S. A. Amin, R. N. G. Naguib, A. Todman, and H.  
499 Al-Omishy, “Breast self examination training through the use of multi-  
500 media: A prototype multimedia application,” *IEEE Eng. Med. Biol.*  
501 *Soc.*, vol. 2, no. 21, pp. 295–298, 2003.
- 502 [10] B. V. Marcela, “The system does work,” *J. Am. College Radiol.*, vol. 1,  
503 no. 6, pp. 438–440, 2004.
- 504 [11] L. Wyld, “Mammographic Breast Screening in Elderly Women,” in *Man-*  
505 *agement of Breast Cancer in Older Women*, part 3, M. W. Reed and R.  
506 A. Audisio, Eds. London, U.K.: Springer, 2010, ch. 9, pp. 127–142.
- 507 [12] R. G. Blanks, S. M. Moss, C. E. McGahan, M. J. Quinn, and P. J. Babb,  
508 “Effect of NHS breast screening programme on mortality from breast  
509 cancer in England and Wales, 1990–1998: Comparison of observed with  
510 predicted mortality,” *BMJ*, vol. 321, no. 7262, pp. 665–669, 2000.
- 511 [13] S. S. Epstein, *The Politics of Cancer*. New York: Doubleday, 1979,  
512 pp. 537.
- 513 [14] G. Burton, *Alternative Medicine*. Washington, DC: Future Medicine  
514 Publishing, 1997.
- 515 [15] Cancer Research U.K. (2007, Jul. 14). *Cancer Incidence—U.K. Statis-*  
516 *tics* [Online]. Available: <http://info.cancerresearchuk.org/cancerstats/incidence/index.htm>  
517
- 518 [16] P. Forest, *Breast Cancer Screening—A Report to the Health Ministers of*  
519 *England, Scotland, Wales and Northern Ireland*. London, U.K.: HMSO,  
520 1986.
- 521 [17] Medicine net (2010 Feb. 18). *Breast Cancer* [Online]. Available:  
522 [http://www.medicinenet.com/breast\\_cancer/page3.htm](http://www.medicinenet.com/breast_cancer/page3.htm)
- 523 [18] I. Pirjo, L. Kauhava, I. Parvinen, H. Helenius, and P. Klemi, “Customer  
524 fee and participation in breast cancer screening,” *The Lancet*, vol. 358,  
525 p. 1425, 2001.
- 526 [19] S. H. Woolf, “The 2009 Breast Cancer Screening Recommendations of the  
527 US Preventive Services Task Force,” *JAMA*, vol. 303, no. 2, pp. 162–163,  
528 2010.
- 529 [20] American Cancer Society Inc., (2010, Feb. 18) *Cancer Reference*  
530 *Information* [Online]. Available: [http://www.cancer.org/docroot/CRI/CRI\\_2\\_5x.asp?dt=5](http://www.cancer.org/docroot/CRI/CRI_2_5x.asp?dt=5)  
531
- 532 [21] D. P. Weller and C. Campbell, “Uptake in cancer screening programmes:  
533 A priority in cancer control,” *Brit. J. Cancer*, vol. 101, pp. 55–59, 2009.
- 534 [22] Y. Zheng, “Breast cancer detection with gabor features from digital mam-  
535 mograms,” *Algorithms*, vol. 3, pp. 44–62, 2010.
- 536 [23] K. W. Eilbert, K. Carroll, J. Peach, S. Khatoun, I. Basnett, and N. Mc-  
537 Culloch, “Approaches to improving breast screening uptake: Evidence  
538 and experience from Tower Hamlets,” *Brit. J. Cancer*, vol. 101, no. 2,  
539 pp. 64–67, 2009.
- 540 [24] D. Schopper and C. de Wolf, “How effective are breast cancer screening  
541 programmes by mammography? Review of the current evidence,” *Eur. J.*  
542 *Cancer*, vol. 45, no. 11, pp. 1916–1923, Jul. 2009.
- 543 [25] E. S. Cassandra, “Breast cancer screening: Cultural beliefs and diverse  
544 populations,” *Health Soc. work*, vol. 31, no. 1, pp. 36–43, 2006.
- 545 [26] NHS Cancer Screening Programmes. (2007, Apr.) *Disclosure of Audit*  
546 *Results in Cancer Screening Advice on Best Practice* (Cancer  
547 Screening Series 3), J. Patnick, Ed. [Online]. Available: <http://www.cancerscreening.nhs.uk/publications/cs3.pdf>  
548
- 549 [27] K. Okane. (2005, Apr. 20). *Mumps Language Bioinformatic*  
550 *Database Resources* [Online]. Available: [http://bioinformatics.org/forums/forum.php?forum\\_id=1035](http://bioinformatics.org/forums/forum.php?forum_id=1035)  
551
- 552 [28] V. Baskaran, R. K. Bali, R. N. G. Naguib, and H. Arochena, “A Knowl-  
553 edge Management approach to increase uptake in a breast screening pro-  
554 gramme,” presented at the IEEE 2nd Humanoid, Nanotechnology, In-  
555 formation Technology, Communication and Control, Environment and  
556 Management (HNICEM) Int. Conf., Philippines, Mar. 2005.
- 557 [29] S. Tarver, K. Cronin-Cowan, and M. E. Wheaton, “A pilot’s life for us,”  
558 *Breast Cancer Res.*, vol. 6, suppl. 1, p. 52, 2004.
- 559 [30] H. E. Arochena, “Modelling and prediction of parameters affecting atten-  
560 dance to the NHS breast cancer screening programme,” Ph.D. dissertation,  
561 Dept. Comp. Sci., Coventry Univ., Coventry, U.K., 2003.
- 562 [31] C. Bankhead, S. H. Richards, T. Peters, D. Sharp, R. Hobbs, J. Brown,  
563 L. Roberts, C. Tydeman, V. Redman, J. Formby, S. Wilson, and J. Austoker,  
564 “Improving attendance for breast screening among recent non-attenders:  
565 A randomised controlled trial of two interventions in primary care,” *J.*  
566 *Med. Screening*, vol. 8, no. 2, pp. 99–105, 2001.



**Vikraman Baskaran** is currently an Assistant Professor at the School of Information Technology Management of Ryerson University, Toronto, ON, Canada. His research interests include finding a viable application of the KM paradigm in healthcare application. His special interest in developing HL7 messaging and health informatics has provided opportunities in excelling in these fields. His current activities include KM, e-health, artificial intelligence, and healthcare informatics.

He is a member of the HL7 U.K. and Canada.



**Aziz Guergachi (M'xx)** is currently an Associate Professor at the Ted Rogers School of Information Technology Management of Ryerson University, Toronto, ON, Canada. Prior to becoming part of the Ryerson community, he was involved in the development of a large software system for trade promotion management and collaborative sales forecasting. His current research interests include advanced system modeling and machine learning with applications to business management and engineering systems.

He is the recipient of the New Opportunities Award of the Canada Foundation for Innovation and currently runs a research laboratory for advanced systems modeling.



**Rajeev K. Bali (SM'xx)** is currently a Reader in Healthcare Knowledge Management at Coventry University, U.K. His main research interests include clinical and healthcare knowledge management (from both technical and organisational perspectives). He has published peer-reviewed journals and is the an author/editor of several textbooks on healthcare knowledge management.

He serves on various editorial boards and conference committees and is regularly invited to deliver presentations and speeches internationally.



**Raouf N. G. Naguib (SM'xx)** is currently a Professor of Biomedical Computing and Head of BIOCORE, Coventry, U.K. Prior to this appointment, he was a Lecturer at Newcastle University, Newcastle Upon Tyne, U.K. He has published more than 240 journals and conference papers and reports in many aspects of biomedical and digital signal processing, image processing, artificial intelligence, and evolutionary computation in cancer research.

He was awarded the Fulbright Cancer Fellowship in 1995–1996 when he carried out research at the University of Hawaii, Mānoa, on the applications of artificial neural networks in breast cancer diagnosis and prognosis. He is a member of several national and international research committees and boards.

## QUERIES

620

- Q1. Author: Please check whether the edits made in the sentence “This large percentage of nonattendance not only . . .” retain your intended sense. 621
- Q2. Author: Refs. [5], [6], [8], [12], [14], [15], [16], and [23] are not cited in the text. Please check and provide citations. 622
- Q3. Author: Please provide the expansion of KM. 623
- Q4. Author: Please provide the educational details of all the authors. 624
- Q5. Author: Please provide the year in which Aziz Guergachi became “Member” of the IEEE. 625
- Q6. Author: Please provide the year in which Rajeev K Bali became “Senior Member” of the IEEE. 626
- Q7. Author: Please provide the year in which Raouf N. G. Naguib became “Senior Member” of the IEEE. 627
- 628

IEEE  
Proof

# Predicting Breast Screening Attendance Using Machine Learning Techniques

Vikraman Baskaran, Aziz Guergachi, *Member, IEEE*, Rajeev K. Bali, *Senior Member, IEEE*,  
and Raouf N. G. Naguib, *Senior Member, IEEE*

**Abstract**—Machine learning-based prediction has been effectively applied for many healthcare applications. Predicting breast screening attendance using machine learning (prior to the actual mammogram) is a new field. This paper presents new predictor attributes for such an algorithm. It describes a new hybrid algorithm that relies on back-propagation and radial basis function-based neural networks for prediction. The algorithm has been developed in an open source-based environment. The algorithm was tested on a 13-year dataset (1995–2008). This paper compares the algorithm and validates its accuracy and efficiency with different platforms. Nearly 80% accuracy and 88% positive predictive value and sensitivity were recorded for the algorithm. The results were encouraging; 40–50% of negative predictive value and specificity warrant further work. Preliminary results were promising and provided ample amount of reasons for testing the algorithm on a larger scale.

**Index Terms**—Breast screening, cancer, machine learning, neural networks, prediction, screening attendance.

## I. INTRODUCTION

**B**REAST cancer is the most common cancer for women in North America [1]. In the U.K., over 40 000 women are being diagnosed with breast cancer each year [2], [3]. Mortality due to breast cancer is also one of the highest in the world [1], [4], and is the second highest of all cancers in the Canada [7]. Breast cancer should ideally be diagnosed at the earlier stages of its development to considerably reduce mortality. Possible treatments include removing or destroying the cancer cells to avoid the spread of the affected cells. Breast self-examination is an effective and noninvasive type of checking for any lumps in the breast tissue. Unfortunately, this greatly depends on the size of the lump, technique, and experience in carrying out a self-examination procedure by a woman [9]. An ultrasound test, examining breast tissue using sound waves, can be utilized to detect lumps but this is usually suited for women aged below 35

Manuscript received April 11, 2010; revised October 12, 2010; accepted December 23, 2010. Date of publication; date of current version. This work was supported in part by the NHS Cancer Screening Programs, U.K.

V. Baskaran and A. Guergachi are with Ryerson University, TRSM, Toronto, ON M5B 2K3, Canada (e-mail: vikraman@ryerson.ca; a2guerga@ryerson.ca).

R. K. Bali is with the KARMAH Group, Health Design and Technology Institute (HDTI), Coventry University, Coventry University Technology Park, Coventry CV1 2TT, U.K. (e-mail: r.bali@coventry.ac.uk).

R. N. G. Naguib is with the Biomedical Computing and BIOCORE, HDTI, Coventry University, Coventry University Technology Park, Coventry CV1 2TT, U.K. (e-mail: r.naguib@coventry.ac.uk).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TITB.2010.2103954

owing to the higher density of breast tissue [1]. Having a tissue biopsy via a fine needle aspiration or an excision is often used to examine the cells histopathologically and to diagnose if the growth, lump, is benign or cancerous. These investigations are mostly employed in treatments or post-treatment examination and as second rung diagnostic confirmation methods [10]. Performing a computed tomography or an MRI scan would result in a thorough examination of the breast tissue but these techniques are not favored due to reasons which include cost, needs preparation, noise, time, and images that may not be clear [10].

Mammography is a technique for detecting breast tissue lumps using a low dosage of X-ray. This technique can even detect a 3-mm-sized lump. The X-ray image of the breast tissue is captured and the image is thoroughly read by experienced radiologists and specialist mammogram readers [10]. Preliminary research suggests that women aged 50 and above are more susceptible to breast cancer; mammography is more suited to women in this age range due to the lower density of breast tissue [11]. Even though mammography has its critics—mainly due to its high rate of false positives and false negatives [13]—it has become the standard procedure for screening women by the NHS National Breast Screening Program in the U.K. [3], [15]. Mammography is the best and most viable tool for mass screening to detect cancer in the breast at an early stage [17]; however, the effectiveness of diagnosis through screening is directly dependent on the percentage of women attending the screening program [18]–[20]. The NHS Breast Screening Program, catering to the entire eligible women population, is funded by the Department of Health in the U.K. It covers 2.5 million women every year and detected nearly 16 500 cancers in the screened population for the year 2007–2008 [3]. Currently, the screening program routinely screens women between the ages 50 and 70.

Early breast cancer detection through screening is fundamental for increasing the efficacy of cancer treatment [11], [21]. Mammography has been accepted as the best and most economically viable tool for population screening [22]. Maximizing coverage for the target population is crucial for the success of such screening programs [11]. Currently, the breast cancer screening attendance rates are below expectations in many countries that have publicly funded healthcare programs [24]. This paper proposes a set of protocols to increase breast screening attendance for the U.K.'s NHS breast screening program. Based on this protocol, a new software prototype was created and tested. The prototype tests the prediction algorithm and shares the prediction results with multiple healthcare stakeholders for initiating opportunistic interventions on nonattendees. This prototype is a radical new idea that uses machine learning techniques for

86 predicting screening attendance and shares this knowledge by  
87 adopting the health informatics initiative of the NHS.

## 88 II. CHALLENGE

89 The NHS Breast Screening Program Annual Review (2008)  
90 states that, out of invited women, only 74% attend the screen-  
91 ing program [3]. This sizeable nonattendance could result in  
92 missed cancer detection for nearly 4 000 women (based on the  
93 cancer detection rate within screened women) [3]. This large  
94 percentage of nonattendance not only result in loss of life due  
95 to breast cancer but also result in loss of screening resources  
96 through costly imaging equipment laying idle, underutilization  
97 of specialist-imaging expertise, wasted screening slots, and so  
98 forth. Screening units are unable to arrange buffered attendees  
99 for the idle slots since the units do not know *a priori* which  
100 women will attend and which will not. In addition, there is a  
101 sizeable cost factor involved in sending repeat screening ap-  
102 pointments letters to nonattending women.

103 Reasons for nonattendance may be largely attributed to dis-  
104 interest in attending a mammography session, prior or current  
105 medical problems, and fear of X-rays [11], [24]. These rea-  
106 sons can be negated by proper education provided to women.  
107 Education has to be directed at explaining the advantages and  
108 importance of screening and assist in removing the sociocultural  
109 and personal barriers [25]. Other possible options include con-  
110 venience in terms of time, place, and dates provided to women  
111 for encouraging their attendance.

112 In spite of the expedient measures provided to the women,  
113 nonattendance has been a grave concern for the NHS—National  
114 Screening Program. This scenario can be properly addressed if  
115 those women who may probably not attend a screening appoint-  
116 ment can be identified in advance so that additional resources  
117 can be directed at interventions that can increase screening  
118 attendance.

119 A proposal enumerating the complete software solution is  
120 summarized at the end of Section IV. The National Screening  
121 Program has been constantly striving to provide better services  
122 to the public and one of the new enhancements offered by the  
123 screening services is to increase the screening age limit from  
124 64 to 70 [26]. This effectively increases the number of screen-  
125 ing episodes and results in augmenting the need for effective  
126 use of the already stretched NHS resources. All the aforemen-  
127 tioned factors underline the need to increase the breast screening  
128 attendance.

## 129 III. SOLUTION PROPOSED

130 To address these challenges, a set of protocols were devel-  
131 oped as part of the ongoing research. The protocols are based on  
132 two components: 1) machine learning algorithms for knowledge  
133 creation; and 2) health informatics for knowledge sharing. This  
134 paper elaborates on how the prediction-based knowledge was  
135 created through a machine learning algorithm. Machine learning  
136 [Artificial Intelligence (AI)-based algorithm] was implemented  
137 through the creation of a prototype software based on open

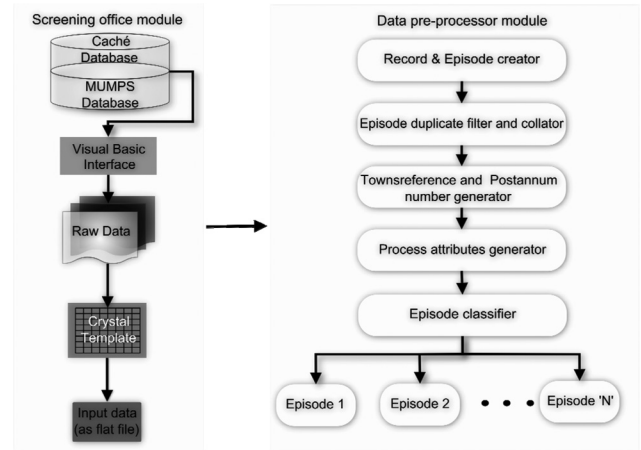


Fig. 1. Data filtering, preparation, and preprocessing.

138 source technologies. The prototype software was automated to  
139 produce the preprocessed data and eventually normalize the  
140 data for neural network (AI) assimilation. These activities were  
141 performed sequentially without human involvement for repeata-  
142 bility, reliability, and accuracy.

143 The AI-based neural network incorporates all additional  
144 transformations that occurred within the screening process (in-  
145 cluding the change in the screening upper age limit). The pro-  
146 totype framework was called JAABS—Java-based attendance  
147 prediction by AI for breast screening. The prototype combines  
148 the demographic data pertaining to the nonattending women  
149 and information related to their family physician as a package.  
150 This package then triggers the generation of an electronic mes-  
151 sage based on the Health Level 7 (HL7) standards and utilizes  
152 web services as the message delivering technology. This paper  
153 focuses on the machine learning techniques used within the pro-  
154 totype and subsequent testing of the algorithm for its prediction  
155 accuracy.

### 156 A. Data Preprocessing Module

157 The prototype was constructed using two main modules: 1)  
158 data preprocessing module; and 2) AI module. The data prepro-  
159 cessing module (see Fig. 1) consists of “Screening office mod-  
160 ule” that accomplishes data extraction from the screening unit’s  
161 database. The demography details for the three-year call/recall  
162 were downloaded (extraction date—Jan 2008) from the local  
163 health care authority’s database. The downloading is affected  
164 via the health link network onto a standalone system within  
165 the breast screening unit. The historical data related to screen-  
166 ing, appointments, and results pertaining to screening women  
167 are retained within the screening unit’s “Massachusetts Gen-  
168 eral Hospital Utility Multi-Programming System” (MUMPS)  
169 database. MUMPS, also known as the Oxford system, is one of  
170 the earliest programming languages used since the 1960s [27].  
171 This language was extensively employed to write database ap-  
172 plications explicitly for the healthcare domain.



Generate input data as flat file from "Crystal Report"  
template  
For every record  
  Separate records for each woman  
    Remove duplicate episodes  
    Collate episodes into one record  
  Generate townsend reference and post annum numbers  
  Generate attributes  
Classify and save record into their respective episode  
groups  
End

Pseudo-code 1. Pseudo-code for filtering raw data and preprocessing it to generate predictor attributes and classify them based on their episode details.

TABLE I  
THIRTEEN-YEAR DATASET DETAILS

Description	Number of records
Total valid women's record	159,412
Number of records deleted due to multiple entries	15,778
Records with missing values	9,799
CR template output records	540,539

173 The MUMPS database is based on the disk operating system  
174 (DOS) and employs character-based user interface for database  
175 interrogation [27]. The cumbersome DOS-based system is prone  
176 to erroneous data entry and hence warranted a change in the  
177 system. A new software package, the National Breast Screen-  
178 ing Computer System (NBSS), was developed in 2002–2003  
179 to address these issues [28]. This NBSS consists of a Visual  
180 Basic (VB) front end connected to a "Caché" database which  
181 is seamlessly integrated with the MUMPS database [29]. Due  
182 to the aforementioned factors, an unstable environment, thus,  
183 resulted in considerable complexities during data extraction for  
184 the current research. The screening office module (see Fig. 1)  
185 is executed with the existing software programs available in the  
186 breast screening office.

187 The VB front end made data extraction straightforward from  
188 the MUMPS database through Structured Query Language  
189 (SQL) queries directed at the Caché database. Currently, the  
190 breast screening office is employing "Crystal Report" (CR) as  
191 part of the NBSS to generate reports for all the screening activi-  
192 ties, including screening, administration, invitation, etc. Part of  
193 the data preprocessing was implemented through the CR soft-  
194 ware. The screening unit had earlier indicated that the routine  
195 functioning of the screening office should not be affected during  
196 the data extraction process.

197 Hence, prior to data extraction, a CR template was created to  
198 reflect the format of the data to be exported (see pseudo-code  
199 1). This template was used to export the data as a flat file to  
200 negate any system instability. All the screening units around the  
201 country were expected to have some form of minimum facility  
202 for creating datasets in a flat file format. Coupled with this, a  
203 need for a low overhead on the existing IT system and minimum  
204 additional complexities was considered as fundamental for the  
205 prototype. All the aforementioned rationale strengthened the  
206 need for adopting a compromised strategy that exports data as  
207 a flat file, so that the mode of data transfer can be standardized  
208 across the country with minimum or no interrogation with the  
209 screening database.

210 The SQL query generated details for all the women in as  
211 many records, pertaining to the demography and episodes. The  
212 demographic data were incomplete and only the first record of  
213 a particular woman had the complete dataset and the remaining  
214 records of the women corresponded to the historical episode  
215 details (see Table I). The women's address and name were ex-  
216 cluded from the study to address data protection and maintain

anonymity. In spite of its necessity for the messaging module, 217  
the complete dataset was generated without the personal infor- 218  
mation of the screening women. The post code of the women 219  
is indispensable for the current study, as it generates the im- 220  
portant predictor variable in the form of Townsend's reference 221  
(Townsend deprivation score denotes the socioeconomic status 222  
of a given postcode) and post annum number. 223

To address this without compromising the research work, 224  
variables related to postcode, such as the Townsend score, post 225  
annum (post annum is an arbitrary number associated with the 226  
women's postcode) and screening distance, were all processed to 227  
generate categorical variables within the screening unit and then 228  
the data were ported to the AI module. The individual women 229  
were identified by their SX number (pseudo-anonymised unique 230  
identifier). The AI module generated the attendance prediction, 231  
which formed the core of the knowledge transfer. The recipient 232  
of the knowledge transfer is the woman's family physician; 233  
hence, family physician information in the form of surname, 234  
surgery address, and postcode was later collated for sending the 235  
HL7-based message. 236

For each episode group  
  Normalize data for AI module  
  Generate networks (BPNN and RBFN) and train  
  For each network  
    Validate data  
  Test data  
  Generate screening attendance prediction  
  Collate the best and save output with women's detail  
End

Pseudo-code 2. Pseudo-code for the AI module and results collation for the final output

One "Record" object was associated with one or more 237  
"Episode" objects (see Fig. 2). The gaps in the demographic 238  
record have to be filled and the episode details were associ- 239  
ated with the women's demographic data. Exhaustive analyses 240  
of the data indicated that the CR report had duplicate episode 241  
details and are to be removed before further processing can be 242  
implemented (see Table I). Each record read from the CR re- 243  
port has to be first partitioned into episode details and stored 244  
as "Episode" objects. They are finally collated and associated 245  
with the women's demographic details (represented as "Record" 246  
object). In addition to this, all the records have to be automat- 247  
ically validated. The earlier work by Arochena had identified 248  
all the contributing predictor attributes through comprehensive 249

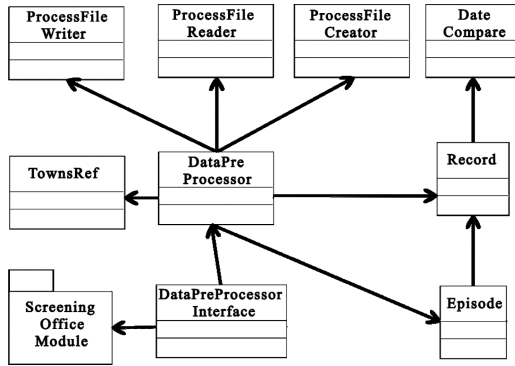


Fig. 2. UML class diagram for data preprocessing module (with I/O processing submodule).

TABLE II  
DATASET SPREAD ACROSS THE EPISODES AND ITS TRI-FURCATED DATA

Episode number	Total records	Train set	Valid set	Test set
Episode 1	23,277	4653	4708	13916
Episode 2	33,765	6838	6734	20193
Episode 3	29497	5868	5891	17738
Episode 4	43584	8792	8839	25953
Episode 5	26669	5340	5338	15991
Episode 6	2366	473	485	1408
Episode 7	238	36	39	163
Episode 8	16	3	3	10

250 statistical analyses [30]. After generating the required attributes,  
 251 the preprocessor module classifies the “Record” objects based  
 252 on the number of “Episode” objects it contains (see Fig. 2). This  
 253 dataset was then written as an in-process flat file for reference.  
 254 All errors generated during the execution of the preprocessing  
 255 module are written in a log (error) and is also saved as a flat file  
 256 for future reference.

257 The data preprocessing module identified episodes with miss-  
 258 ing data and removed them from the study. In total 2% (9 799)  
 259 were removed as records with missing data (see Table I). It fur-  
 260 ther deleted almost 3% (15 778) of the total records due to dupli-  
 261 cate entries. The valid records constituted 86% (159 412) of the  
 262 extracted dataset; on an average, each record had 3.2 episodes.  
 263 Table II depicts the spread of data for each episode. The highest  
 264 number of records was reached for the fourth episode. The first  
 265 to fifth episodes had an average of 31 000 records. For the re-  
 266 maining episodes (sixth, seventh, and eighth) the average is only  
 267 800 records. This might have a significant impact on the actual  
 268 prediction capacity of the JAABS algorithm for these episodes.

## 269 B. AI Module

270 JAABS is the new algorithm designed and developed in a  
 271 JAVA environment. As the design process was based on more  
 272 of an evolutionary type, a modular design strategy was selected.  
 273 This assists in parallel development of the implementation and  
 274 also enables testing as modules rather than as one single mono-  
 275 lithic program. The modular design also ensured that any addi-  
 276 tions or changes happening within the screening unit’s business

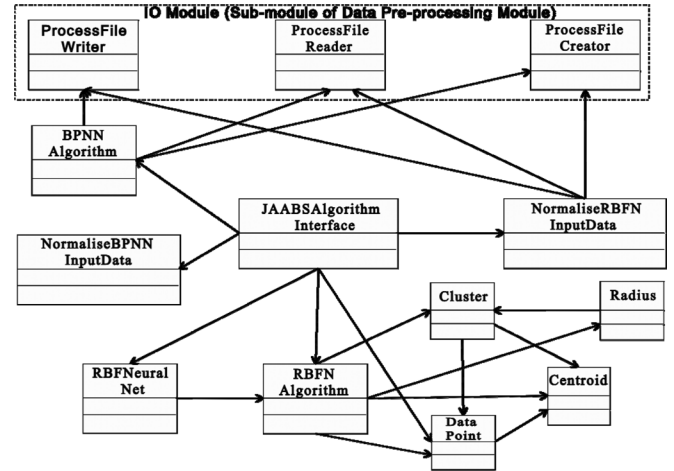


Fig. 3. UML class diagram of JAABS algorithm showing back propagation-based neural network and radial-basis function-based neural.

logic can be implemented without affecting the other modules 277  
 (see pseudo-code 2.). The “AI Module” encompasses the data 278  
 normalizer; the neural networks; and the results collator (see 279  
 Fig. 3). The Java-based algorithm implements two different 280  
 neural networks: feed-forward back-propagation neural network 281  
 (BPNN) and radial basis function neural network (RBFN). 282

The neural network algorithm requires the input data vector 283  
 classified as binary values; hence, the input data are normalized. 284  
 The input data in the RBFN are first passed through a radial basis 285  
 function algorithm, to identify the clusters and assign a radius 286  
 for cluster classification. These cluster centers are calculated 287  
 and the real-time data are checked against these established 288  
 cluster centers. Once the distance is calculated, the input dataset 289  
 is then associated with its nearest cluster. These data then trigger 290  
 a neural network for performing the prediction on attendance. 291  
 Each episode has a different set of predictor attributes; hence, 292  
 each episode is fed through separate neural networks that were 293  
 trained with their respective training dataset. 294

The results module collects the collated prediction for each 295  
 episode and submits it to a “Pooler” based classifier (see Fig. 4). 296  
 The “Pooler” finds the best prediction for the given episode 297  
 and generates the final prediction output based on the confi- 298  
 dence value of the prediction. This is fed into the prediction 299  
 result collator for all the input (women) based on each episode. 300  
 The consolidated result is used to generate the nonattendance 301  
 list and written as a flat file for processing by the “messaging 302  
 module” for message generation. The final output is associated 303  
 with the women’s SX number so that general physician details 304  
 can be added for knowledge sharing and to initiate physician 305  
 intervention. 306

## 307 IV. ANALYSES

The predictor attributes (PA: post annum is an arbitrary num- 308  
 ber associated with the women’s postcode, TS: townsend depri- 309  
 vation score denotes the socioeconomic status of a given post- 310  
 code, AttBin: previous episode’s attendance, NumTest: number 311  
 of tests in the previous episodes, Cancer: denotes if cancer was 312  
 diagnosed in previous episodes, FP: false positive in previous 313

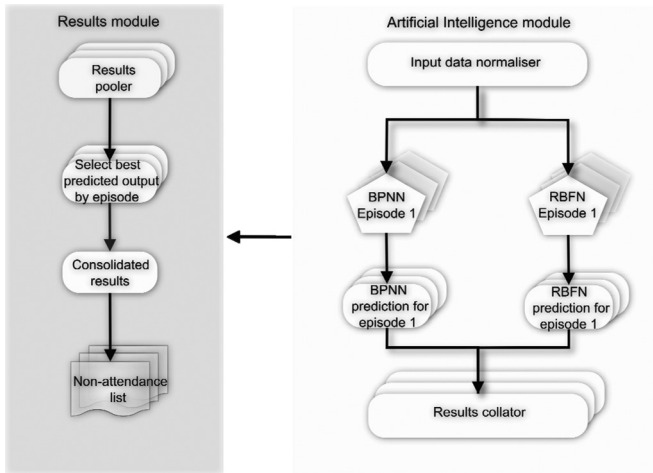


Fig. 4. Machine learning algorithm containing artificial intelligence and results module.

TABLE III  
PREDICTOR ATTRIBUTES AND THEIR ASSOCIATION TO THE SCREENING ATTENDANCE EPISODE WISE

Independent variables	Epi1	Epi2	Epi3	Epi4	Epi5	Epi6	Epi7	Epi8
PA	❖	❖	❖	❖	❖	✓	✓	✓
TS	❖	❖	❖	❖	❖	❖	❖	✓
AttBin		✓	✓	✓	✓	✓	✓	✓
NumTest		✓	✓	✓	✓	✓	✓	✓
Cancer		●	●	❖	●	●	●	
FP		●	●	●	●	●	●	
HFP			●	●	●	●	●	✓
HC			●	●	●	●	●	
AttTypeBin	✓	✓	✓	✓	✓	✓	✓	✓
AgeBand	❖	❖	●	●	●	●	●	✓
Slip	✓	✓	✓	✓	✓	✓	✓	✓
ScrDist	●	●	●	●	●	❖	✓	

✓ Association more than 0.2  
 ❖ Association more than 0.1 and less than 0.2  
 ● Association more than 0.0 and less than 0.1  
 No association is left blank

314 episodes, HFP: history of false positive, HC: history of cancer,  
 315 AttTypeBin: type of attendance like first or later episodes, Age-  
 316 Band: age categories, Slip: difference in days between screening  
 317 appointment and actual screening date, ScrDist: distance traveled  
 318 by the women for getting a mammogram) were initially  
 319 verified for their association with the screening attendance (see  
 320 Table III). The variables, being categorical, were analyzed  
 321 through parameters such as Lambda, Uncertainty, Phi (), Cram-  
 322 mer's V, and Contingency (confidence level at 95%).

323 These tests for association were conducted for establishing  
 324 some kind of linear relationship between the dependent and in-  
 325 dependent variables. Even though an association was not strong,  
 326 it was used only to establish some form of relationship between  
 327 the variables. This was used as an indication and as a first step  
 328 for resolving the real problem space which is multispatial. This  
 329 strategy assisted in filtering out the nonparticipating attributes  
 330 and to reduce the introduction of background noise.

331 Episode 1 lacked the historical variables and had to rely  
 332 only on demographic details. The rest of the episodes have

TABLE IV  
ROC FOR ALL EPISODES—AIATT AND JAABS (JAVA AND CLEMENTINE)

AI-ATT- Clementine (version 5)						
AIATT	ACC	NPV	PPV	SPC	SEN	
Episode 1	67.01	20.45	87.48	41.81	71.43	
Episode 2	87.76	56.1	92.85	58.91	93.14	
Episode 3	86.49	50.54	92.91	55.99	91.32	
Episode 4	81.65	41.26	92.51	64.59	85.42	
Avg. for 4						
Episodes	80.73	42.09	91.44	55.33	85.33	
JAABS- Java						
JAABS	ACC	NPV	PPV	SPC	SEN	
Episode 1	67.29	42.07	76.71	40.22	78.05	
Episode 2	69.38	47.65	77.87	45.66	79.22	
Episode 3	69.95	39.45	76.46	26.29	85.59	
Episode 4	79.17	39.25	87.06	37.37	87.93	
Episode 5	76.23	51.61	83.84	49.64	84.89	
Episode 6	57.79	46.51	64.77	44.92	66.21	
Episode 7	51.39	30.02	76.53	60.05	48.18	
Avg. for 4						
Episodes	71.45	42.11	79.53	37.39	82.7	
Average	67.31	42.37	77.61	43.45	75.72	
JAABS-Clementine (version 12)						
JAABS	ACC	NPV	PPV	SPC	SEN	
Episode 1	68.16	52.58	69.35	11.57	95.04	
Episode 2	79.61	74.59	81.33	57.93	90.28	
Episode 3	81.24	72.56	83.86	57.63	90.99	
Episode 4	85.73	74.91	88.45	62	93.34	
Episode 5	80.81	74.43	82.56	53.88	92.18	
Episode 6	67.88	63.8	70.36	56.7	76.16	
Episode 7	78.99	86.49	77.61	41.56	96.89	
Avg. for 4						
Episodes	78.68	68.66	80.75	47.28	92.41	
Average	77.49	71.34	79.08	48.75	90.7	

333 both the demographic and historical attributes as predictors; es-  
 334 pecially the new attribute in the form of screening distance  
 335 was found to increase the prediction efficiency for all the  
 336 episodes. The JAABS algorithm and its predictor attributes  
 337 were compared with its predecessor [AI-based attendance pre-  
 338 diction algorithm(AI-ATT)] for validation [30]. The AI-ATT  
 339 algorithm was developed in a visual modeling environment—  
 340 Clementine [30]. This off-the-shelf software assisted in design-  
 341 ing and implementing the algorithm rapidly, but created new  
 342 functional challenges such as the need for licensing the software  
 343 for all the screening units, specialist requirement for running the  
 344 algorithm, as it was not automated, and is based on outdated data  
 345 and semantics (1989–2001) to name just a few.

346 AI-ATT provided a base line for comparison and a reference  
 347 for validating the JAABS algorithm. To make the validation  
 348 more up-to-date, the same dataset that was applied to the JAABS  
 349 algorithm was also tested on Clementine (version 12.0). The  
 350 dataset was trifurcated into training, validating, and test sets (see  
 351 Table II). The training set contained equal numbers of women  
 352 categorized as attendees and nonattendees. The validating set  
 353 contained data that were never exposed during the training and  
 354 contained an equal number of attendees and nonattendees. The  
 355 test set contained skewed data, where nonattendees were only a  
 356 small proportion. This ensures that the test set reflects the real-  
 357 time dataset that would also be skewed (less nonattendees). The  
 358 JAABS algorithm was tested with the complete set of episodes  
 359 after appropriate training and validation.

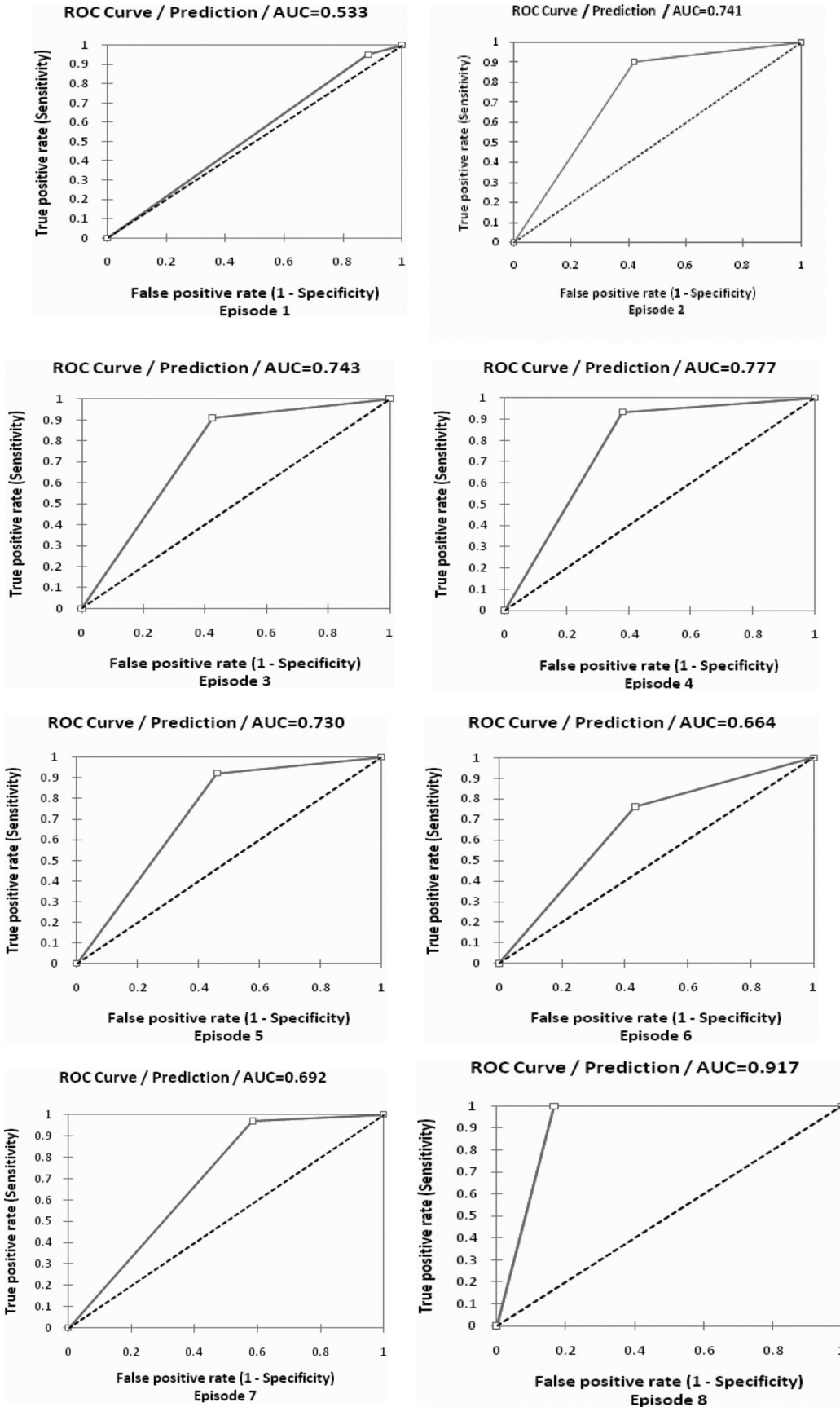


Fig. 5. ROC curve for Episodes one to eight for the machine learning algorithm.



The receiver operator characteristics (ROC) are summarized in Table IV (ACC: accuracy, NPV: negative predictive value, PPV: positive predictive value, SPC: specificity, SEN: sensitivity). The algorithm's final prediction of the screening attendance was based on a polling strategy that relies on the prediction confidence. The accuracy of the algorithm was around 68% for the first three episodes. Episode 4 had the maximum accuracy at 79%, closely followed by the fifth episode. The accuracies of the sixth and seventh episodes were lowest (57% and 51%, respectively). The NPV was the maximum at 51% for the fifth episode. The rest of the episodes had NPV values between 41% and 47%.

Episode 7 had the lowest NPV (30%). These lower NPVs were expected as the proportion of nonattendees was lesser in the test set (unbalanced). The PPVs for the fourth and fifth episodes were higher between 83% and 87%. The remaining episodes had values in the seventies range, except for the sixth episode where it was 64%. Specificity was highest for the seventh episode at 60%, but this may not be a true indicator as this episode had only 238 records in total. The next highest value was in the fifth episode at 49%. Episodes 1, 2, and 6 had values between 40% and 45%. Episodes 3 and 4 had lower values at 26% and 37%, respectively. The sensitivity was around 80% for the first four episodes, peaking at 85% for Episode 3. The higher the training set of records, the higher the sensitivity values. Since the previous algorithm (AI-ATT) had only four episodes, the averages for the first four episodes were used for comparing the JAABS and AI-ATT algorithms. The same set of attributes, when presented to commercial software (Clementine), generated improved results (see Table IV).

The first three episodes show an almost 10% increase in accuracy. Similarly, the later episodes (Episodes 4 and 5) when predicted by the JAABS–Clementine model, on average, do 6% better than the JAABS–Java algorithm, whereas Episodes 6 and 7 illustrated the maximum difference in accuracy (10–27%); this shows that the commercial software performed better even with a reduced training dataset. The NPV was lowest for the first episode, but was double when compared to AI-ATT and nearly 10% more than JAABS (Java). The NPV for the rest of the episodes (second to fifth) was around 73%. The remainder (sixth and seventh) were at 63% and 86%, respectively. The NPV is the metric that corresponds to the prediction of nonattendance and this was much better than that was achieved by the AI-ATT. Specificity is the next important measure and tests on Clementine showed promising results for all the episodes except for the first one.

The ROC curves for JAABS (Clementine) showed good prediction characteristics for all episodes except for Episode 1 (see Fig. 5). From the model's performance perspective, all these prediction characteristics were positive. The AI model proposed (JAABS—implemented in both Java and Clementine) was consistent and even outperformed the earlier model (AI-ATT) in many aspects. This could be attributed to the larger database and more complete attribute set and even the new predictor variable (screening distance) assisting in improving the algorithm's efficiency. The knowledge creation by applying AI (JAABS) is not only consistent, repeatable, and economical, but also ensures

minimal human intervention. This is ideal for automating the whole process.

The proposed AI network (JAABS) for predicting screening nonattendance would be incorporated in a new breast screening software model that connects to the screening database to generate the screening batch. Based on the prediction, an automated message would be sent to the women's healthcare stakeholders (GPs, nurses, and other clinical specialists). These messages would be assimilated by the clinical system used by the stakeholders and would eventually flag the women as a nonattender. When a woman's clinical record is opened, a flag/pop-up window would trigger opportunistic interventions that are aimed at educating the woman. This knowledge transfer would empower the woman to make an informed decision toward screening. This multistakeholder-based opportunistic intervention strategy would increase the overall breast screening attendance.

## V. CONCLUSION

This paper discussed the details of how a machine learning-based prediction tool can be effectively applied to increase the breast cancer screening attendance. The need for a high degree of automation was highlighted to simplify the algorithm's adoption; such automation would also reduce overheads and make integration as seamless as possible [31]. From the model's performance perspective, all the prediction characteristics were positive. The machine learning-based AI model (JAABS—implemented in both Java and Clementine) proposed was consistent and even outperformed the earlier model (AI-ATT) in many aspects. The performance improvement could be attributed to the larger database, more complete attribute set and even the new predictor variable (screening distance). The knowledge creation by applying AI (JAABS) is not only reliable, repeatable, and economical, but also ensures minimal human intervention. There is still scope for improving the prediction efficiency and this can be achieved through better predictor attributes and/or improved machine learning techniques. The former would be difficult to achieve as the data source itself may not be available but the latter would be possible as better AI models, such as support vector machines, fuzzy logic, and genetic algorithms or a combination of these, would enable further investigation for increasing the efficiency.

## ACKNOWLEDGMENT

The authors would like to thank J. Patnick CBE, Director, NHS Cancer Screening Programs (U.K.), for funding this research, Dr. M. Wallis, Consultant Radiologist, Cambridge Breast Unit team, and Margot Wheaton, Program Manager for the Warwickshire, Solihull and Coventry Breast Screening Service at Coventry and Warwickshire Hospital, for their excellent support and guidance throughout this research.

## REFERENCES

- [1] American Cancer Society. (2010, Feb. 10). *Breast Cancer Facts & Figures 2009–2010* [Online]. Available: [http://www.acsevents.org/downloads/STT/F861009\\_final%209-08-09.pdf](http://www.acsevents.org/downloads/STT/F861009_final%209-08-09.pdf).

- 469 [2] Cancer Research U.K. (2010, Feb. 10). *Breast Cancer—U.K.*  
470 *Mortality Statistics*. [Online]. Available: <http://info.cancerresearchuk.org/cancerstats/types/breast/mortality/index.htm>.  
471
- 472 [3] NHS Breast Screening Programme—Cancer Screening Programmes  
473 Annual Review 2009. (2010, Feb. 10). [Online]. Available: <http://www.cancerscreening.nhs.uk/breastscreen/publications/nhsbsp-annualreview2009.pdf>.  
474  
475
- 476 [4] K. Turner, J. Wilson, and J. Gilbert, “Improving breast screening uptake:  
477 Persuading initial non-attenders to attend,” *J. Med. Screening*, vol. 1,  
478 pp. 199–202, 1994.
- 479 [5] A. Majeed, R. Given-Wilson, and E. Smith, “Impact of follow up letters  
480 on non-attenders for breast screening: A general practice based study,” *J.*  
481 *Med. Screening*, vol. 4, pp. 19–20, 1997.
- 482 [6] J. P. Sin and A. S. Leger, “Interventions to increase breast screening  
483 uptake: Do they make any difference?,” *J. Med. Screening*, vol. 6, no. 1,  
484 pp. 170–181, 1999.
- 485 [7] Canadian Cancer Society. (2006). *Canadian Researchers Find Common*  
486 *Breast Cancer Chemotherapy Regime Inferior at Preventing Disease*  
487 *Recurrence* [Online]. Available: [http://www.cancer.ca/Canadawide/About%20us/Media%20centre/CW-Media%20releases/CW2006/Canadian%20Researchers%20Find%20Common%20Breast%20Cancer%20Chemotherapy%20Regime%20Inferior%20at%20Preventing%20Disease%20Recurrence.aspx?sc\\_lang=en](http://www.cancer.ca/Canadawide/About%20us/Media%20centre/CW-Media%20releases/CW2006/Canadian%20Researchers%20Find%20Common%20Breast%20Cancer%20Chemotherapy%20Regime%20Inferior%20at%20Preventing%20Disease%20Recurrence.aspx?sc_lang=en).  
488  
489  
490  
491
- 492 [8] Canadian Cancer Society. (2008, Mar. 22). *Canadian Cancer Statistics*  
493 *2008* [Online]. Available: [http://www.cancer.ca/Canada-wide/About%20cancer/Cancer%20statistics/~media/CCS/Canada%20wide/Files%20List/English%20files%20heading/pdf%20not%20in%20publications%20section/Canadian%20Cancer%20Society%20Statistics%20PDF%202008\\_614137951.ashx](http://www.cancer.ca/Canada-wide/About%20cancer/Cancer%20statistics/~media/CCS/Canada%20wide/Files%20List/English%20files%20heading/pdf%20not%20in%20publications%20section/Canadian%20Cancer%20Society%20Statistics%20PDF%202008_614137951.ashx).  
494  
495  
496  
497
- 498 [9] A. Oikonomou, S. A. Amin, R. N. G. Naguib, A. Todman, and H.  
499 Al-Omishy, “Breast self examination training through the use of multi-  
500 media: A prototype multimedia application,” *IEEE Eng. Med. Biol.*  
501 *Soc.*, vol. 2, no. 21, pp. 295–298, 2003.
- 502 [10] B. V. Marcela, “The system does work,” *J. Am. College Radiol.*, vol. 1,  
503 no. 6, pp. 438–440, 2004.
- 504 [11] L. Wyld, “Mammographic Breast Screening in Elderly Women,” in *Man-*  
505 *agement of Breast Cancer in Older Women*, part 3, M. W. Reed and R.  
506 A. Audisio, Eds. London, U.K.: Springer, 2010, ch. 9, pp. 127–142.
- 507 [12] R. G. Blanks, S. M. Moss, C. E. McGahan, M. J. Quinn, and P. J. Babb,  
508 “Effect of NHS breast screening programme on mortality from breast  
509 cancer in England and Wales, 1990–1998: Comparison of observed with  
510 predicted mortality,” *BMJ*, vol. 321, no. 7262, pp. 665–669, 2000.
- 511 [13] S. S. Epstein, *The Politics of Cancer*. New York: Doubleday, 1979,  
512 pp. 537.
- 513 [14] G. Burton, *Alternative Medicine*. Washington, DC: Future Medicine  
514 Publishing, 1997.
- 515 [15] Cancer Research U.K. (2007, Jul. 14). *Cancer Incidence—U.K. Statis-*  
516 *tics* [Online]. Available: <http://info.cancerresearchuk.org/cancerstats/incidence/index.htm>  
517
- 518 [16] P. Forest, *Breast Cancer Screening—A Report to the Health Ministers of*  
519 *England, Scotland, Wales and Northern Ireland*. London, U.K.: HMSO,  
520 1986.
- 521 [17] Medicine net (2010 Feb. 18). *Breast Cancer* [Online]. Available:  
522 [http://www.medicinenet.com/breast\\_cancer/page3.htm](http://www.medicinenet.com/breast_cancer/page3.htm)
- 523 [18] I. Pirjo, L. Kauhava, I. Parvinen, H. Helenius, and P. Klemi, “Customer  
524 fee and participation in breast cancer screening,” *The Lancet*, vol. 358,  
525 p. 1425, 2001.
- 526 [19] S. H. Woolf, “The 2009 Breast Cancer Screening Recommendations of the  
527 US Preventive Services Task Force,” *JAMA*, vol. 303, no. 2, pp. 162–163,  
528 2010.
- 529 [20] American Cancer Society Inc., (2010, Feb. 18) *Cancer Reference*  
530 *Information* [Online]. Available: [http://www.cancer.org/docroot/CRI/CRI\\_2\\_5x.asp?dt=5](http://www.cancer.org/docroot/CRI/CRI_2_5x.asp?dt=5)  
531
- 532 [21] D. P. Weller and C. Campbell, “Uptake in cancer screening programmes:  
533 A priority in cancer control,” *Brit. J. Cancer*, vol. 101, pp. 55–59, 2009.
- 534 [22] Y. Zheng, “Breast cancer detection with gabor features from digital mam-  
535 mograms,” *Algorithms*, vol. 3, pp. 44–62, 2010.
- 536 [23] K. W. Eilbert, K. Carroll, J. Peach, S. Khatoun, I. Basnett, and N. Mc-  
537 Culloch, “Approaches to improving breast screening uptake: Evidence  
538 and experience from Tower Hamlets,” *Brit. J. Cancer*, vol. 101, no. 2,  
539 pp. 64–67, 2009.
- 540 [24] D. Schopper and C. de Wolf, “How effective are breast cancer screening  
541 programmes by mammography? Review of the current evidence,” *Eur. J.*  
542 *Cancer*, vol. 45, no. 11, pp. 1916–1923, Jul. 2009.
- 543 [25] E. S. Cassandra, “Breast cancer screening: Cultural beliefs and diverse  
544 populations,” *Health Soc. work*, vol. 31, no. 1, pp. 36–43, 2006.
- 545 [26] NHS Cancer Screening Programmes. (2007, Apr.) *Disclosure of Audit*  
546 *Results in Cancer Screening Advice on Best Practice* (Cancer  
547 Screening Series 3), J. Patnick, Ed. [Online]. Available: <http://www.cancerscreening.nhs.uk/publications/cs3.pdf>  
548
- 549 [27] K. Okane. (2005, Apr. 20). *Mumps Language Bioinformatic*  
550 *Database Resources* [Online]. Available: [http://bioinformatics.org/forums/forum.php?forum\\_id=1035](http://bioinformatics.org/forums/forum.php?forum_id=1035)  
551
- 552 [28] V. Baskaran, R. K. Bali, R. N. G. Naguib, and H. Arochena, “A Knowl-  
553 edge Management approach to increase uptake in a breast screening pro-  
554 gramme,” presented at the IEEE 2nd Humanoid, Nanotechnology, In-  
555 formation Technology, Communication and Control, Environment and  
556 Management (HNICEM) Int. Conf., Philippines, Mar. 2005.
- 557 [29] S. Tarver, K. Cronin-Cowan, and M. E. Wheaton, “A pilot’s life for us,”  
558 *Breast Cancer Res.*, vol. 6, suppl. 1, p. 52, 2004.
- 559 [30] H. E. Arochena, “Modelling and prediction of parameters affecting atten-  
560 dance to the NHS breast cancer screening programme,” Ph.D. dissertation,  
561 Dept. Comp. Sci., Coventry Univ., Coventry, U.K., 2003.
- 562 [31] C. Bankhead, S. H. Richards, T. Peters, D. Sharp, R. Hobbs, J. Brown,  
563 L. Roberts, C. Tydeman, V. Redman, J. Formby, S. Wilson, and J. Austoker,  
564 “Improving attendance for breast screening among recent non-attenders:  
565 A randomised controlled trial of two interventions in primary care,” *J.*  
566 *Med. Screening*, vol. 8, no. 2, pp. 99–105, 2001.



**Vikraman Baskaran** is currently an Assistant Professor at the School of Information Technology Management of Ryerson University, Toronto, ON, Canada. His research interests include finding a viable application of the KM paradigm in healthcare application. His special interest in developing HL7 messaging and health informatics has provided opportunities in excelling in these fields. His current activities include KM, e-health, artificial intelligence, and healthcare informatics.

He is a member of the HL7 U.K. and Canada.



**Aziz Guergachi (M'xx)** is currently an Associate Professor at the Ted Rogers School of Information Technology Management of Ryerson University, Toronto, ON, Canada. Prior to becoming part of the Ryerson community, he was involved in the development of a large software system for trade promotion management and collaborative sales forecasting. His current research interests include advanced system modeling and machine learning with applications to business management and engineering systems.

He is the recipient of the New Opportunities Award of the Canada Foundation for Innovation and currently runs a research laboratory for advanced systems modeling.



**Rajeev K. Bali (SM'xx)** is currently a Reader in Healthcare Knowledge Management at Coventry University, U.K. His main research interests include clinical and healthcare knowledge management (from both technical and organisational perspectives). He has published peer-reviewed journals and is the an author/editor of several textbooks on healthcare knowledge management.

He serves on various editorial boards and conference committees and is regularly invited to deliver presentations and speeches internationally.



**Raouf N. G. Naguib (SM'xx)** is currently a Professor of Biomedical Computing and Head of BIOCORE, Coventry, U.K. Prior to this appointment, he was a Lecturer at Newcastle University, Newcastle Upon Tyne, U.K. He has published more than 240 journals and conference papers and reports in many aspects of biomedical and digital signal processing, image processing, artificial intelligence, and evolutionary computation in cancer research.

He was awarded the Fulbright Cancer Fellowship in 1995–1996 when he carried out research at the University of Hawaii, Mānoa, on the applications of artificial neural networks in breast cancer diagnosis and prognosis. He is a member of several national and international research committees and boards.

## QUERIES

620

- Q1. Author: Please check whether the edits made in the sentence “This large percentage of nonattendance not only . . .” retain your intended sense. 621
- Q2. Author: Refs. [5], [6], [8], [12], [14], [15], [16], and [23] are not cited in the text. Please check and provide citations. 622
- Q3. Author: Please provide the expansion of KM. 623
- Q4. Author: Please provide the educational details of all the authors. 624
- Q5. Author: Please provide the year in which Aziz Guergachi became “Member” of the IEEE. 625
- Q6. Author: Please provide the year in which Rajeev K Bali became “Senior Member” of the IEEE. 626
- Q7. Author: Please provide the year in which Raouf N. G. Naguib became “Senior Member” of the IEEE. 627
- 628

IEEE  
Proof