

1 A Machine Learning Approach to Predict Health- 2 care Cost of Breast Cancer Patients

3 Pratyusha Rakshit, Onintze Zaballa-Larumbe, Aritz Pérez, Elisa Gómez-
4 Inhiesto, Maria T. Acaiturri-Ayesta and Jose A. Lozano

5 Basque Center for Applied Mathematics, Spain

6 {prakshit, ozabella, aperez, jlozano}@bcamath.org

7 Osakidetza, Spain

8 {mariaelisa.gomezinhiesto, mariateresa.acaiturriayesta}@osakidetza.eus

9 **Abstract.** This paper presents a novel machine learning approach to per-
10 form an early prediction of the healthcare cost of breast cancer patients. The
11 learning phase of our prediction method considers the following two steps: i) in
12 the first step, the patients are clustered taking into account the sequences of ac-
13 tions undergoing similar clinical activities and ensuring similar healthcare costs,
14 and ii) a Markov chain is then learned for each group to describe the action-
15 sequences of the patients in the cluster. A two step procedure is undertaken in
16 the prediction phase: i) first, the healthcare cost of a new patient's treatment
17 is estimated based on the average healthcare cost of its k -nearest neighbors in
18 each group, and ii) finally, an aggregate measure of the healthcare cost estimated
19 by each group is used as the final predicted cost. Experiments undertaken reveal
20 a mean absolute percentage error as small as 6%, even when half of the clinical
21 records of a patient is available, substantiating the early prediction capability of
22 the proposed method. Comparative analysis substantiates the superiority of the
23 proposed algorithm over the state-of-the-art techniques.

24
25 **Keywords:** Healthcare Cost, Clustering, Markov Chain, k Nearest
26 Neighbor.

27 1 Introduction

28 An *electronic health record* (EHR) is an electronic version of a patient's clinical
29 history over time. It comprises all administrative clinical data of a patient in a
30 healthcare organization, including his/her demographics, diagnosis, medications,
31 laboratory data, and associated costs, and so on. The plethora of longitudinal
32 patients' data of an EHR can be utilized for developing patient-centered per-
33 sonalized healthcare solutions, including cost. It is however worth mentioning
34 that the healthcare costs, ranging from clinician's fees to the cost of hospital
35 stays and medicines, are escalating at a rapid rate around the world [1] [2]. It
36 has motivated the researchers to take keen interest in controlling this upsurge

37 in the healthcare costs. The crucial step to control the healthcare cost is to
38 enable the healthcare organizations to predict the possible future cost of indi-
39 vidual patients. It in turn helps to identify the individuals at the highest risk of
40 enduring the significant costs in future. It thus helps to prioritize the allocation
41 of scarce resources among the patients in a healthcare organization for efficient
42 care management.

43 Moreover, a report from The Commonwealth Fund (2012) emphasizes the
44 need to identify high-cost patients as the first step towards achieving “rapid
45 improvements in the value of services provided” [22]. A proactive approach to
46 address this problem is to identify patients who are at risk of becoming high-cost
47 patients accurately before substantial unnecessary costs have been incurred and
48 health condition has deteriorated further. Eventually, this calls for prediction
49 of possible total healthcare cost of a patient as early as possible when a limited
50 volume of clinical records of the given patient is provided. In other words,
51 another important aspect in the context of healthcare cost prediction is to devise
52 a model using a training set of complete clinical records of some patients to
53 predict the total healthcare cost of a new patient as accurately and also as
54 early as possible, preferably before the availability of the patient’s full-length
55 clinical record. Such early prediction of future healthcare cost can be used to
56 judiciously identify high-risk high-cost patients and prevent crises in healthcare
57 organizations. It is obvious that the earliness of the prediction may affect the
58 accuracy. It has motivated the researchers to build a model to predict healthcare
59 cost as early as possible while maintaining an appropriate level of accuracy.

60 Nevertheless, healthcare cost prediction based on individual patient’s char-
61 acteristics is a challenging issue from the data mining perspective due to the
62 non-Gaussian skewed distribution of the cost data of the patients [5]. Studies in
63 [6], [7] reveal dubious efficacy of the statistical methods to predict the healthcare
64 cost. Furthermore, the traces of linear regression and rule-based approaches are
65 also found in literature [2], [7] for the cost prediction. But the requirement of
66 a lot of domain knowledge has restricted their applications for most of the real
67 world economic data of the patients [8]. Now-a-days, machine learning algo-
68 rithms, including clustering and classification techniques, have emerged as an
69 alternative effective tool for this purpose [9], [10].

70 This paper proposes a machine learning based novel approach for healthcare
71 cost prediction of individual patient’s treatments based on their clinical actions,
72 jointly including the clinical activities and the respective cost over time. The
73 activity here represents diagnosis, medication, pharmacy and the like. A two-
74 step procedure is employed in the learning phase: i) in the first step, the ordered
75 sequences of clinical actions of the patients’ treatments are clustered using the
76 hierarchical DBSCAN [15] with an aim to identify the group of patients under-
77 taking similar clinical activities and incurring similar healthcare costs, and ii)
78 each group is then modelled by means of a Markov chain [11] delineating the
79 probability distributions of transitions between different clinical actions. A new

80 distance measure is also proposed to measure the similarity of the treatment
81 patterns of the patients during clustering.

82 The prediction phase, concerned with prediction of the healthcare cost of the
83 sequence of clinical actions of a new patient’s treatment, also encompasses two
84 steps: i) first, for each group, we compute a tentative cost of the new sequence
85 by averaging the cost of its k -nearest neighbor [12] sequences in the group, ii)
86 the final cost is obtained as a weighted sum of the cost estimated by each of the
87 groups. The weights for each group are the likelihood of the new sequence to
88 the respective group as assigned by the corresponding Markov chain.

89 The performance of the proposed healthcare cost prediction algorithm is eval-
90 uated with the economic information together with information of the clinical
91 activities of the breast cancer patients obtained from the health administrative
92 department of the public health care system of the Basque Country, Spain. A 10-
93 fold cross validation is employed with the training dataset resulting the optimal
94 value of k of k -NN as three in the present application with respect to the mean
95 absolute percentage error ($MAPE$) [2]. Moreover, the proposed method results
96 in an $MAPE$ measure of less than 6% when half of the clinical records of a new
97 patient is available, irrespective of the value of k . It substantiates the capability
98 of the proposed stratagem for early prediction of healthcare cost. Experiments
99 undertaken also reveal that the proposed algorithm outperforms its state-of-the-
100 art contenders with respect to $MAPE$ metric. The comparative analyses verify
101 the significance of jointly considering the clinical activity and the associated cost
102 data to effectively capture the clinical records of patients for accurate healthcare
103 cost prediction as early as possible.

104 The paper is divided into following sections. Section 2 delineates the proposed
105 method of healthcare cost prediction. Experiments undertaken and the results
106 are reported in Section 3. Section 4 concludes the paper.

107 2 Method

108 2.1 Data Transformation

109 This section refers to transforming the database of individual patient’s treat-
110 ments into a series of actions, sorted by time. Here, we provide some definitions
111 which will be used throughout the paper to develop a solution to the healthcare
112 cost prediction problem.

113 **Definition-1: Action.** Let \mathbf{X} be the set of all clinical activities, including
114 diagnosis, procedure, medicine and the like, $\mathbf{Y} \in \mathbb{R}$ be the set of all possible
115 incurred healthcare cost as recorded in the database and \mathbf{T} be the set of visiting

116 times of the patients to the hospital. An *action*, say a , is then expressed as a
 117 three-tuple, given by

$$a = \{(x, y, t) \mid \forall x \in \mathbf{X}, \forall y \in \mathbf{Y}, \forall t \in \mathbf{T}\}. \quad (1)$$

118 **Definition-2: Patient's treatment.** A *patient's treatment* is defined by a
 119 sequence of its corresponding actions, sorted by the visiting time. Symbolically,
 120 a patient's treatment P is represented by

$$P = (a_1, a_2, \dots, a_n) \quad (2)$$

121 where $a_i = (x_i, y_i, t_i)$ represents the action encompassing the clinical activity
 122 $x_i \in \mathbf{X}$ and its respective healthcare cost $y_i \in \mathbf{Y}$ incurred during visiting time
 123 $t_i \in \mathbf{T}$ of the specific patient. For sake of simplicity of readers, we drop the
 124 notion of visiting time and hence a_i now can be simplified as

$$a_i = \{(x_i, y_i) \mid x \in \mathbf{X}, y \in \mathbf{Y}\}. \quad (3)$$

125 The clinical actions of P in (2) are chronologically ordered. Evidently, if
 126 $i < j$, a_i occurs before a_j . A sequence of actions of a patient's treatment is used
 127 to jointly track the progression of its activity-outcome and the corresponding
 128 healthcare cost over time. The length of the sequence varies across patients
 129 because of the diversity in their treatments over time.

130 **Definition-3: Modified cost.** Intuitively, the number of possible actions
 131 for all patients in the database is huge due to infinite number of healthcare cost
 132 elements in \mathbf{Y} . For the sake of simplicity, \mathbf{Y} is reduced to a finite set in a two
 133 step procedure described below.

134 1) *Discretization*: First, the entire range of \mathbf{Y} is discretized into n_s segments
 135 defined by the n_s -quantiles of \mathbf{Y} . In other words, we set the lower and the
 136 upper limit of the i -th segment respectively to the $(i - 1)$ -th quantile and the
 137 i -th quantile of the healthcare cost elements for all possible clinical activities,
 138 recorded in the database.

139 2) *Quantization*: Then a real healthcare cost element, lying in the i -th seg-
 140 ment is replaced by the mean value of all cost elements of the i -th segment.

141 The strategy is pictorially demonstrated in Fig. 1 for the healthcare cost
 142 information of two patients only with $n_s = 8$. The setting of $n_s = 8$ and the cost
 143 values used here are illustrative examples only. The healthcare cost, referred
 144 henceforth, denotes the modified cost.

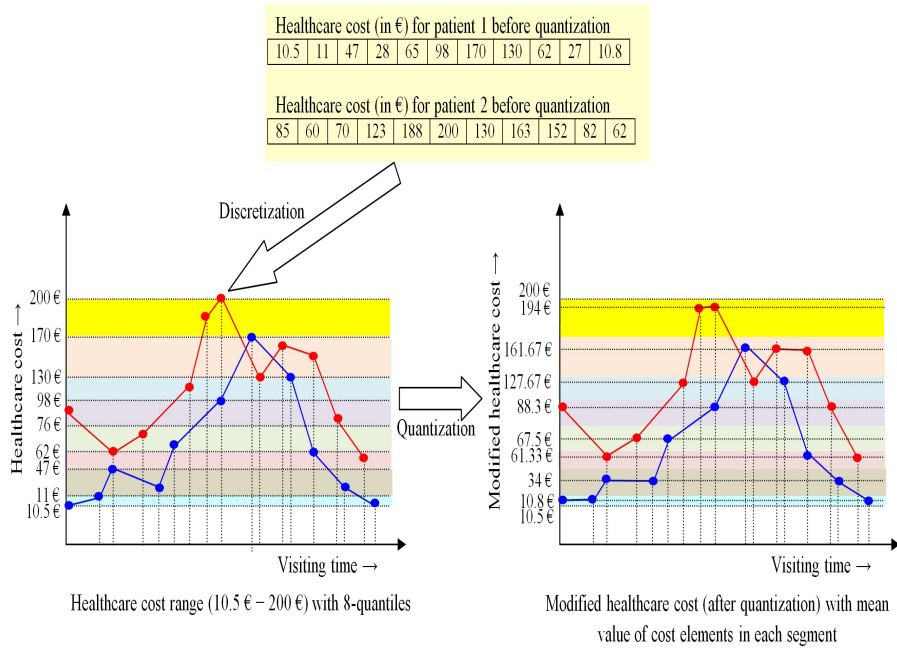


Figure 1: Calculation of modified healthcare cost of two patients with 8-quantiles

145 2.2 Clustering Patients' Action-Sequences

146 It is noteworthy that patients undergoing various clinical activities reveal consid-
 147 erable diversity of their corresponding cost information. Hence, prior to predict
 148 cost of a new action-sequence, we cluster the action-sequences of the existing
 149 patients into groups. We then consult the cost information of the specific group
 150 of patients providing the maximum similarity with the action-sequence of the
 151 new patient to predict the respective possible future cost.

152 Two significant issues to categorize the patients based on their action se-
 153 quences include: i) design of an appropriate distance measure to capture the
 154 similarity between action-sequences of varying length, and ii) selection of an ef-
 155 ficient clustering algorithm to ensure that action-sequences within a group are

156 similar to each other than those in other groups.

157 **Design of distance measure:** There exists plethora of literature on us-
158 ing *edit distance* [13] to measure the dissimilarity of two strings of characters
159 (or words). Given two strings S_1 and S_2 over a finite alphabet, an edit dis-
160 tance $ED(S_1, S_2)$ between S_1 and S_2 can be defined as the minimum cost of
161 transforming S_1 to S_2 through a sequence of weighted edit operations. These
162 operations primarily include insertion, deletion, and substitution of one symbol
163 by another. Usually, the edit operations are assigned with equal weights of unity.
164 Nevertheless, the string in this paper denotes the action-sequences.

165 However, there is a major limitation of using the conventional ED directly
166 in the present context. The conventional ED compares two strings of characters
167 (or words) only. In the present work, the components of the string (or action-
168 sequence) is not only representing character (symbolizing a clinical activity)
169 but an activity-cost pair. Hence, application of the conventional ED in the
170 present scenario captures the difference between two action-sequences based on
171 their respective clinical activities only, ignoring the corresponding healthcare cost
172 information. It thus loses the cost information and the temporal relationship of
173 the activity-cost pairs over time.

174 Consequently, the clusters of patients based on the conventional ED measures
175 identify patients ensuring similar clinical activities only. Evidently, the accuracy
176 of the healthcare cost prediction based on the clusters, thus formed, is reduced to
177 great extent. It has motivated us to design an appropriate distance measure to
178 jointly capture the dissimilarity of two clinical activities (of two different action
179 sequences) and their respective healthcare costs.

180 The proposed distance measure, referred to as *treatment pattern difference*
181 (TPD) is an extended version of the conventional ED . In case of the conventional
182 ED , all possible edit operations are associated with equal cost of unity. In TPD ,
183 the edit costs are modified as follows to consider the healthcare cost components
184 of two action-sequences.

185 Let P_1 and P_2 be two different action-sequences. The cost of insertion of a
186 clinical activity x_i (or a character) to convert P_2 to P_1 is given by

$$C_1 = y_i \tag{4}$$

187 where y_i denotes the healthcare cost of the clinical activity x_i at the visiting
188 time t_i in the action-sequence P_1 . Similarly, the cost of deleting an action x_j
189 from P_1 to covert it to P_2 is given by

$$C_2 = y_j \tag{5}$$

190 where the symbols carry their usual meanings. If the clinical activity x_i of P_1
 191 is substituted with a different clinical activity x_j of P_2 , the corresponding edit
 192 cost is given by

$$C_3 = |y_i - y_j + \epsilon|. \quad (6)$$

193 Here ϵ is a small positive constant. It is used to ensure that even when $y_i = y_j$
 194 for $x_i \neq x_j$, at least $C_3 = \epsilon$ is used as the edit cost for substitution of x_i by x_j .

195 It is noteworthy that if $x_i = x_j$, the conventional *ED* gives a zero penalty.
 196 However, there are instances of different healthcare costs for the same clinical
 197 activity of two different patients. To capture this, *TPD* uses an additional edit
 198 cost, given by

$$C_4 = |y_i - y_j|. \quad (7)$$

199 Hence, the total edit cost to convert an action-sequence P_1 to another action-
 200 sequence P_2 is given by


$$TPD(P_1, P_2) = w_1 \times \left(\sum_{\forall ins.} C_1 + \sum_{\forall del.} C_2 + \sum_{\forall sub.} C_3 \right) + w_2 \times \sum_{\forall match} C_4. \quad (8)$$

201 Here, w_1 and w_2 denote the weight for the edit operations respectively for
 202 different and similar activities. Intuitively, $w_2 < w_1$ as it corresponds to the
 203 penalty corresponding to similar activities with different healthcare cost. After
 204 a wide experimentation, we set $w_1 = 0.7$ and $w_2 = 0.3$. An example of evaluating
 205 the dissimilarity of two action-sequences based on the *TPD* measure is presented
 206 in Fig. 2.

Action-sequence $P_1 = (\{G, 10 \text{ €}\}, \{C, 69 \text{ €}\}, \{A, 25 \text{ €}\}, \{T, 53 \text{ €}\}, \{G, 10 \text{ €}\}, \{C, 97 \text{ €}\}, \{U, 25 \text{ €}\})$

Action-sequence $P_2 = (\{G, 10 \text{ €}\}, \{A, 25 \text{ €}\}, \{T, 53 \text{ €}\}, \{T, 53 \text{ €}\}, \{A, 30 \text{ €}\}, \{C, 69 \text{ €}\}, \{A, 25 \text{ €}\})$

Edit operations to convert P_1 to P_2 based on clinical activities only

P_1 to		{G, 10 €}	{C, 69 €}	{A, 25 €}	{T, 53 €}	-	{G, 10 €}	{C, 97 €}	{U, 25 €}
P_2		{G, 10 €}	-	{A, 25 €}	{T, 53 €}	{T, 53 €}	{A, 30 €}	{C, 69 €}	{A, 25 €}

Position	Edit operations	Edit cost
1	Match (G)	$C_4 = 10-10 = 0$
2	Deletion (C)	$C_2 = 69$
3	Match (A)	$C_4 = 25-25 = 0$
4	Match (T)	$C_4 = 53-53 = 0$
5	Insert (T)	$C_1 = 53$
6	Substitution (G by A)	$C_3 = 10-30+1 = 19$
7	Match (C)	$C_4 = 97-69 = 28$
8	Substitution (U by A)	$C_3 = 25-25+1 = 1$

$$TPD(P_1, P_2) = 0.7 \times (69 + 53 + 19 + 1) + 0.3 \times 28 = 107.8$$

Figure 2: Calculation of TPD of two action sequences

207 **Selection of clustering algorithm:** The TPD measures of each pair of pa-
 208 tients' treatments in the given record are used to cluster the similar sequences in
 209 the same subgroups. The *hierarchical density-based spatial clustering of applica-*
 210 *tions with noise* (hierarchical DBSCAN) algorithm [15] is employed to identify
 211 the groups of patients' treatments. The selection of DBSCAN in the present
 212 context is justified because of its merit of clustering similar data points (here,
 213 the action-sequences of patients) into same groups based on the density (number
 214 of nearby neighbors) without prior setting of the number of clusters. Moreover,
 215 unlike the traditional partitioning algorithms, DBSCAN can be applied for clus-
 216 ters of arbitrary shape, even when the data may be contaminated with noise
 217 [16].

218 It is however worth mentioning that the huge economic database includes
 219 clusters of records of patients characterized at different density levels. The tra-
 220 ditional DBSCAN algorithm with a single global density threshold often fails to
 221 effectively identify such clusters. This impasse is overcome here by using the
 222 hierarchical DBSCAN, proposed in [15], which discovers all DBSCAN-identified
 223 clusters for an infinite range of density thresholds. Finally, it identifies a simpli-
 224 fied hierarchical structure of significant clusters only.

2.3 Markov Chain Representation of a Cluster

This step is concerned with representing each cluster of patients' action-sequences by a Markov chain [11]. The crux of such representation is founded on the underlying premise that the medical practitioners take their decision based on the previous clinical activities. Again, our cost prediction algorithm greatly relies on the recorded action-sequence of a patient.

A first order Markov chain exhibits memoryless property where the current state only depends on the previous state. Let N be the possible number of actions (activity-cost pairs) in the database. The Markov chain model of a group of patients, say G_l , is then demonstrated by a state-transition probability distribution, which is denoted as:

$$M_l = [m_{i,j,l}] \text{ for } i, j = 1, 2, \dots, N \quad (9)$$

$$\text{where } m_{i,j,l} = p_l(x_{t+1} = s_j | x_t = s_i) = \frac{q_{i,j,l}}{\sum_{k=1}^N q_{i,k,l}}. \quad (10)$$

Here $q_{i,j,l}$ and $p_l(x_{t+1} = a_j | x_t = a_i)$ respectively denote the number of cases and the probability of transition from the current action $x_t = a_i$ to the immediate next action $x_{t+1} = a_j$ in the specific group G_l of action-sequences. Evidently, it satisfies

$$m_{i,j,l} \geq 0 \text{ and } \sum_{j=1}^N m_{i,k,l} = 1. \quad (11)$$

In addition to M_l , we also evaluate the initial probability $p_l(a_i)$ of action a_i considering all the action-sequences in the group G_l for $i = 1, 2, \dots, N$ as follows.

$$p_l(a_i) = \frac{s_{i,l}}{\sum_{k=1}^n s_{k,l}} \quad (12)$$

Here $s_{i,l}$ denotes the number of action-sequences initiated with the action a_i in G_l for $i = 1, 2, \dots, N$. This entire process is repeated for all groups identified by the hierarchical DBSCAN.

245 **2.4 Cost Prediction of a Patient's Treatment from Action**
 246 **Sequence**

247 The aim of this step is to predict the possible total cost of a patient from the
 248 respective action-sequence. The action-sequence of the patient is formed follow-
 249 ing the principle given in Section 2.1. Let the ordered sequence of actions of the
 250 new patient's treatment be denoted by $P = (a_1, a_2, \dots, a_n)$ where the action a_i
 251 represents the activity-cost pair at the visiting time instant t_i . The prediction
 252 of future cost based on P is undertaken in three phases.

253 **Phase-1: Cost estimation of P based on a specific group.** We employ
 254 k -nearest neighbor (k -NN) to identify k action-sequences from a group, say G_l ,
 255 that offer maximum similarity with P based on TPD measure as given in (8).
 256 First, we compute the TPD values between P and each member sequence of
 257 the group G_l . The member sequences are then sorted in ascending order of their
 258 TPD measures thus evaluated. The first k members are selected as the k nearest
 259 neighbors of P . Next, each of the k members is assigned a weight $w_{j,l}$, inversely
 260 proportional to its TPD measure from P for $j = 1, 2, \dots, k$. Consequently, the
 261 total cost $\hat{c}_l(P)$ of the new action-sequence P estimated by the group G_l is given
 262 by

$$\hat{c}_l(P) = \frac{\sum_{j=1}^k w_{j,l} \times c_{j,l}}{\sum_{j=1}^k w_{j,l}}. \quad (13)$$

263 Here $c_{j,l}$ denotes the total cost incurred by the j -th nearest neighbor of P in
 264 G_l for $j = 1, 2, \dots, k$. $\hat{c}_l(P)$ is computed for all clusters of patients identified by
 265 the hierarchical DBSCAN.

266 **Phase-2: Evaluation of the likelihood of P to patients' groups.** This
 267 step is concerned with evaluating the likelihood of P to each subgroup of patients
 268 based on the respective Markov chain model. The likelihood of the ordered
 269 sequence of actions $P = (a_1, a_2, \dots, a_n)$ to a specific group G_l is given by

$$\lambda_l(P) = p_l(a_1) \times \prod_{i=1}^{n-1} p_l(a_{i+1}|a_i). \quad (14)$$

270 Here a_1 denotes the initial action of P and a_i represents the action of P
 271 occurred at visiting time t_i for $i = 1, 2, \dots, n$. Evidently, $p_l(a_1)$ and $p_l(a_{i+1}|a_i)$
 272 respectively symbolize the initial probability of action a_1 and the probability of
 273 transition from the current action a_i to the immediate next action a_{i+1} of P
 274 as described by the group G_l . Expression (14) is evaluated using the Markov chain
 275 model \mathbf{M}_l representing the group G_l .

276 After evaluating $\lambda_l(P)$ for all groups, the normalized likelihood of P to each
 277 subgroup is computed using

$$\hat{\lambda}_l(P) = \frac{\lambda_l(P)}{\sum_{\forall k} \lambda_k(P)}. \quad (15)$$

278 **Phase-3: Cost prediction based on all groups.** After evaluating the
 279 estimated cost and the normalized likelihood of P to all groups, the total cost
 280 of P is finally predicted following

$$\bar{c}(P) = \sum_{\forall l} \hat{\lambda}_l(P) \times \hat{c}_l(P). \quad (16)$$

281 3 Results

282 3.1 Database

283 The study is performed on the economic data, along with the clinical activi-
 284 ties of the patients obtained from the health administrative department of the
 285 public health care system (OSAKIDETZA) of the Basque Country, Spain. The
 286 database includes medical history of 579798 patients treated in different lev-
 287 els of healthcare organizations (including 1 hospital, 11 outpatients clinics and
 288 emergency care) from January 1, 2017 to December 31, 2019. The clinical data
 289 of the patients primarily consists of their clinical assistance and the respective
 290 healthcare cost information.

291 To validate the proposed method of cost prediction, the present work con-
 292 siders the pool of breast cancer patients only. The selection of breast cancer
 293 patients from the database conforms the International Statistical Classification
 294 of Diseases and Related Health Problems (10-th revision) [17], stating that ev-
 295 ery code starting by C50 corresponds to breast cancer diagnosis. A few filtering
 296 steps are then carried out following [21] to judiciously select the pool of patients
 297 of interest. The filtering process affirms that the selected patients have their
 298 complete treatment in the above-mentioned time period of two years. Following
 299 the medical guideline, a final set of 972 patients is identified. 70% of the entire
 300 database is ultimately used as the training dataset, while the remaining as the
 301 test data. A 10-fold cross validation is employed on the training dataset for
 302 judicious selection of the value of k for k -NN.

303 **3.2 Identification and Representation of Patients' Action-**
 304 **Sequences**

305 The final record of the 464 patients consists of 23 unique clinical activities as
 306 described in Table 1. The healthcare cost is next discretized into n_s segments.
 307 In Fig. 3, we present a plot of normalized quantization error values for different
 308 settings of the number of quantiles n_s , varied from 2 to 12 to check a significant
 309 improvement in performance. The normalized quantization error (NQE) is given
 310 by (17).

$$NQE = \frac{\frac{1}{N_c} \sum_{i=1}^{N_c} |c(i) - c_m(i)|}{\max_{i=1}^{N_c} c(i) - \min_{i=1}^{N_c} c(i)} \quad (17)$$

311 Here $c(i)$ and $c_m(i)$ respectively denote the true and the modified i -th health-
 312 care cost (after discretization) of the database with N_c cost elements for $i = 1,$
 313 $2, \dots, N_c$. Fig. 3 reveals that the quantization error is reduced with an increase
 314 in the number of segments n_s . However, it is also observed that there is no
 315 significant change in the error for $n_s \geq 8$. We have thus fixed $n_s = 8$. It is worth
 316 mentioning that the setting of n_s here is biased to the healthcare cost values
 317 of the present database. The quantization of the healthcare cost range of the
 318 present database using 8-quantiles ensures a balanced number of healthcare cost
 319 elements in each of the eight cost-segments.

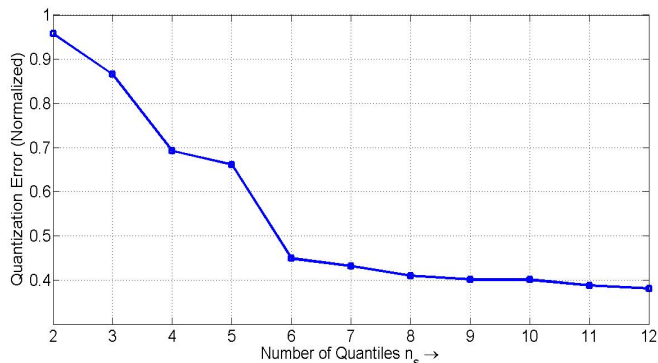


Figure 3: Normalized quantization error for different values of n_s

320 Next, the healthcare cost of all clinical activities of 464 patients is discretized
 321 in eight segments based on 8-quantiles of the healthcare cost range, as demon-
 322 strated in Fig. 1. Let the segments (sorted in ascending order) be denoted as
 323 *very-very-low (VVL)*, *very-low (VL)*, *low (L)*, *medium-low (ML)*, *medium-high*
 324 *(MH)*, *high (H)*, *very-high (VH)* and *very-very-high (VVH)*. Eventually, there
 325 exist $22 \times 8 = 176$ actions to jointly represent a pair of clinical activity and the

326 corresponding healthcare cost. However, a close scrutiny of the final record re-
 327 veals only 63 possible pairs from the recorded medical history of the 464 patients,
 328 as reported in Table 2.

329 The hierarchical DBSCAN algorithm is then employed on the training dataset
 330 to cluster the sequences using *TPD* values. The algorithm results in eight clus-
 331 ters. The clusters thus identified are pictorially represented in Fig. 4. The
 332 descriptions of the actions of the sequences, shown in different colors, are tabu-
 333 lated in Table 2. Each cluster is then described by a Markov chain following
 334 Section 2.3.

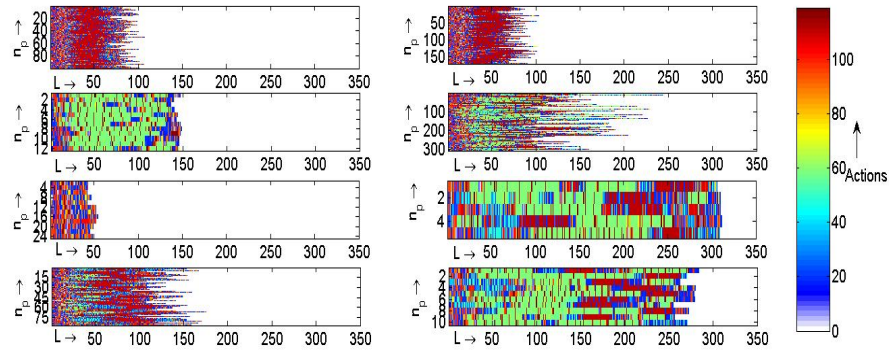


Figure 4: Cluster of sequences of visit records (activity-cost pairs) of patients with n_p as number of patients and L as the length of the sequence

335 3.3 Performance Evaluation of Proposed Healthcare Cost 336 Prediction Method

337 **Performance metric:** The performance of the proposed cost prediction algo-
 338 rithm is evaluated with respect to *mean absolute percentage error (MAPE)* with
 339 a lower error indicating a better performance.

$$MAPE = \frac{\frac{1}{N_t} \sum_{i=1}^{N_t} |c(P_i) - \bar{c}(P_i)|}{\frac{1}{N_t} \sum_{i=1}^{N_t} c(P_i)} \times 100. \quad (18)$$

340 Here $c(P_i)$ and $\bar{c}(P_i)$ (evaluated using (16)) respectively represent the true
 341 and the predicted cost of the i -th patient's treatment P_i in the test dataset with
 342 N_t records for $i = 1, 2, \dots, N_t$.

Table 1. Description of the clinical activities

Activity	Abbreviated form	Full form
1	ANES	Aesthesia
2	APAT	Pathological Anatomy
3	CEXT	External Consultation
4	CONS	Consultation
5	FAMB	Hospital Pharmacy Services
6	FAMR	Pharmacy
7	HCRI	Critical Care Hospitalization
8	HDIA	Day Hospital
9	HDOM	Home Hospitalization
10	HOSP	Hospitalization
11	INCO	Interconsultation
12	LABO	Laboratory
13	MNUC	Nuclear Medicine
14	OSAT	Osatek (Magnetic Resonance Service)
15	PFUN	Functional Testing
16	QUIR	Surgery Unit
17	RADI	Radiology
18	REHA	Rehabilitation
19	RTER	Radiotherapy
20	UCRI	Nursing Critical Care Unit
21	UCSI	Surgery without Hospitalization
22	UENF	Nursing Unit
23	URP	Post Anesthesia Care Unit

344 **Validation of Earliness Prediction and Selection of k of k -NN:** The
345 capability of the proposed algorithm to predict the possible total healthcare cost
346 of patients is verified by varying the length of sequence of the recorded treat-
347 ments of the patients from 20% to 100%. The appropriate selection of k (of
348 k -NN) for the optimal performance is undertaken using 10-fold cross validation
349 on the training dataset. The *MAPE* values for different settings of k and per-
350 centage of length of sequence of the recorded treatments of the patients using
351 10-fold cross validation are tabulated in Table 3 (for the training data). Table
352 3 reveals that the longer the length of the sequence, the better is the prediction
353 accuracy with smaller *MAPE* measures, irrespective of the setting of k . The
354 optimal performance of the method is obtained for $k = 3$ with the entire se-

355 quence information. It is also noted that an *MAPE* smaller than 6% is obtained
 356 even when 50% of a visit sequence is utilized. It proves the effectiveness of the
 357 proposed method for an early prediction of the healthcare cost.

358 **Table 2A.** Description of the clinical actions (activity-cost pairs)

Action	Activity	Cost	Action	Activity	Cost	Action	Activity	Cost
1	ANES	VVL	26	FAMB	L	51	HDIA	VVH
2	ANES	VL	27	FAMB	ML	52	HDOM	VVL
3	ANES	L	28	FAMB	H	53	HDOM	VL
4	ANES	ML	29	FAMB	VH	54	HDOM	L
5	ANES	MH	30	FAMB	VVH	55	HDOM	ML
6	ANES	H	31	FAMR	VVL	56	HDOM	H
7	ANES	VH	32	FAMR	VL	57	HDOM	VH
8	ANES	VVH	33	FAMR	L	58	HDOM	VVH
9	APAT	VVL	34	FAMR	ML	59	HOSP	VVL
10	APAT	L	35	FAMR	MH	60	HOSP	VL
11	APAT	VH	36	FAMR	H	61	HOSP	L
12	APAT	VVH	37	FAMR	VH	62	HOSP	ML
13	CEXT	VL	38	FAMR	VVH	63	HOSP	MH
14	CEXT	L	39	HCRI	VVL	64	HOSP	H
15	CEXT	ML	40	HCRI	VL	65	HOSP	VH
16	CEXT	H	41	HCRI	L	66	HOSP	VVH
17	CONS	VVL	42	HCRI	VH	67	INCO	L
18	CONS	L	43	HCRI	VVH	68	INCO	ML
19	CONS	ML	44	HDIA	VVL	69	INCO	MH
20	CONS	MH	45	HDIA	VL	70	INCO	VH
21	CONS	H	46	HDIA	L	71	INCO	VVH
22	CONS	VH	47	HDIA	ML	72	LABO	L
23	CONS	VVH	48	HDIA	MH	73	LABO	MH
24	FAMB	VVL	49	HDIA	H	74	LABO	VH
25	FAMB	VL	50	HDIA	VH	75	LABO	VVH

Table 2B. Description of the clinical actions (activity-cost pairs)

Action	Activity	Cost	Action	Activity	Cost	Action	Activity	Cost
76	MNUC	L	91	QUIR	MH	106	REHA	ML
77	MNUC	ML	92	QUIR	H	107	REHA	MH
78	MNUC	H	93	QUIR	VH	108	RTER	VVL
79	MNUC	VH	94	QUIR	VVH	109	RTER	MH
80	MNUC	VVH	95	RADI	VVL	110	RTER	H
81	OSAT	L	96	RADI	VL	111	RTER	VH
82	OSAT	H	97	RADI	L	112	RTER	VVH
83	OSAT	VH	98	RADI	ML	113	UCRI	VVH
84	OSAT	VVH	99	RADI	MH	114	UCSI	VH
85	PFUN	VVL	100	RADI	H	115	UENF	MH
86	PFUN	VL	101	RADI	VH	116	UENF	H
87	PFUN	L	102	RADI	VVH	117	UENF	VH
88	QUIR	VL	103	REHA	VVL	118	UENF	VVH
89	QUIR	L	104	REHA	VL	119	URP	ML
90	QUIR	ML	105	REHA	L			

360 **Comparative performance analysis:** The next experiment aims at com-
361 parative performance analysis of our proposed algorithm. Three state-of-the-
362 art techniques are considered in the comparative framework, including *gradient*
363 *boosting* (GB) [18], *artificial neural net* (ANN) [19] and *lasso* [20]. These existing
364 methods have utilized the healthcare cost data only to predict the future cost
365 [2]. The *MAPE* measures for these algorithms are tabulated in Table 4. The re-
366 ported results substantiate that our proposed method overcomes its contenders
367 with GB acquiring the second rank. It in turn validates the efficiency of jointly
368 considering the clinical activity and the associated cost data for the healthcare
369 cost prediction.

370 **Table 3.** *MAPE* values (training error during 10-fold cross validation) for dif-
371 ferent values of k and length of sequence (in percentage *per*)

<i>per</i> k	20%	30%	40%	50%	60%	70%	80%	90%	100%
1	9.25	7.43	6.01	6.50	5.85	5.53	4.68	4.15	3.76
2	8.83	8.08	6.95	6.54	6.04	5.80	5.15	4.10	3.65
3	8.87	7.74	5.89	5.39	4.98	4.62	4.38	4.07	3.49
4	9.36	7.82	6.04	5.47	4.86	4.65	3.94	3.77	3.63
5	8.90	7.14	5.69	5.13	4.84	4.49	4.24	3.77	4.03
6	9.01	7.35	5.76	5.32	5.17	5.15	4.51	4.26	4.33
7	9.29	7.58	5.77	5.27	5.58	4.97	4.62	4.46	4.35
8	9.13	7.39	5.72	5.64	5.18	5.28	4.46	4.13	4.21
9	9.57	7.66	5.92	5.42	5.57	5.14	4.64	4.11	4.08
10	9.68	8.51	6.27	6.26	6.11	5.93	5.22	4.60	4.41

372 **Table 4.** *MAPE* values (with test data) for different competitive methods for
373 different length of sequence (in percentage *per*)

<i>per</i>	20%	30%	40%	50%	60%	70%	80%	90%	100%
Prop. method	8.89	7.84	6.27	6.05	5.94	5.53	5.29	4.17	3.59
GB	9.70	8.63	6.41	6.37	6.57	5.91	5.55	4.94	4.43
ANN	10.89	10.69	9.83	9.45	8.58	7.98	7.55	6.95	6.61
LASSO	12.08	11.93	11.86	10.80	9.85	9.32	8.85	8.12	7.65

374 4 Conclusion

375 The paper presents a novel method to predict healthcare cost of breast cancer
376 patients as early and accurately as possible. The early prediction capability of
377 the proposed method is used for identifying patients at risk of becoming high-
378 cost healthcare users, before incurring substantial avoidable costs. The merit of
379 the paper lies in the following counts. First, it considers the clinical activity and
380 the associated healthcare cost data jointly to model the treatment of a patient.
381 Second, it recommends a novel distance measure to capture the dissimilarity of
382 two treatment patterns, encompassing both clinical activities and healthcare cost
383 information. Third, it employs the hierarchical DBSCAN to categorize patients
384 into different clusters with an aim to effectively identify the high-need and/or
385 high-cost patients. Fourth, each cluster of patients is depicted by a Markov chain
386 model to mathematically represent the treatment patterns. Finally, the Markov

387 chain models of all the clusters are used to predict the possible future (total)
388 cost of a patient's treatment. The performance of the proposed algorithm is
389 compared for different length of sequence of the recorded treatments of patients.
390 The experimental results reveal that the method achieves an *MAPE* value, as
391 small as 6% even with half of the clinical records of a patient. Experiments
392 undertaken also substantiate the superiority of the proposed algorithm to three
393 state-of-the-art techniques which utilize only the healthcare cost data of the
394 patients for prediction.

395 As a continuation of the present work, we first plan to test our method on
396 different databases from different healthcare organizations for patients suffering
397 from different diseases. More experiments on different databases could help to
398 take a deeper dive into the data and explore ways to obtain more solid evidence
399 on the performance of the proposed method, irrespective of databases. Second,
400 we may consider the socio-demographic information of the patients along with
401 the clinical actions with an aim to be utilize their joint explanatory power to
402 understand the root causes of patient' costs. Third, we have not exploited time
403 feature in the present work. Intuitively, inclusion of time feature may effec-
404 tively capture the differences of treatment patterns of patients and thus may
405 enhance the prediction performance of the proposed method. Finally, appropri-
406 ate stratagem needs to be developed to effectively balance the trade-off between
407 the accuracy and earliness of the healthcare cost prediction.

408 Acknowledgement

409 This work is supported by the Basque Government under the grant "Artificial
410 Intelligence in BCAM number EXP. 2019/00432", BERC 2018-2021 program
411 and through the ELKARTEK program, and by the Ministry of Science, In-
412 novation and Universities: BCAM Severo Ochoa accreditation SEV-2017-0718.
413 Jose A. Lozano is partially supported by the Basque Government through the
414 BERC 2018-2021 program, IT1244-19 and grant "Artificial Intelligence in BCAM
415 number EXP. 2019/00432" and by the Spanish Ministry of Science, Innovation
416 and Universities: BCAM Severo Ochoa accreditation SEV-2017-0718, TIN2016-
417 78365-R and PID2019-104966GB-I00. Aritz Pérez is also supported by Spanish
418 Ministry of Economy and Competitiveness MINECO through TIN2017-82626-
419 R funded by (AEI/FEDER, UE). Onintze Zaballa holds a grant of the Basque
420 Government EJ-GV 2019.

421 References

- 422 1. The Centers for Medicare and Medicaid Services (CMS) DoHaHS, United
423 States, National Health Expenditure Data, 2016.

- 424 2. Morid M. A., Kawamoto K., Ault T., Dorius J., and Abdelrahman S.:
425 Supervised learning methods for predicting healthcare costs: systematic
426 literature review and empirical evaluation. In: AMIA Annual Symposium
427 Proceedings, pp. 1312-1321. American Medical Informatics Association
428 (2017).
- 429 3. Stacey D., Légaré F., Lewis K., Barr, M. J., Bennett C. L., Eden K. B., and
430 Trevena L.: Decision aids for people facing health treatment or screening
431 decisions. *Cochrane database of systematic reviews*, 4 (2017).
- 432 4. Zhang Y., and Padman R.: Data-driven clinical and cost pathways for
433 chronic care delivery. *The American journal of managed care*, 22(12), 816-
434 820 (2016).
- 435 5. Jones A.: Models for health care. Technical report, HEDG, c/o Depart-
436 ment of Economics, University of York (2010).
- 437 6. Gregori D., Petrinco M., Bo S., Desideri A., Merletti F., and Pagano
438 E.: Regression models for analyzing costs and their determinants in health
439 care: an introductory review. *International Journal of Quality Health Care*
440 23(3), pp. 331-341 (2011).
- 441 7. Diehr P., Yanez D., Ash A., Hornbrook M., and Lin D. Y.: Methods
442 for analyzing health care utilization and costs. *Annual Review of Public*
443 *Health*, 20(1), pp. 125-44 (2007).
- 444 8. Sushmita S., Newman S., Marquardt J., Ram P., Prasad V., Cock, M. D.,
445 and Teredesai A.: Population cost prediction on public healthcare datasets.
446 In: *Proceedings of the 5th International Conference on Digital Health*, pp.
447 87-94(2015).
- 448 9. Bertsimas D., Bjarnadóttir M. V., Kane, M. A., Kryder J. C., Pandey R.,
449 Vempala S., and Wang, G.: Algorithmic prediction of health-care costs.
450 *Operations Research*, 56(6), pp. 1382-1392(2008).
- 451 10. Lahiri C. B., and Agarwal N.: Predicting healthcare expenditure increase
452 for an individual from medicare data. In: *Proceedings of the ACM SIGKDD*
453 *Workshop on Health Informatics* (2014).
- 454 11. Brooks S., Gelman A., Jones, G., and Meng XL.: eds. *Handbook of Markov*
455 *chain Monte Carlo*. CRC press (2011).
- 456 12. Dudani S. A.: The distance-weighted k-nearest-neighbor rule. *IEEE Trans-*
457 *actions on Systems, Man, and Cybernetics*, 4, pp. 325-327 (1976).
- 458 13. Ristad E. S., and Yianilos P. N.: Learning string-edit distance. *IEEE*
459 *Transactions on Pattern Analysis and Machine Intelligence*, 20(5), pp. 522-
460 532 (1998).

- 461 14. Marzal A, and Vidal E.: Computation of normalized edit distance and
462 applications. *IEEE Transactions on Pattern Analysis and Machine Intelli-*
463 *gence*, 15(9), pp. 926-932 (1993).
- 464 15. Campello R. J., Moulavi D., and Sander J.: Density-based clustering based
465 on hierarchical density estimates. In: *Pacific-Asia conference on knowledge*
466 *discovery and data mining*, pp. 160-172. Springer, Berlin, Heidelberg
467 (2013).
- 468 16. Schubert E., Sander J., Ester M., Kriegel H. P., and Xu, X.: DBSCAN
469 revisited, revisited: why and how you should (still) use DBSCAN. *ACM*
470 *Transactions on Database Systems (TODS)*, 42(3), pp. 1-21(2017).
- 471 17. Organization WH. ICD-10: international statistical classification of dis-
472 eases and related health problems: tenth revision. 2nd ed. World Health
473 Organization; 2004.
- 474 18. Sutton C. D.: Classification and regression trees, bagging, and boosting.
475 *Handbook of statistics*, 24, pp. 303-329(2005).
- 476 19. Yegnanarayana B.: *Artificial neural networks*: PHI Learning Pvt. Ltd.
477 (2009).
- 478 20. Tibshirani R.: Regression shrinkage and selection via the lasso. *Journal of*
479 *the Royal Statistical Society Series B (Methodological)*. pp. 267-88 (1996).
- 480 21. Zaballa O., Pérez A., Inhiesto E. G., Ayesta T. A., and Lozano J. A.:
481 Identifying common treatments from electronic health records with missing
482 information: An application to breast cancer. *PLOS ONE* (accepted, to
483 be published) (2020).
- 484 22. Billings J., Dixon J., Mijanovich T., Wennberg D.: Case finding for patients
485 at risk of readmission to hospital: Development of algorithm to identify
486 high-risk patients. *British Medical Journal* 333(7563): 327. (2006).