

MedClick Health Recommendation Algorithm

Recommending Healthcare Professionals Handling Patient Preferences and Medical Specialties

Rui Miguel Dos Santos Patornilho¹ and André Vasconcelos^{1,2}

¹*Instituto Superior Técnico, Department of Computer Science and Engineering, Lisboa, Portugal*

²*Instituto de Engenharia de Sistemas e Computadores, Investigação e Desenvolvimento, Lisboa, Portugal*

Keywords: Diagnosis, Recommendation, Health, Interactivity, Reliability.

Abstract: Today's health has a determinant role and it is a subject of concern by society. Diagnosing a disease or obtaining a medical specialty, given a set of symptoms, is not a trivial task and different decisions and approaches can be adopted to solve and handle this problem. Expert systems advise patients about a possible diagnosis, associated diseases, treatments and more concrete information about a disease considering simple symptoms. However, most systems don't have the recommendation component of a medical doctor, which will be the differentiating factor of this research. The aim of this paper is to develop an algorithm capable of determining the medical specialties associated with a set of symptoms and diseases, and based on the medical specialties obtained, recommend the most suitable specialists. The algorithm is divided into two phases: Health Screening and Healthcare Professional Recommendation. Health Screening has the purpose of determining and computing all the medical specialties probabilities, given a set of patient symptoms and applying a statistical model based on all the relations symptom \rightarrow disease and disease \rightarrow medical specialty. Healthcare Professional Recommendation has the purpose of recommending the best healthcare professionals, given a set of patient preferences, applying a weighted mean average, where each weight of a healthcare professional feature is given by a patient according to his preferences. This algorithm was evaluated through a set of test cases, having a database with information about symptoms, diseases and medical specialties. This algorithm was later compared to other systems that have the same purpose, to access its quality. The comparison result between the algorithm and WebMD system indicates that the diseases found by the solution are in 80% of all the cases equal to the diseases found and pointed by WebMD system.

1 INTRODUCTION

Health today has been more and more often a subject of continuous research and improvements. These improvements are ranging from new treatments, new medicines, new improvements into how to diagnose a patient, among others. Added to these improvements, an easier way of exchanging information about diseases, symptoms, treatments, among other health topics have been performed. Numerous platforms as web health blogs and health applications allow scheduling of medical services as appointments, medical diagnosis based on symptoms and biometric data, among many other functionalities available.

This research has the goal of developing a health algorithm that based on a patients symptom set must determine and compute all the medical specialties probabilities, having into account all the relations

between: symptom \rightarrow disease and disease \rightarrow medical specialty. After computing all the medical specialties probabilities, the algorithm must recommend a set of healthcare professionals most suitable to deal with the medical specialty set obtained previously. A set of features such as Rating, Price and Distance, beside the medical specialty probability can influence the recommendation process.

A variety of different techniques and a set of systems that are able to perform diagnosis have been analyzed in **Background** Section. These techniques and systems are the research and the basis of the solution definition and respective implementation, described in **Health Recommendation Algorithm** Section. This solution has been tested with a dataset, a set of different test cases and compared with other similar systems, in order to measure its accuracy and reliability. These results are presented in **Results And Dis-**

Discussion Section. This article ends with conclusions and future work that can be done to improve the solution.

The major contribution of this work is to develop a new system capable of recommending the most suitable healthcare professionals given a set of patient symptoms and that fulfils as most as possible all the patient preferences.

2 BACKGROUND

This section presents the state of the art of the methodologies and techniques used to perform diagnosis and systems able of performing diagnosis and/or recommending healthcare professionals according to a set of features and considering all the patient's preferences. A comparative analysis to all the different methodologies and technologies is performed, in order to point all the advantages and disadvantages associated. A comparative analysis to all the different systems analyzed is also performed, in order to point the main differences and decide which systems are the most suitable to be a base for the solution.

2.1 Machine Learning

Machine learning is an approach that has the objective of learning from data and making predictions from data. Machine Learning is characterized for making predictions or decisions with all the inputs received, building a knowledge model and applying the knowledge model to a new entry. Based on the previous experience, one algorithm developed according to a Machine Learning approach can learn by its own and it will be able to react with a new experience, in order to reach a solution or predict something (Mannilat, 1996; Munoz, 2014). Applying Machine Learning techniques to self-diagnosis is possible, for example if there is information available about diseases, relations symptom(s) \rightarrow disease(s) able to build a knowledge model, which can learn and infer conclusions based on a set of inputs given. Afterwards the algorithm will be able to apply the knowledge model built to a new input in order to infer a conclusion (diagnosis).

2.2 Bayesian Network

A Bayesian network is a probabilistic model that represents a set of variables and their conditional dependencies via a directed acyclic graph. For example, a Bayesian Network could represent the probabilistic relationships between diseases and symptoms. Given

all the symptoms, a network can be used to compute the probabilities of the presence of various diseases. A directed acyclic graph is composed by **nodes** and **nodes** are linked by **edges**. **Nodes** represent random variables. Meanwhile **edges** represent the conditional dependencies between nodes, but two nodes can be **conditionally independent** (no path that link the two nodes). A probability function is associated with each node, that receive as input a set of values from the node's parent variables, and gives as output the probability of the variable represented by the node (Bayesianas and Fred, 1993; Leung, 2007).

2.3 Decision Trees

A decision tree is a decision support tool that uses a tree-like graph or model of decisions and their possible consequences, including chance event outcomes, resource costs, and utility. A decision tree can be linearized into decision rules, for example:

if Fever and Cough and Migraine then Flu

These decision rules represent conditions derived by association rules that are methods implemented for discovering relations between variables. For example if we have Fever, Cough and Migraine it is a possible sign of Flu. So starting with a main decision node the tree is traversed by a set of decisions until it's reached a leaf node that corresponds to a derived conclusion based on all the decisions made (Kami and Jakubczyk, 2017; Quinlan, 1987). Applying a decision tree to a self-diagnosis algorithm is possible if we have a set of relations symptom(s) \rightarrow disease(s), the tree is traversed using symptoms and reached an outcome that corresponds to a disease.

2.4 Fuzzy

Fuzzy logic is a form of many-valued logic where variables can have values between 0 and 1. A value is considered as completely true if equals 1 and completely false if equals 0. Fuzzy logic is used in cases where it is not linear to identify if a variable is absolutely true or false, due to the variable's value and respective interpretation is ambiguous (different viewpoints) to people (Zadeh, 1996).

An example of a set of rules is:

IF cough and migraine then Flu is moderate
IF cough, migraine and fever then Flu is high

With these rules it is possible to obtain outputs that correspond, in this case, to the degree of expressiveness of Flu given a set of symptoms. Fuzzy logic can be applied to self-diagnosis with a fuzzy set of each disease probability given a set of symptoms.

2.5 Comparative Analysis

Machine Learning offers the ability of discovering data patterns and possible associations between the input data. It is also possible to apply different techniques, from artificial neural networks to learning classifier systems. Machine Learning approaches can suffer from performance issues, due to processing a large input dataset (mandatory) or indexing data updates. It is necessary to train the system with several examples (large dataset), in order to reduce the risk of bad learning, incorrect calculations and to minimize system performance issues.

Decision Trees can become very big structures if the number of characteristics and relations between the data is very large. As such this poses a performance issue as computations will become more complex, the greater the tree is. Operations like recomputing a tree, in order to represent a new variable or value will have a greater delay, diminishing the performance of the overall system. However, this approach is intuitive and easy to learn and use. It is also possible to represent data such as a hierarchical structure with levels. It is also possible to combine decision trees with other different techniques.

Fuzzy Inference is a process which is both flexible and intuitive. It represents a natural way of expressing uncertain information (possibility of incomplete data). This system however can have poor performance due to a large amount of fuzzy inference rules. The system is also prone to error if the fuzzy inference rules given are too generic or too specific. These process can be used in combination with other techniques, like machine learning.

Bayesian network is a graphical model that represents probabilistic relationships among variables and that applies the Bayesian probability. With a Bayesian network it is possible to handle incomplete data sets, learn and infer knowledge about a set of casual data relationships. It allows the representation and information extraction from two factors that are believed to be correlated (a priori knowledge). It is also possible to update the weight of the directed edges, based on new data. There are disadvantages associated as: the computational difficulty of exploring a previously unknown network, due to the need of calculating all the branches of the network to obtain a value of one branch. The second disadvantage is that the prior knowledge must be reliable to have a useful network. An excessively optimistic or pessimistic expectation of the quality of these prior beliefs, will distort the entire network and invalidate the results.

2.6 WebMD Symptom Checker

This application has the purpose of inferring a list of possible conditions, indicated by a set of questions performed to a user such as: age, gender, followed by a series of questions including a part of the body where symptoms occur. A user must answer the following questions to obtain a more precise diagnosis, however the system is able to determine a possible diagnosis in cases of lack of information. After presenting the list with all the possible conditions, if a user clicks on a condition will be presented to him more information about that condition,(how common is the disease, the degree of severity, among others). This application is also able to point the most suitable medical doctors (specialists) to deal with that condition. A user is able to choose the feature that will be responsible for ordering a set of medical doctors, ranging from name, years of experience to distance. Here the user's preferences are taken into account at the time of decision (WebMD, 2017a; WebMD, 2017b; Whysel, 2012).

2.7 Isabel Symptom Checker

This application is responsible for promoting the search of medical knowledge to all people, using the professional Isabel Diagnosis Check-list System, used by doctors around the world when they're unsure of a diagnosis. To obtain diagnosis information it is only necessary for a user to indicate his symptoms (unlike most symptom checkers, a user can put in as many symptoms as he wishes) and it will be provided a list of the most possible diagnoses that are related to those symptoms. Each diagnosis has medical information associated that complement and explain diseases, treatments and other symptoms. It can also be an auxiliary way to understand, obtain more information about a health care status, in order to discuss a possible health topic with a medical doctor. It is also possible for a user to find a doctor, offering a functionality composed by a set of links to various web resources that offer this functionality. These links take into account the medical specialty obtained and location, in order to return the most relevant medical doctors (Isabel, 2017).

2.8 Mayo Clinic Symptom Checker

This application is responsible for providing information, not diagnosing, a given symptom. The idea is for the patient to initially choose a symptom. After choosing the symptom, a user needs to point at least one related factor, in order to complement the information

that the application will need to determine the possible related diseases. After a user points his symptoms and associated causes, a set of diseases that match at least one cause is presented to a user. If the user clicks on one of the presented diseases, information about itself as all the symptoms related and a description, for example the expressiveness on a population (age, gender, severity, among others) will be presented to a user. To each disease it is also present another link that describes all the associated factors, highlighting in Bold the user's selection. This application doesn't take into account information related with user's profile as (age, gender, location, among others). It is also possible for a user to request a medical appointment, however it doesn't have any recommendation tool that points the best medical doctors to deal with the diseases obtained (MayoClinic, 2017).

2.9 Results and Comparative Analysis

In terms of the number of symptoms insertion, only "WebMD" and "Isabel" allow an insertion of any number of symptoms as a user want, differing from "Mayoclinic" that only allows a insertion of one symptom. In terms of the profile information, only "Mayoclinic" does not take into account the biometric information about a user (age, gender, weight, among others). All the other systems, except "Isabel", perform a set of questions, in order to obtain more information that will increase and improve the diagnosis. In other words, the answers to the set of questions will provide complementary information that will result into a more precise and accurate diagnosis (not to generic). All the systems, if a user requests a more detailed information, are able to provide a more concrete and detailed explanation about all the diseases obtained with information about risk factors, symptoms associated and treatments.

In terms of medical doctors recommendation, only "WebMD" and "Isabel" are able to provide a clear way of recommending the most suitable medical doctors to deal with a specific diagnosis. Isabel differs from "WebMD", due to the recommendation being done with external applications that can take or not into account patient preferences such as location, years of experience, rating, among others. The most complete system, considering all the functionalities is "WebMD", due to being able to recommend a medical doctor, ask a user for complementary information with a set of questions, introduce as many symptoms as a user needs (no restrictions) and by providing a detailed and complete information about a health topic (for example, disease). This information will dot the user with more medical knowledge and provide to

him the basic knowledge to discuss the results obtained with his medical doctor.

3 HEALTH RECOMMENDATION ALGORITHM

The methodology adopted to develop the solution was to divide **Health Algorithm** into two different parts: **Health Screening** and **Healthcare Professional Recommendation** identifying all the dependencies between them.

3.1 Solution Overview

In Figure 1 it is represented a diagram that presents all the components that interact directly with the solution component, including all the interactions between them:

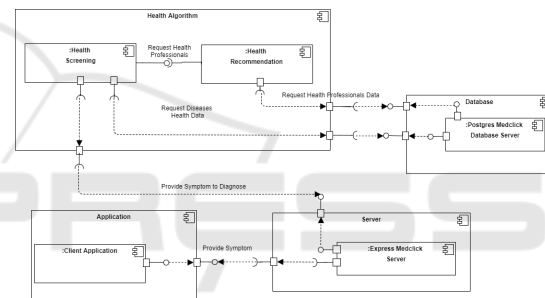


Figure 1: Solution Component Algorithm Diagram.

Figure 1 represents all the interactions between the **Health Algorithm** component and all the remaining components **Database**, **Server** and **Application** components. It is possible to see that the **Health Application** components interacts with the **Database** component, in order to request data, healthcare professional data or the health diseases data. Otherwise, this component connects with the **Server** component, in order to obtain the most significant symptom of a patient. That symptom will feed the algorithm and being responsible for starting the diagnosis process. Finally, the last interaction represented is between the **Application** and **Server** component, which is represented by the symptom exchange between the two components.

3.2 Health Screening

Health Screening starts by interpreting the most significant patient symptom, in other words, the most important symptom or more expressive. The result from the symptom interpretation is a set of diseases

that have the symptom and all the remaining symptoms. This result is going to be used to compute each symptom relevance, in order to obtain complementary information (if possible). Each symptom relevance corresponds to the number of diseases where the symptom is present, higher the number of diseases present, higher the probability of having the symptom.

The relevance of each symptom is computed by the Equation 1:

$$\text{symptom relevance} = \frac{\text{number of diseases}}{\text{total diseases}} \quad (1)$$

where:

- number of diseases, corresponds to the number of diseases that have the most important symptom.
- total diseases, corresponds to the total number of diseases existing.

Equation 1 computes the degree of how a symptom is presented into all the diseases set (including the diseases that don't have the most significant symptom). After computing all the symptoms relevances, the algorithm will ask a set of questions to a patient, in order to look for complementary information (presence of a symptom), improving the accuracy of a diagnosis. Having all the patient answers and the respective information, the algorithm is now able to compute the diseases weight.

The disease weight is computed by the Equation 2:

$$\text{disease weight} = \frac{\sum \text{relevance symptom}}{\text{total relevance symptom}} \quad (2)$$

where:

- relevance symptom, corresponds to the relevance of a patient symptom present in the disease symptoms set.
- total relevance symptom, corresponds to the total relevance value of all the symptoms that a patient has.

Equation 2 computes the relation between the sum of the relevance associated with the set of symptoms presented by a patient and that are present in a disease, divided by the total relevance of all the symptoms presented by a patient. Computed all the diseases weights the algorithm will compute all the medical specialties probabilities associated with all the diseases.

The medical specialty probability is computed by the Equation 3:

$$\text{specialty probability} = \frac{\sum \text{disease weight}}{\text{total specialties mean}} \quad (3)$$

where:

- disease weight, corresponds to the weight value of a disease that belongs to the specialty.
- number diseases specialty, number of total diseases of specialty
- total specialties mean, corresponds to the total value of the mean diseases weights of specialty values passed to a variable before being updated.

Equation 3 computes the relation between the mean diseases weight of all the diseases of one medical specialty divided by the total sum of all the means. After obtaining all the medical specialties' probabilities, the algorithm begins the **Healthcare Professionals Recommendation** phase.

3.3 Healthcare Professionals Recommendation

Healthcare Professionals Recommendation starts by considering a set of healthcare professionals features, that will be used for recommending a healthcare professional by applying a features set combination. Several features can be a target of concern.

In our research, considering the data available, we propose to use for recommend a healthcare professionals the following features:

- **Rating**
- **Distance**
- **Price of Medical Appointment**

Despite all the concerns related to the availability of the data, these features have been chosen by their importance into performing a recommendation. **Rating** can be used to obtain information about a healthcare professional quality, measure expressed by all the patients. **Distance** is important, in terms of determining the nearest healthcare professionals to a patient, given his location and **Price of Medical Appointment** is important to choose the cheapest healthcare professionals for scheduling a medical appointment, among other factors.

Entering **Healthcare Professionals Recommendation** and depending of the specialties obtained it is provided a healthcare professionals set, composed by all the healthcare professionals able to handle the medical specialties (have at least one medical specialty). Having the healthcare professionals set, it is necessary to convert all the healthcare professionals features to a homogeneous scale, in order to obtain better results and work with homogeneous values. **Rating** is by default into a range from 1 to 5 so it is not necessary to convert this feature. Converting the **Distance** and the **Price of Medical Appointment** to a

scale from 1 to 5 require dealing with the maximum value and the minimum value.

Having these two values it is necessary to set the minimum value to 5 and the maximum value to 1 applying the Equation 4:

$$y = mx + b \tag{4}$$

where:

- m, corresponds to the straight slope
- x, corresponds to the independent variable of the function $y = f(x)$.
- b, corresponds to the y-intercept of the line

Equation 4 is a straight line equation to compute a price between the minimum and maximum that has points point1 (x_1, y_1) and point2 (x_2, y_2) , where $x_1 =$ minimum, $y_1=5$, $x_2 =$ maximum and $y_2=1$. Equation 4 has a set of variables that need to be computed, so it is necessary to follow a set of steps.

The first step necessary is computing the m variable, applying the expression:

$$m = \frac{\Delta y}{\Delta x} = \frac{y_1 - y_2}{x_1 - x_2} \tag{5}$$

The second step corresponds to the computation of the b variable, applying the expression:

$$b = y_1 - (m \times x_1) \tag{6}$$

After having the two computations it is only necessary to obtain the respective y, that is the scaled value.

Completing this set of steps and having all the scaled features of each healthcare professional, the algorithm will ask the patient about his preferences, in terms of **Rating**, **Distance** and **Price**, attributing a value from 1 to 5, where 1 is less important and 5 the most important. These values are weights that will be used a **Weighted Mean Average**. In this step, it is necessary to be considered the following scenarios: patient points his preferences, patient doesn't point his preferences and the past medical history of a patient.

If a patient points his preferences these values are the weights that will be used a **Weighted Mean Average**. If a patient doesn't point his preferences, the algorithm will search for his past medical history (past medical appointments), inferring and computing the respective weights based on the patient's medical history data. Finally, in case of none of the previous two scenarios, the weights will be equal to each feature, except the medical specialty probability that is always assigned a value of 5.

Having all the weights computed, it is necessary to apply the Equation 7:

$$\text{weight average} = \frac{\sum_{i=1}^n w_i * x_i}{\sum_{i=1}^n w_i} \tag{7}$$

where:

- x, set of healthcare professional features
- w_i , weight of feature x at ith position
- n, total number of features

The final step necessary is ordering all the results obtained, applying the Equation 7 to each healthcare professional.

3.4 Technological Architecture

This algorithm has been developed with the following technological architecture represented in Figure 2:

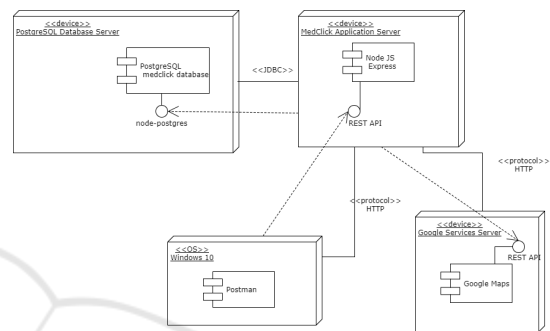


Figure 2: Technological Architecture Diagram.

In Figure 2 it is represented the technological architecture associated with the algorithm. It is represented a MedClick Application Server that has relations with all the three other components. These relations are represented by all the HTTP requests performed to the server component, and by all the requests performed from the server to the database component. **Google Maps API** is a component responsible for handling all the requests performed by the MedClick Application Server, ranging from geocoordinates, addresses to compute distance, among others.

Node JS modules are a set of components that are responsible for "feeding" our solution with data structures, methods, among other relevant libraries that will allow our solution to work in a set of different situations. In this figure an example of a Node JS component is **Node JS Express** responsible for creating a server responsible for answering and handling requests from other devices.

The **Postman** component is an external component responsible for making HTTP request, in order to obtain a given answer. The last external component that was used is **PostgreSQL**, responsible for handling all the connections and requests between the database and the application server. This module for example, will load the data that a application will need and that will be requested by the MedClick Application Server.

3.5 Solution Assumptions and Approach

Different approaches, ranging from machine learning to inference rules are possible ways, differing into some aspects. Our algorithm measure and points a set of medical specialty probabilities, allowing a friendly interactive way of communication between a patient and the system. The approach chosen is based on Bayesian Network, due to being possible to represent all the relations symptom → disease in an interconnected network and taking into account all the relations and dependencies between the data. The reasons that led to choose this approach are: the genericness of the algorithm (several medical specialties, symptoms and diseases), due to not be suitable to create datasets and test cases (large volume) that are mandatory to ensure reliability and accuracy to the algorithm by a Machine Learning approach. Our solution points and calculates probabilities (more than one possible medical specialty), being in this case Decision Tree a bad approach, due to it possibly being a huge and sparse structure, being hard to transverse and in the most of cases it only is able to point one output.

To the **Healthcare Professional Recommendation** the initial approach is to use a correlation methodology, in order to find similarities between the healthcare professionals data and the patient characteristics (past medical appointment records). However this can have some issues such as, for example, if we have two professionals with same characteristics only differing in rating. Assuming that a patient visited only healthcare professionals with low rating, the tendency is to first recommend the healthcare professionals with the lower rating, however the healthcare professional with the higher rating should be recommended first. The chosen approach is a weighted average, after being asked all the patient preferences, features that the patient consider more relevant.

4 RESULTS AND DISCUSSION

The algorithm proposed has been evaluated considering the two components in separate and after with a global evaluation. The **Health Screening** was evaluated by a set of internal tests and comparing the results of this algorithm with similar systems that perform diagnosis. The **Healthcare professionals Recommendation** was evaluated by a set of internal tests, which is composed by different cases of patient preferences, ranging from preference for one feature in demand of the others, more that one feature preference to equal feature preference.

4.1 Dataset

Starting with the **Health Screening**, this algorithm has a dataset composed by: 50 diseases, 50 diseases associated with 10 medical specialties different, and 388 relations symptom → disease. The dataset available for performing the **Healthcare Professionals Recommendation** is composed by: 7 healthcare professionals with the respective information about their rating, appointment price, health care provider and medical specialty(ies). Also, there are 4 health care providers present in the dataset.

4.2 Health Screening Results

This subsection presents results of a set of medical specialties probabilities obtained (output), given a set of patient symptoms (input).

Figure 3 represents a set of results expressing the variation of medical specialties probabilities, given a symptoms set.

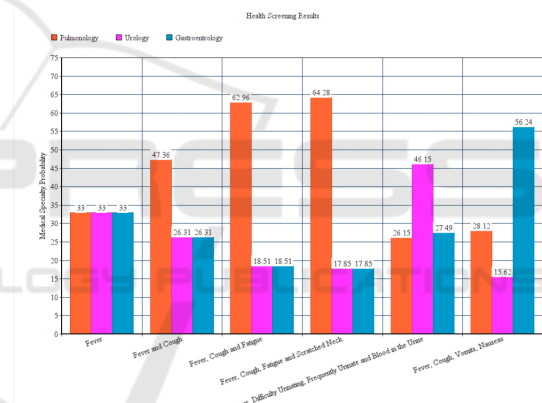


Figure 3: Health Screening Results.

From Figure 3 it is possible to see the set of respective symptoms inputs (represented by the x axis) and the corresponding algorithm outputs medical specialties probabilities (represented by y axis). It is expected that with different symptoms the medical specialties probabilities are different, indicating that, for example one medical specialty is more expressive given a set of symptoms and that these symptoms are more relevant on diseases that belong to that medical specialty.

The different symptoms inputs are composed by symptoms ranging from one symptom to four, however more symptoms can be added to the set. The input characterized by only a single symptom presents an equal probability to each medical specialty, due to this symptom be present in all diseases with equal weights and how each disease has a similar weight, all the medical specialties probabilities will be similar.

The dataset contains diseases that only have one symptom or two, among other different possibilities, being the results obtained according all the symptoms. The symptoms set with more than one symptom will have a medical specialty with a higher value as represented in Figure 3. The set that has Fever and Cough present a higher result in Pulmonology, due to these symptoms being present with more expressiveness on diseases of that medical specialty. Otherwise, adding Vomits and Nauseas to the previous symptoms set will increase the probability of Gastroentology, due to the previous symptoms and the new ones be more expressive on diseases that belong to the Gastroentology medical specialty.

These results are expressive that the algorithm with more than one symptom is able to identify in a more clear way all the medical specialties and computing their probabilities on a more precise way. More information provided, more than one symptom, corresponds to be possible to not having all diseases with the same symptoms, due to the existence of symptoms that are more relevant into a restrict set of diseases, that will allow different medical specialties probabilities.

4.3 Healthcare Professionals Recommendation Results

This section presents the results of applying patient preferences on the recommendation phase to different healthcare professionals features.

In this case all healthcare professionals have a **Rating** associated that can be different, their medical specialties can be **Cardiology** or **Pulmonology** and they work in different places expressed at:

- **Provider 1**, Distance: 11.874 Km
- **Provider 2**, Distance: 20.053 Km
- **Provider 3**, Distance: 333.792 Km

The price of medical appointments is expressed at:

- **Cardiology**, 55 euros
- **Pulmonology**, 60 euros

To simplify it is implicit that the medical specialties probabilities obtained are equal (0.5).

The Healthcare Professionals Features are available in:

- Professional 1: Cardiology, Rating: 5, Provider: Provider 3
- Professional 2: Cardiology, Rating: 3, Provider: Provider 1
- Professional 3 : Pulmonology, Rating: 2, Provider: Provider 1

- Professional 4: Pulmonology, Rating: 4, Provider: Provider 2

Figure 4 represents a set of results expressing the variation of the healthcare profession weighted mean average, given a set of healthcare professional features.

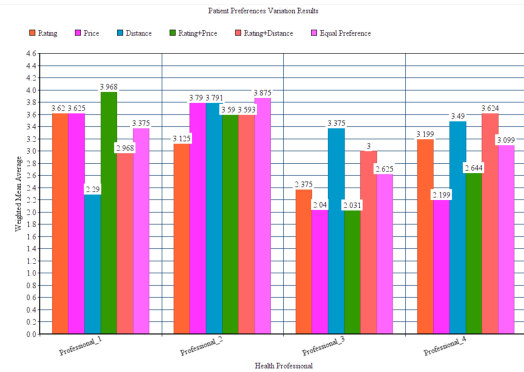


Figure 4: Healthcare Professionals Recommendation Results.

From Figure 4 it is possible to see a set of respective healthcare professionals (represented by the x axis) and the corresponding algorithm outputs (represented by y axis) Weighted Average result according to a patient preference. Each bar is representative of a patient preference that could be only one preference pointed, equal preference or a combination of features preferences. This figure represents the respective classification to all the healthcare professionals considering all the patient’s preferences, scaled from 1 to 5 as the healthcare professionals features.

To an individual preference if there aren’t significant differences on the remaining features it is expected that the healthcare professionals with the higher value on the patient feature preference have a higher value. For Distance preference, it is expected that the healthcare professional with the highest value is Professional 2, due to this professional has a lowest distance to the patient address. Figure is representative of this fact, given the presence of 3.791 for Professional 2, followed by Professional 4 with a value of 3.49. The same expectations are applied to the Price and Rating features. To the equal preferences it is expected that the professionals with a higher mean features values (closer to 5) have the highest results. It is expected that Professional 2 has a highest value, given having the best medical appointment price, working in the closest health provider and having a good rating value (3), fact once more represented by the results obtained.

Combining preferences will indicate better results in healthcare professionals that present a highest value on the features that are combined. In terms of

combining Rating and Price it is expected that the highest value is presented by Professional 1, given the highest values of these two features. The obtained results show that this fact is reached, due the highest value be associated with Professional 1. To the combination of Rating and Distance it is expected that the professional with the highest value is Professional 4, due to despite not having the highest values as possible (Professional 1 has a higher value of Rating) this professional has the highest combination values of these features. The results present that this fact is addressed.

4.4 Health Systems Comparative Analysis

The system was compared with similar systems by a set of 19 test cases, where each test case is identified by a set of systems (provided as input) and a set of diseases (provided by each algorithm as output). The results present in Figure 5 are the respective percentage of the same outputs identified by the system with the remaining systems output.

In Figure 5 it is represented a comparative analysis between our system outputs and the other systems that have the same behavior:

System	Systems								
	Isabel			Mayoclinic			WebMD		
	High	Mean	Low	High	Mean	Low	High	Mean	Low
MedClick	1	0.5268	1/8	1	0.3342	0	1	0.8004	1/3

Figure 5: Comparative Systems Analysis.

In Figure 5 are represented the respect high probability value between all the diseases output matches (our system and other systems), the Mean probability value of all the examples used to compare our solution with the remaining systems and the Lowest probability value obtained. From this figure it is possible to see that all the systems have a High matching probability value of 1, indicating that in the presence of that symptoms the output of our system is the same of all the remaining systems. Note that the datasets are different and the respective relations symptom → disease are different from ours. In terms of the Low matching probability value, the "Mayoclinic" system has a value of 0, indicating that no disease is pointed by our system, given the symptoms presented.

Finally, in terms of the Mean matching probability value, the highest probability value is presented by the "WebMD" system with a value of 0.8, closer to 1 (representing a perfect match). Comparing our solution to all the other systems it is possible to see that, despite some diseases matching probabilities being low, with "WebMD" system these probabilities

values are higher. The obtained probability result is of 0.8, which indicates that our system is able to identify the same diseases in 80% cases as the "WebMD" system, despite the different approaches and different datasets.

5 CONCLUSIONS

Medclick **Health Recommendation Algorithm** is able to diagnose a patient given a symptom and with the symptoms suggested to the patient is able to perform and reach a sustainable conclusion about the probability of each medical specialty. The algorithm is also able to recommend an ordered list of medical doctors able to deal with the health specialties obtained at the health screening phase and that fulfill all the patient preferences, in terms of distance, rating and price. Our algorithm is able to deal with clinical, patient and healthcare professional data and perform a decision that is the most suitable as possible, in order to facilitate the patients work when it has to choose and decide where he should go given a set of symptoms.

The obtained results indicate that comparing our solution with other similar approaches already in use show that MedClick achieves acceptable results. The comparison measures are a mean value of 0.8 with "WebMD", 0.52 with "Isabel" and 0.33 with "Mayoclinic". Despite the higher value be 0.8, and considering that all these systems have different datasets, including different relations symptom → disease, it is possible to conclude that our solution is able to identify and point the majority of the diseases correctly as any other similar system. Otherwise, all the tests performed to our solution point that it is able to identify and point diseases and respective medical specialties in the presence of few and large symptoms. The results obtain show that comparing the probabilities variation with one symptom and more than one, the solution increases the accuracy in the presence of more than one symptom. Results show that our solution with Fever as a symptom points an equal probability to each medical specialty and with Fever and Cough as symptoms, the medical specialty probability of Pulmonology increases facing the decrease of the other medical specialties (accuracy increases).

Finally, to the healthcare professional recommendation the results indicate that our solution is able to recommend the "best" medical doctor according to all the patient preferences. In our solution, a set of patient preferences is a determinant factor that influence the recommendation process, in order to reach, such as possible, all the patient preferences and ex-

pectations. Results show that if a patient prefers healthcare professionals with lower Price, the first healthcare professional to be pointed is the healthcare professional with the lowest Price. Otherwise, it is also possible to combine features, Price and Distance. In this case the results show that the healthcare professionals will be present by ordering all the healthcare professionals by the Price and Distance. The first healthcare professional to be recommended to a patient is the professional with the lowest Price and the lowest Distance.

5.1 Future Work

This system can be target of improvements. One possible improvement is to decrease its genericness, choosing carefully which medical specialties it will be able to diagnose, in order to improve its accuracy and reliability applying techniques as Machine Learning, providing concrete text cases and examples that will feed the system and building a knowledge model to be applied to increase the precision and accuracy of a diagnose process. This system is also able to handle more features as age of healthcare professional, gender, information about his Curriculum Vitae (CV), statistics about the probability of a given population have a certain disease, among others. These set of features added will improve the accuracy and reliability, in terms of diagnosing and healthcare professional recommendation. For example, there are diseases that are more expressive in males than in females. Otherwise having more features available for healthcare professionals, will allow that a patient can have more options to point his preferences and that the recommendation will be done with a more diverse set of features and can be a more efficient and diverse process. The last possible improvement can be, providing the system with an interface available for user testing, in order to obtain the degree of satisfaction and respectively measuring the systems user quality. A similar approach to healthcare professionals can be also performed.

ACKNOWLEDGEMENTS

This work was supported by national funds through Fundação para a Ciência e a Tecnologia (FCT) with reference UID/CEC/50021/2013.

REFERENCES

- Bayesianas, R. and Fred, A. L. N. (1993). Redes Bayesianas. *Instituto Superior Tecnico da Universidade de Lisboa*.
- Isabel (2017). Isabel Symptom Checker Info. <https://symptomchecker.isabelhealthcare.com/home/main>.
- Kami, B. and Jakubczyk, M. (2017). A framework for sensitivity analysis of decision trees. *SGH Warsaw School of Economics, Warsaw Poland*.
- Leung, K. M. (2007). Naive Bayesian Classifier. *POLYTECHNIC UNIVERSITY*.
- Mannilat, H. (1996). Data mining : machine learning , statistics , and databases. *University of Helsinki*, pages 2–9.
- MayoClinic (2017). MayoClinic Symptom Checker Info. <http://www.mayoclinic.org/about-this-site/about-symptom-checker>.
- Munoz, A. (2014). Machine Learning and Optimization. Technical report, Courant Institute of Mathematical Sciences, New York.
- Quinlan, J. (1987). Simplifying decision trees. *International Journal of Man-Machine Studies*, 27(3):221 – 234.
- WebMD (2017a). WebMD Symptom Checker. <http://symptoms.webmd.com/%0A>.
- WebMD (2017b). WebMD Symptom Checker. <http://symptoms.webmd.com/default.htm#introView>.
- Whysel, N. (2012). A User Test of WebMD Symptom Checker A User Test of WebMD Symptom Checker Table of Contents. *WebMD*, pages 1–29.
- Zadeh, L. A. (1996). *Selected Papers by Lotfi A . Zadeh*. World Scientific.