

Medical Decision Data Mining using Naïve Bayes & K-Means Clustering

PATHALLAPALLI RAKESH VARMA¹, PUVVALA MANI KRISHNA² & UMMADI JANARDHAN REDDY³ ¹B-Tech Dept. Information Technology Vignan's University Vadlamudi-522213 Guntur Andhra Pradesh India Mail Id: - <u>varma.283@gmail.com</u>

²B-Tech Dept. Information Technology Vignan's University Vadlamudi-522213 Guntur Andhra Pradesh India Mail Id: - <u>pmanibvm@gmail.com</u>

³Assistant Professor Dept. Information Technology Vignan's University Vadlamudi-522213 Guntur Andhra Pradesh India Mail Id: - <u>ummadi.janardhan@gmail.com</u>

Abstract

In this paper, by utilizing data mining we can evaluate many patterns which will be used in future to make keenly intellective systems and decisions By data mining refers to sundry methods of identifying information or the adoption of solutions predicated on cognizance and data extraction of these data so that they can be utilized in sundry areas such as decision-making, the presage value for the presage and calculation. In our days the health industry has amassed astronomical amounts of patient data, which, infelicitously, is not "engendered" in order to give some obnubilated information, and thus to make efficacious decisions, which are connected with the base of the patient's data and are subject to data mining. This research work has developed a Decision Support in Heart Disease Presage System (HDPS) utilizing data mining modelling technique, namely, Naïve Bayes and Kmeans clustering algorithms that are one of the most popular clustering techniques; however, where the initial cull of the centroid vigorously influences the final result. Utilizing of medical data, such as age, sex, blood pressure and blood sugar levels, chest pain, electrocardiogram, analyzes of different study patient, etc. graphics can presage the likelihood of the patient. This paper shows the efficacy of unsupervised learning techniques, which is a k-betokens clustering to ameliorate edifying methods controlled, which is ingenuous Bayes. It explores the integration of K-designates clustering with verdant Bayes in the diagnosis of disease patients. It withal investigates different methods of initial centroid cull of the K-designates clustering such as range, inlier, outlier, arbitrary attribute values, and desultory row methods in the diagnosis of heart disease patients. The results designate that the integration of the K-betokens clustering with naïve Bayes with different initial centroid culling naïve Bayesian amend precision in diagnosis of the patient.

Keywords: Data Mining, Naïve Bayes, K-Means Clustering.



p-ISSN: 2348-6848 e-ISSN: 2348-795X Volume 04 Issue 06 May 2017

1. INTRODUCTION

Data mining this is revelation process in the raw antecedently unknown, data nonfrivolous, virtually utilizable. the interpretation of the available erudition indispensable for decision-making in the sundry spheres of human activity. This search for relationship with subsisting astronomically immense associated data that are obnubilated among astronomically immense amounts of data and refers to the "mining" cognizance from sizably voluminous amounts of data. Subsisting systems are habituated to avail in decisionmaking, referred to as data mining. These systems represent an iterative sequence of pre-processing as cleaning, data integration, and data cull is veridical the pattern identification of data mining and erudition representation. Data mining is the search for relationships and ecumenical patterns that subsist in astronomically immense databases, but obnubilated among the plethoras of data. Computer diagnosis of diseases is the medico for the same instrument, the calculations for an engineer: design diagnostics does not supersede the medico, but it avails. The practice of examining immensely colossal preexisting data bases in order to engender incipient information. It coverts raw data into

subsidiary information. It analyzes the data for relationships that have not antecedently been discovered. The steps of data mining are: Data cleaning, data integration, data cull, data transformation, data mining, pattern evaluation and cognizance representation. Medical data mining is a of imprecision domain of lot and skepticality. The clinical decisions are conventionally predicated on the medicos intuition. Consequently this may lead to disastrous consequences. Due to this there are many errors in the clinical decisions and it results in extortionate medical costs. Serialization is withal utilized in this system. It converts the data objects into streams of bytes and stores it into database.

2. RELATED WORK

Many hospital information systems are designed to fortify patient billing, inventory management and generation of simple statistics. Some hospitals use decision support systems, but they are largely constrained. They can answer simple queries like "What is the average age of patients who have heart disease?", "How many surgeries had resulted in hospital stays longer than 10 days?", "Identify the female patients who are single, above 30 years old, and who have been treated for cancer." However, they cannot answer intricate



p-ISSN: 2348-6848 e-ISSN: 2348-795X Volume 04 Issue 06 May 2017

queries like "Identify the paramount Preoperative prognosticators that increase the length of hospital stay", "Given patient records on cancer, should treatment include chemotherapy alone, radiation alone, or both chemotherapy and radiation?", and "Given patient records, soothsay the probability of patients getting a heart disease." Clinical decisions are often made predicated on doctors" intuition and experience rather than on the erudition- opulent data obnubilated in the database. This practice leads to unwanted biases, errors and extortionate medical costs which affects the quality of accommodation provided to patients. Wu, et al proposed that integration of clinical decision support with computerbased patient records could reduce medical errors, enhance patient safety, decrement unwanted practice variation, and ameliorate patient outcome. This suggestion is promising as data modeling and analysis implements, e.g., data mining, have the potential to engender a cognizance-opulent environment which can avail to significantly amend the quality of clinical decisions

3. IMPLEMENTATION

Naïve Bayes Algorithm

Ingenuous Bayes classifier can be trained in supervised learning setting. It utilizes the method of maximum kindred attribute. It has

been worked in involute authentic world situation. It requires iota of training data. It estimates parameters for relegation. Only the variance of variable need to be tenacious for each class not the entire matrix. Naïve bayes is mainly used when the inputs are high. It gives output in more sophisticated form. The probability of each input attribute is shown from the prognosticable state. Machine learning and data mining methods are predicated on naïve bayes relegation. Naïve bayes will rudimentally soothsay the output whether the patient will have chances of getting the heart disease or not. The model dataset which we get after applying K-Betokens algorithm will compared the values of dataset with a trained dataset. It will apply the bayes theorem and the probability will be obtained whether the patient will have heart disease or not

Algorithm Steps

Outlook	Temp	Humidity	Windy	Play Golf	
Rainy	Hot	High	False	No	
Rainy	Hot	High	True	No	
Overcast	Hot	High	False	Yes	
Sunny	Mild	High	False	Yes	
Sunny	Cool	Normal	False	Yes	
Sunny	Cool	Normal	True	No	
Overcast	Cool	Normal	True	Yes	
Rainy	Mild	High	False	No	
Rainy	Cool	Normal	False	Yes	
Sunny	Mild	Normal	False	Yes	
Rainy	Mild	Normal	True	Yes	
Overcast	Mild	High	True	Yes	
Overcast	Hot	Normal	False	Yes	
Sunny	Mild	High	True	No	

Fig 1 Sample Medical Data Set values



https://edupediapublications.org/journals

-		Play G	olf			Play G	olf
7		Yes	No			Yes	No
	Sunny	3 3/9	2 2/5	Temp.	Hot	2 2/9	2 2/5
Outlook	Overcast	4 4/9	0 0/5		Mild	44/9	2 2/5
	Rainy	2 2/9	33/5		Cool	33/9	11/5
		Play G	olf			Play G	olf
		Play G Yes	No			Play G Yes	olf
	High	Play G Yes 3 3/9	No 44/9		False	Play G Yes 6 6/9	olf No 2 2/5

T

Fig 2 Frequency Tables from Data Set

values

 The posterior probability can be calculated by first, constructing a frequency table for each attribute against the target

Then, transforming the freq. tables to likelihood tables and finally using the

Naive Bayesian equation to calculate the posterior probability for each class • The class with the highest posterior probability is the outcome of prediction



Fig 3 Posterior Probability Calculation



Fig 4 Result

4. K-MEANS CLUSTERING

K-denotes is simplest learning algorithm to solve the clustering quandaries. The process is simple and facile, it relegates given data set into certain number of clusters. It defines

k centroids for each cluster. They must be placed as much as possible far away from each other. Then take each point belonging to given data set and relate into the most proximate centroid. If no point is pending then an group age is done. Then we recalculate k incipient centroid for the cluster resulting from anterior steps. When we get the k centroid an incipient binding is to be done between lucid data points and most proximate centroid. A loop is been engendered because of this loop key centroid transmute the location step by step until no more changes are done.

Algorithm Steps



Fig 5 K-means clustering algorithm

5. EXPERIMENTAL RESULTS



International Journal of Research

Available at https://edupediapublications.org/journals

p-ISSN: 2348-6848 e-ISSN: 2348-795X Volume 04 Issue 06 May 2017

Heart Disease Dataset Logout No Age Gender Smoling Heart Ret Chest Pain Chest Chalestereit Blood pressure Sugar Blood Sugar Blood Sugar Blood Persure Sugar Persure Sugar	Home											
Logout Sno Age Genderr Smoking Heart Refe Cluest Pain Cluestered Blood pressure Blood Sugar 1 33 M Y 45 6 200 65 80 M 2 55 F N 66 9 256 88 99 M 3 77 M Y 87 5 222 142 151 M 4 55 M Y 55 2 155 121 200 M 6 89 M N 88 5 240 120 222 M 7 78 M Y 77 6 355 91 99 M	Heart Disease Dataset											
1 33 M Y 45 6 200 65 80 M 2 55 F N 66 9 256 88 99 M 3 77 M Y 87 5 222 142 151 M 4 55 M Y 55 2 155 121 200 M 5 66 M Y 56 8 239 139 122 M 6 89 M N 88 5 240 120 222 M 7 78 M Y 77 6 355 91 99 M	Logout		Sno	Age	Gender	Smoking	Heart Rate	Chest Pain	Cholesterol	Blood pressure	Blood Sugar	Heart Attack
2 55 F N 66 9 256 88 99 1 3 77 M Y 87 5 222 142 151 1 ↓ 4 55 M Y 55 2 155 121 200 1 5 66 M Y 56 8 239 139 122 1 6 89 M N 88 5 240 120 222 1 7 78 M Y 77 6 355 91 99 N			1	33	М	Y	45	6	200	65	80	¥
3 77 M Y 87 5 222 142 151 M ↓ 4 55 M Y 55 2 155 121 200 M 5 66 M Y 56 8 239 139 122 M 6 89 M N 88 5 240 120 222 M 7 78 M Y 77 6 355 91 99 M			2	55	F	N	66	9	256	88	99	N
↓ 4 55 M Y 55 2 155 121 200 Y 5 66 M Y 56 8 239 139 122 1 6 89 M N 88 5 240 120 222 Y 7 78 M Y 77 6 355 91 99 Y			3	77	М	Y	87	5	222	142	151	N
5 66 M Y 56 8 239 139 122 1 6 89 M N 88 5 240 120 222 1 7 78 M Y 77 6 355 91 99 N		Δ	4	55	М	Ŷ	55	2	155	121	200	Y
6 89 M N 88 5 240 120 222 1 7 78 M Y 77 6 355 91 99 1		м	5	66	М	Ŷ	56	8	239	139	122	N
7 78 M Y 77 6 355 91 99 N			6	89	М	N	88	5	240	120	222	Ŷ
			7	78	М	Y	77	6	355	91	99	Y

Fig 6 Medical Data Set

Heart Disease Dataset						
Logout Frequen	Frequency Table for age data					
Age Catego	ry Yes Score	No Score				
180	2	3				
31-60	304	196				
61-90	53	40				
91-120	0	1				

Fig 7 Frequency Table for Age Data



Enter Patient Details

Age	22 I 🗘
Gender	Male ~
Smoker	Yes ~
Heart Rate	78
Chest Pain	\$3
Cholesterol	
Blood Pressur	e
Blood Sugar	
Submit	

Fig 8 Input values

Туре	Set	Yes Score	No Score
Age	31-60	305.0	197.0
Gender	F	60.0	121.0
Smoker	N	120.0	61.0
Heart Rate	81-100	121.0	79.0
ChestPain	0-3	180.0	1.0
Cholesterol	100-200	166.0	1.0
Bloodpressure	91-120	121.0	1.0
BloodSugar	80-150	181.0	181.0
(Yes_Score_of_bp/Tot of Ye p(yes)=3.777237407132538	s Score)*(Yes_Score_of_bs 86E-4	Viot of Yes Score)"(Yes	Tot/Total)
p(no)=(No_Score_of_age/To (No_Score_of_hr/Tot of No (No_Score_of_hp/Tot of No 10	ot of No Score)*(No_Score Score)*(No_Score_of_cp/T Score)*(No_Score_of_bs/	e_of_gen/Tot of No Scor Tot of No Score)*(No_Sc Tot of No Score)*(Tot N	e)*(No_Score_of_emoker/Tot of No Score)* ore_of_ch/Tot of No Score)* o Score/Total) p(no)=7.568008013928807E
p(yes)+p(no)=3.777244975	1405525E-4		
p(yes)/(p(yes)+p(no))=0.99	99979964211843		

p(no)/(p(yes)+p(no))=2.0035788157073925E-6

RESULT: POSSITIVE

Fig 9 Result after Mining

6. CONCLUSION

In this paper we are proposing heart disease prognostication system utilizing naïve bayes and k-designates clustering. We are utilizing k-betokens clustering for incrementing the efficiency of the output. This is the most efficacious model to prognosticate patients with heart disease. This model could answer intricate queries, each with its own vigor with deference to facilitate of model interpretation, access to detailed information and precision

7. REFERENCES

[1] Sellappan Palaniappan, Rafiah Awang "Intelligent Heart Disease Prediction System Using Data Mining Techniques"Department of Information Technology Malaysia University of Science and Technology Block C, Kelana Square, Jalan SS7/26 Kelana Jaya, 47301 Petaling Jaya, Selangor, Malaysia .



[2] "CSV File Reading and Writing" (http://docs. python. org/library/csv. html).
Retrieved July 24, 2011. "is no "CSV standard""

[3] Y. Shafranovich. "Common Format and MIME Type for Comma- Separated Values (CSV) Files" (http://tools.ietf.org/html/ rfc4180) Retrieved September 12, 2011.

[4]

home.deib.polimi.it/matteucc/Clustering/tut orial_html/kmeans.html "A tutorial on clustering algorithms".

[5] Shadab Adam Pattekari and Asma
Parveen "Prediction System For Heart
Disease Using Naïve Bayes" International
Journal of Advanced Computer and
Mathematical Sciences ISSN 2230-9624.
Vol 3, Issue 3, 2012, pp 290-294.

[6] Mrs.G.Subbalakshmi (M.Tech), Mr. K.
Ramesh M.Tech, Asst. Professor Mr. M.
Chinna Rao M.Tech,(Ph.D.) Asst. Professor,
"Decision Support in Heart Disease
Prediction System using Naïve Bayes"
G.Subbalakshmi et al. / Indian Journal of
Computer Science and Engineering
(IJCSE)2011.

[7] Jesmin Nahar, Tasadduq Imama, Kevin S. Tickle, Yi-Ping Phoebe Chen "Association rule mining to detect factors which contribute to heart disease in males and females" Expert Systems with Applications 40 (2013) 1086–1093. [8] Oleg Yu. Atkov (MD, PhD), Svetlana G. Gorokhova (MD, PhD), Alexandr G. Sboev (PhD), Eduard V. Generozov (PhD), Elena V. Muraseyeva (MD, PhD), Svetlana Y. Moroshkina, Nadezhda N. Cherniy "Coronary heart disease diagnosis by artificial neural networks including genetic polymorphisms and clinical parameters" Journal of Cardiology (2012) 59, 190-194. [9] Shantakumar B.Patil Y.S.Kumaraswamy "Intelligent and Effective Heart Attack Prediction System Using Data Mining and Artificial Neural Network" European Journal of Scientific Research ISSN 1450-

216X Vol.31 No.4 (2009), pp.642-656.

[10] Sivagowry, Dr. Durairaj. M2 and Persia. "An Empirical Study on applying Data Mining Techniques for the Analysis and Prediction of Heart Disease" 2013.