

MINING OF NUTRITIONAL INGREDIENTS IN FOOD FOR DISEASE ANALYSIS

S.Raju¹, P.V.Sarath Chand², Dr. M Bal Raju³, Dr.G.S.S.Rao⁴,
 Assistant Professor^{1,2}, Professor^{3,4},
 Department of CSE,
 Pallavi Engineering College^{1,2,4}, Swamy Vivekananda Institute of Technology³,
 Mail ID:srajunayak@gmail.com, Mail ID:chandsarath70@gmail.com, Mail ID:drarajucse@gmail.com
 Kuntloor(V), Hayathnagar(M), Hyderabad, R.R. Dist. -501505.

ABSTRACT

In the prevention and treatment of noncommunicable illnesses, such as cancer, it has long been recognised that a well-balanced, nutritious diet is essential (NCDs). Research has been conducted on the nutritional components of food that are beneficial in the rehabilitation of noncommunicable diseases, on the other hand, but only a small amount has been done. Because of the use of data mining technologies, we were able to conduct a thorough investigation into the association between food components and illnesses. In order to get started, we compiled a list of more than 7,000 disorders, after which we decided which foods were recommended for each condition and which foods were strictly forbidden. Using the China Food Nutrition as a reference, we went on to predict which nutritional ingredients are most likely to have beneficial impacts on disease using noise-intensity and information entropy.

At the conclusion of the research, we proposed an improved technique called CVNDA Red, which is based on rough sets and is used to select the necessary core ingredients from among the most favourable nutritional components. CVNDA Red is based on rough sets and is used to select the necessary core ingredients from among the most favourable nutritional components. A contraction of the phrases CVNDA and Red, which translates as "CVNDA Red." CVNDA Red is a trademark of the CVNDA Corporation. According to our knowledge, this is the first research in China to analyse the association between nutritious elements in food and illnesses via the use of data mining techniques based on rough set theory, which we believe is the case. We have shown via experiments carried out on real-world information that our data mining technique outperforms the conventional statistical approach, with accuracy 1.682 times greater than the conventional statistical methodology. By way of aside, our research has been beneficial in uncovering the first two to three nutritional components contained within foods that may be used to aid in the rehabilitation of a range of common conditions such as type 2 diabetes, hypertension, and cardiovascular disease. These experimental findings indicate the utility of using data mining to choose nutritional components in food for illness analysis when choosing nutritional ingredients in food when selecting nutritional elements in food when selecting nutritional components in food.

1. INTRODUCTION

As defined by the National Council on Chronic Illnesses (NCDS), chronic illnesses are those that are primarily caused by occupational and environmental factors, as well as lifestyle and behavioural variables. According to the organisation, chronic illnesses include obesity and diabetes as well as hypertension and tumours, among other diseases. Global Health Organization's (WHO) Global Status Report on Noncommunicable Diseases (Global Status Report on NCDs) states that the number of people who die each year from NCDs is increasing, resulting in a significant economic burden for the whole world's population. Noncommunicable diseases (NCDs) are responsible for over 40 million deaths per year worldwide, accounting for approximately 70% of all mortality on the globe. Chinese chronic disease and nutrition statistics show that the number of patients suffering from noncommunicable diseases (NCDs) in the country

outnumbers those in any other country on the planet, and that China's current prevalence rate has risen far above that found in any other country on the planet. According to government statistics, the number of people aged 60 and over in China has surpassed 230 million, with noncommunicable diseases (NCDs) accounting for around two-thirds of those affected (NCDs). As a consequence, relevant departments in each nation, particularly in China, such as medical schools, hospitals, and disease research organisations, are all worried about noncommunicable diseases (NCDs), which are a result of this (NCDs). It is vital to eat nutritious meals (NCDs) in order to maintain health and avoid the advent of noncommunicable diseases (NCDs) (NCDs). As a result of the rising adoption of this paradigm in China, the nation has also re-configured the relationship between food and health. However, data mining is still considered to be a novel method of investigating the nutritional

features of food that might be used to help in the rehabilitation of people suffering from ailments in China, according to experts. China is only getting started when it comes to developing IT (Information Technology) infrastructure for smart health-care delivery. Most research into the relationship between nutritional components in food and illness are still carried out using expensive precision tools or long-term clinical trials, as is still the case today. Beyond that, further preventative studies have been published; however, most of them only focused at one or a few disorders at a time.

When it comes to using data mining to investigate the association between dietary elements and illnesses, Chinese academics are just in their infancy. When it comes to patients suffering from noncommunicable diseases (NCDs), physicians tend to prescribe just certain meals, neglecting to provide any further crucial nutrition counselling to these patients, particularly on the nutritional elements found in foods. Answers to NCDs must be based on the use of a wide range of competencies. In light of the massive amounts of data being gathered today, data mining has emerged as a critical method of unearthing new information across a broad variety of businesses, particularly in the fields of illness prediction and precision healthcare delivery (AHC) (AHC). In recent years, it has emerged as a major source of funding for researchers in a variety of fields, including preventive medicine, fundamental medical, and clinical medicine. Our key contributions to the field of disease analysis are as follows in terms of sickness analysis via the mining of nutritional components in food: To evaluate which nutritional elements in food may have a positive impact on illness, we used the noise-intensity and information entropy measures. In addition, the data in this research is continuous and does not contain any choice features. We feel that you will find this document to be of great use. As a consequence, we developed an improved algorithm, dubbed CVNDA Red, that is based on rough set theory and is capable of selecting more closely linked core components from among the positive nutritional compounds

found in food. CVNDA Red is a mixture of the drugs CVNDA and CVNDA Red. CVNDA Red is the name of the approach used in this case. In accordance with the following, the following is the organisational structure of this work: This section offers readers a high-level summary of current research in the disciplines of disease analysis and data mining that is relevant to the subject. There is a thorough description of the specific data mining methods utilised in this work, as well as the rationale behind the selection of the algorithms and the two assessment indexes that we employed. An in-depth examination of the data, experimental results, and statistical analysis will be provided in this section. This section provides an examination of the differences and similarities between various approaches. Some of the findings are presented, as well as some potential future study subjects that have been examined.

2. LITERATURE SURVEY

It is the technique for doing a literature review that is the most critical phase in the software development process. Before the tool can be developed, many aspects must be taken into account, including time limits, budgetary concerns, and the general strength of the business. Once these requirements have been met, the following ten steps will be utilised to identify which operating system and programming language will be used to construct the tool in order to finish the project. Once these phases have been completed, the project will be completed. The programmers will need a significant amount of assistance from other sources in order to accomplish their task once they begin working on the tool. For example, assistance may be gained from senior programmers, books, or websites, among other sources of information. Preliminary considerations are given to the components listed above prior to the building of the system, which will allow for the development of the proposed system. During a large meeting in India, researchers conducted a retrospective study on hypertension screening and discovered that The significance of these findings for

noncommunicable disease prevention and control programmes are examined in further depth further down this page. Cardiovascular disease is the most common kind of non-communicable disease in India, and it is responsible for the great majority of non-communicable deaths in the country. National Program for the Prevention and Control of Cancer, Diabetes, Cardiovascular Disease, and Stroke is a government-sponsored initiative that aims to prevent and control illnesses such as cancer as well as diabetes, cardiovascular disease, and stroke. With this initiative, the government of India hopes to expand capacity development for NCD screening, referral, and management across the country. The project includes outreach and screening programmes in local communities, as well as government-sponsored research and other activities. During religious holidays that bring large throngs of people to the nation, the Indian government is typically responsible for providing basic medical care.

During the 2015 KumbhMela, which took place in Nashik and Trimbakeshwar, the state government expanded its services to include a hypertension screening programme, which was established and implemented by the Maharashtra government. Over the course of this essay, we will analyse the merits and disadvantages of opportunistic screening for large groups of individuals. The number of people who elected to have their blood pressure checked at the KumbhMela was 5760, with each participant obtaining a single blood pressure test to record their decision. 1783 people (33.6 percent) tested positive in total, with 1580 of those participants being entirely oblivious of their condition prior to being subjected to testing. Prescription drugs were given to patients with previously diagnosed hypertension, with 240 (79 percent) of those who got treatment following their doctors' orders (that is, 52.8 percent under treatment). The blood pressure of 55 individuals (18 percent) in this research was within normal ranges, according to the findings (BP under control).

According to the results of the study, cigarette smokers had a greater prevalence of hypertension (39

percent) as compared to non-users (28 percent), which was statistically significant ($P=0.001$). Because of a failure to collect phone numbers (0.01 percent), any phone-based follow-up was unable to take place as a result of this. Due to poor levels of knowledge, treatment, and control of the condition, India continues to have a significant challenge in terms of both screening and management of hypertension, and this will continue in the future.

3. SYSTEM ANALYSIS

3.1 Existing system:

Because they would not be recommended unless there was a compelling reason to do so, all recommended meals must have a high degree of stability in order to be recommended in the first place. The use of specific dietary ingredients within the context of other advised meals if they are not designated PNIs for a certain condition are examples of alternatives. On the other hand, this is not always the case (poor stability). Specifically, the SA considers just the quantity of nutritional component values included in a product when deciding whether or not a product includes PNIs.

Disadvantages:

One significant change is that the overall level of performance has been reduced.

The one that is most cost-effective in the long run. As an example, the system that has been proposed is as follows: We may soon be able to make meal suggestions based on the amount of Creatinine that is present in the body, according to a number of theories. However, the research described above is mostly carried out through long-term clinical trials that only propose diets for specific disorders, and it rarely explores the association between nutritional elements and diseases through the use of data mining techniques, as opposed to the research described above.

Advantages:

1. Overall performance has improved as a result of these changes.

Second, they are less cost-effective than other solutions.

4. ALGORITHM**4.1 Linear Regression**

According to one definition, linear regression is a statistical model that evaluates the linear connection that exists between a dependent variable and a collection of independent variables. It is also known as a correlation analysis. According to this relationship between variables, when the value of one or more independent variables changes (increases or decreases), the value of the dependent variable will change in proportion to the increase or decrease in the value of the independent variables. This relationship between variables can be expressed as (increase or decrease).

With the help of the following equation, we can quantitatively define the relationship between the variables:

$$Y = mX + b$$

When we are trying to make predictions about the dependent variable, Y, we are referring to the situation in which we are now working.

In order to make predictions, it is important to have a dependent variable, which is represented by the letter X in the equation.

As a result, in regression analysis, M is the slope of a regression line, which indicates how well an X-Y relationship holds up under various conditions (e.g., under different conditions).

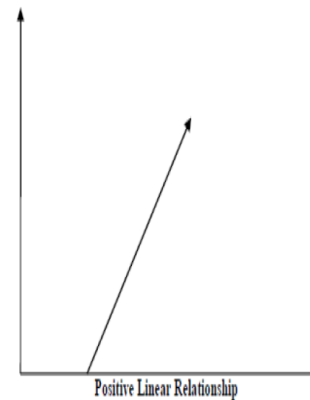
It is a constant that is symbolised by the letter b in the equation, which is the Y-intercept. Suppose that X

= 0 and Y = 0, in which case the value of Y equals the value of b, as shown in the diagram.

Positive linearity exists in a connection when the two parties are in agreement.

When both the independent and dependent variables are positive, this is referred to as a positive linear connection, which is described as follows:

As a consequence of this, the value of the variable increases. To help you better comprehend the problem, the accompanying diagram may be of use to you.

**4.2 K-nearest neighbor's algorithm (k-NN)**

The k-nearest neighbours algorithm (k-NN) is a non-parametric classification and regression approach that may be used to any kind of data collection and is referred to as the k-NN algorithm. As is true in both circumstances, the input is comprised of the k training examples that are the closest to each other in the feature space, which is identical to the input in both cases. When the k-NN approach is used to classification or regression issues, the following results are obtained as a result of the algorithm:

As an example, in the case of k-NN classification, the outcome is a class membership that may be utilised to improve the classification even more. An item is

categorised according to the votes of its neighbours who have cast a majority vote in a classification system, with the object being allocated to the class that is most often represented among its k nearest neighbours (k is a positive integer, typically small). The object is allocated to the class of its closest neighbour if k is equal to 1, in which case it is assigned to the class of the object's single nearest neighbour, and so on.

This is the value of a property associated with the object that is the outcome of a K-NN regression in this instance. In this circumstance, the average of the values of the k closest neighbours is utilised to get the result. Slow learning is a kind of instance-based learning that differs from other types of learning in that it approximates the function only locally and that all computation is postponed until after classification is completed. It is an effective approach for both classification and regression when weights are applied to the contributions of the neighbours, since the neighbours who are closer to the average contribute more to the average than the neighbours who are farther away. For example, a common weighting technique would be to assign a weight of $1/d$ to each neighbour in a pair of neighbours, where d signifies the distance between the neighbours in a pair of neighbours and $1/d$ denotes the weight of the neighbour in the pair of neighbours. An object collection with known classes (in the case of k-NN classification) or an object collection with known attribute values of objects (in the case of k-NN regression) is used to generate neighbours for classification and regression in both classification and regression. There is a possibility that this data will serve as the algorithm's training set; nevertheless, there is no need for an explicit training phase to be performed on this data.

4.3 Working of KNN Algorithm

In the K-nearest neighbours (KNN) algorithm, values of new datapoints are predicted based on their 'feature similarity,' which means that a new data point will be assigned a value based on how closely it matches the points in the training set, which is further

defined as how closely it matches the points in the training set, and so on. If we study the stages stated below, we may have a better grasp of how it operates:

First, we'll take a look at the options available. Any algorithm's implementation is difficult to do without the usage of a dataset of some kind. As a result, during the first phase of the method, we must enter both the training and test data into the KNN simultaneously.

We must pick the value of K that corresponds to those data points which are the closest to it from the available possibilities in order to proceed to the second step. K may be any integer, positive or negative, and it does not have to be an even number.

The following processes should be followed throughout the third stage, at each point in time when the test data is collected:

3.1 Calculate the distance between each row of test data and each row of training data using one of the distance algorithms listed below: the Euclidean distance, the Manhattan distance, or the Hamming distance, depending on your preference. While there are other approaches to determine distance, the Euclidean distance computation technique is the most often utilised. It's also the most accurate of the three.

In 3.2, arrange them in ascending order depending on the distance value that exists between them at this moment in time. 3.3

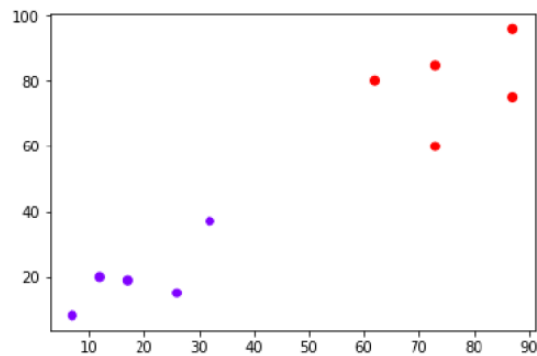
In the next stage, the first K rows of data from the sorted array will be selected as the starting point for the analysis.

Each test point will be classified based on the class that appears the most often in these rows, with the classification decided as a result of the most frequently appearing class in these rows.

This is the fourth step. We have reached the conclusion of the process.

Example

Using the following example, you will be able to better grasp the notion of K as well as the operation of the KNN algorithm. For example, suppose we have a dataset that can be shown in the following way:



Now, we need to classify new data point with black dot (at point 60,60) into blue or red class. We are assuming $K = 3$ i.e. it would find three nearest data points. It is shown in the

5. Results:

View Data

| Labels | Food Type | Disease | Minerals | Grains |
|--------|------------------|----------------|-------------|--------|
| A | vegetables | Angina | Vitamin A | 400 |
| B | meat | Acne | C | 0.8 |
| C | fruits | cardiovascular | E | 400 |
| D | Dairy Foods | cholan | Vitamin B12 | |
| E | Grains | Stroke | ماغنسيوم | 3 |
| F | Beans and Nuts | tooth decay | potassium | m |
| G | Fish and Seafood | Psoriasis | iron | p |
| H | liquid drinks | liver disease | copper | NaN |
| I | tobacco food | oral cancers | Vitamin A | 2 |
| J | potato chips | Hypertension | Sodium | 100 |
| K | vegetables | Kidney stone | calcium | 400 |
| L | meat | NaN | NaN | 0.8 |
| M | fruits | Angina | Vitamin A | 400 |

Mining of Nutritional Ingredients in Food for Disease Analysis

Disease Predicted : ANGINA

Food Type: Enter between (0-4)

Minerals: Enter upto (0-10)

Grains: Enter upto (400)

Submit

All rights reserved ©

Mining of Nutritional Ingredients in Food for Disease Analysis

Disease Predicted : ANGINA

Food Type: 5

Minerals: 1

Grains: 400

Submit

All rights reserved ©

Mining of Nutritional Ingredients in Food for Disease Analysis

Disease Predicted : CARDIO VASCULAR

Food Type: Enter between (0-4)

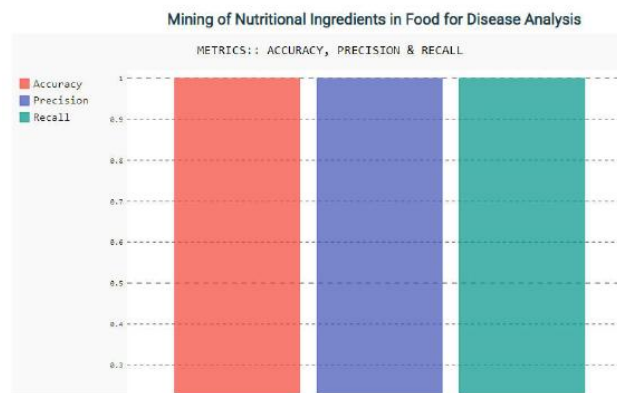
Minerals: Enter upto (0-10)

Grains: Enter upto (400)

Submit

All rights reserved ©

Graph:



CONCLUSION

Not to mention, the project is primarily concerned with the construction of a system for analysing the performance of university students as a whole. As a result of this study, it is possible to predict students' performance in the course "TMC1013 System Analysis and Design" by using a data mining technique and a classification system. While in the classroom, one of the most important contributions of the SPAS is to assist teachers in doing student performance analysis, which is the SPAS's most significant contribution. The system's aid in identifying students who are most likely to fail the course TMC1013 System Analysis and Design will be beneficial to professors teaching the subject. Furthermore, by offering a centralised database, SPAS supports lecturers in gaining access to information regarding their students' performance throughout the semesters. In the first part, we gathered and classified more than seven thousand Chinese diseases from medical and government websites, as well as the foods that were recommended and forbidden in the context of each disease; in the second, we discussed the relationship between nutritional ingredients and diseases, with the primary goal of identifying which ingredients are beneficial in the rehabilitation of diseases. According to our knowledge, this is the first research in China to examine the association between nutritious elements

in food and illnesses via the use of computer-assisted data mining methods. Although the trials demonstrated that data mining approaches were unsuccessful in discovering and choosing all of the beneficial dietary elements for illnesses, the first two or three items were accurately identified and picked, indicating that the methodology was effective. Furthermore, if our point of view can be combined with foods that are prohibited, the results are more likely to be better and more in line with reality, which is the direction in which our future work will be directed. First and foremost, executing this activity has enormous time savings implications. The identification of excellent nutritional ingredients that are beneficial to the rehabilitation of illnesses as exactly as possible may be aided by us, and we can also help doctors and sickness researchers in the prevention of disease. To get information about diseases, suggested and forbidden foods, and related nutrition information that is relevant to this article, you may check our website², which is accessible through the link provided above. Some information is not yet accessible because it is presently undergoing medical verification at the time of writing, which means that it is not yet available. We also hope that, since our knowledge base is always expanding, researchers would inform us if they come across anything in our work that is inaccurate, allowing us to make our study even more accurate.

REFERENCES

- [1] CNS, "2016 Global Nutrition Report," in *Chinese Nutrition Society*, 2016.
- [2] WHO, "Global Status Report on Noncommunicable Diseases (2014)," in *World Health Organization*, 2014.
- [3] S. Balsari, P. Vemulapalli, M. Gofine et al., "A Retrospective Analysis of Hypertension Screening at a Mass Gathering in India: Implications for Non-communicable Disease Control Strategies," *Journal of Human Hypertension*, vol. 31, no. 11, pp. 750–753, 2017.
- [4] DNHFPC of PRC, "Chinese Resident's Chronic Disease and Nutrition (2015)," in *National Health and Family Planning Commission of the People's Republic of China*, 2015.
- [5] S. Tellier, A. Kiaby Lars, P. Nissen et al., "Basic Concepts and Current Challenges of Public

Health in Humanitarian Action,” International Humanitarian Action, pp. 229–317, 2017.

[6] F. Ara1, F. Saleh, S. J. Mumu, F. Afnan and L. Ali, “Awareness Among Bangladeshi Type 2

Diabetic Subjects Regarding Diabetes and Risk Factors of Non-communicable Diseases,”

Diabetologia, pp. S379, 2011. DOI:10.1007/s00125-011-2276-4.

[7] QIANZHAN, “Report of Market Prospective and Investment Strategy Planning on China

Intelligent medical construction industry (2017-2022),” in *Qianzhan Intelligence CO.LTD*, 2017.

[8] W. H. Ling, “Progress of Nutritional Prevention and Control on Noncommunicable Chronic

Diseases in China,” *China J Dis Control Prev*, vol. 21, no. 3, pp. 215–218, 2017.

[9] M. B. Margaret, B. K. Barbara and D. Colette, “Developing Health Promotion Workforce

Capacity for Addressing Non-communicable Diseases Globally,” *Global Handbook on*

Noncommunicable Diseases and Health Promotion, pp. 417–439, 2013.

[10] M. Williams and H. Moore, “Lumping Versus Splitting: the Need for Biological Data Mining in

Precision Medicine,” *BioData Mining*, vol. 8, no. 16, pp. 1–3, 2015.

[11] G. M. Oppenheimer, “Framingham Heart Study: The First 20 Years,” *Progress in Cardiovascular*

Diseases, vol. 53, no. 1, pp. 55–61, 2010.

[12] W. Y. Jiao, Y. Xue, T. C. He, Y. M. Zhang and P. Y. Wang, “Association Between South

Korean Dietary Pattern and Health,” *Food and Nutrition in China*, vol. 23, no. 5, pp. 81–84, 2017.

[13] K. W. Lee and M. S. Cho, “The Traditional Korean Dietary Pattern Is

Associated with Decreased Risk of Metabolic Syndrome: Findings from the Korean National Health

and Nutrition Examination Survey 1998–2009,” *Journal of Medicinal Food*, vol. 17, no. 1, pp. 43–56, 2014.