



SIES (NERUL) COLLEGE OF ARTS, SCIENCE AND COMMERCE

PROJECT PROPOSAL

on

DIET RECOMMENDED SYSTEM

**PROJECT WORK SUBMITTED IN PARTIAL FULFILLMENT
OF THE REQUIREMENT FOR THE AWARD OF THE
DEGREE OF MSc. (COMPUTER SCIENCE)**

SUBMITTED BY

Rohan Thakur

PROJECT GUIDE

Asst. Professor Manasvi Sharma



SIES (NERUL) COLLEGE OF ARTS SCIENCE AND COMMERCE
NAAC ACCREDITED 'A' GRADE COLLEGE
(ISO 9001:2015 CERTIFIED INSTITUTION)
NERUL, NAVI MUMBAI - 400706

Certificate

THIS IS TO CERTIFY THAT THE PROJECT TITLED

DIET RECOMMENDED SYSTEM

IS UNDERTAKEN BY

ROHAN KIRAN THAKUR

Seat No: 20 —

In partial fulfilment of MSc - IT / CS Degree (Semester — **III** — Examination

in the academic year **2021-2022** and has not been submitted for any other examination and does not form part of any other course undergone by the candidate. It is further certified that he/she has completed all the required phases of the project.

Project Guide

External Examiner

Head of Department

Principal

ACKNOWLEDGMENT

I extend my heartfelt gratitude and thanks to Asst. Professor Mansvi Sharma sir for providing me excellent guidance to work on this project and for their understanding and assistance by providing all the necessary information needed for my project topic.

I would also like to acknowledge all the staffs for providing a helping hand to us in times of queries & problems. The project is a result of the efforts of all the peoples who are associated with the project directly or indirectly, who helped us to work to complete the project within the specified time frame. They motivated me in the project and gave a feedback on it to improve my adroitness.

Thanks to all my teachers, who were a part of the project in numerous ways and for the help and inspiration they extended to me and for providing the needed motivation.

With all Respects & Gratitude, I would like to thanks to all the people, who have helped for the development of the Project.

ROHAN THAKUR

MSc. Computer Science (Part-II)

SIES (Nerul) College of Arts, Science, and Commerce.

Table of Content

1) Title	5
2) Introduction	6
3) Related Works	8
4) Objective	15
5) Methodology	16
6) Conclusion	20
7) Reference	22

Project On:
Diet Recommended System

Introduction:

Recommender systems (RS) suggest items of interest to users of information systems or e-business systems and have evolved in recent decades. A typical and well-known example is Amazon's suggest service for products. We believe the idea behind recommender systems can be adapted to cope with the special requirements of the health domain. Recommendations based on single foods or food groups are easier to implement when only a few foods serve as the major sources of an essential dietary component (e.g., dairy products, which are the primary source of calcium in the Indian diet).

During the last decades huge amounts of data have been collected in clinical databases representing patients' health states (e.g., as laboratory results, treatment plans, medical reports, diet plans). Hence, digital information available for patient-oriented decision making has increased drastically but is often scattered across different sites. As a solution, personal diet recommender systems (DRS) are meant to centralize an individual's health data and to allow access for the owner as well as for authorized health professionals

Nutrients are Essential compounds that the body can't make or can't make insufficient quantity. According to the World Health Organization (WHO), these nutrients must come from food, and they're vital for disease prevention, growth, and good health. Macronutrients are eaten in large amounts and include the primary building blocks of your diet protein, carbohydrates, and fat which provide your body with energy. Vitamins and minerals are micronutrients, and small doses go a long way. Most of disease occurred due to efficiency of nutrients. To fill these nutrients, we can suggest natural diet (that have no side effects), and precautions to user.

Nutrient-based food recommendations (e.g., due to vitamin D deficiency increase the risk skin disease to reduce this we suggest oily fish, meat etc.) might be easy for public health personnel to interpret and implement, recommendations pertaining to nutrient intake would usually need to be translated by professionals into guidance about food choices for the public.

Nutrient-based recommendations must be derived often from the epidemiologic (Distribution of health-related data) data on dietary patterns. For example, the statement that diets with a high plant food and low-fat content are associated with reduced rates of certain cancers more accurately reflects present knowledge than do conclusions that diets high in selenium or isothiocyanates are likely to reduce cancer risk. The latter requires an inference about cause and effect that is not yet justified by the data.

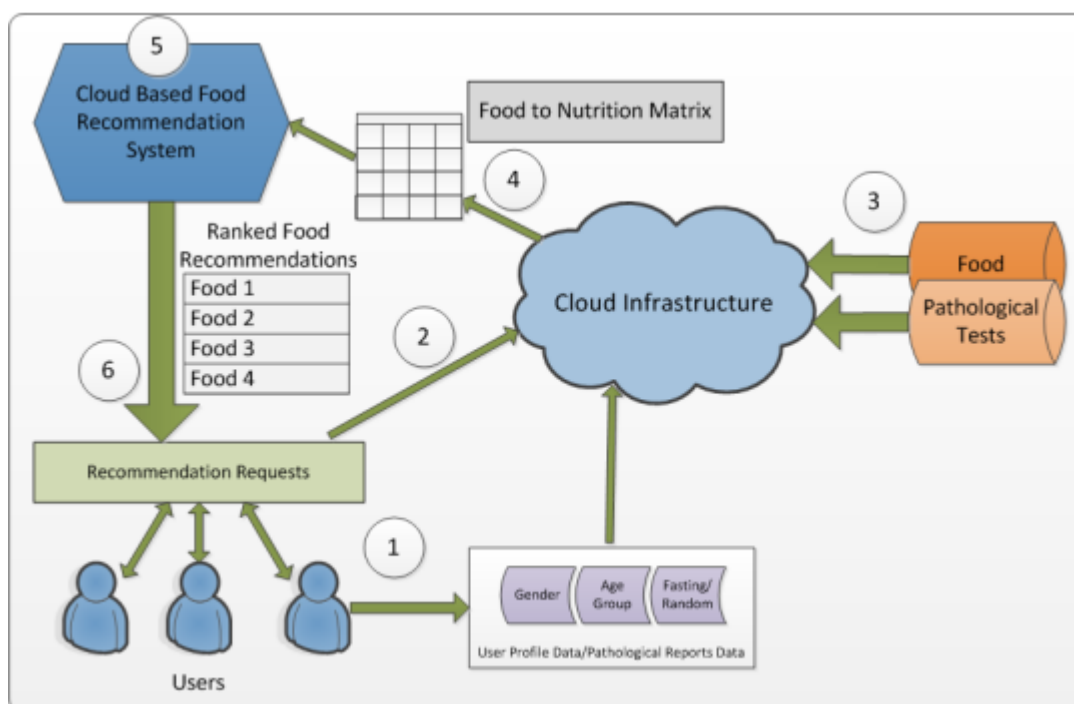
Similar work: Electronic Medical Records (EMRs) and Electronic Health Records (EHRs) provide the technology for the electronic storage of medical data and enables hospitals and other healthcare players to share such data electronically among authorized caregivers.

Related Work:

Several works have been proposed for different recommendation systems related to diet and food. These systems are used for food recommendations, menu recommendations, diet plan recommendations, health recommendations for specific diseases, and recipe recommendations. Majority of these recommendation systems extract users' preferences from different sources like users' ratings.

The main focus of this work is to provide dietary assistance to different people who are suffering from common diseases. The proposed model recommends various foods and nutrition to the people based on their pathological test reports. Every pathological report has some indicators that are calculated based on the nature of the tests. For instance, if a doctor advised a patient to take pathological test of blood, then the common test entries include the values of haemoglobin, red blood cells (RBC), white blood cells (WBC), plasma, and sugar. Normal ranges of the aforementioned indicators are usually given in the test reports. In this way, the patient can identify the abnormalities after comparing with the normal ranges. In our proposed system, a user is provided with the complete list of the test parameters to make selection from. The user inputs the specific values of test report in the selected parameters. We gathered the data of normal ranges for tests including blood (plasma and serum), urine, stool, cerebrospinal fluid (CSF), and gastric and secretion tests. A matrix of 345 entries was constructed. Each individual components of a test (e.g., blood test) have normal ranges with lower and upper bound. The ranges of the same component may differ on the basis of gender, age groups, and fasting or no fasting. Our system is trained on various types of age groups and their respective ranges of parameters. This allows the system to suggest diets as per needs of the users.

3.1 Diet-Right Architecture In majority of the existing food recommendation systems, centralized architecture is used [8-15]. The main disadvantage of using such systems is scalability, when dealing with the massive amount of data. We propose a cloud-based solution to offer the scalability and pervasiveness, where the smart phone users can conveniently access the recommendation system (see Fig. 1). The model takes the input values as a first step. User enters the demographic data including gender and age as well as the value of the pathological test reports. These values are sent to cloud infrastructure in second step and are compared with the normal ranges that are stored in our database. In the third step, the abnormality level of the pathological test reports is computed. In the next step, the weight assignments and matrix generation process are carried out. In the fifth step, ranks are calculated for each food item and are sorted in descending order. In the sixth step, the user is provided the recommended list of food items.



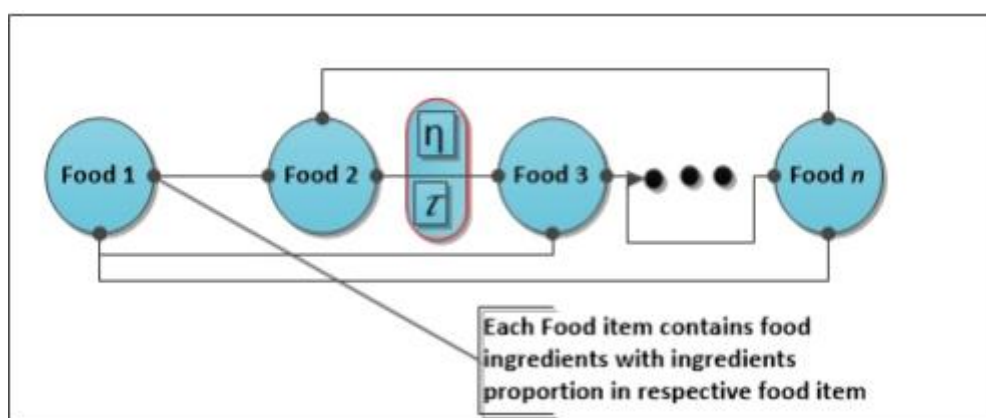
3.2 Proposed Algorithm based on Ant Colony Optimization (ACO) In this subsection, we present the food recommendation process using variant

of ant colony approach on a graph of foods to generate the optimal food set for the users. In Diet-Right, we have used Ant Colony Optimization technique. ACO metaheuristic is a constructive and population based-approach which relies on the social behavior of ants. It is recognized as a most powerful approach for the solution of combinatorial optimization problems [33]. The main steps used in our proposed Algorithm are explained as follows:

Each food item is placed on nodes and a strongly connected graph is generated as shown in Fig. 2. Each link of graph has associated η and τ values, where τ is the randomly initialized pheromone, and η is the heuristic information initialized as the inverse of squared sum of difference of all the ingredients I . In Equation (1), k represents index for food ingredient, i and j represent i th and j th food item, l represents a single ingredient in a certain food item, and m is the total number of ingredients in a certain food item.

$$\eta_{ij} = \sum_{k=1}^m (I_{ik} - I_{jk})^2.$$

Where, η is used to control exploration and exploitation of ACO and the values of $\eta \in (0,1)$.



After initialization, each ant constructs its local solution by visiting nodes which provide best cost in terms of low error compared to targets. Target

vector represents the amount of food ingredients required against the particular disease. Target vector is predefined based on pathological reports, for instance, target vector for user with calcium deficiency may ranges from 9 to 10.5. The different nodes or food items are selected using transition rule which selects a path with highest transition probability. Transition probability is given by Equation (2):

$$p_i^k(t) = \begin{cases} \frac{[\tau_i(t)]^\alpha \times [\eta_i(t)]^\beta}{\sum_{r \in j^k} [\tau_r(t)]^\alpha \times [\eta_r(t)]^\beta}, & \text{if } i \in j^k \\ 0 & \text{Otherwise} \end{cases} \quad (2)$$

Where, $\tau_i(t)$ represents the pheromone level at time t , $\eta_i(t)$ is the heuristic information at time t , and α, β are the hyper parameters in the model used to weight heuristic information

and pheromone level (used for fine tuning). Moreover, k represents ant, i represents initial node, j is the target node, r is index of current selected path, and jk represents the solution. When an ant selects a path among all existing paths (excluding the path in the solution), it updates the pheromone level locally as depicted in Equation (3).

$$\tau_{ij}^k(t+1) = \begin{cases} \tau_{ij}^k(t) \times \rho + \delta \tau^k, & \text{if } ij \in S^k(t) \\ (1 - \rho) \tau_{ij}^k(t) & \text{otherwise} \end{cases} \quad (3)$$

Where, $\tau_{ij}^k(t+1)$ is new pheromone level that is increased by amount $\delta \tau^k$, and evaporation is governed by multiplication of pheromone decaying parameter ρ . Also, S is the solution of the jk ant at time t . Each ant provides locally optimized food set based on the nutrition expert recommendations, but here, we are interested in the globally optimized solution. As we are using supervised approach, we use Root Mean Square Error (RMSE) for the selection of globally best solution. To do so, we initialize global solution to EMPTY set and solve for each ant. Initially, solution returned by the first ant is

considered as the global solution. Afterwards, when rest of the ants return with a solution, their RMSE is compared with the current global solution replacing the global solution with the solution having minimum RMSE value. For fast convergence of the solution, we update the pheromone level again using the same formula, but the update is only for the path that is globally optimal solution as depicted in Equation (4).

$$\tau_{ij}(t+1) = \begin{cases} \tau_{ij}^g(t) \times \rho + \delta\tau^g, & \text{if } ij \in S^g(t) \\ (1 - \rho)\tau_{ij}^g(t) & \text{otherwise} \end{cases} \quad (4)$$

In food selection process, there is a need to select diversified foods to enhance the acceptance of foods among different people. We manage and update the heuristic information in such a way that the diversity among foods is maximized. For heuristic information update, we use Equation (5).

$$\eta_i = \frac{1}{m_i} \sum_{k=1}^n \gamma(S^k(t)) \left(1 + \phi_i e^{\frac{-|S^k(t)|}{n}} \right), \text{ if } i \in S^t(t) \quad (5)$$

Where, m is the selected number of foods, ϕ is the number of times a food is selected in whole iteration, and m , ϕ are the parameters used to balance the solution in terms of local and global perspective. The used heuristic information facilitates in the selection of foods with minimal redundancy. Algorithm 1 presents food recommendation using ACO.

Content based food recommender system [3] is proposed which recommend food recipes according to the preferences already given by the user. The preferred recipes of the user are fragmented into ingredients which are assigned ratings according to the stored users' preferences. The recipes with the matching ingredient are recommended. The authors do not consider the nutrition factors and the balance in the diet. Moreover, chances of identical recommendation are also present because the preference of the user may not change on daily basis. The above-mentioned diet recommendation systems are specifically dealing with some diseases or related to balance the diet plans. In case of food recommendation for specific diseases, the systems recommend different foods for patients without knowing the level of disease which may vary in different cases and cause severe effects on patients. Similarly, in case of food recommendations to balance the diet, nutrition factors are ignored which are very much important to recommend food and balance diet.

Methodology Seeing that the requirements of the application were clear, the waterfall model was used to develop the application.

Algorithm used Content Based Filtering Algorithm: In a content-based recommender system, keywords or attributes are used to describe items. Calculate the weight of each feature, namely Calories, that has the lowest value. This step is in order to make an accurate list of food items. If the category i has the lowest value C_i , but it has many contents' features j in specific category i , the application needs to compute the weight W_i of each feature content value C_i within that category. The equation given below is adopted from below.

$$W_i = \frac{C_i * \sum_{j=0}^n C_j}{n}$$

After computing the weight of each feature content, the highest weight is seen as having the lowest value. Only the food items which have the lowest calories value is shown in the list to the user. The list is sorted in value of the calories from lowest to highest

Objectives:

The aim of this project was to build a food recommendation system for disease. This included to pre-process nutrition-based food and nutrition-based disease dataset for diet retrieval. To optimize the vocabulary of food to match them in the disease. To train, evaluate and test a K means clustering model and Random Forest able to predict the Food that disease belongs, by considering its set of nutrition. Predict the probability food according to disease

This project study is to consider various important aspects of the user's lifestyle and make sure that these factors are incorporated while the system works on a solution to build and recommend a healthy and nutritious diet for the user. A good nutritious healthy diet and a moderate amount of physical activity can help in maintaining a healthy weight. But the benefits of good nutrition have a lot more to do than just managing the weight. Being fit is all about the 70/30 rule. Here's how it goes, for a person to stay healthy he/she must focus 70% on his dietary intake and 30% on his physical activity/exercise.

DRS is project to provide users with healthy diet and individual nutritional recommendation. Health dietary was based on user disease. Furthermore, DRS contains catalogues of typical foods from an Asian region. Several Calabrian foods have been inserted because of their nutraceutical properties widely reported in several quality studies. DRS includes disease-based precautions.

Methodology:

Research methodology is a process that includes a number of activities to be performed. These are then arranged in proper sequence for conducting research. It is a master plan specifying the method and procedures for collecting and analysing needed information. Descriptive Research is used in this study as the main aim is to describe characteristics of the phenomenon or a situation. As the issue is well understood, it focuses on the development of in-depth knowledge the facts will be used to analyse and evaluate the data.

How it will be done:

The Data:

- Collecting dataset: I have found a dataset from Kaggle also I am going to collect the data from WHO.
- I will have to manually collect data using web scrapping.

Data Pre-Processing:

- At this stage removing punctuation, stop-words, special characters, making the table heading lowercase etc. will be done.
- Convert Integer values into Float.
- Null values in dataset can be removed or filled.

Splitting of the dataset:

- The dataset will be divided into training set and testing set.
- My train dataset size(population) is 6692 so my test (sample) data size will be 364 with margin of error is 5%.

Topic Modelling:

- Performance Evaluation: Here the models used will be evaluated using evaluation techniques like f1- score, accuracy, etc.
- Complaints:
 - I am use normal diet as required to normal Human.
 - If disease become detected serious then that become excluded.
 - False-Negative result will be included.

Implementation:

1. User's will enter the necessary information like their disease name, weight etc. on the form.
2. The information will then go through the ML model in following manner:
 - i. **K-Means** is used for clustering to cluster the food according to nutrition's.
 - ii. **Random Forest Classifier** is used to classify the food items and predict the food items based on input
3. The System will then recommend diet to the users based on input
4. The Users can choose from multiple recommended items and make their diet plan.

Classification Model:

We use the training dataset to get better boundary conditions that could be used to determine each target class. Once the boundary conditions are determined, the next task is to predict the target class. The whole process is known as classification. Classification will be based on various factors that are quality, quantity, search process and payment etc. To find out answers of above questions, here I am using some algorithms base on that I can easily classify my data.

Random Forests and k-Means falls under different categories of algorithms. k-Means is an unsupervised clustering algorithm wherein you group/cluster records. And Random Forests is a supervised learning algorithm used generally for classification and regression problems.

K Means Clustering:

K-Means clustering is an **unsupervised learning algorithm**. There is no labelled data for this clustering, unlike in supervised learning. we often use classification or regression algorithms in supervised learning methods to predict categories or values, we still often encounter situations where we need to use unsupervised learning methods to obtain a set of data categories. When the amount of data is large, you can consider using clustering algorithms to get different data categories. Clustering is subordinate to unsupervised learning, which does not rely on the defined classes and training examples of class

labels. Among them, K-means clustering is a very classic clustering method [8]

Given a set of elements, where each element has observable attributes, use a certain algorithm to divide into subsets, and require the degree of difference between the elements within each subset as much as possible low, and the element dissimilarity of different subsets is as high as possible. Concentration, each subset is called a cluster. Different from classification, classification is exemplary learning, which requires that each category be clarified before classification and that each element is mapped to a category, while clustering is observational learning, and the category may not be known or even the number of categories may not be known before clustering.

K-means tries to find the natural category of the data. The user sets the number of categories to find a good category centre. The algorithm flow is as follows:

- (1) Enter the number of data sets and categories K
- (2) Randomly assign the centre point of the category
- (3) Put each point into the set of the category centre point closest to it
- (4) Move the category centre point to the set where it is
- (5) Go to step 3 until convergence

After a number of cycles, the best classification effect can be obtained. Different from marine shale reservoirs, the relationship between food content of coal reservoirs and nutrition's of the coal reservoirs is relatively poor, and the laws are inconsistent, which also leads to the unreliability of the final prediction model. This is because coal reservoirs are more complex than shale reservoirs and have worse continuity, which causes the logging of coal seams to be affected by multiple factors. Using the clustering method to obtain multiple categories and establishing corresponding prediction models based on different categories can greatly improve the prediction results.

Random Forest Algorithm:

Random forest is a **supervised learning algorithm**. The "forest" it builds, is an ensemble of decision trees, usually trained with the

“bagging” method. The general idea of the bagging method is that a combination of learning models increases the overall result. Random forest is a highly flexible machine learning algorithm that has just emerged in the 21st century. It refers to a classifier that contains multiple decision trees. The thinking behind it is similar to group wisdom. In the 1980s, Breiman et al. invented an algorithm for classification trees, which performed classification or regression through repeated dichotomy of data, which greatly reduced the amount of calculation. In 2001, Breiman combined the classification trees into a random forest, that is, randomized the use of variables and the use of data, generated many classification trees, and then summarized the results of the classification trees [8].

Random forest improves the prediction accuracy without a significant increase in the amount of calculation. Random forest is not sensitive to multivariate collinearity, and the results are relatively robust to missing data and unbalanced data and can well predict the effect of thousands of explanatory variables.

Random forest uses a random method to build a forest. There are many decision trees in the forest, and there is no correlation between each decision tree in the random forest. After obtaining the forest, when a new input sample enters, let each decision tree in the forest make a judgment separately to see which category the sample belongs to. The class with the most classification times is the predicted class. Random forest can handle quantities whose attributes are discrete values. The construction process of random forest is as follows:

- (1) If there are **N** samples, samples are randomly selected for replacement (one sample is randomly selected each time and then returned to continue selection). Use the selected **N** samples to train a decision tree as the sample at the root node of the decision tree
- (2) When each sample has **M** attributes, when each node of the decision tree needs to be split, then **m** attributes are selected from these **M** attributes, and the condition $m < M$ is satisfied. Then, from these **m** attributes, strategies such as information gain are used to select one attribute as the split attribute of the node
- (3) In the process of decision tree formation, each node must be split according to step 2 until it can no longer be split. Note that there is no pruning during the entire decision tree formation process.

(4) Follow steps 1-3 to build a large number of decision trees to form a random forest

In the process of building each decision tree, attention should be paid to the impact of sampling and complete splitting. The first is two random sampling processes. Random forest samples the input data in rows and columns. For line sampling, a replacement method is used, that is, in the sample set obtained by sampling, there may be duplicate samples.

Assuming that there are N input samples, there are also N samples sampled. In this way, when training, the input samples of each tree are not all samples, making it relatively difficult to overfitting. Then, perform column sampling, from M features, select $m(m < M)$.

After that, a decision tree is built using a completely split method for the sampled data, so that a certain leaf node of the decision tree cannot continue to split, or all the samples in it point to the same category.

Generally, many decision tree algorithms have an important step-pruning, but this is not done here. Since the previous two random sampling processes ensure randomness, even if pruning is not performed, overfitting will not occur. Using a random forest method to predict food content should be able to achieve better results.

Combination Method of K-Means Clustering and Random Forest:

It is difficult to evaluate the food content of nutrition's, because the result has been affected by various factors, resulting in a poor relationship between food content and nutrition's. Only by using clustering and other methods to truly combine nutrition's for classification, different types of data are affected differently, and the relationship between nutrition's and food content in different categories is closer.

Therefore, K-means clustering is performed first, and then based on the results of the clustering, a random forest model of different types is established for final application. In fact, the inherent meaning of this model is similar to that of random forests. It uses K-means clustering combined with random forests to form a "forest group" to predict food content more accurately. The modelling and forecasting process is as follows:

- i. Use K-means clustering to divide the data into several categories.
The measurement method usually used to compare the results of

different **K** values is the average distance between a data point and its cluster centroid. Since increasing the number of clusters will always reduce the distance to the data point, when is the same as the number of data points, increasing **K** will always reduce the metric to zero. Therefore, this indicator cannot be used as the sole target. Conversely, the average distance to the centre of mass is plotted as a function of **K**, and the “elbow point” at which the reduction rate changes sharply can be used to roughly determine the **K** value.

- ii. Use **K** sets of data and random forest algorithm to train **K** models. After determining the category of the new data, the corresponding model can be used to calculate the food content.
- iii. When predicting new data, first determine the category of the new data by calculating the Euclidean distance between the sample data and the centroids of multiple classes of data. The new data belongs to the category corresponding to the centroid with the smallest Euclidean distance. After the category is determined, the corresponding model is used for prediction, and the predicted value of the food content of the sample point is obtained, and the reliability of the algorithm is determined by comparing with the real value

Requirements:

The experiment would be implemented completely in python language specification of it may be as follows:

- Python 3.5 and above
- Python IDE and Visual Code will be used as package manager for python environment.
- Google Collab will be used as online python environment.

Conclusion:

This project satisfying need will help to put patients in control of their own health data and therefore increase patients' autonomy. An approach of integrating recommender systems into personal Diet recommender system (DRS) was outlined. we can suggest natural diet (that have no side effects), precautions to user. The proposed system builds a user's health profile and, accordingly, provides individualized nutritional recommendations, also with attention to food geographical origin. The importance of nutritional guidance is increasing day by day to lead a healthy and fit life and by accepting the user's preferences and a user's profile in the system a healthy diet plan is generated.

Reference:

1. <https://ieeexplore.ieee.org/document/8765311>
2. <https://www.hindawi.com/journals/geofluids/2021/9321565/>
3. <https://www.sciencedirect.com/science/article/abs/pii/S0169260716306927#:~:text=DIETOS%20is%20a%20novel%20food,attention%20to%20food%20geographical%20origin>
4. <https://www.sciencedirect.com/science/article/abs/pii/S0169260716306927#:~:text=DIETOS%20is%20a%20novel%20food,attention%20to%20food%20geographical%20origin>.
5. Ge, M., Elahi, M., Fernaández-Tobías, I., Ricci, F., & Massimo, D., "Using tags and latent factors in a food recommender system," in Proc. of the 5th International Conference on Digital Health, pp. 105-112, ACM., May 2015.
6. Freyne, J., & Berkovsky, S., "Evaluating recommender systems for supportive technologies," User Modeling and Adaptation for Daily Routines, pp. 195-217, Springer London, 2013.
7. Recommender systems in the health care domain: state-of-the-art and research issues by ThiNgocTrangTran¹ · Alexander Felfernig¹ · Christoph Trattner² · Andreas Holzinger³
8. Atas, M., Tran, T.N.T., Felfernig, A., Polat-Erdeniz, S., Samer, R., Stettinger, M. (2019). Towards similarity aware constraint-based recommendation. In Advances and trends in artificial intelligence, lecture notes in computer science, Springer.
9. Bankhele, S., Mhaske, A., Bhat, S. (2017). V., s.: a diabetic health care recommendation system. International Journal of Computer Applications,