

Case-based Reasoning for Skin Diseases Diagnose using Minkowski Distance

Mihuandayani¹, Yufika Sari Bagi¹ and Theofani Christi Irene Momongan¹

¹Department of Information System, STMIK Multicom Bolaang Mongondow, Kotamobagu, Indonesia

Keywords: Case-based Reasoning, Skin Disease, Minkowski Distance.

Abstract: The skin is one of the crucial organs in the human body because it functions to receive stimuli such as touch, pain and other influences from the outside. According to the data, skin disease is the third out of the ten most diseases in Indonesia, and the skin specialists carry out the treatment of patients. However, due to the limitations of the experts resulting in the slow handling of patients so that needed a tool that can help diagnose patients with skin diseases. Case-Based Reasoning is one of the problem-solving techniques to build a system by making decisions from new cases based on the solutions of the old cases that are closest to the new case. The process of diagnosis is to enter new problems, compared to the old case, and then calculate the value of proximity using Minkowski distance. This study produced a case-based reasoning system to diagnose skin diseases due to the virus and bacterial infections based on selecting the symptoms suffered by patients and providing treatment solutions. The testing is done by comparing the expert diagnosis data and system diagnostic results with 92.10 % accuracy.

1 INTRODUCTION

The skin is a vital organ in humans located in the outer layer of the body, which functions to receive stimuli from the outside. Human skin condition based on data (Karimkhani et al., 2017) in 2013, contributed around 1.79% of the world disease burden. Damage to the skin barrier is one of the common problems related to nursing. Dermatological diseases can affect significantly in physical, psychological and social (Hay et al., 2014). Unhealthy skin can cause various skin diseases, so it needs to maintain skin health earlier to avoid diseases. A person's body skin that is affected by the disease can interfere with the appearance and activities of the person. Skin disease is often underestimated because it tends to be harmless and not cause death. The 2010 Indonesian Health Profile data shows that skin disease is ranked as the third out of the ten most diseases in outpatients in Indonesian hospitals. The incidence of skin diseases in Indonesia is still relatively high and becomes a significant problem. Various skin diseases can be caused by several factors, such as the environment and bad daily habits, climate change, viruses, bacteria, fungi, allergies, endurance, and others.

Skin diseases due to virus infections are classified into several types, including Verruca Vulgaris, Herpes Zoster, Herpes Genital, and Varicella. Verruca Vulgaris clinically in the form of solid papules/plaque and the surface is in the form of small papules measuring 1-3 mm, and the most frequent of Human Papilloma Virus (HPV) infection. Herpes Zoster is a disease caused by the Varicella-Zoster virus, mainly affecting adults with the characteristics of radicular, unilateral pain and hordes of vesicles scattered according to the dermatome that is conserved by a sensory nerve ganglion. Herpes Genital becomes an infection of the skin disease caused by the Herpes Simplex Virus (HSV) transmitted through sexual contact. Varicella is a skin disease with scattered vesicles, primarily affecting children, which is easily transmitted by the Varicella-Zoster virus.

Then, for skin diseases due to bacterial infections classified into several types including Acne vulgaris is a chronic inflammation of the pilosebaceous follicles which is characterized by blackheads, papules, cysts and pustules in predilection areas. Furuncle is an acute infection of one hair follicle, which usually experiences necrosis caused by Staphylococcus aureus. Leprosy is a disease caused by Mycobacterium leprae and belongs to a chronic

disease that attacks the peripheral nerves, skin, and other body parts (Job, 1994). Impetigo is a superficial and infectious pyogenic skin disease caused by *Staphylococcus* or *Streptococcus*.

To diagnose patients suffering from skin diseases can be known from the symptoms that appear or experience by the patients. Handling in patients with skin diseases is carried out by experts, namely skin specialists, but due to limited expert expertise resulting in slow handling of patients, so we need a tool that can help diagnose patients with skin diseases. Cases stored in medical records regarding the diagnosis of skin diseases by experts to determine the type of skin disease suffered by patients can be reused as a reference to determine the type of skin disease when there are new cases. Utilization of cases that have occurred before or old cases is generally known by the term Case-Based Reasoning (CBR).

CBR is a psychological theory of human cognition that overcomes problems in memory, learning, planning, and problem-solving (Slade, 1991). CBR is a computer reasoning method that utilizes old knowledge to solve new problems. Ancient knowledge in the form of documentation of problems that already have solutions. The solution can be used to solve similar new problems. The CBR method was chosen because of its advantages, namely the problem-solving process using existing cases and having a revision process used to correct diagnostic errors or inaccurate solutions, the results of the revision can be stored as a new knowledge base on the system, so that the system can continue to develop. Then, distance metrics are widely used in the estimation of similarity. In this study, the method of calculating distances between cases uses Minkowski distance. Minkowski distance is a metric in vector space where the $n-1$ distance is also called the distance of a city block and is often considered a generalization of two distance, the Euclidean distance and the distance of Manhattan (John, 1995). Besides, in the study of Distance (Sreedevi and Padmavathamma, 2015) and also (Karakoc et al., 2006) explained that the Minkowski distance method is relatively better than other distance methods.

The system that was built only covered a few types of skin diseases, namely skin diseases due to viruses and bacteria in adults, then the process of revising the results of the diagnosis manually by experts. The testing process was carried out to measure the performance of Minkowski Distance with data based on medical records for five years of work, obtained from the Clinic dr. Giana Sugeha, Sp.KK. Based on the background problem, this study produced a system to solve diagnosis in cases of skin

diseases using CBR and the Minkowski distance calculation method. This research can help health workers to diagnose patients with the virus and bacterial skin diseases, which can be used as a reference before providing treatment therapy recommendations.

2 RELATED WORKS

Research that applies CBR to cases of skin disease has been done before. While in this study CBR to diagnose skin diseases due to viruses and bacteria with several types of diseases such as *Verruca Vulgaris*, Herpes, Varicella, Furuncle, Leprosy, Acne Vulgaris, Impetigo using the Minkowski distance calculation method. By entering the symptoms suffered by the patient will produce the final result that is the illness.

Other studies related to the computational efficiency of the k-means algorithm with distance metrics find similar data objects that lead to the development of robust algorithms for data mining functionality. The research experiment applied the K-means algorithm with three metrics, namely Euclidean distance, Manhattan distance and Minkowski distance. The results obtained that the selection of distance metrics has a vital role in clustering (Singh et al., 2013).

Other studies discuss the development of Minkowski distance generalizations using IOWA, and it is referred to as the Minkowski OWA operator (IMOWAD) or induced long-distance OWA operator (IGOWAD) by obtaining various distances. Development of new application approaches in the problem of decision making about investment selection (José and Montserrat, 2011).

Research that discusses early detection of diabetes that helps prevention measures using genetic algorithms with the Multimodal Evolution algorithm. The research that discusses the diagnosis of diabetes is carried out using Genetic Algorithm and distance calculations, including the distance of Minkowski on the PIMA Indian dataset. Through the calculation of the distance can be proven better accuracy using Minkowski Distance (Sreedevi and Padmavathamma, 2015).

The research discussion is related to image processing of various types of skin cancer by proposing a method for detecting two types of skin, one is cancerous skin, and the other is affected but not cancerous skin based on several feature values. Some values extracted from the Gray Level Co-occurrence Matrix (GLCM) also include in Euclidean distance,

Manhattan Distance, Minkowski Distance, and Hamming Distance. This process is a secure system rather than a doctor's biopsy procedure. The system consumes less time and gets better results than conventional systems. This study shows that the combination of co-occurrence of matrix and neural network provides a technique for detecting cancer cells and non-cancer cells (Bhuiyan et al., 2015).

Research analyzing and comparing images based on twelve distance calculations including the distance of Minkowski, Euclidean, Mahalanobis, Manhattan, Chebychev, and others for diagnosis of 176 images of dermoscopic skin lesions. The result shows that CBDIR performance can be significantly improved by using Canberra and Bray-Curtis Distance compared to conventional measures (Khadidja and Mostefai, 2015).

Research that develops a theoretical basis for the analysis and construction of shape sizes of time series association. Some general methods of construction of such sizes are suitable for measuring the time of the equation of the series form and proposed form association. Size associations of time series forms based on Minkowski distance and data standardization methods were considered. The cosine similarity and Pearson correlation coefficient were obtained as a part of the case from the proposed general method which can also be used for the construction of new association steps in data analysis (Batyrrshin, 2013).

Research that discusses the search for structural similarities between small molecules focuses on the method of the closest K-neighbor. The research shows an optimal computation with weighted Minkowski distance to maximize discrimination between active and inactive compounds, then shows the structure of KNN-based pruning data for the distance of wL_p which minimizes the time for finding similarities. The result of the experiment shows that the classification of KNN with Minkowski wL_1 distance gets better accuracy than LDA and MLR (Karakoc et al., 2006).

The discussion presents the Minkowski type distance, including the Hamming, Euclidean, and Chebyshev distances, for the fuzzy orthopair r -rung set. It introduced Minkowski-type distance for orthopair values, based on orthopairs rankings. Then propose some distance in the fuzzy orthopair set of r -rung and discuss it for multi-distribution decision-making problems. Each element was expressed as an ordered value pair. The first showed support for membership and the last supports membership. The study presents a method based on Minkowski

distance for a discussion of its application in multi-attribute decision making (Du, 2018).

In other studies, the weighted geographical regression model GWR was adapted to benefit from a variety of distance metrics, where it was shown that a well-chosen distance metric could improve the performance of the model. Minkowski's approach is proposed, which allows the selection of optimal distance metrics for a given GWR model. This approach was evaluated in a simulation experiment consisting of three scenarios. This result suggests that the Minkowski approach can be useful when there is no knowledge, understanding or insight into the 'true' distance metric. There is a significant drawback to the Minkowski approach, where it is often difficult to describe how certain Minkowski distance functions can be measured, except for some general cases, such as Manhattan ($p = 1$) or Euclidean ($p = 2$). This deficiency tends to make the Minkowski approach more suitable for predictive purposes with GWR than for exploration or inferential purposes with GWR (Binbin, 2015).

Research that develops new decision-making models using induction ordered weighted average operators and Minkowski distance fuzzy linguistic variables. Several examples of the new method obtain the results. Generalization of an IOWA operator that uses order trigger variables is carried out to assess the complex rearrangement process, blurred linguistic information and fuzzy linguistic Minkowski distance (Xian et al., 2014).

Research that discusses how to measure perceptual similarities between two objects using Minkowski type metrics. In the Minkowski metric there is no substantial similarity to the same object, through mining a broad set of visual data, the study has found the perception of the function of distance by calling a function that is found to be a partial dynamic function. When compared empirically with dynamic partial functions, Minkowski type distance functions in shooting and image capture transition detection, DPF performs significantly better (Li et al., 2003). Research that presents cellular-based medical assistance aimed at diagnosing skin diseases using CBR and image processing to increase awareness of disease prevention. CBR is used to determine the symptoms of skin diseases based on image data as setting a new knowledge base (Aruta et al., 2015).

In a study comparing 14 distance measurements, including Minkowski Distance and modification between feature vectors with the performance of the primary component, PCA method proposed a modification of sum square error (SSE) -based distance. The experiments showed that the proposed

distance measure is the first three best steps for the characteristics of different biometric systems. Besides, it shows that using an algorithmic combination of distance measurements gets better results than using distances separately (Perlibakas, 2004). There are various uses of distance calculation and its advantages in previous studies. This study proposes calculations using Minkowski Distance in the case of skin disease caused by viruses and bacteria to be able to diagnose the skin disease and produce a system that can help patients who have symptoms of skin disease.

3 RESEARCH METHODOLOGY

This study has several stages or methods used, ranging from data collection to calculation of results for concluding. In the initial stages, observations are made on the condition and procedures of a clinic. At present, the process that occurs is that the patient registers at the clinic then the nurse records the patient’s medical record information from the examination of the patient’s condition related to symptoms, and blood pressure. The patient waits according to the queue number. Next, the nurse is taken to the doctor’s office. The doctor will make a diagnosis, and further examination then gives a prescription by the illness. Furthermore, the patient will complete the stages to get the medicine according to the doctor’s prescription.

The next stage is the data collection is done by conducting interviews directly to those who have the capacity and information needed in this research, namely dr. Giana Sugeha, Sp.KK, as an expert in skin and genital specialist. Furthermore, literature studies are carried out on textbooks, research journals, and other references. Documentation was carried out for taking medical records of patients on sample cases related to skin diseases caused by viruses and bacteria for five years of work obtained from the Clinic dr. Giana Sugeha, Sp.KK.

Data analysis was obtained to determine data processing needs, especially in the application of CBR, which is a problem-solving method that uses knowledge of previous experience to solve new problems (Pal and Shiu, 2004). A problem is solved by looking for the same old case, if found, then the solution of the two is also identical. However, if it is not found, the system looks for old cases that have the highest similarity and needs to be adapted so that the problem finds a solution. This solution is done based on similar situations, referred to as cases, which have previously been stored in the system.

The stages of the CBR process are new problems that are matched with cases that are in the case store database and find one or more similar cases (retrieve). Solutions that are recommended through case matching are then reused for similar cases, the solutions offered may be changed and adopted (revised). If new cases do not match in the case storage database, CBR will store the new cases (retain) in the knowledge base. The Retrieval process in CBR is based on the hypothesis that the solution to previous problems can help resolve current problems, as long as there are similarities between them. The following is an overview of the process flow of implementing CBR in cases of diagnosing skin diseases shown in Figure 1.

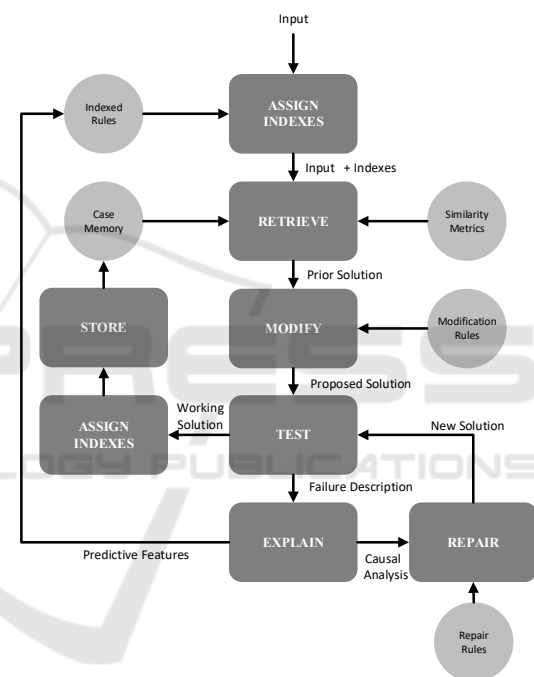


Figure 1: CBR Processing Pipeline.

The similarity of the two points can be calculated by distance. The more similar the two points, the smaller the distance and vice versa. The more different two points, the higher the distance. The Minkowski distance calculation is used to calculate the similarity. Minkowski distance is a metric in Euclidean space which is a generalization of Euclidean’s distance and Manhattan’s distance. The following is the Minkowski distance formula.

$$Dist_{XY} = \left(\sum_{k=1}^d |X_{ik} - X_{jk}|^{\frac{1}{p}} \right)^p \tag{1}$$

In the formula, it is explained that when $p = 2$, the distance is called the Euclidean distance. Then, when $p = 1$, it becomes the distance of the city block. Furthermore, the Chebyshev distance is one of the Minkowski distance variants where $p = \infty$ at the threshold.

Furthermore, to design a system that is capable of diagnosing, modelling is needed in the form of use case diagrams. System analysis is intended to determine the functional and non-functional requirements of the newly designed system. Functional requirements related to functions must be provided by the system to meet the primary needs and supporting needs in running the system. Analysis of the functional requirements of the system is modelled using a use case diagram can be seen in Figure 2.

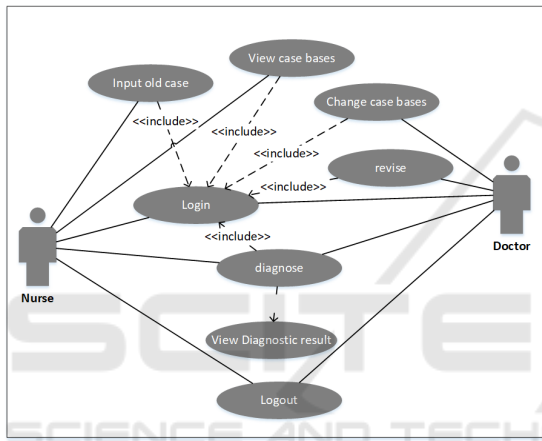


Figure 2: Use Case Diagram of System.

Based on the use case diagram, the nurse can input the old case, make a diagnosis and be able to see the diagnosis, see the basis of the case and change the basis of the case. Whereas the doctor can diagnose and can see the results of the diagnosis, can see the basis of the case and make revisions in cases where the value of the symptoms does not match the threshold. In processing the system, the user enters symptoms to diagnose the disease, then the system takes the value of the selected symptoms and compares the new case with the old case.

After that, the system calculates the proximity between cases, namely the value and takes the case with the lowest distance value then adjusts to the threshold. If it matches the threshold, the system will display the diagnostic results, namely the case code, disease, and solution. If it does not match the threshold, the system will save the case for revision. Next, to find out the data processed in the design of this study, the following is shown the relation table in Figure. 3.

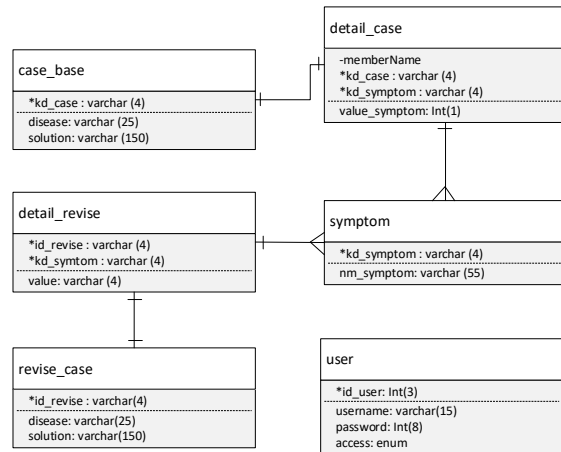


Figure 3: Data Table Relations.

The relation table above describes the relationship between data processing related to the system. Then, a system that can diagnose skin diseases and manage the similarity calculation with Minkowski Distance is done so that the data can be evaluated in the system test and compared with data from the clinic of dr. Giana Sugeha, Sp.KK.

4 RESULT AND DISCUSSION

4.1 Algorithm Calculation

Data samples used in the system of diagnosing skin diseases due to the virus and bacterial infections were as many as 103 case data taken from the patient’s medical record in the skin and genital specialist clinic of dr. Giana Sugeha, Sp.KK. The data were taken in the form of medical records of patients who have skin diseases due to the virus and bacterial infections. The data is then collected into case data with various types of diseases that have been diagnosed by doctors and solutions or therapies given. The data is divided into 65 training data and 38 testing data. Samples of the old case data are shown in Table 1.

4.2 Evaluation

The result of system testing aims to determine the accuracy of the system in diagnosing the diseases, also to test the system work process whether it is following the design that has been made using a threshold with a predetermined value. If the threshold value is set at 0.80, and the diagnosis is obtained more than the threshold value of 0.80, it is necessary to revise it. Whereas if the results obtained are the same or less than the threshold value, then there is no need for a revision and it can be determined that the patient has A disease. Two tests have been compared to the results of manual calculations with expert calculation results. This test is conducted to find out whether the manual calculation has a right level of accuracy by the results of expert calculations or not. At this phase, system testing will be carried out to check the accuracy of the system produced whether the system can be run according to specific standards. The accuracy calculations used the following formula.

$$Accuracy = \frac{Correct\ Test\ Data}{Amount\ of\ Test\ Data} \times 100\ \% \quad (2)$$

The following test results on 38 test data using predetermined threshold values are shown in Table 3.

Table 3: Testing Using Threshold Values.

Threshold	Correct Test Data	Amount of Test Data	Accuracy
Threshold 1	32	38	84.21 %
Threshold 1.5	35	38	92.10 %
Threshold 2	35	38	92.10 %

The results of the accuracy calculation are represented in the diagram, which is shown in Figure 4.

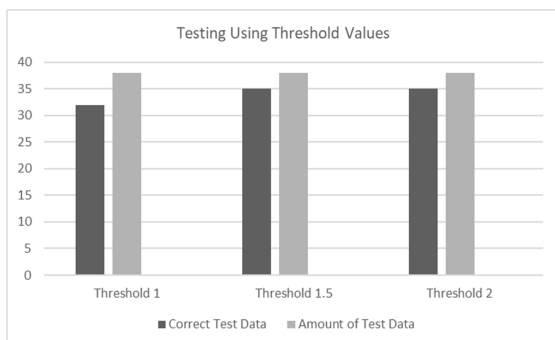


Figure 4: Testing Using Threshold 1, Threshold 1.5, and Threshold 2.

Based on Table 3 and Figure. 4, it was obtained that the calculation with the highest accuracy value for diagnoses of skin diseases due to viruses and infections was obtained from the threshold 1.5 and threshold 2 with 92.10% accuracy.

5 CONCLUSIONS

Based on the research and results of system testing, it can be concluded that this study produced a CBR system for the diagnosis of skin and genital diseases due to the virus and bacterial infections by calculating the proximity between new problems and old cases based on symptoms by accommodating case feature values and confidence levels. The system can diagnose diseases based on symptoms and display the results of patient diagnosis to provide treatment solutions. The system provides disease diagnosis based on the similarity between old cases and new cases. The diagnosis is considered correct if the distance value is ≥ 1.5 . The test results on skin and genital disease test data due to the virus and bacterial infections indicate that the system can recognize skin and genital diseases using the Minkowski distance method correctly with an accuracy rate of 92.10%.

ACKNOWLEDGEMENTS

STMIK Multicom Bolaang Mongondow for the support in providing the facility which is needed along with the study. Especially for dr. Giana Sugeha, Sp.KK. She has given the clinic data and documentation for the research. This good cooperation among all sector to provide a paper that hopes can be useful for others.

REFERENCES

- Aruta CL, Calaguas CR, Gameng JK, Prudention MV, and Lubaton AACJ. Mobile-based Medical Assistance for Diagnosing Different Types of Skin Diseases Using Case-based Reasoning with Image Processing. International Journal of Conceptions on Computing and Information Technology Vol. 3, Issue. 3, October' 2015; ISSN: 2345 - 9808.
- Batyrshin, I. Constructing Time Series Shape Association Measures: Minkowski Distance and Data Standardization. 2013 BRICS Congress on Computational Intelligence & 11th Brazilian Congress on Computational Intelligence.

- Bhuiyan MA, Miah MBA, and Mia MR. Detection of Cancerous and Non-cancerous Skin by Using GLCM Matrix and Neural Network Classifier. *International Journal of Computer Applications* (0975 - 8887). Volume 132 - No.8, December 2015.
- Binbin Lu, Martin Charlton, Chris Brunson & Paul Harris. (2015). The Minkowski approach for choosing the distance metric in geographically weighted regression. *International Journal of Geographical Information Science*, DOI: 10.1080/13658816.2015.1087001.
- Du, WS. Minkowski-type distance measures for generalized orthopair fuzzy sets. *Int J Intell Syst.* 2018;1-16. <https://doi.org/10.1002/int.21968>.
- Hay RJ, Johns NE, Williams HC, Bolliger IW, Dellavalle RP, Margolis DJ, Marks R, Naldi L, Weinstock MA, Wulf SK, Michaud C, Murray CJL, & Naghavi M. (2014). The global burden of skin disease in 2010: An analysis of the prevalence and impact of skin conditions. *Journal of Investigative Dermatology*, 134, 1527-1534.
- Job CK. (1994). Pathology of leprosy. In: Hasting RC (ed). *Leprosy*, 2nd ed., Edinburgh, Churchill Livingstone, p 193-224.
- John P. Van de Geer. Some Aspects of Minkowski Distances. Leiden University, Department of Data Theory, 1995.
- José M. Merigó & Montserrat Casanovas. (2011). A New Minkowski Distance Based on Induced Aggregation Operators. *International Journal of Computational Intelligence Systems*, 4:2, 123-133, DOI: 10.1080/18756891.2011.9727769.
- Karakoc E, Cherkasov A, and Sahinalp SC. Distance-based algorithms for small biomolecule classification and structural similarity search. *Bioinformatics*. Vol. 22 No. 14 2006, pages e243 - e251.
- Karimkhani C, Dellavalle RP, Coffeng LE, Flohr C, Hay RJ, Langan SM, Nsoesie EO, Ferrari AJ, Erskine HE, Silverberg JI, Vos T, & Naghavi M. (2017). Global Skin Disease Morbidity and Mortality: An Update From the Global Burden of Disease Study 2013. *JAMA Dermatol*, 153, 406-412.
- Khadija B, & Mostefai Sihem. (2015). Similarity measures for Content-Based Dermoscopic Image Retrieval: A comparative study. 1-6. 10.1109/NTIC.2015.7368761.
- Li B, Chang E, and Wu Y. Discovery of a perceptual distance function for measuring image similarity. *Multimedia Systems* 8: 512 - 522 (2003). DOI: 10.1007/s00530-002-0069-9.
- Pal SK, and Shiu SCK. (2004). *Foundations of Soft Case-Based Reasoning*. John Wiley & Sons Inc, New Jersey.
- Perlibakas, V. Distance Measures for PCA-based face recognition. *Pattern Recognition Letters* 25 (2004) 711 - 724.
- Singh A, Yadav A, and Rana A. K-means with Three different Distance Metrics. *International Journal of Computer Applications* (0975 - 8887). Vol. 67 - Number 10, April 2013.
- Slade, Stephen. (1991). Case-Based Reasoning: A Research Paradigm. *AI Magazine*. 12. 42-55. 10.1609/aimag.v12i1.883.
- Sreedevi E, Padmavathamma MA. Threshold genetic algorithm for diagnosis of diabetes using Minkowski distance method. *International Journal of Innovative Research in Science, Engineering and Technology*, Vol. 4, Issue 7, July 2015.
- Xian S, Sun W, Xu S, and Gao Y. (2014). Fuzzy linguistic induced OWA Minkowski distance operator and its application in group decision making. *Pattern Anal Applic*, DOI: 10.1007/s10044-014-0397-3.