# Final Report: KNN And Decision Tree

CSE-0408 Summer 2021

Rubyat Jesmin Shammi

*Department of Computer Science and Engineering*
*State University of Bangladesh (SUB)*
Dhaka, Bangladesh
rubyatshammiss@gmail.com

*Abstract*——**This paper introduced for KNN problem and Decision Tree problem.**
*Index Terms*—

## I. INTRODUCTION FOR KNN

KNN is a lazy learning, non-parametric algorithm. It uses data with several classes to predict the classification of the new sample point. KNN is non-parametric since it doesn't make any assumptions on the data being studied, the model is distributed from the data.

## II. ADVANTAGES FOR KNN

1. Simple to implement and intuitive to understand
2. Can learn non-linear decision boundaries when used for classfication and regression. Can came up with a highly flexible decision boundary adjusting the value of K.

## III. DISADVANTAGES FOR KNN

1. Performance Issue with large data-set: The time required to calculate the distance between the new point and each existing points is huge. Which then degrades the performance of the algorithm.
2. Value of K: It is really crucial to determine what value to assign to k. with different value of K you get different results

## IV. LITERATURE REVIEW FOR KNN

Author's introduction: Zhongheng Zhang, MMed. Department of Critical Care Medicine, Jinhua Municipal Central Hospital, Jinhua Hospital of Zhejiang University. Dr. Zhongheng Zhang is a fellow physician of the Jinhua Municipal Central Hospital. He graduated from School of Medicine, Zhejiang University in 2009, receiving Master Degree. He has published more than 35 academic papers (science citation indexed) that have been cited for over 200 times. He has been appointed as reviewer for 10 journals, including Journal of Cardiovascular Medicine, Hemodialysis International, Journal of Translational Medicine, Critical Care, International Journal of Clinical Practice, Journal of Critical Care.

## V. CONCLUSION FOR KNN

KNN is an effective machine learning algorithm that can be used in credit scoring, prediction of cancer cells, image recognition, and many other applications. The main importance of using KNN is that it's easy to implement and works well with small data sets.

## VI. INTRODUCTION FOR DECISION TREE

In a decision tree, the algorithm starts with a root node of a tree then compares the value of different attributes and follows the next branch until it reaches the end leaf node. It uses different algorithms to check about the split and variable that allow the best homogeneous sets of population.

## VII. ADVANTAGES FOR DECISION TREE

1. When using Decision tree algorithm it is not necessary to normalize the data.
2. Decision tree algorithm implementation can be done without scaling the data as well.
3. When using Decision tree algorithm it is not necessary to impute the missing values.
4. The data pre-processing step for decision trees requires less code and analysis.
5. The data pre-processing step for decision trees requires less time.

## VIII. DISADVANTAGES FOR DECISION TREE

1. The mathematical calculation of decision tree mostly require more memory.
2. The mathematical calculation of decision tree mostly require more time.
3. The space and time complexity of decision tree model is relatively higher.
4. Decision tree model training time is relatively more as complexity is high.

## IX. LITERATURE REVIEW FOR DECISION TREE

Angel Insua, Alberto Monje, Hom-Lay Wang, Marita Inglehart, Patient-Centered Perspectives and Understanding of Peri-Implantitis, Journal of Periodontology, 10.1902/jop.2017.160796, 88, 11, (1153-1162), (2017).

Wiley Online Library Nima D. Sarmast,Howard H. Wang,Nikolaos K. Soldatos,Nikola Angelov,Samuel Dorn,Raymond Yukna,Vincent J. Iacono, First published: 01 December 2016
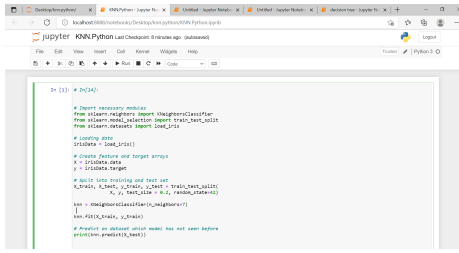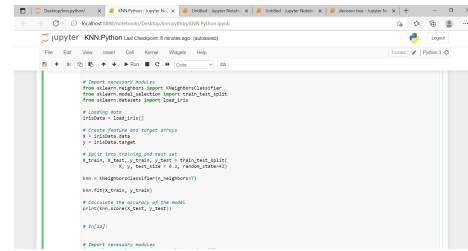
Fig. 1.



Fig. 2.

## X. CONCLUSION FOR DECISION TREE

Decision trees assist analysts in evaluating upcoming choices. The tree creates a visual representation of all possible outcomes, rewards and follow-up decisions in one document. Each subsequent decision resulting from the original choice is also depicted on the tree, so you can see the overall effect of any one decision. As you go through the tree and make choices, you will see a specific path from one node to another and the impact a decision made now could have down the road.



Fig. 3.

## ACKNOWLEDGMENT

Fig. 4.

## REFERENCES

[1] Zhou Q, Yang D, Ombrello AK, et al. Early-onset stroke and vasculopathy associated with mutations in ADA2. N Engl J Med. 2014;370:911–20.

[2] Navon Elkan P, Pierce SB, Segel R, et al. Mutant adenosine deaminase 2 in a polyarteritis nodosa vasculopathy. N Engl J Med. 2014;370:921–31.

[3] Oda H, Kastner DL. Genomics, biology, and human illness: advances in the monogenic autoinflammatory diseases. Rheum Dis Clin N Am. 2017;43:327–45.

[4] Gonzalez Santiago TM, Zavialov A, Saarela J, et al. Dermatologic features of ADA2 deficiency in cutaneous polyarteritis nodosa. JAMA Dermatol. 2015;151:1230–4.

[5] Short RD, Fukunaga K. The optimal distance measure for nearest neighbor classification. IEEE Transactions on Information Theory 1981;27:622-7. 10.1109/TIT.1981.1056403 [CrossRef] [Google Scholar]

[6] Weinberger KQ, Saul LK. Distance metric learning for large margin nearest neighbor classification. The Journal of Machine Learning Research 2009;10:207-44. [Google Scholar]

[7] Cost S, Salzberg S. A weighted nearest neighbor algorithm for learning with symbolic features. Machine Learning 1993;10:57-78. 10.1007/BF00993481 [CrossRef] [Google Scholar]
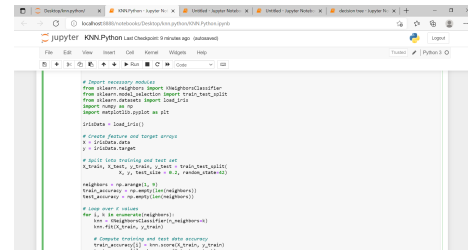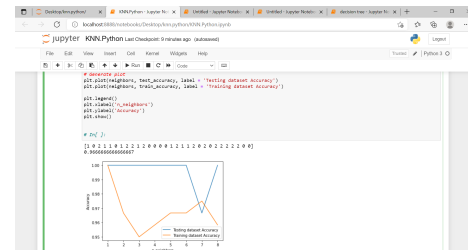
## XI. CODE



Fig. 5.



Fig. 6.

Fig. 7.