

1.

- On hepatitis dataset, using KNN model can generally predict more precise results with appropriate choice of value K. On the other hand, Decision Tree algorithm implemented with misclassification function can perform as approximately as KNN with appropriate K.
 - On messidor features dataset, KNN can predict more accurately than DT does.
2. As K increases to a certain threshold, KNN can obtain a best prediction on its test accuracy even it is not necessarily the K value that gives smallest gap between test accuracy and train accuracy. In addition, as K increments, train accuracy gradually decreases.
3. Maximum tree depth can barely affect the performance of Decision Tree when it is large enough. However, when we start from max depth = 1, there are some rapid change of test accuracy. It stabilizes after reaching a certain depth. In addition, the train accuracy will fix at 100% after max depth meets a threshold, which implies that the model is overfitting the training dataset.
4. i. For hepatitis dataset, both distance functions make roughly same predictions in KNN model. The entropy test of DT algorithm makes the best prediction among all three cost functions.
- ii. For messidor features dataset, manhattan function predicts better than euclidean function in KNN and cost entropy generates the most precise prediction among all cost functions
5. We found the pair of features with highest correlation and neither of them should be a binary feature (either yes or no). The plots of such pairs of features are key features decision boundary plot.