

Ques 1:

I have implemented the decision tree algorithm with both Entropy and gini index as criteria for splitting decision tree. Multi-way split, and multivariate split strategy has been used. Accuracy has been defined as the average of 10-fold accuracies.

- a) The classification accuracy for binary split with entropy as measure = **0.7927**
- b) The classification accuracy for binary split with gini index as measure = **0.7932**
- c) The classification accuracy for multiway and multivariate split = **0.8182**
- d) The classification accuracy by using pruning after splitting = **0.8283**

Gini index measure performs slightly better as it takes into account the probability of misclassification.

The accuracy of Decision trees in which numeric attributes are split several ways are more comprehensible than the usual binary trees because attributes rarely appear more than once in any path from root to leaf. Hence, it may prove to be a better approach than binary split. But multiway split is complicated and takes a lot of time for execution.

Pruning reduces the size of decision trees by removing sections of the tree that provide little power to classify instances. It reduces the complexity of the final classifier, and hence improves predictive accuracy by the reduction of overfitting.

Ques 2:

I used decision tree algorithm for prediction of cuisines given the ingredients.

Following are the Kaggle prediction scores:

- Prediction score with entropy as splitting criterion in decision tree: **0.56556**
- Prediction score with gini index as splitting criterion in decision tree: **0.57602**

Entropy is a way to measure impurity, whereas **Gini index** is a criterion to minimize the probability of misclassification. **Entropy** reaches maximum value when all classes in the node have equal probability. The **Gini index** is maximal if the classes are perfectly mixed, for example, in a binary class.

For this dataset, Gini index criterion performs better than the entropy measure for information gain as the prediction score for Gini Index is higher. This is because the cuisines can be misclassified based on a few changes in ingredients, and Gini index minimizes these misclassifications and pruning the tree accordingly.