Liver cancer, particularly hepatocellular carcinoma (HCC), presents significant diagnostic challenges due to its asymptomatic nature in its early stages and complex pathogenesis. This study introduces a novel approach to liver cancer detection, leveraging the diagnostic potential of microRNA (miRNA) expression profiles and advanced machine-learning techniques. Recent research in molecular oncology has spotlighted miRNAs as critical players in carcinogenesis. These small, non-coding RNA sequences are instrumental in gene regulation and have shown distinct expression patterns in various cancer types, including HCC. The dataset used in this study contains miRNA expression levels (denoted by 'MIMAT' codes) and the health statuses of subjects. The preprocessing of this dataset involves isolating relevant miRNA expressions and categorizing subjects into 'Healthy' or 'Cancer' groups based on their health status. The core of our methodological framework employs the Random-Forest-Classifier, a robust ensemble learning method suitable for handling high-dimensional and non-linear data typical in genomic studies. The dataset is partitioned into training and testing sets, ensuring an unbiased evaluation of the model's performance. Hyperparameter tuning, crucial for optimizing the RandomForestClassifier's performance, is conducted through a RandomizedSearchCV, exploring a range of values for parameters like 'n_estimators', 'max_features', and 'max_depth'. The best parameters obtained from this search are then used to train the final model. The developed random forest model serves two purposes. Firstly, it offers a 'predict_cancer' function, which allows for the prediction of liver cancer based on the expression levels of the microRNA. The second purpose is to demonstrate high accuracy in the process of distinguishing between healthy individuals and those at risk of liver cancer. This is quantified by the accuracy metric, which shows a 97% accuracy, a strong score. In addition to assessing the model's accuracy, we have implemented an evaluation of the logistic loss, or log loss, which provides a more nuanced insight into the performance of our RandomForestClassifier model. The log loss is particularly informative as it penalizes false classifications, offering a perspective on the model's confidence in its predictions. Lower log loss values indicate better model performance, with a perfect model having a log loss of 0. Upon training the model with the best parameters obtained from RandomizedSearchCV, we conduct predictions on the test set. These predictions are twofold: the binary classification (y_pred) and the probability estimates of each class (y_pred_proba). The latter is crucial for the calculation of log loss. The accuracy score, derived from y_pred, provides a straightforward measure of how often the model correctly predicts the health status.  However, the log loss, computed using y_pred_proba, delves deeper by considering the probability distributions of the output. It penalizes incorrect classifications more severely when the model is confident about a false prediction, thus providing a more stringent assessment of the model's performance. We also compute a confusion matrix, which offers a comprehensive view of the model's performance across different classes, highlighting true positives, false positives, true negatives, and false negatives. This effort leads to a log loss of .164, which is a remarkably good result. The graph of both the false positive rate to the true positive rate will be plotted below, as well as the confusion matrix and feature importance.