

# DIC Phase 2

Abdul Wasi Lone

50609995

**Introduction:** This study aims to identify factors that impact the **length of hospitalization** for patients and explore the relationship between **patient characteristics and ASA ratings**.

We employed two machine learning algorithms, **Decision Tree and Random Forest Classifiers**, to analyze a comprehensive medical dataset and derive insights.

**Dataset Overview:** Our analysis utilized a dataset containing various medical features, including:

- BIRTH\_DATE (patient age)
- GENDER
- ASA\_RATING\_C (ASA physical status classification)
- AN\_TYPE (type of anesthesia)
- ICU\_ADMIN\_FLAG (ICU admission status)
- AN\_LOS\_HOURS (length of stay in hours)

## **Methodology:**

**Data Preprocessing:** We implemented the following preprocessing steps:

Categorical variable encoding using OneHotEncoder

Numerical feature scaling using StandardScaler

Data splitting: 80% training, 20% testing

## **Model Selection and Justification :**

**Decision Tree Classifier:** We chose the Decision Tree Classifier for its:

- Interpretability, aligning well with medical decision-making processes
- Ability to handle both numerical and categorical data effectively
- Capacity to capture non-linear relationships between features

**Random Forest Classifier:** We selected the Random Forest Classifier to:

- Reduce overfitting and improve generalization as an ensemble method
- Provide robust feature importance rankings
- Effectively handle high-dimensional data Model Tuning

We used **GridSearchCV** for **hyperparameter** optimization:

**Decision Tree Parameters tuned:**

- max\_depth: [3, 5, 7, 9]
- min\_samples\_split: [2, 5, 10]
- min\_samples\_leaf: [1, 2, 5, 10]

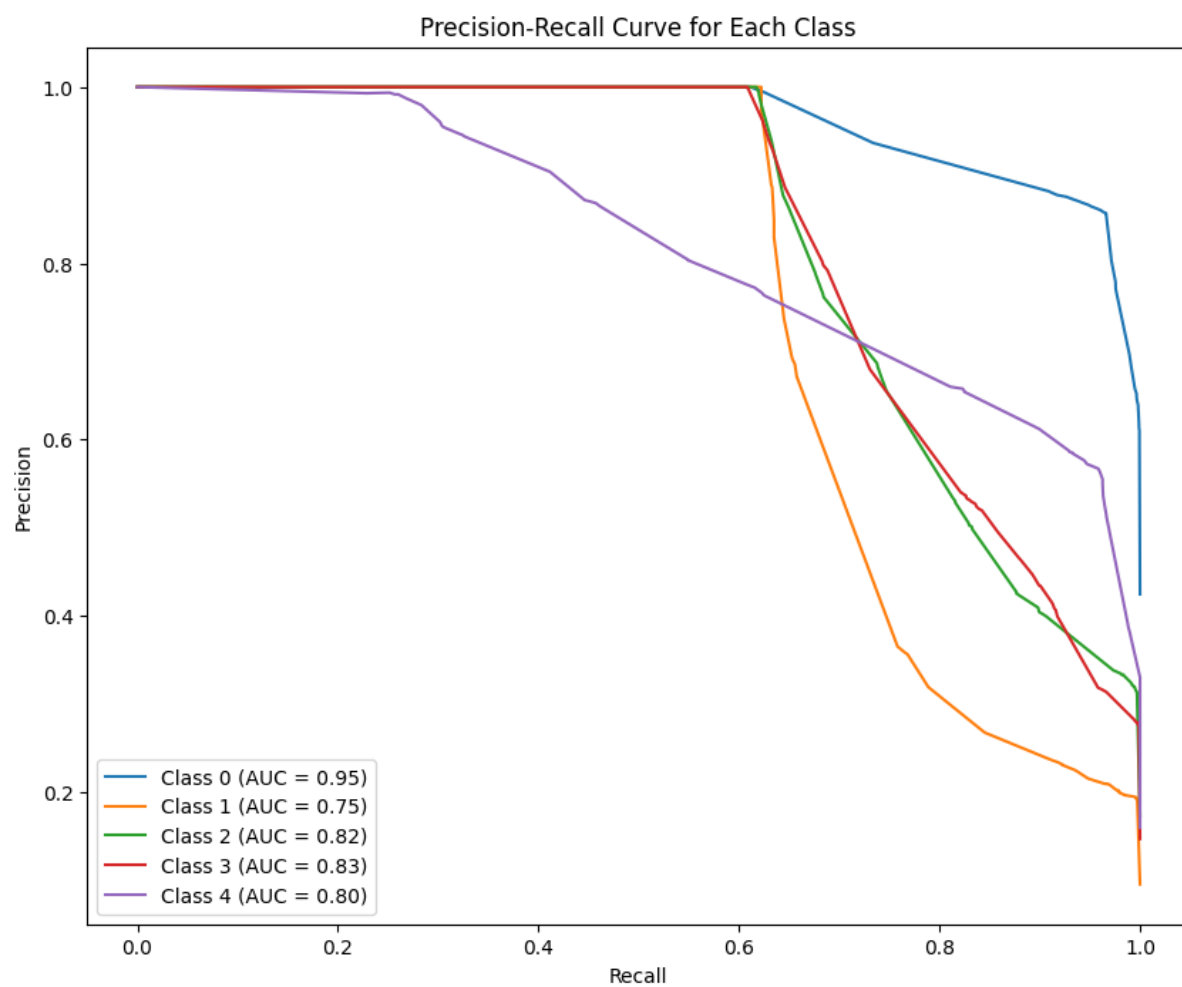
Best parameters: max\_depth=5, min\_samples\_split=2, min\_samples\_leaf=5

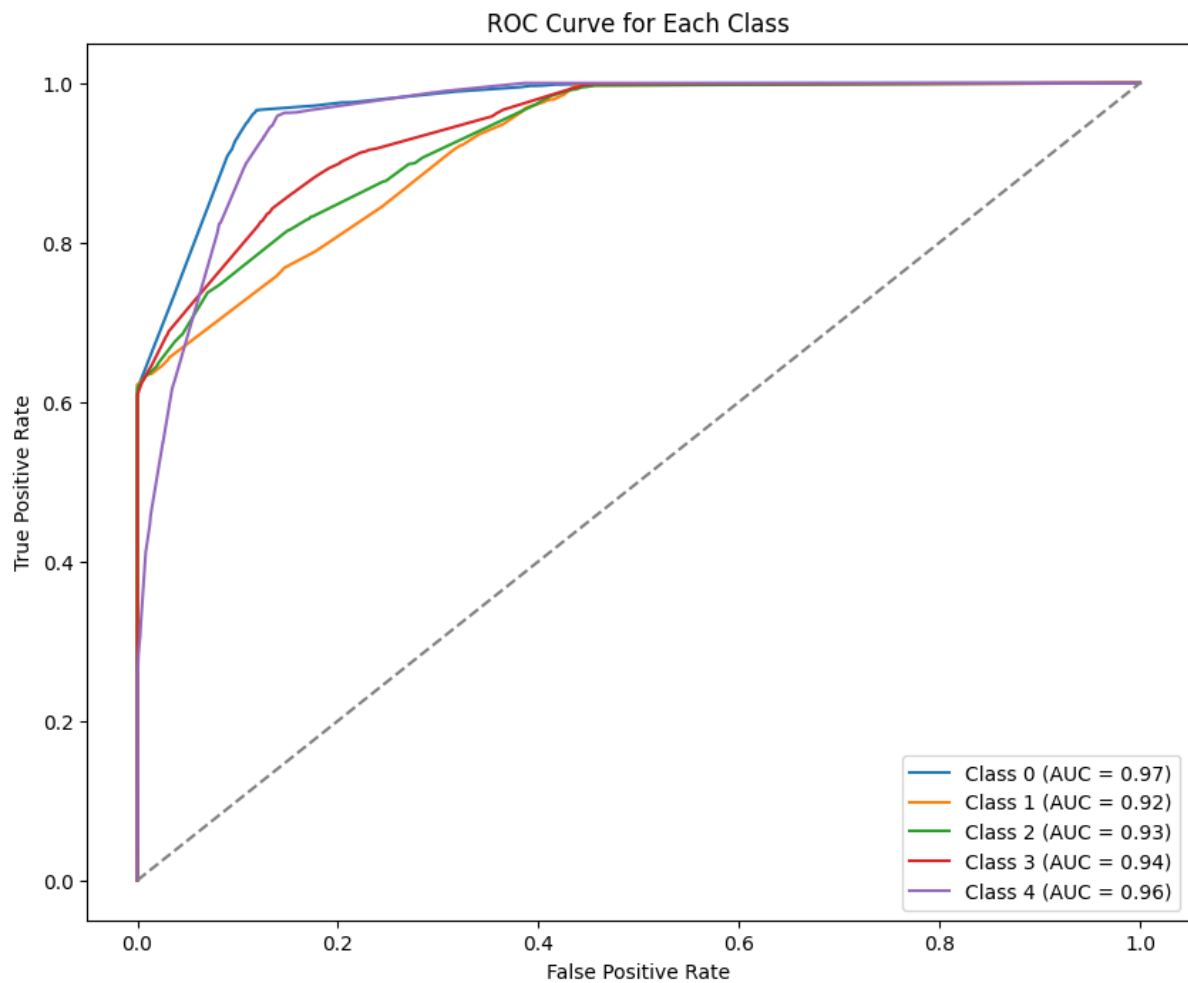
**Random Forest Parameters tuned:**

- n\_estimators: [100, 200]
- max\_features: ['auto', 'sqrt']
- max\_depth: [10, 20, 30, None]

Best parameters: n\_estimators=200, max\_features='sqrt', max\_depth=10

Results and Analysis Model Performance Decision Tree:

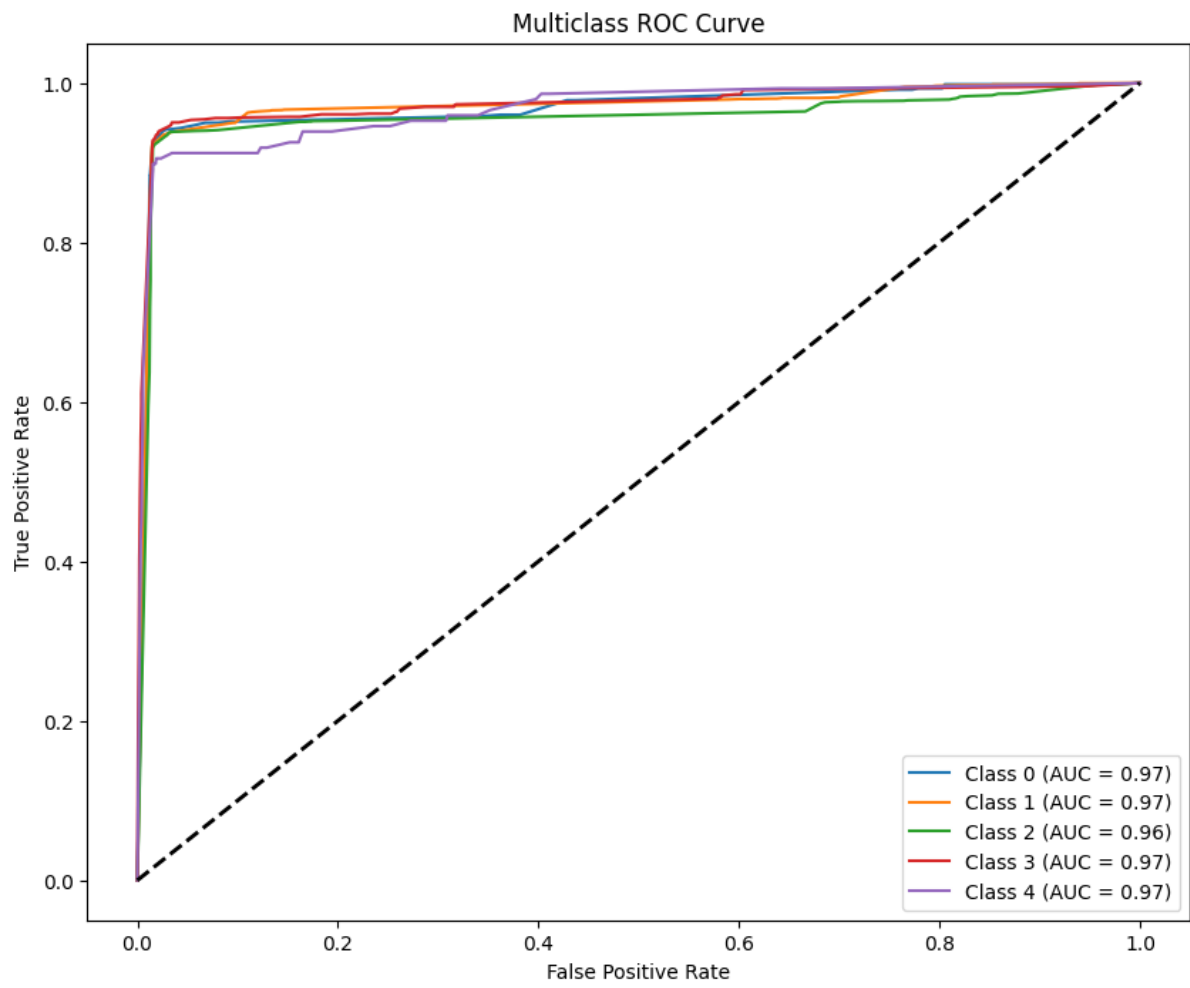


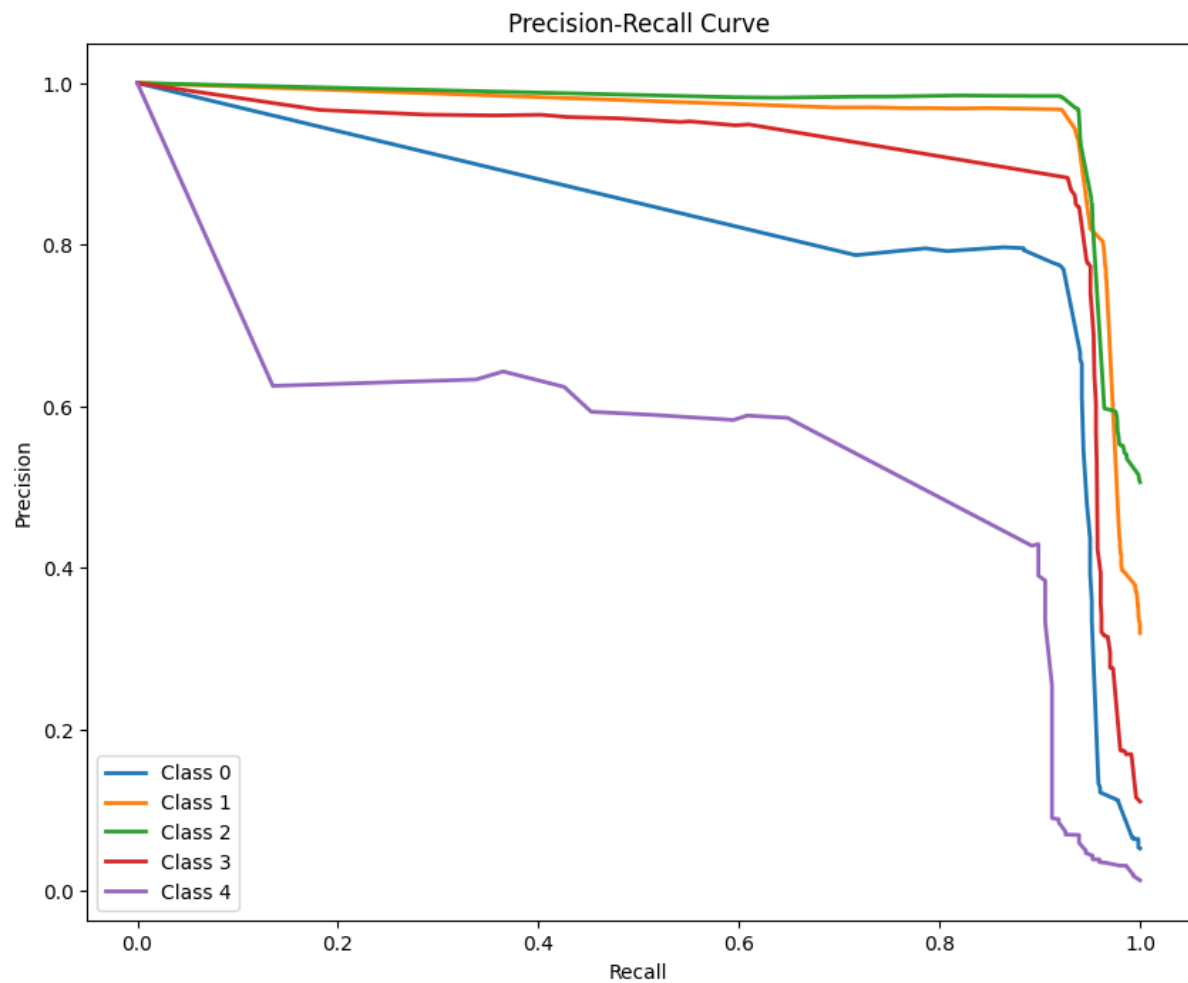


```
Best hyperparameters: {'classifier__max_depth': 10, 'classifier__min_samples_leaf': 1, 'classifier__min_samples_split': 2}
Accuracy: 0.818378471921779
Classification Report:
```

	precision	recall	f1-score	support
0	0.86	0.97	0.91	4907
1	1.00	0.62	0.77	1104
2	1.00	0.62	0.76	1997
3	0.92	0.63	0.75	1700
4	0.59	0.93	0.72	1849
accuracy			0.82	11557
macro avg	0.87	0.75	0.78	11557
weighted avg	0.86	0.82	0.82	11557

Results and Analysis Model Performance Random Forest Tree:





```
Best hyperparameters: {'Classifier__max_depth': 5, 'Classifier__min_samples_split': 2, 'Classifier__n_estimators': 100}
Accuracy: 0.926060606060606
Classification Report:
      precision    recall  f1-score   support

     0       0.77     0.92     0.84         603
     1       0.95     0.93     0.94        3683
     2       0.98     0.92     0.95        5842
     3       0.85     0.94     0.89        1274
     4       0.42     0.90     0.57         148

 accuracy          0.93        11550
 macro avg         0.79     0.92     0.84        11550
 weighted avg         0.94     0.93     0.93        11550
```