**Titanic Survival Prediction Final Report**

## 1. Main Objective of the Analysis

The primary goal of this analysis is to build a **predictive model** to estimate the survival probability of passengers on the Titanic. Through this model, we aim to:

- **Optimize Rescue Resource Allocation**: Provide data-driven insights for prioritizing rescue efforts in future similar events.

- **Identify Key Survival Factors**: Understand how demographic features (e.g., gender, class) influence survival probability.

- **Enhance Model Interpretability**: Offer transparent insights into the drivers of survival for business decision-making.
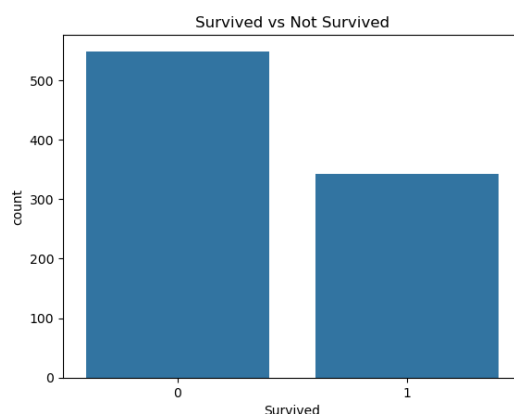
## 2. Dataset Description

### Dataset Source

The Titanic dataset from Kaggle (train.csv and test.csv) includes training data for 891 passengers and test data for 418 passengers.

### Key Features

- **Target Variable**: Survived (0 = Did Not Survive, 1 = Survived).


Survived vs Not Survived

- **Input Features**:

  - Demographic: Sex (gender), Age (age).

  - Socioeconomic Status: Pclass (ticket class, 1st/2nd/3rd class).

  - Family Structure: SibSp (number of siblings/spouses), Parch (number of parents/children).

  - Other: Fare (ticket fare), Embarked (port of embarkation).

**Analysis Goal**

To predict passenger survival probability using machine learning models and identify key factors influencing survival.

---

## 3. Data Exploration and Cleaning

**Data Exploration**

- **Survival Rate**: Approximately 38% of passengers survived (see Figure 1).

- **Key Feature Distributions**:

  - Females had a significantly higher survival rate than males (74% vs. 19%).

  - First-class passengers had a much higher survival rate (63%) compared to third-class passengers (24%).

**Data Cleaning and Feature Engineering**

1. **Handling Missing Values**:

   - Age: Filled with the median value.

   - Embarked: Filled with the most frequent value (S port).

   - Cabin: Dropped due to a high percentage of missing values.

2. **Creating New Features**:

o   FamilySize: Family size (SibSp + Parch + 1).

o   IsAlone: Whether the passenger was traveling alone (FamilySize == 1).
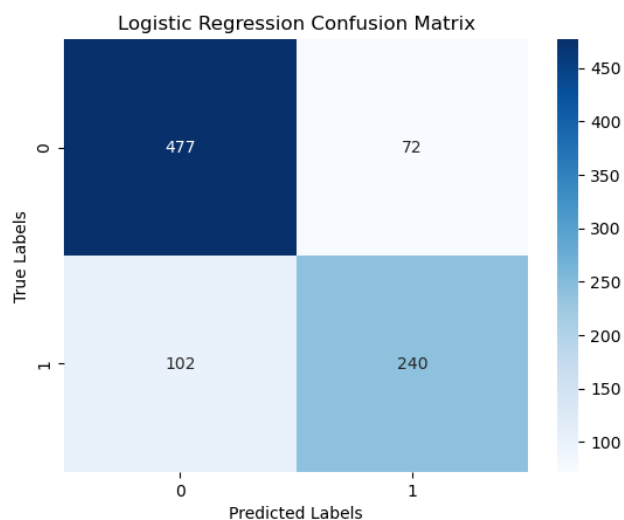
3.  **Removing Redundant Features**: Dropped irrelevant features such as PassengerId, Name, and Ticket.
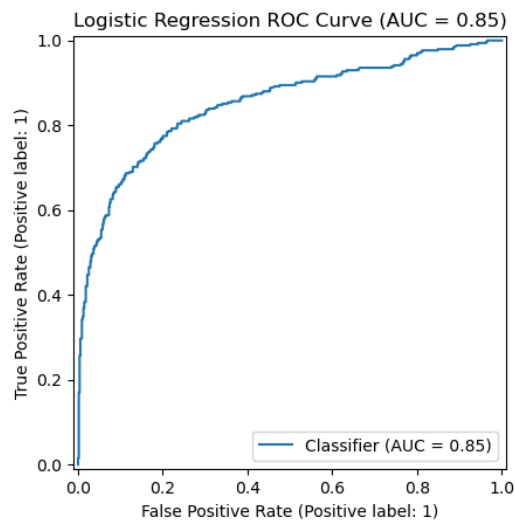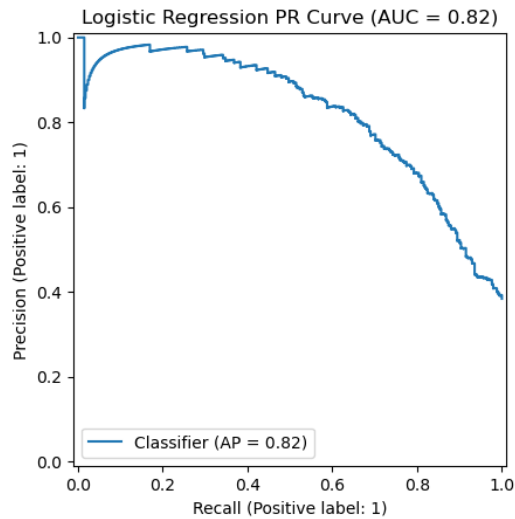
---

## 4. Model Training and Evaluation

We trained three different classifiers, ensuring consistency through **5-Fold Cross-Validation**:
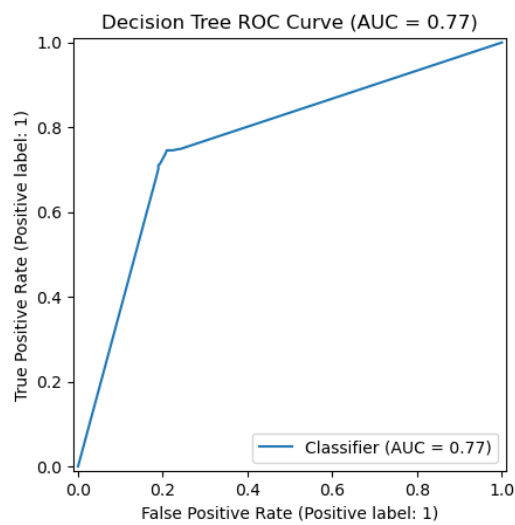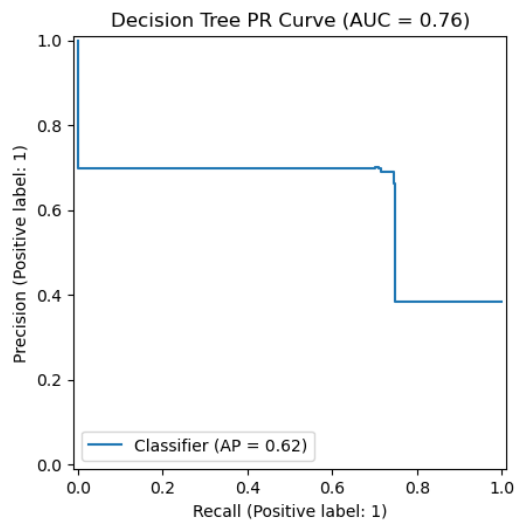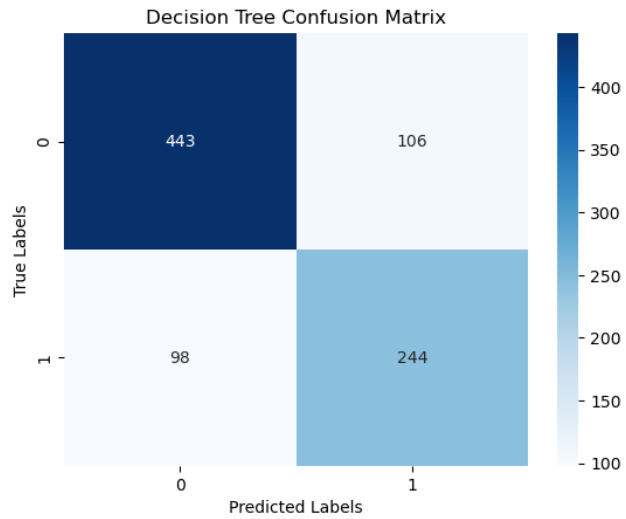
**Model 1: Logistic Regression**

- **Characteristics**: High interpretability, serves as a strong baseline model.

- **Strengths**: Fast computation, effective for linearly separable data, and provides clear insights into feature importance.

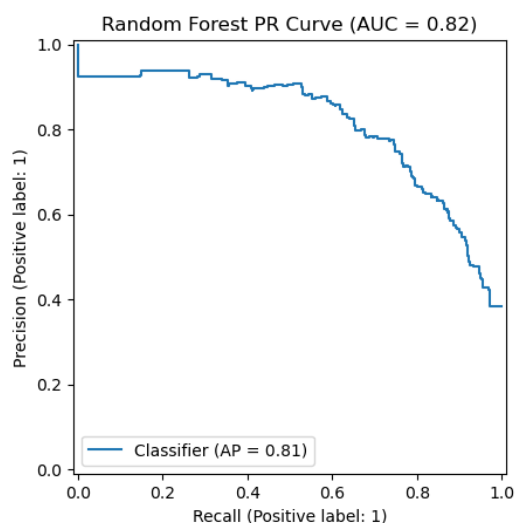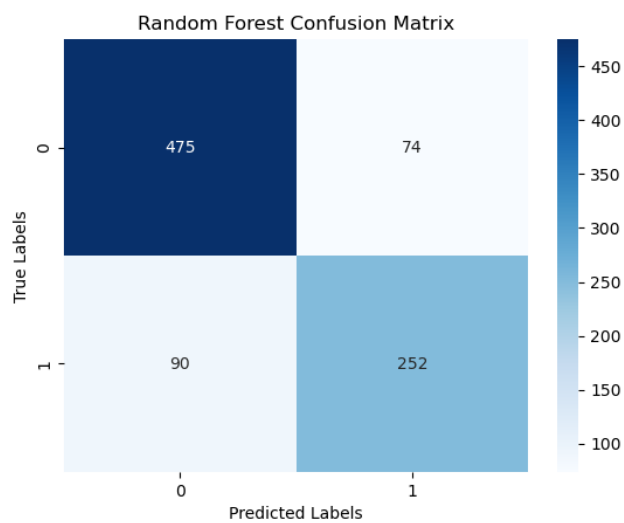- **Weaknesses**: Struggles with capturing nonlinear relationships and complex interactions.



Logistic Regression Confusion Matrix

Logistic Regression PR Curve (AUC = 0.82)


Logistic Regression ROC Curve (AUC = 0.85)
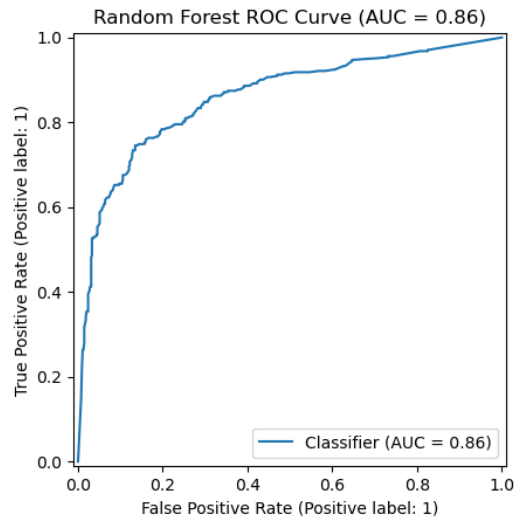
**Model 2: Decision Tree**

- **Characteristics**: Simple, interpretable, and capable of handling nonlinear relationships.

- **Strengths**: Automatically selects key features and is not sensitive to feature scaling.

- **Weaknesses**: Prone to overfitting, leading to poor generalization.

Decision Tree Confusion Matrix



Decision Tree PR Curve (AUC = 0.76)



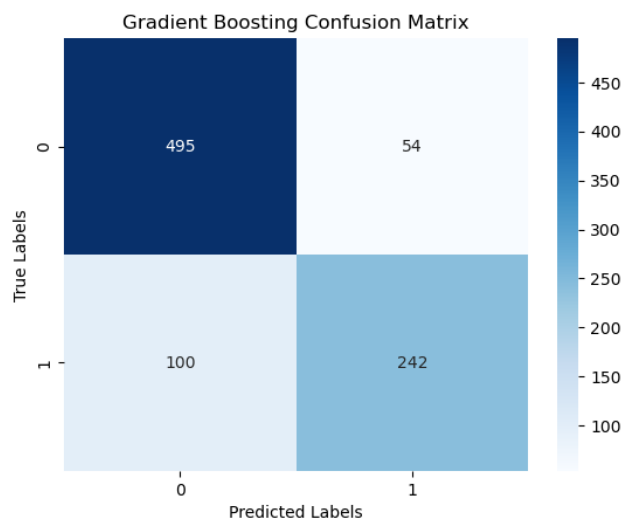Decision Tree ROC Curve (AUC = 0.77)

**Model 3: Random Forest**

- **Characteristics**: An ensemble method that balances predictive performance and interpretability.

- **Strengths**: Reduces overfitting compared to a single decision tree, provides feature importance scores.

- **Weaknesses**: Higher computational cost and less interpretable than logistic regression.
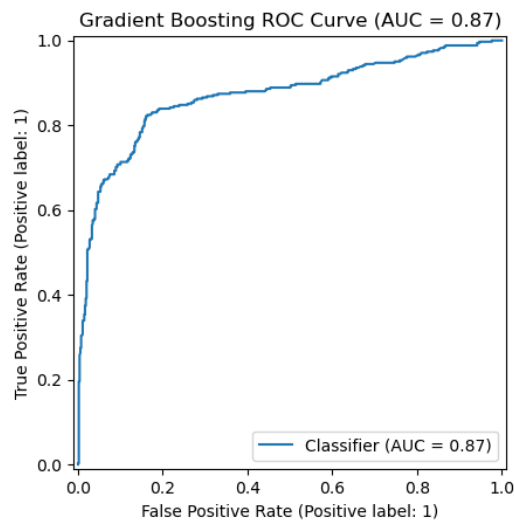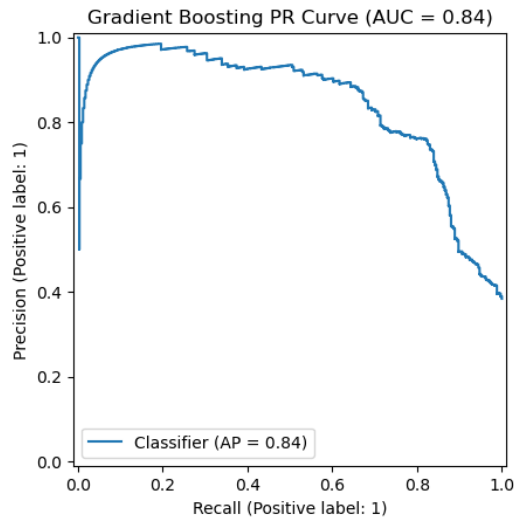


Random Forest Confusion Matrix



Random Forest PR Curve (AUC = 0.82)

Random Forest ROC Curve (AUC = 0.86)

**Model 4: Gradient Boosting**

- **Characteristics: An ensemble method that builds models sequentially, improving prediction accuracy at each step.**

- **Strengths: High predictive performance, captures complex patterns, and reduces bias by focusing on hard-to-predict instances.**

- **Weaknesses: Can be prone to overfitting with noisy data, requires careful tuning of hyperparameters, and has higher computational costs.**
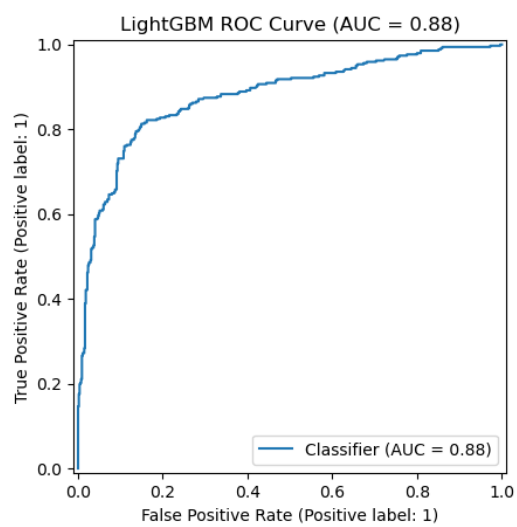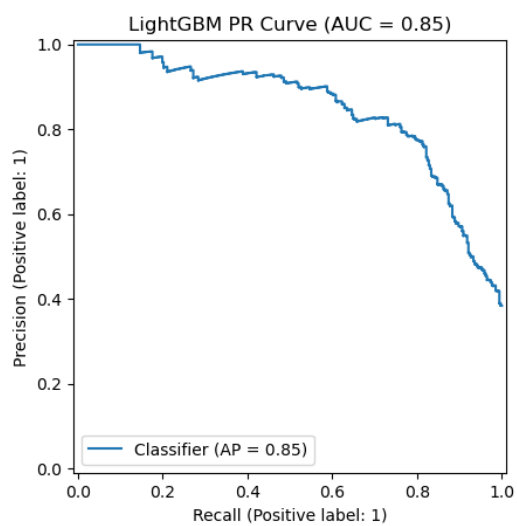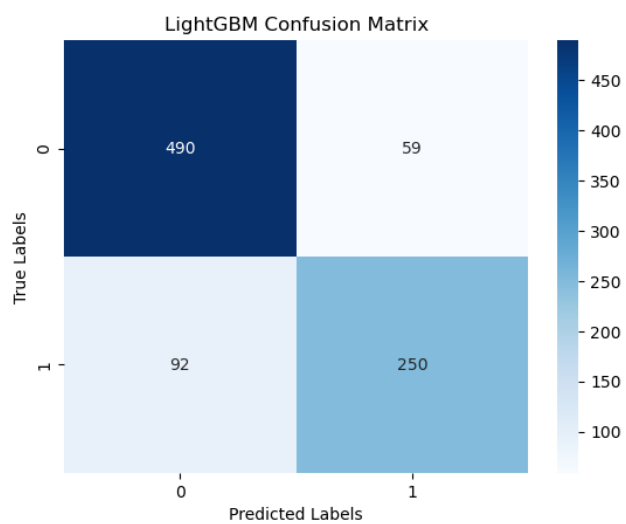


Gradient Boosting Confusion Matrix

Gradient Boosting PR Curve (AUC = 0.84)


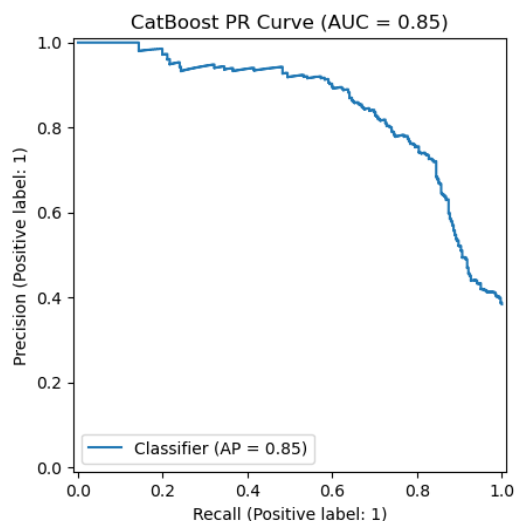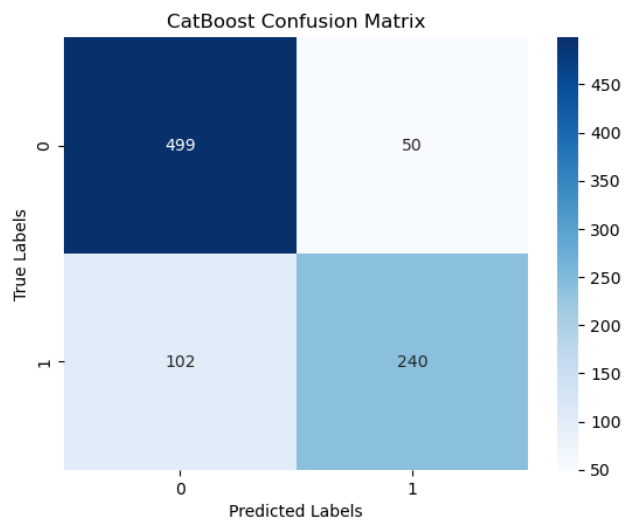Gradient Boosting ROC Curve (AUC = 0.87)

**Model 5: LightGBM**

- **Characteristics: A gradient boosting variant optimized for efficiency.**

- **Strengths: Fast training, effective for large datasets, and captures complex patterns.**

- **Weaknesses: Can overfit on small datasets, sensitive to hyperparameter tuning.**
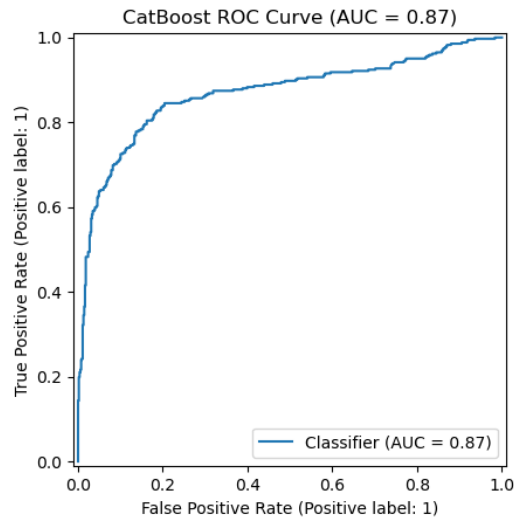
LightGBM Confusion Matrix



LightGBM PR Curve (AUC = 0.85)



LightGBM ROC Curve (AUC = 0.88)
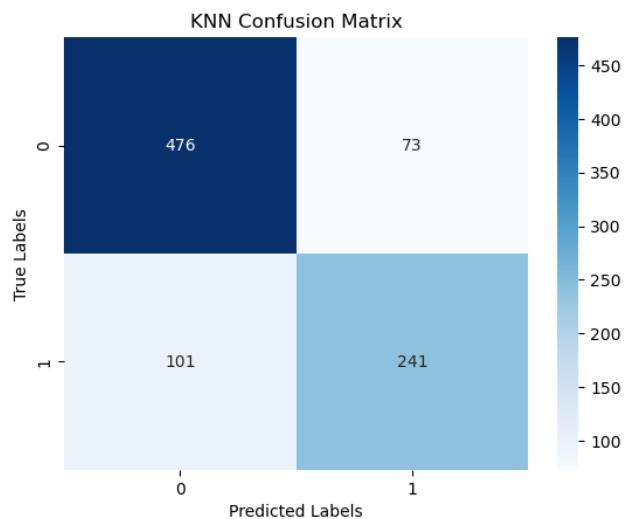
**Model 6: CatBoost**

- **Characteristics**: A gradient boosting model optimized for categorical features.

- **Strengths**: Handles categorical data efficiently, reduces preprocessing needs, and is robust to missing values.

- **Weaknesses**: Slower training time, may not significantly outperform other boosting models on small datasets.
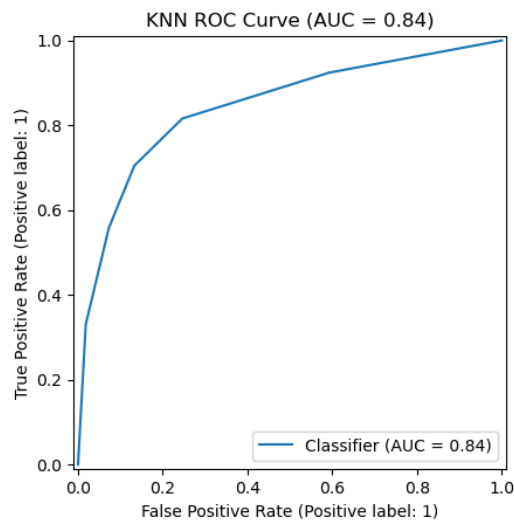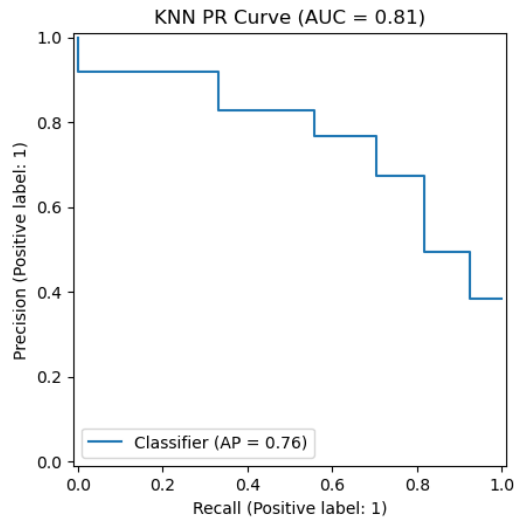


CatBoost Confusion Matrix



CatBoost PR Curve (AUC = 0.85)

CatBoost ROC Curve (AUC = 0.87)
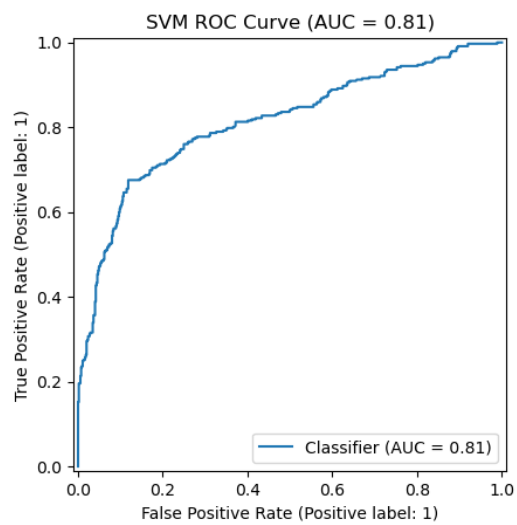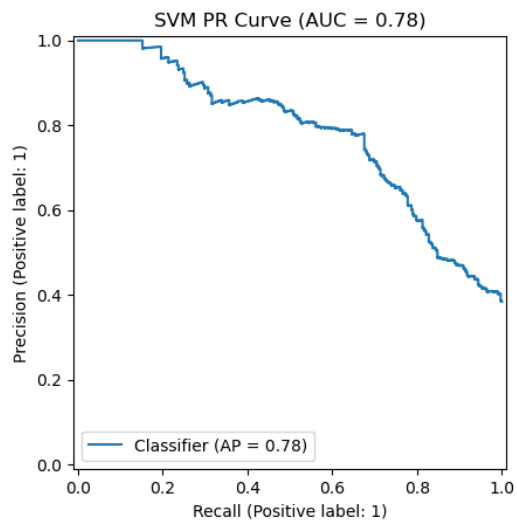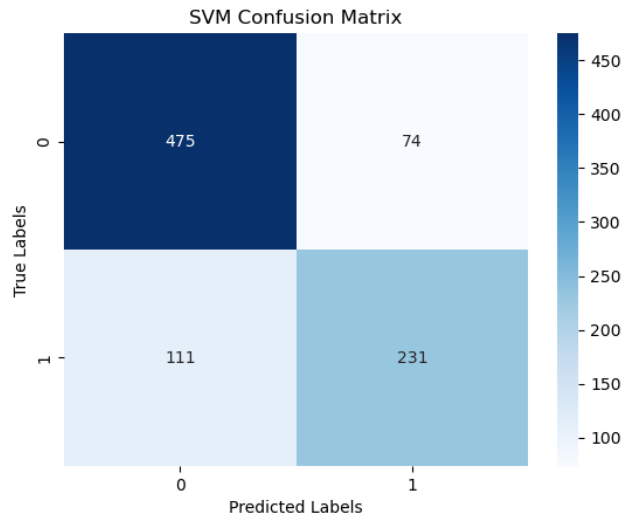
**Model 7: K-Nearest Neighbors (KNN)**

- **Characteristics**: A non-parametric model based on proximity to labeled samples.

- **Strengths**: Simple and intuitive, requires no explicit training phase.

- **Weaknesses**: Computationally expensive, struggles with high-dimensional and imbalanced data.
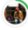


KNN Confusion Matrix

KNN PR Curve (AUC = 0.81)


KNN ROC Curve (AUC = 0.84)

**Model 8: Support Vector Machine (SVM)**

- **Characteristics**: A powerful classification algorithm that constructs an optimal hyperplane.

- **Strengths**: Works well with high-dimensional data, effective in small datasets.

- **Weaknesses**: Computationally expensive, especially with large datasets, and requires careful tuning of kernel parameters.

## SVM Confusion Matrix



## SVM PR Curve (AUC = 0.78)



## SVM ROC Curve (AUC = 0.81)

**5. final Kaggle score**



| 1490 | thomas #2 | | 0.78708 | 4 | 2h |
| 1491 | **Tsai Cheng Hsung** | | 0.78708 | 14 | 35m |

**Your Best Entry!**
Your most recent submission scored 0.78708, which is an improvement of your previous score of 0.77751. Great job!

Tweet this

| 1492 | Jonathan Peteza | | 0.78468 | 10 | 2mo |

**6. Key Findings and Insights**

1. **Gender and Class Are Core Survival Drivers**:

   o Females were 3.9 times more likely to survive than males.

   o First-class passengers were 2.6 times more likely to survive than third-class passengers.

2. **Impact of Family Structure**:

   o Passengers traveling alone (IsAlone=1) had a lower survival rate (30% vs. 50%).

3. **Non-linear Impact of Age**:

   o Children (<10 years) had a higher survival rate, while adult males had the lowest survival rate.

**7. Future**

1. **Data Enhancement**:

   o Collect detailed Cabin information to improve feature granularity.

   o Add passenger occupation or social status data (e.g., titles like Mr/Mrs).

2. **Model Optimization**:

   o Experiment with neural networks to capture more complex interactions.