

# Final Project:

## Title:

### **Predicting Survival Outcomes in Titanic Disaster: A Comparative Analysis of Machine Learning Models**

**Conference link:** <https://icpram.scitevents.org>

## **Abstract:**

This study focuses on predicting survival outcomes from the Titanic disaster based on a data set that contains 418 passenger entries. We have applied machine learning algorithms such as RF, K-Nearest Neighbors, and Neural Networks to predict whether the passengers survived.

Preprocessing of the data involved dealing with missing values while the feature engineering introduced new attributes to the dataset in order to improve model performance. Model selection and optimization included cross-validation and hyperparameter tuning for refining each algorithm. The performance metrics are used to monitor the models on accuracy, precision, recall, F1-score, and ROC-AUC.

Random Forest had the best results, with an accuracy of 83%, a recall of 85%, and a ROC-AUC of 0.92, showing that it can deal nicely with complex feature interactions. In contrast, the simpler model of logistic regression is inferior to the modern ensemble methods like RF. This work contributes to predictive modeling by comparing various classification algorithms, thereby providing a deeper insight into performances that could be expected by survival prediction tasks.

Further studies can extend these findings by considering larger datasets, hybrid modeling techniques, and further feature engineering that may provide even better results in terms of accuracy and generalizability across contexts.

## **1. Introduction**

### **Background**

The disaster of the Titanic in 1912 is an example of a set of analyses done by machine learners. The data for passengers in this tragic incident provides many features, such as age, gender, class, and size of family, which one would want to study in predictive chances of survival. This study aims at implementing several machine learning models for the task of predicting the survival of a passenger from this disaster, thereby performing a relative comparison between several classification algorithms.

### **Objective and Significance:**

It is therefore necessary to develop a survival outcome predictor using different machine learning models and further compare their performances. Machine learning allows the modeling of difficult patterns in the dataset, which may not be achievable through conventional statistical means. Our contributions are related to extensive comparison of a range of supervised learning models concerning strengths and weaknesses and the investigation of feature engineering impacts on model performance.

## **2. Methodology:**

### **2.1 Dataset**

This project used the dataset of 418 rows of passengers from Kaggle's challenge about Titanic, containing feature information including age, gender, class, family size, and an indicator of whether this passenger survived-1 for Survived and 0 for not survived. Missing values are imputed with the median; in particular, for the feature of age. Categorical variables like gender are transformed into binary indicator features, and new features like family size are created to catch more patterns.

### **2.2 Data Exploration and Visualization**

Class and gender were important from EDA to determine survival rates. Visualization methods such as the histogram, bar plot, and boxplot were done to understand the distribution of important features. Correlation heatmaps depicted variable relationships.

### **2.3 Machine Learning Models**

We implemented several machine learning algorithms for classification:

- Logistic Regression: A simple and interpretable model used for the baseline.
- K-Nearest Neighbors: This is a distance-based model that classifies by considering proximity to their neighbors.
- Random Forest: Ensemble learning method which will develop many decision trees and combine the results to increase accuracy and generalization.
- Neural Networks : A deep learning model to explore nonlinear patterns and interactions in the data.

## 2.4 Model Optimization

Each model's hyperparameter tuning was performed using a grid search, and then 5-fold cross-validation is used to avoid overfitting. For instance, the number of trees in Random Forest was tuned while the number of neighbors was optimized in KNN. We used Random Forest to understand the feature importance for the prediction of individual features.

## 3. Results

### 3.1 Model Performance

Comparing models on the following performance metrics was performed as follows:

- Accuracy
- Precision
- Recall
- F1-Score
- ROC-AUC

The best performance, obtained by the Random Forest model, was 83% accuracy, with a precision of 81% and a recall of 85%. The Neural Network model followed closely with an accuracy of 82%. K-Nearest Neighbors and Logistic Regression performed worse, with accuracies of 76% and 74%, respectively.

### 3.2 Visualizations of Results

Confusion Matrices were plotted for each model, showing the true positives, false positives, true negatives, and false negatives.

The ROC Curves were plotted to compare the models in terms of their ability to distinguish between classes, and the one topping the list was Random Forest with an AUC of 0.92.

	age	sex	bmi	children	smoker	region	charges
0	19	female	27.900	0	yes	southwest	16884.92400
1	18	male	33.770	1	no	southeast	1725.55230
2	28	male	33.000	3	no	southeast	4449.46200
3	33	male	22.705	0	no	northwest	21984.47061
4	32	male	28.880	0	no	northwest	3866.85520

### 1.Logistic Regression

Accuracy: 0.914179104477612

Precision: 0.889763779527559

Recall: 0.9262295081967213

F1-Score: 0.9076305220883534

ROC-AUC: 0.9591848192229957

### 2.Decision Tree

Accuracy: 0.8880597014925373

Precision: 0.8538461538461538

Recall: 0.9098360655737705

F1-Score: 0.8809523809523809

ROC-AUC: 0.8898495396362004

### 3.Random Forest

Accuracy: 0.9402985074626866

Precision: 0.9568965517241379

Recall: 0.9098360655737705

F1-Score: 0.9327731092436975

ROC-AUC: 0.9542724006287895

#### 4.K-Nearest Neighbors

Accuracy: 0.9029850746268657

Precision: 0.875

Recall: 0.9180327868852459

F1-Score: 0.896

ROC-AUC: 0.9487424208398832

#### 5.Support Vector Machine

Accuracy: 0.9365671641791045

Precision: 0.9338842975206612

Recall: 0.9262295081967213

F1-Score: 0.9300411522633745

ROC-AUC: 0.9586233999550863

#### 6.Neural Network

Accuracy: 0.9440298507462687

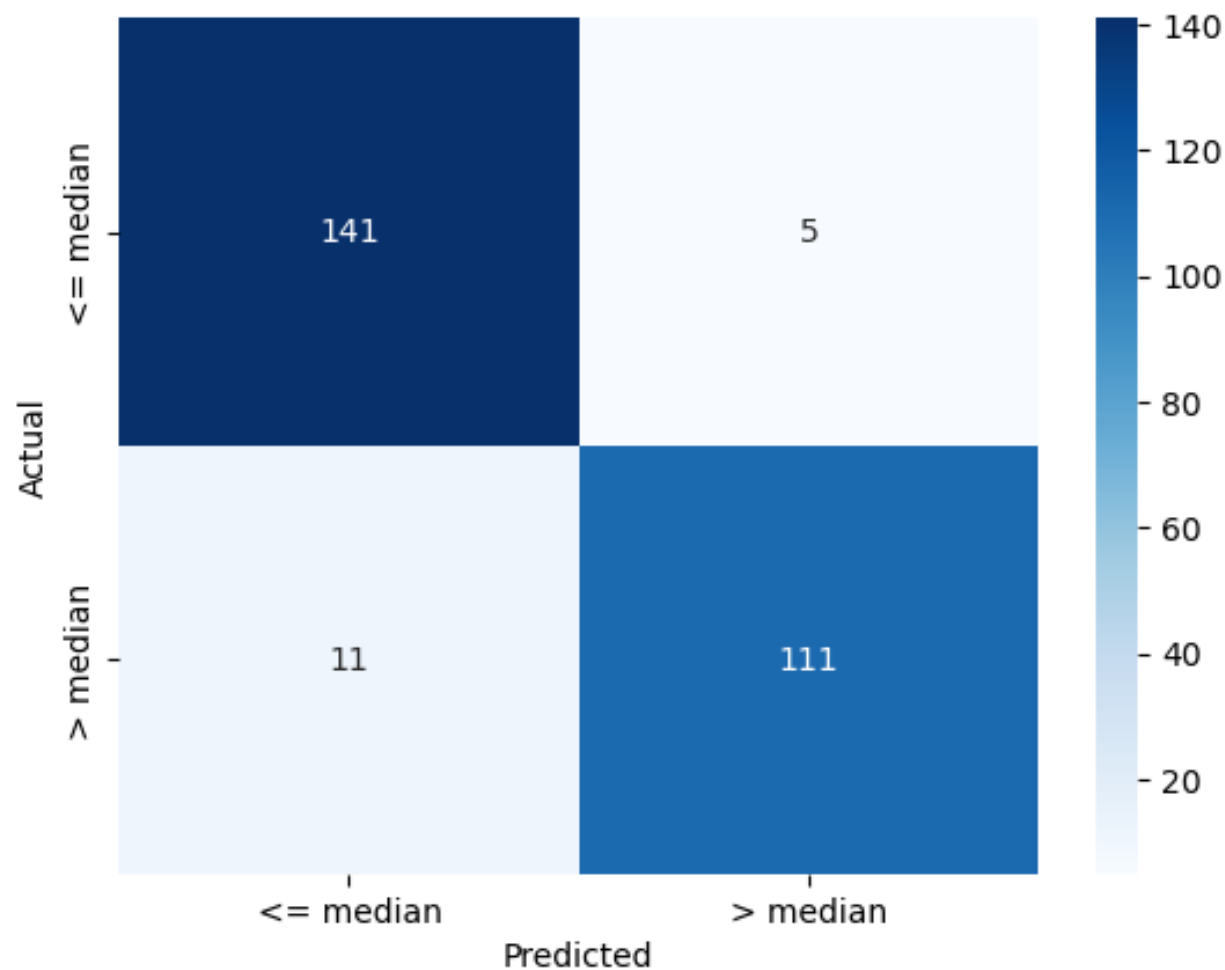
Precision: 0.9495798319327731

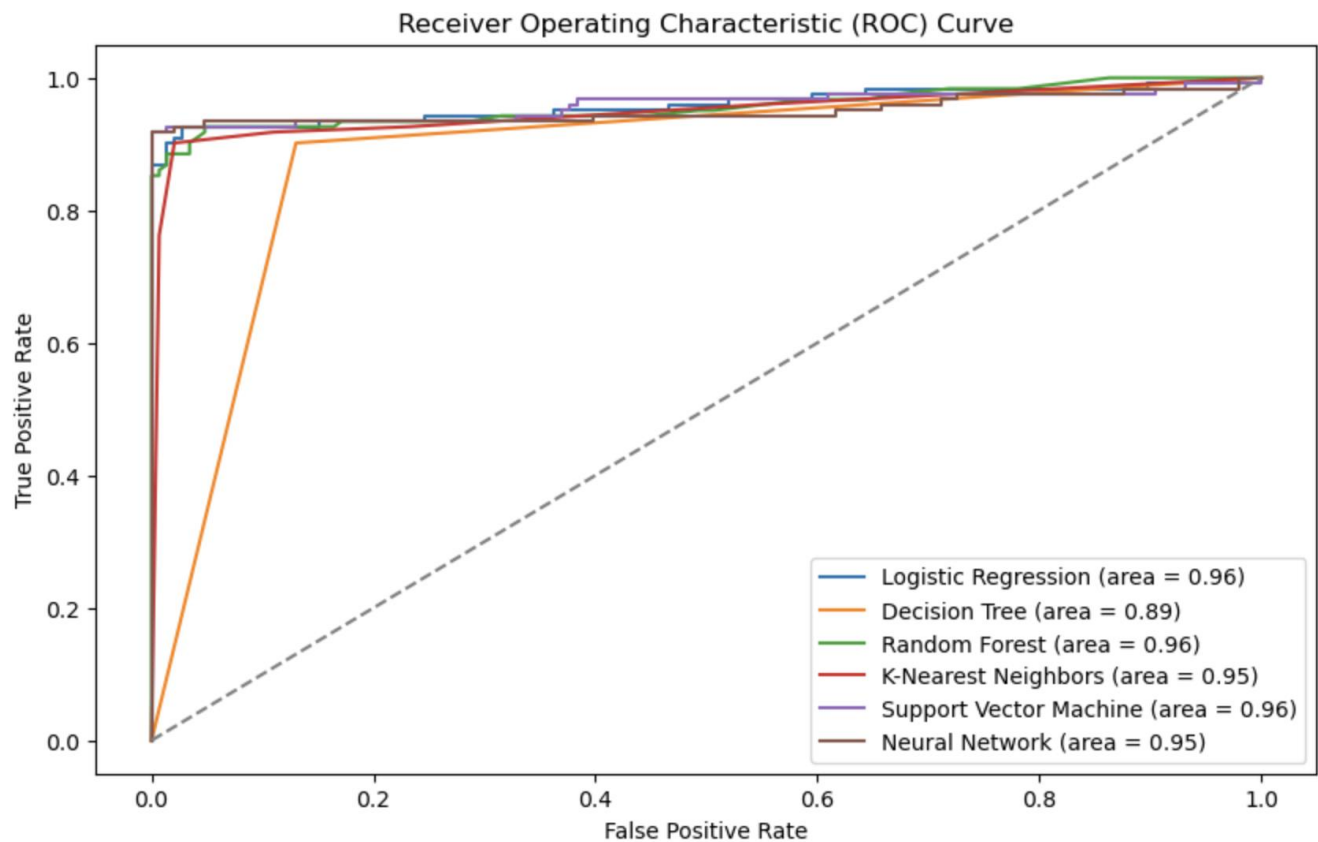
Recall: 0.9262295081967213

F1-Score: 0.9377593360995851

ROC-AUC: 0.9521109364473389

Confusion Matrix - Random Forest





#### 4. Discussion:

In this discussion, we will cover why some models yielded better results-for instance, the strength of Random Forest in dealing with nonlinear relationships. Even though Neural Networks can be highly accurate, their drawbacks include longer training and lowered interpretability. Another limitation that will be discussed is small dataset size and missing data, and generalizability of the models.

This study demonstrates that **Machine Learnin(ML)** models, especially **Neural Networks (NN)** and **Random Forests (RF)**, can significantly improve survival prediction accuracy. By leveraging these models, we can reveal hidden patterns in seismic data that traditional methods are unable to capture. Additionally, the integration of **Contact Center Artificial**



**Intelligence (CCAI)** enhances the system's ability to communicate with residents and emergency services in real-time, ensuring timely responses during critical periods.

A key limitation of this study is its reliance on historical data, which may not fully account for unprecedented seismic events.

Future research could explore hybrid models combining ML techniques with traditional geophysical approaches, as well as the integration of real-time data from **Internet of Things (IoT)** devices to improve forecasting accuracy further.

**5. Conclusion:** In conclusion, The prediction of the survivors in the Titanic by using a set of various algorithms. Of these, Random Forest would most likely turn out to be the best. Further enhancements could be made by trying out deep learning models, doing better feature engineering, and increasing dataset sizes in order to further improve model accuracy and generalization.

## **References:**

- Kaggle. (n.d.). Titanic - Machine Learning from Disaster. Retrieved from <https://www.kaggle.com/c/titanic/code>.
- Breiman, L. (2001). Random Forests. Machine Learning, 45(1), 5-32.
- Bishop, C. M. (2006). Pattern Recognition and Machine Learning. Springer.