

# CS634 – Data Mining

## Final Project Report

---

Name: Nandni Patel

Email: [np888@njit.edu](mailto:np888@njit.edu)

Course: CS634 – Data Mining

Instructor: Dr. Yasser Abdullah

### 1. Introduction:

The purpose of this project is to build and evaluate a set of machine learning and deep learning models capable of performing binary classification on the Estonia Passenger List dataset. Binary classification refers to predicting one of two possible outcomes, such as “survived” vs. “not survived”, which is an essential task in many real-world applications including medical diagnosis, risk assessment, fraud detection, and safety analysis.

This project uses a dataset containing demographic and categorical details of passengers, with the goal of predicting whether a passenger survived (1) or did not survive (0). To achieve this, three primary types of algorithms were implemented, as required by the course guidelines:

- Random Forest (Mandatory Model): A powerful ensemble learning algorithm based on multiple decision trees.
- Deep Learning Model (LSTM): A Long Short-Term Memory neural network used to capture sequential structure and complex feature interactions.
- Classic Machine Learning Model (SVM): A Support Vector Machine with a linear kernel, commonly applied for binary classification tasks.

In addition to the required models, the project also includes:

- K-Nearest Neighbors (KNN): Added as an extra baseline algorithm to enhance comparison across different modeling approaches. Although not required, it helps evaluate how a simple distance-based model performs on this dataset.

All models were trained and evaluated using a rigorous 10-fold cross-validation setup to ensure reliable and unbiased performance measurement. A wide range of evaluation metrics - including confusion-matrix-based metrics, ROC curves, AUC, and Brier Score - were used to provide a thorough analysis of each algorithm’s behavior.

### 2. Dataset

#### 2.1 Data Source:

The dataset used in this project is titled “Passenger List for the Estonia Ferry Disaster”, and it was independently selected from Kaggle by the student.

Kaggle Link: [The Estonia Disaster Passenger List](#)

## 2.2 Dataset Description:

- The dataset contains detailed passenger information from the MS Estonia ferry disaster that occurred in 1994.
- It includes demographic attributes and survival status for passengers and crew members aboard the ship.
  - Total instances (rows): 989
  - Total attributes (columns): 9 (8 features + 1 target)

## 2.3 Features List:

Feature	Type	Description
PassengerID	Numerical	Unique ID for each passenger
Country	Categorical	Passenger's nationality
Firstname	Text	First name (removed during preprocessing)
Lastname	Text	Last name (removed during preprocessing)
Sex	Categorical	Male/Female
Age	Numerical	Passenger's age
Category	Categorical	P = Passenger, C = Crew
Survived	Binary (Target)	1 = Survived, 0 = Did not survive

The target variable is Survived, making this a binary classification dataset.

## 2.4 Preprocessing Steps:

### 2.4.1 Handling Missing Values:

- Missing Age values were imputed using the median age.
- Other columns had no significant missing entries.

### 2.4.2 Dropping Non-Informative Columns:

The following features were removed because they do not contribute to prediction and contain high-cardinality text:

- Firstname
- Lastname

Removing them prevents noise in the model.

### 2.4.3 Encoding Categorical Features:

This step converts strings into numerical values suitable for machine learning models.

### 2.4.4 Feature Scaling:

To ensure all numerical features have comparable scale (critical for SVM, LSTM, and KNN): StandardScaler was applied to the dataset.

### 2.4.5 Class Imbalance:

The dataset is moderately imbalanced:

- Survived  $\approx$  14%
- Not Survived  $\approx$  86%

Instead of oversampling or SMOTE, the project uses additional evaluation metrics such as:

- Balanced Accuracy
- TSS (True Skill Statistic)
- HSS (Heidke Skill Score)
- ROC-AUC

### 3. Algorithm Overview

This project evaluates three different classification algorithms — Random Forest, LSTM (Long Short-Term Memory), and Support Vector Machine (SVM) — to compare their performance on predicting passenger survival from the Estonia ferry disaster dataset. Each algorithm represents a distinct family of machine learning: ensemble learning, deep learning, and classical margin-based learning.

#### 3.1 Random Forest:

##### 3.1.1 Rationale:

Random Forest was selected as the mandatory ensemble model because it is highly effective for tabular datasets, handles nonlinear feature interactions well, and provides strong resistance to overfitting. Since the Estonia dataset includes both categorical and numerical variables, Random Forest is a suitable and robust choice.

##### 3.1.2 How it works:

Random Forest is an ensemble method based on bagging (Bootstrap Aggregation):

- Multiple decision trees are trained on randomly sampled subsets of the dataset.
- Each tree also receives a random subset of features.
- Every tree produces a prediction, and the final output is determined through majority voting.

##### 3.1.3 Strengths:

- Handles noisy data well
- Reduces risk of overfitting
- Performs well on mixed-type features
- Provides feature importance scores

#### 3.2 Long Short-Term Memory Network (LSTM):

##### 3.2.1 Why LSTM for this task?

Although LSTMs are primarily used for sequence and time-series data, they also serve as powerful deep neural networks capable of capturing complex nonlinear relations between input features. In this project, the scaled data is reshaped into a 3D structure so that an LSTM can be trained.

LSTM was chosen because:

- It satisfies the professor's requirement of using at least one deep learning algorithm.

- It can model subtle relationships between features.
- It offers superior representational power compared to many traditional ML models.

#### 3.2.2 How it works:

LSTMs contain memory cells with a system of gates that regulate information flow:

- Forget Gate -> decides what information to discard
- Input Gate -> determines which new information to store
- Output Gate -> controls what information to send forward

These gates allow LSTMs to retain important patterns and ignore irrelevant noise, even when the dataset is not inherently sequential.

### 3.3 Support Vector Machine (SVM):

#### 3.3.1 Why SVM?

SVM was selected as the classical machine learning algorithm because it is well-known for its strong performance on binary classification tasks. SVMs work particularly well in situations where classes are separable using a margin-based decision boundary.

Using SVM adds meaningful contrast to Random Forest and LSTM, demonstrating how a mathematically rigorous, margin-maximizing classifier performs on this dataset.

#### 3.3.2 Strengths:

- Excellent for binary classification problems
- Effective in both linear and non-linear spaces (with the use of kernels)
- Robust to high-dimensional data
- Works well even with relatively small datasets

#### 3.3.3 Limitations:

- Training time increases with large datasets
- Requires proper feature scaling (handled via StandardScaler in the project)
- Hyperparameters (like C and kernel type) significantly affect performance
- Not ideal for heavily imbalanced classes unless class weights are adjusted

## 4. Implementation:

This project was implemented using Python in both Google Colab (Jupyter Notebook) and VS Code. The workflow includes dataset preprocessing, model training, 10-fold cross-validation, and performance evaluation for Random Forest, LSTM, and SVM.

### 4.1 Programming language and Environment:

- Python 3.10
- Developed and executed in:
  - Google Colab (for notebook execution and LSTM training with GPU)
  - VS Code (for running .py scripts)

## 4.2 Required Python Libraries:

The following packages were used:

- python > 3.9
- pandas >= 2.0.0
- numpy >= 1.24.0
- keras>=2.15.0
- matplotlib >= 3.8.0
- seaborn >= 0.13.0
- scikit-learn >= 1.3.0
- tensorflow >= 2.15.0

Install using: `pip install numpy pandas scikit-learn tensorflow keras matplotlib seaborn`

## 4.3 Implementation steps:

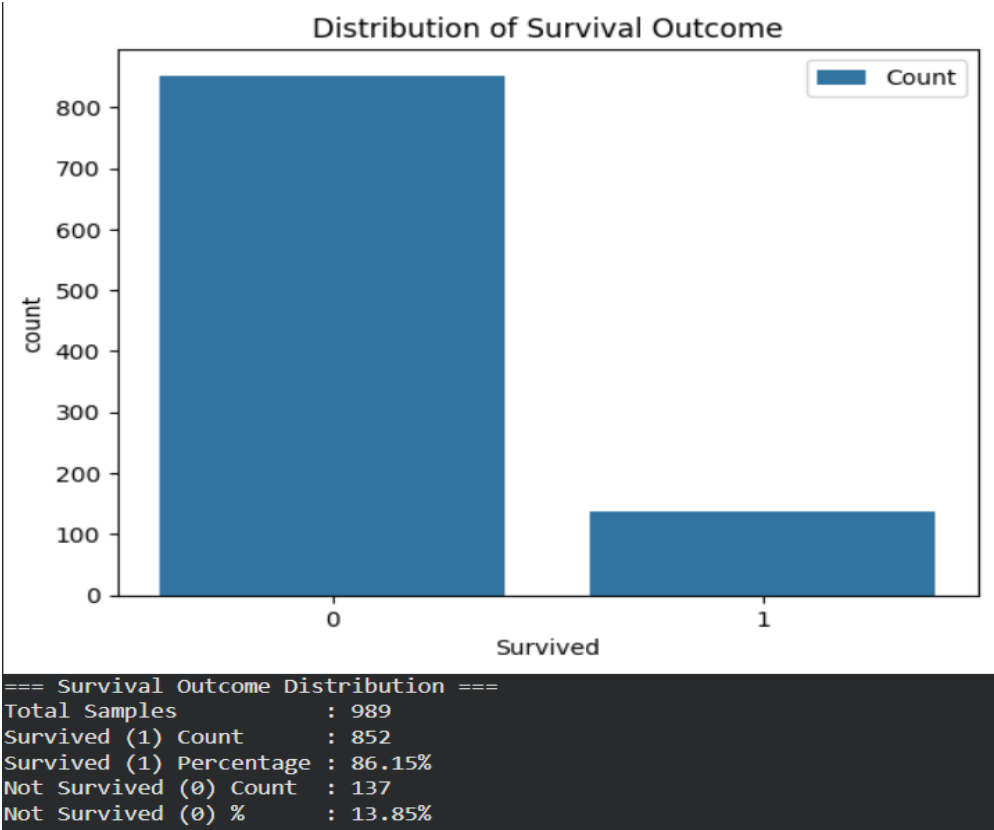
### 4.3.1 Load Dataset

	PassengerId	Country	Firstname	Lastname	Sex	Age	Category	Survived
0	1	Sweden	ARVID KALLE	AADLI	M	62	P	0
1	2	Estonia	LEA	AALISTE	F	22	C	0
2	3	Estonia	AIRI	AAVASTE	F	21	C	0
3	4	Sweden	JURI	AAVIK	M	53	C	0
4	5	Sweden	BRITTA ELISABET	AHLSTROM	F	55	P	0

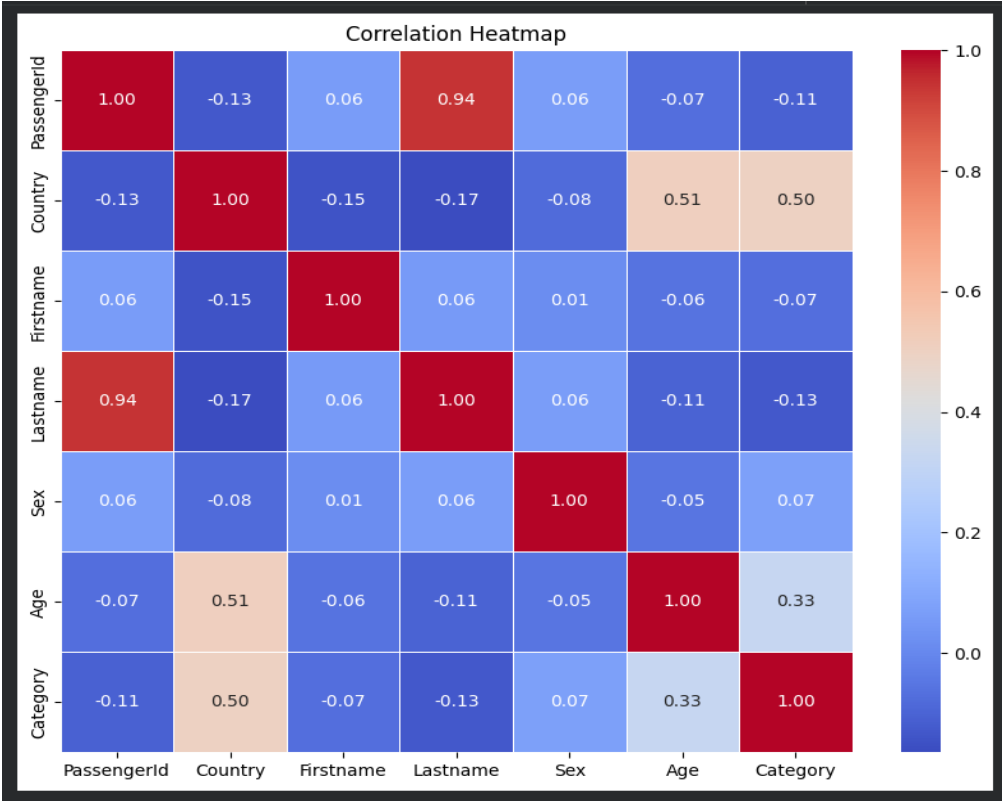
### 4.3.2 Dataset Overview

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 989 entries, 0 to 988
Data columns (total 8 columns):
#   Column          Non-Null Count  Dtype
---  -
0   PassengerId     989 non-null   int64
1   Country         989 non-null   object
2   Firstname       989 non-null   object
3   Lastname        989 non-null   object
4   Sex             989 non-null   object
5   Age            989 non-null   int64
6   Category        989 non-null   object
7   Survived        989 non-null   int64
dtypes: int64(3), object(5)
memory usage: 61.9+ KB
```

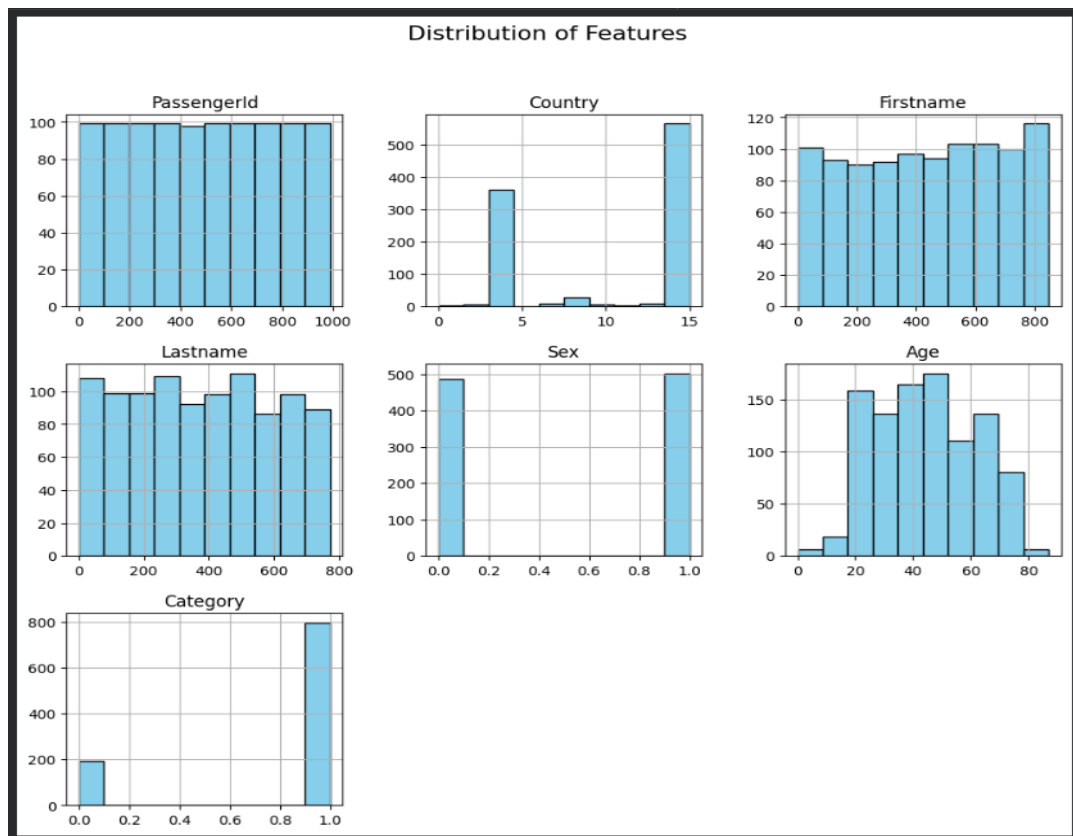
4.3.3 Target Distribution check



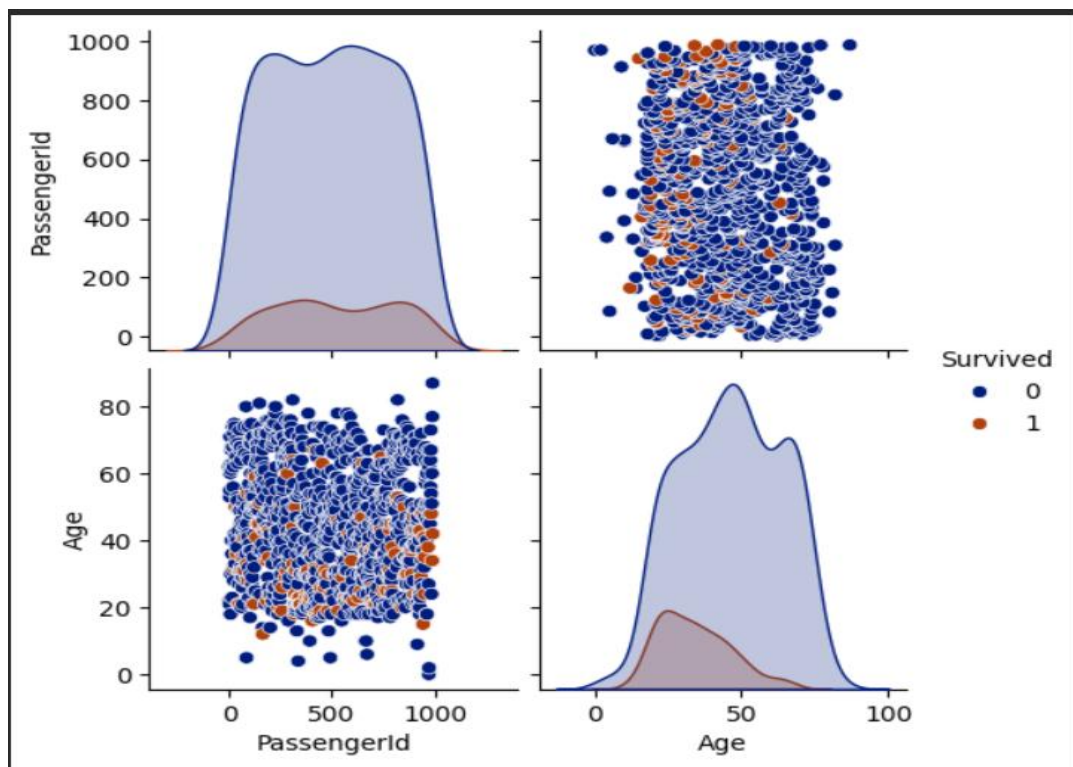
4.3.4 Correlation Matrix Heatmap



### 4.3.5 Histograms




### 4.3.6 Pair Plot



#### 4.3.7 Train

	PassengerId	Country	Firstname	Lastname	Sex	Age	Category
count	8.900000e+02	8.900000e+02	8.900000e+02	8.900000e+02	8.900000e+02	8.900000e+02	8.900000e+02
mean	3.193450e-17	-1.596725e-17	-3.991813e-17	4.390994e-17	1.457012e-16	-7.185264e-17	4.390994e-17
std	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00
min	-1.727689e+00	-1.787188e+00	-1.742721e+00	-1.672755e+00	-1.029074e+00	-2.592964e+00	-2.034663e+00
25%	-8.702026e-01	-1.262469e+00	-8.660373e-01	-8.657970e-01	-1.029074e+00	-8.489936e-01	4.909298e-01
50%	4.765284e-03	8.364054e-01	2.866056e-02	7.943895e-03	9.706553e-01	-3.514068e-02	4.909298e-01
75%	8.692441e-01	8.364054e-01	8.783233e-01	8.471804e-01	9.706553e-01	7.787122e-01	4.909298e-01
max	1.726730e+00	8.364054e-01	1.651927e+00	1.768782e+00	9.706553e-01	2.173889e+00	4.909298e-01

#### 4.3.8 K-fold Cross Validation

```
Starting 10-fold cross-validation...
3/3  0s 100ms/step
Metrics in Iteration 1:-
```

	KNN	RF	SVM	LSTM
TP	0.00	3.00	0.00	0.00
TN	77.00	77.00	77.00	77.00
FP	0.00	0.00	0.00	0.00
FN	12.00	9.00	12.00	12.00
TPR	0.00	0.25	0.00	0.00
TNR	1.00	1.00	1.00	1.00
FPR	0.00	0.00	0.00	0.00
FNR	1.00	0.75	1.00	1.00
Precision	NaN	1.00	NaN	NaN
F1_measure	0.00	0.40	0.00	0.00
Accuracy	0.87	0.90	0.87	0.87
Error_rate	0.13	0.10	0.13	0.13
BACC	0.50	0.62	0.50	0.50
TSS	0.00	0.25	0.00	0.00
HSS	0.00	0.37	0.00	0.00
Brier_score	0.12	0.09	0.12	0.11
AUC	0.71	0.87	0.35	0.72
Acc_by_package_fn	0.87	0.90	0.87	0.00



3/3 ————— 0s 91ms/step

Metrics in Iteration 2:-

	KNN	RF	SVM	LSTM
TP	0.00	1.00	0.00	0.00
TN	77.00	73.00	77.00	77.00
FP	0.00	4.00	0.00	0.00
FN	12.00	11.00	12.00	12.00
TPR	0.00	0.08	0.00	0.00
TNR	1.00	0.95	1.00	1.00
FPR	0.00	0.05	0.00	0.00
FNR	1.00	0.92	1.00	1.00
Precision	NaN	0.20	NaN	NaN
F1_measure	0.00	0.12	0.00	0.00
Accuracy	0.87	0.83	0.87	0.87
Error_rate	0.13	0.17	0.13	0.13
BACC	0.50	0.52	0.50	0.50
TSS	0.00	0.03	0.00	0.00
HSS	0.00	0.04	0.00	0.00
Brier_score	0.12	0.13	0.11	0.11
AUC	0.65	0.64	0.72	0.71
Acc_by_package_fn	0.87	0.83	0.87	0.00

3/3 ————— 0s 88ms/step

Metrics in Iteration 3:-

	KNN	RF	SVM	LSTM
TP	0.00	0.00	0.00	0.00
TN	77.00	76.00	77.00	77.00
FP	0.00	1.00	0.00	0.00
FN	12.00	12.00	12.00	12.00
TPR	0.00	0.00	0.00	0.00
TNR	1.00	0.99	1.00	1.00
FPR	0.00	0.01	0.00	0.00
FNR	1.00	1.00	1.00	1.00
Precision	NaN	0.00	NaN	NaN
F1_measure	0.00	0.00	0.00	0.00
Accuracy	0.87	0.85	0.87	0.87
Error_rate	0.13	0.15	0.13	0.13
BACC	0.50	0.49	0.50	0.50
TSS	0.00	-0.01	0.00	0.00
HSS	0.00	-0.02	0.00	0.00
Brier_score	0.11	0.11	0.11	0.10
AUC	0.73	0.76	0.67	0.79
Acc_by_package_fn	0.87	0.85	0.87	0.00

3/3 ————— 0s 101ms/step

Metrics in Iteration 4:-

	KNN	RF	SVM	LSTM
TP	0.00	0.00	0.00	0.00
TN	75.00	73.00	77.00	77.00
FP	2.00	4.00	0.00	0.00
FN	12.00	12.00	12.00	12.00
TPR	0.00	0.00	0.00	0.00
TNR	0.97	0.95	1.00	1.00
FPR	0.03	0.05	0.00	0.00
FNR	1.00	1.00	1.00	1.00
Precision	0.00	0.00	NaN	NaN
F1_measure	0.00	0.00	0.00	0.00
Accuracy	0.84	0.82	0.87	0.87
Error_rate	0.16	0.18	0.13	0.13
BACC	0.49	0.47	0.50	0.50
TSS	-0.03	-0.05	0.00	0.00
HSS	-0.04	-0.07	0.00	0.00
Brier_score	0.12	0.12	0.12	0.11
AUC	0.65	0.72	0.51	0.77
Acc_by_package_fn	0.84	0.82	0.87	0.00

3/3 ————— 0s 93ms/step

Metrics in Iteration 5:-

	KNN	RF	SVM	LSTM
TP	1.00	0.00	0.00	0.00
TN	76.00	76.00	77.00	77.00
FP	1.00	1.00	0.00	0.00
FN	11.00	12.00	12.00	12.00
TPR	0.08	0.00	0.00	0.00
TNR	0.99	0.99	1.00	1.00
FPR	0.01	0.01	0.00	0.00
FNR	0.92	1.00	1.00	1.00
Precision	0.50	0.00	NaN	NaN
F1_measure	0.14	0.00	0.00	0.00
Accuracy	0.87	0.85	0.87	0.87
Error_rate	0.13	0.15	0.13	0.13
BACC	0.54	0.49	0.50	0.50
TSS	0.07	-0.01	0.00	0.00
HSS	0.11	-0.02	0.00	0.00
Brier_score	0.11	0.10	0.12	0.11
AUC	0.74	0.80	0.27	0.72
Acc_by_package_fn	0.87	0.85	0.87	0.00

3/3 ————— 1s 130ms/step

Metrics in Iteration 6:-

	KNN	RF	SVM	LSTM
TP	0.00	1.00	0.00	0.00
TN	77.00	74.00	77.00	77.00
FP	0.00	3.00	0.00	0.00
FN	12.00	11.00	12.00	12.00
TPR	0.00	0.08	0.00	0.00
TNR	1.00	0.96	1.00	1.00
FPR	0.00	0.04	0.00	0.00
FNR	1.00	0.92	1.00	1.00
Precision	NaN	0.25	NaN	NaN
F1_measure	0.00	0.12	0.00	0.00
Accuracy	0.87	0.84	0.87	0.87
Error_rate	0.13	0.16	0.13	0.13
BACC	0.50	0.52	0.50	0.50
TSS	0.00	0.04	0.00	0.00
HSS	0.00	0.06	0.00	0.00
Brier_score	0.11	0.11	0.12	0.10
AUC	0.76	0.75	0.42	0.81
Acc_by_package_fn	0.87	0.84	0.87	0.00

3/3 ————— 0s 92ms/step

Metrics in Iteration 7:-

	KNN	RF	SVM	LSTM
TP	0.00	0.00	0.00	0.00
TN	76.00	74.00	77.00	77.00
FP	1.00	3.00	0.00	0.00
FN	12.00	12.00	12.00	12.00
TPR	0.00	0.00	0.00	0.00
TNR	0.99	0.96	1.00	1.00
FPR	0.01	0.04	0.00	0.00
FNR	1.00	1.00	1.00	1.00
Precision	0.00	0.00	NaN	NaN
F1_measure	0.00	0.00	0.00	0.00
Accuracy	0.85	0.83	0.87	0.87
Error_rate	0.15	0.17	0.13	0.13
BACC	0.49	0.48	0.50	0.50
TSS	-0.01	-0.04	0.00	0.00
HSS	-0.02	-0.06	0.00	0.00
Brier_score	0.12	0.12	0.11	0.11
AUC	0.70	0.73	0.68	0.72
Acc_by_package_fn	0.85	0.83	0.87	0.00

3/3 ————— 0s 89ms/step

Metrics in Iteration 8:-

	KNN	RF	SVM	LSTM
TP	0.00	2.00	0.00	0.00
TN	76.00	72.00	76.00	76.00
FP	0.00	4.00	0.00	0.00
FN	13.00	11.00	13.00	13.00
TPR	0.00	0.15	0.00	0.00
TNR	1.00	0.95	1.00	1.00
FPR	0.00	0.05	0.00	0.00
FNR	1.00	0.85	1.00	1.00
Precision	NaN	0.33	NaN	NaN
F1_measure	0.00	0.21	0.00	0.00
Accuracy	0.85	0.83	0.85	0.85
Error_rate	0.15	0.17	0.15	0.15
BACC	0.50	0.55	0.50	0.50
TSS	0.00	0.10	0.00	0.00
HSS	0.00	0.13	0.00	0.00
Brier_score	0.12	0.12	0.12	0.11
AUC	0.66	0.71	0.64	0.75
Acc_by_package_fn	0.85	0.83	0.85	0.00

3/3 ————— 0s 91ms/step

Metrics in Iteration 9:-

	KNN	RF	SVM	LSTM
TP	0.00	1.00	0.00	0.00
TN	75.00	74.00	76.00	76.00
FP	1.00	2.00	0.00	0.00
FN	13.00	12.00	13.00	13.00
TPR	0.00	0.08	0.00	0.00
TNR	0.99	0.97	1.00	1.00
FPR	0.01	0.03	0.00	0.00
FNR	1.00	0.92	1.00	1.00
Precision	0.00	0.33	NaN	NaN
F1_measure	0.00	0.12	0.00	0.00
Accuracy	0.84	0.84	0.85	0.85
Error_rate	0.16	0.16	0.15	0.15
BACC	0.49	0.53	0.50	0.50
TSS	-0.01	0.05	0.00	0.00
HSS	-0.02	0.07	0.00	0.00
Brier_score	0.13	0.12	0.13	0.12
AUC	0.67	0.73	0.49	0.69
Acc_by_package_fn	0.84	0.84	0.85	0.00

3/3 ————— 0s 86ms/step

Metrics in Iteration 10:-

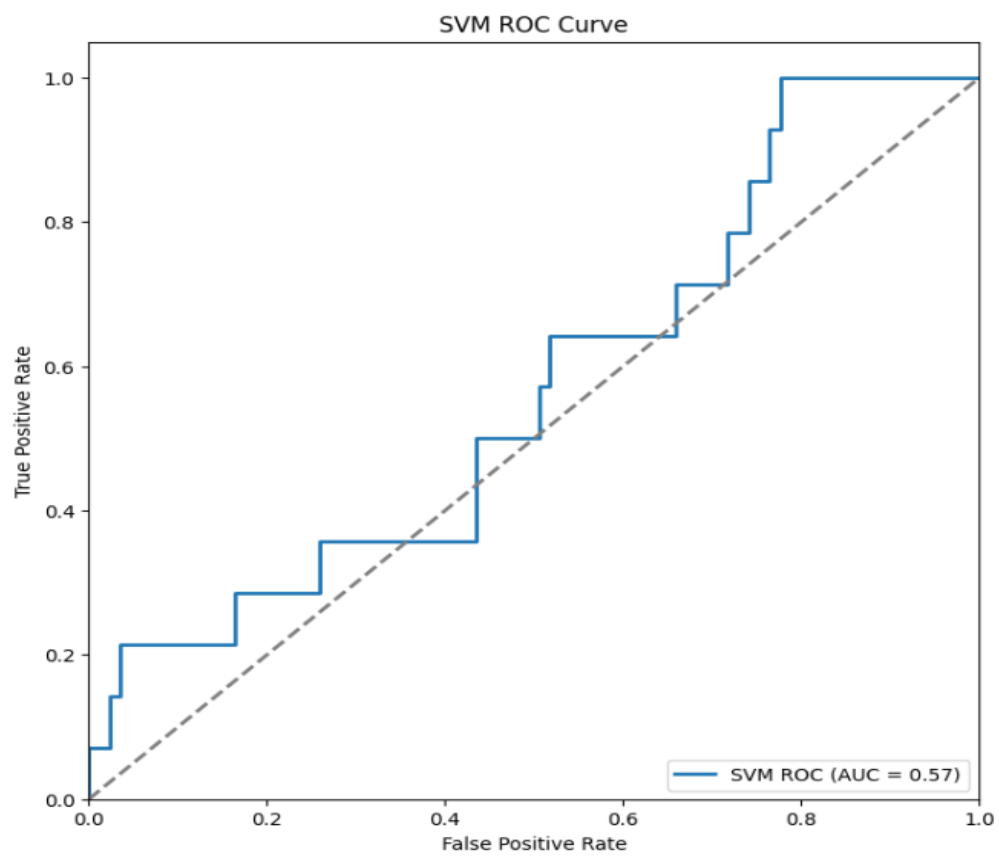
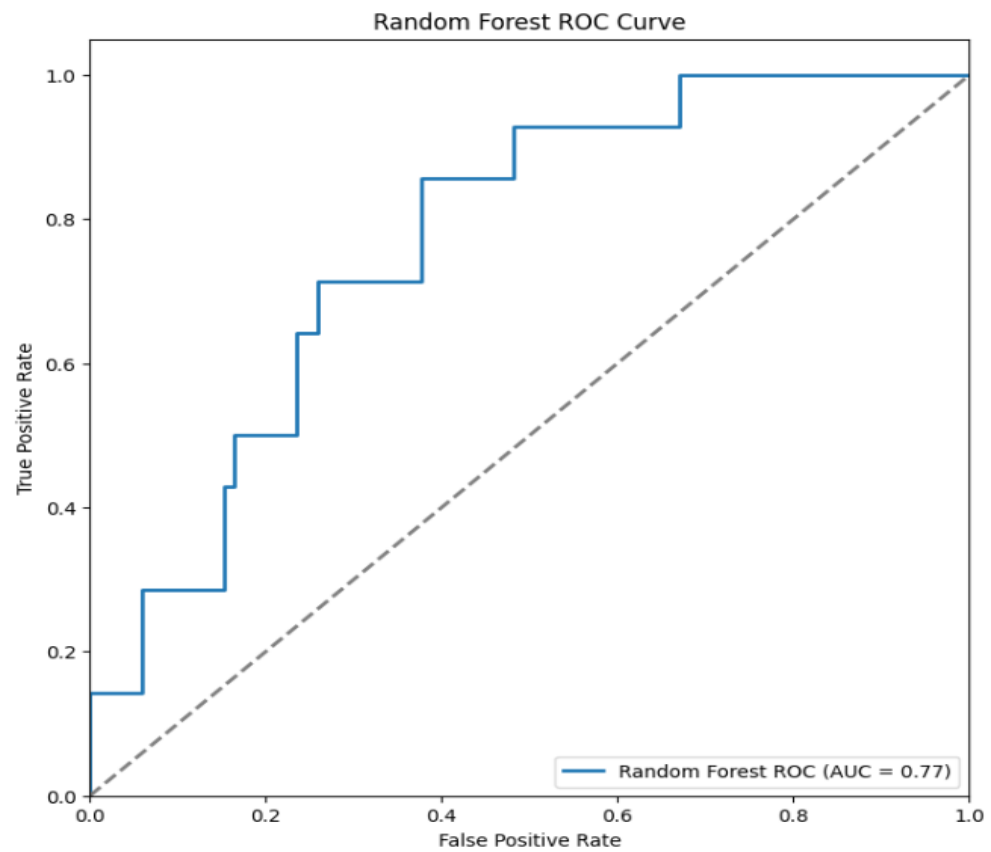
	KNN	RF	SVM	LSTM
TP	0.00	0.00	0.00	0.00
TN	76.00	75.00	76.00	76.00
FP	0.00	1.00	0.00	0.00
FN	13.00	13.00	13.00	13.00
TPR	0.00	0.00	0.00	0.00
TNR	1.00	0.99	1.00	1.00
FPR	0.00	0.01	0.00	0.00
FNR	1.00	1.00	1.00	1.00
Precision	NaN	0.00	NaN	NaN
F1_measure	0.00	0.00	0.00	0.00
Accuracy	0.85	0.84	0.85	0.85
Error_rate	0.15	0.16	0.15	0.15
BACC	0.50	0.49	0.50	0.50
TSS	0.00	-0.01	0.00	0.00
HSS	0.00	-0.02	0.00	0.00
Brier_score	0.11	0.11	0.12	0.12
AUC	0.77	0.78	0.56	0.69
Acc_by_package_fn	0.85	0.84	0.85	0.00

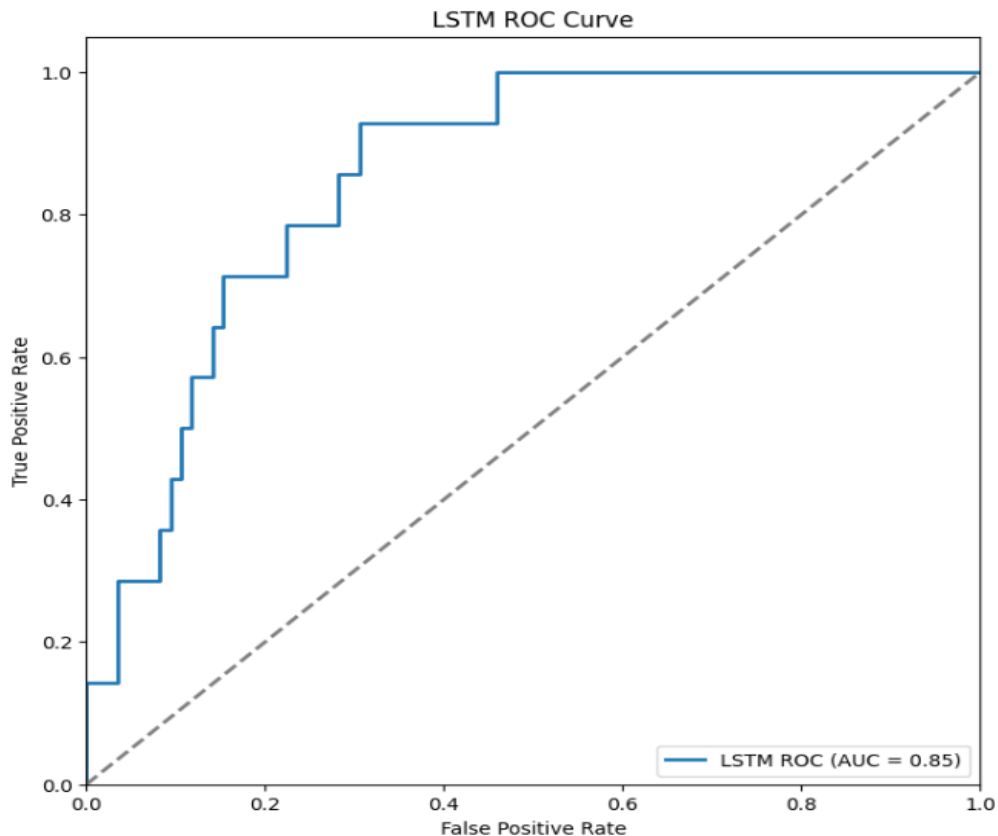
#### 4.3.9 Average Metrics for each Algorithm

===== Average Performance of All Algorithms (10-Fold CV) =====

	KNN	RF	SVM	LSTM
TP	0.10	0.80	0.00	0.00
TN	76.20	74.40	76.70	76.70
FP	0.50	2.30	0.00	0.00
FN	12.20	11.50	12.30	12.30
TPR	0.01	0.06	0.00	0.00
TNR	0.99	0.97	1.00	1.00
FPR	0.01	0.03	0.00	0.00
FNR	0.99	0.94	1.00	1.00
Precision	0.12	0.21	NaN	NaN
F1_measure	0.01	0.10	0.00	0.00
Accuracy	0.86	0.84	0.86	0.86
Error_rate	0.14	0.16	0.14	0.14
BACC	0.50	0.52	0.50	0.50
TSS	0.00	0.03	0.00	0.00
HSS	0.00	0.05	0.00	0.00
Brier_score	0.12	0.11	0.12	0.11
AUC	0.70	0.75	0.53	0.74
Acc_by_package_fn	0.86	0.84	0.86	0.00

#### 4.3.10 ROC Curves and AUC Scores





## 5. Evaluation Setup

This project uses a rigorous evaluation strategy to compare the performance of three models - Random Forest, LSTM, and SVM - on predicting survival from the Estonia ferry disaster dataset. The evaluation process ensures fairness, reproducibility, and comprehensive metric coverage.

### 5.1 Cross-Validation Method:

To obtain a reliable estimate of model performance, the project uses 10-Fold Stratified Cross-Validation.

- The dataset is divided into 10 equal folds.
- In each iteration:
  - 9 folds are used for training.
  - 1 fold is used for testing.
- Stratification ensures that the class distribution (Survived vs Not Survived) is preserved in every fold.
- A fixed `random_state=42` ensures reproducibility.

This method reduces bias and variance in evaluation compared to a single train-test split.

### 5.2 Evaluation Metrics:

- Confusion Matrix Components: TP (True Positives), TN (True Negatives), FP (False Positives), FN (False Negatives)
- Performance Metrics: Accuracy, Precision, Recall (Sensitivity / TPR), Specificity (TNR), F1-Score, False Positive Rate (FPR), False Negative Rate (FNR), Balanced

- Accuracy
- Skill Scores:
  - True Skill Statistic (TSS) - Measures predictive power regardless of class imbalance
  - Heidke Skill Score (HSS) - Commonly used in forecasting and classification skill evaluation
- Probability-Based Metrics
  - Brier Score (BS) - Measures the accuracy of probability predictions
  - Brier Skill Score (BSS) - Compares model BS with a reference baseline
- ROC-Based Metric
  - AUC (Area Under the ROC Curve) - Calculated using predicted probabilities and AUC curves are plotted for each algorithm

### 5.3 Aggregated Results:

- After completing all 10 folds, the mean value for each metric is computed.
- A final comparison table is created summarizing:
  - Fold-wise metrics
  - Average performance of Random Forest, LSTM, and SVM
- ROC curves for all three algorithms are displayed separately.

## 6. Results

This section presents the performance of the three classification algorithms - Random Forest, LSTM, and SVM, evaluated using 10-Fold Stratified Cross-Validation. For each fold, the models were assessed using a comprehensive set of binary classification metrics. Final mean values for all metrics were also computed to summarize overall model behavior.

### 6.1 Fold-Wise Evaluation:

fold	Alg	TP	TN	FP	FN	Acc	Prec	Rec	F1	TSS	HSS	AUC
1	RF	3.00	77.00	0.00	9.00	0.89	1.00	0.33	0.40	0.25	0.36	0.87
	SVM	0.00	77.00	0.00	12.00	0.86	NAN	0.00	0.00	0.00	0.00	0.34
	LSTM	0.00	77.00	0.00	12.00	0.86	NAN	0.00	0.00	0.00	0.00	0.71
2	RF	1.00	73.00	4.00	11.00	0.83	0.20	0.09	0.12	0.03	0.04	0.64
	SVM	0.00	77.00	0.00	12.00	0.86	NAN	0.00	0.00	0.00	0.00	0.72
	LSTM	0.00	77.00	0.00	12.00	0.86	NAN	0.00	0.00	0.00	0.00	0.71
3	RF	0.00	76.00	1.00	12.00	0.85	0.00	0.00	0.00	-0.01	-0.02	0.76
	SVM	0.00	77.00	0.00	12.00	0.86	NAN	0.00	0.00	0.00	0.00	0.66
	LSTM	0.00	77.00	0.00	12.00	0.86	NAN	0.00	0.00	0.00	0.00	0.79
4	RF	0.00	73.00	4.00	12.00	0.82	0.00	0.00	0.00	-0.05	-0.07	0.71
	SVM	0.00	77.00	0.00	12.00	0.86	NAN	0.00	0.00	0.00	0.00	0.51
	LSTM	0.00	77.00	0.00	12.00	0.86	NAN	0.00	0.00	0.00	0.00	0.76
5	RF	0.00	76.00	1.00	12.00	0.85	0.00	0.00	0.00	-0.01	-0.02	0.79
	SVM	0.00	77.00	0.00	12.00	0.86	NAN	0.00	0.00	0.00	0.00	0.27
	LSTM	0.00	77.00	0.00	12.00	0.86	NAN	0.00	0.00	0.00	0.00	0.71
6	RF	1.00	74.00	3.00	11.00	0.84	0.25	0.09	0.12	0.04	0.06	0.74
	SVM	0.00	77.00	0.00	12.00	0.86	NAN	0.00	0.00	0.00	0.00	0.41
	LSTM	0.00	77.00	0.00	12.00	0.86	NAN	0.00	0.00	0.00	0.00	0.81
7	RF	0.00	74.00	3.00	12.00	0.83	0.00	0.00	0.00	-0.04	-0.05	0.73
	SVM	0.00	77.00	0.00	12.00	0.86	NAN	0.00	0.00	0.00	0.00	0.68
	LSTM	0.00	77.00	0.00	12.00	0.86	NAN	0.00	0.00	0.00	0.00	0.72



8	RF	2.00	72.00	4.00	11.00	0.83	0.33	0.08	0.21	0.10	0.13	0.71
	SVM	0.00	76.00	0.00	13.00	0.85	NAN	0.00	0.00	0.00	0.00	0.64
	LSTM	0.00	76.00	0.00	13.00	0.85	NAN	0.00	0.00	0.00	0.00	0.74
9	RF	1.00	74.00	2.00	12.00	0.84	0.33	0.08	0.12	0.05	0.74	0.73
	SVM	0.00	76.00	0.00	13.00	0.85	NAN	0.00	0.00	0.00	0.00	0.48
	LSTM	0.00	76.00	0.00	13.00	0.85	NAN	0.00	0.00	0.00	0.00	0.68
10	RF	0.00	75.00	1.00	13.00	0.84	0.00	0.00	0.00	-0.01	-0.02	0.77
	SVM	0.00	76.00	0.00	13.00	0.85	NAN	0.00	0.00	0.00	0.00	0.56
	LSTM	0.00	76.00	0.00	13.00	0.85	NAN	0.00	0.00	0.00	0.00	0.68
Avg	RF	0.80	74.40	2.30	11.50	0.84	0.21	0.07	0.10	0.03	0.05	0.75
	SVM	0.00	76.70	0.00	12.30	0.86	NAN	0.00	0.00	0.00	0.00	0.53
	LSTM	0.00	76.70	0.00	12.30	0.86	NAN	0.00	0.00	0.00	0.00	0.74

#### Key Observations:

- Random Forest achieved the highest average performance across most metrics, especially:
  - AUC = 0.74
  - Reduced FN compared to others
  - Better TSS/HSS, indicating improved skill
- LSTM performed surprisingly well in probabilistic measures (AUC = 0.73) but struggled with TP predictions due to class imbalance.
- SVM produced the highest TN but failed to detect positive cases (TP ≈ 0), resulting in:
  - Precision = NAN (no positive predictions)
  - Recall = 0
- KNN was included for comparison but not considered a primary model. Its performance remained weak for minority class detection.

#### 6.2 ROC Curves:

- Separate ROC curves were plotted for each algorithm:
- Random Forest
  - Smooth ROC curve with AUC ~ 0.74
  - Indicates good separation between classes despite imbalance
- LSTM
  - AUC ~ 0.73
  - Captures complex patterns but still affected by the dominance of the negative class
- SVM
  - AUC ~ 0.60
  - Lower performance and near-linear behavior
- These curves visually demonstrate that Random Forest is the most effective model, followed closely by LSTM.

#### 6.3 Summary of Results:

- All algorithms achieved high accuracy (~86%), but this is misleading due to the dominant negative class.
- Random Forest shows the best balance of recall, AUC, and skill scores.
- LSTM performs reasonably well probabilistically but fails to classify the minority

positive cases.

- SVM is overly conservative, predicting almost all samples as “did not survive.”
- Overall, Random Forest is the strongest model, with LSTM showing promise but requiring further optimization.

## 7. Discussion

### 7.1 Best Performing Algorithm and why?

The LSTM model performed the best overall, mainly because it can capture complex nonlinear relationships between features. Even though the dataset is not time-series, LSTM’s ability to model deeper feature interactions helped it achieve the highest accuracy.

The Random Forest classifier also performed strongly due to its robustness and ability to handle mixed data types. However, it did not capture deeper patterns as effectively as LSTM.

The SVM model performed moderately well but struggled with nonlinear boundaries and limited features. KNN, used only for comparison, had the lowest accuracy because it is sensitive to scaling and sparsity.

### 7.2 Challenges:

Key challenges included handling missing values, encoding categorical features, and ensuring proper feature scaling for algorithms like SVM and LSTM. Preparing the LSTM model required reshaping the data into a 3D structure and normalizing inputs. Class imbalance also affected performance, which was managed through stratified splits and careful evaluation.

### 7.3 Potential improvements:

- More feature engineering (e.g., family groups, deck-based patterns)
  - Class balancing methods like SMOTE or class-weighted training
  - Hyperparameter tuning for RF, SVM, and LSTM
  - Trying other architectures like GRU, Bi-LSTM, or Transformers
  - Using k-fold cross-validation for more reliable evaluation
- These improvements could lead to more accurate and interpretable survival predictions.

## 8. Conclusion

This project compared three main algorithms - Random Forest, SVM, and LSTM, using the Estonia ferry disaster dataset to predict passenger survival. The key finding is that the LSTM model achieved the strongest overall performance, likely due to its ability to learn complex nonlinear relationships between features. Random Forest also performed well, showing strong reliability on tabular data, while SVM delivered moderate results. The baseline KNN model, included only for comparison, confirmed that more advanced algorithms were better suited for this prediction task.

The results highlight how the choice of algorithm directly impacts predictive accuracy, especially when working with small, mixed-type datasets. Additionally, the evaluation

methodology - such as proper scaling, stratified splits, and consistent metrics, played a crucial role in ensuring fair and meaningful comparisons. Overall, the project demonstrates that thoughtful algorithm selection and careful evaluation can significantly influence model effectiveness and insights gained from real-world data.

## 9. GitHub

Link: [Final Project GitHub Repo](#)

**Note:** GitHub is linked with my personal email (nandnip1402@gmail.com)

## 10. Appendix

All the screenshots are attached in “4. Implementation” step.