

Titanic Survival Prediction

Abstract: About a century ago, one memorable night in April 1912, a world-shattering event happened. The Titanic, the 2,240-passenger luxury cruise ship, sank forever off the coast of Newfoundland in the North Atlantic after extensive damage to its hull by an iceberg on its maiden voyage. Only 705 people survived this disaster. Although nearly a century has passed, the research on Titanic has never stopped, and there are still many studies on it. This study was supposed to predict the survival of passengers on Titanic using different methods based on data from the Kaggle competition "Titanic: Machine Learning from Disaster." It predicted each passenger in the test set who would survive the sinking. The result was the percentage of correct prediction. In the Machine Learning study, the task is to achieve 80% accuracy in predicting the survival distribution of the Titanic disaster based on the demographic data testing notebook by different algorithms models. Using classification is the main point to calculate the efficiency achieved by those models through the test environment. The f-measurement scores obtained from the machine learning technology were in comparison with the f-measurement scores obtained by Kaggle.

Keywords: Machine Learning, Kaggle, Method, Data Analysis, Model Evaluation Introduction.

1. Introduction

On April 14, 1912, at 11 p.m., the Titanic experienced a horrific disaster, and nearly full of its parts got damaged by the catastrophe [1]. Sadly, there were not enough lifeboats to rescue all 2,224 passengers on board [2]. By introducing survivors and officially collected data, such as ferry tickets and corresponding numbers of people records, it can be known that in this disaster, the Age and gender distribution of the dead were uneven [3]. Data research on the disaster has continued to this day. The approach is a death distribution budget based on public data on Kaggle. Using different models will calculate the efficiency achieved through the test environment provided by Kaggle [2].

The requirement for this research is to build a series of machine learning models with an fmeasurement accuracy greater than 80 percent on a given set of demographic information on Kaggle. Predictive modeling is an essential component of predictive analysis [4]. Through the mathematical process, the input data can predict future results. The different attributes of people, such as gender, Age, the category they belong to, and their social class, build the predictive models by providing a good data analysis database [5]. In machine learning, the data has two parts, one is for training, and the other is for testing. The training models are used to fit machine learning models, and the testing models are used to evaluate equipped machine learning models. The train-test split procedure estimates the performance of a machine learning algorithm when it makes predictions based on data that cannot use to train the model [6]. To reach the aim of this research project, the use the construction of Scikit-Learn and related machine learning algorithms to perform data analysis, build models and evaluate algorithm performance on data used to train models in Kaggle [2].

2. Methodology

2.1. Data

There are two groups in the historical data, one is the training set, and the other is the testing set. For the training set, using the testing dataset, which Kaggle provides, can determine whether the passenger survived to build the model for generating prediction patterns.

Loading data from a package requires importing the data into a library and executing instructions. The original library covers a lot of related data, so there is little data to show. After loading the test data content, the study can approach categorizing and recording the data by getting the first ten rows of the dataset to [7]. Table 1 shows the training dataset.

Table 1. Training dataset

	Passengers	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund,Mr.Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cummings,Mrs.JohnBrady(Florence Briggs)	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen,Miss.Laina	female	26.0	0	0	O2.3101282STON/	7.9250	NaN	S
3	4	1	1	Futrelle,Mrs.Jacques Heath(Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen,Mrs.William Henry	male	35.0	0	0	373450	8.0500	NaN	S
5	6	0	3	Moran,Mr.James	male	NaN	0	0	330877	8.4583	NaN	Q
6	7	0	1	McCarthy,Mr.Timothy J	male	54.0	0	0	17463	51.8625	E46	S
7	8	0	3	Palsson,Master.Gosta Leonard	male	2.0	3	1	349909	21.0750	NaN	S
8	9	1	3	Johnson,Mrs.Oscar	female	27.0	0	2	347742	11.1333	NaN	S
9	10	1	2	Nasser,Mrs.Nichola	female	14.0	1	0	237736	30.0708	NaN	C

It is helpful to record each row of passengers on board after receiving the dataset's count of rows and columns. On the other hand, these columns are considered the distribution data for each commuter. Theoretically, the dataset has 891 passengers in the rows and 15 data points in the columns [8]. Mean value, count, standard deviation, and other statistics can also display on the dataset. Subsequently, passengers are classified and analyzed in different categories.

2.2. Section snippets

For this part, an essential process of building machine learning can be implemented, as Figure 1 shown, to fit the collected training data into a data model that can be tested. Under this data model, more accurate data will be determined to compare and analyze the obtained results using the different models. These are essential steps that need to be performed in performing classification and regression analysis, which are summarized in Figure 1.

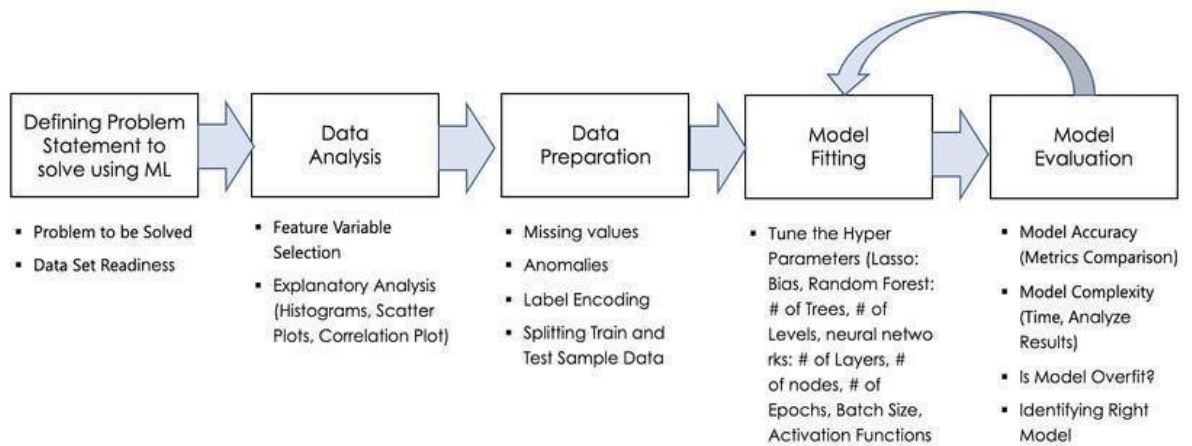


Figure 1. The process of ML

2.3 . Data cleaning

The purpose of data cleaning is to prevent similar errors in data collection. Data cleanup can be divided into these steps identifying the mistakes and either correcting or deleting the data that needs to be modified or manually processing the data, which is an essential part of machine learning to build models. Correct data deletion can save us time and cost and improve data efficiency, and the previous dataset after cleaning is summarized in Table.2.

Table 2. dataset after cleaning

	Survived	Pclass	Sex	Age	Fare	Embarked	IsAlone	Title
0	0	3	1	1	0.0	0	0	1
1	1	1	0	2	3.0	1	0	3
2	1	3	0	1	1.0	0	1	2
3	1	1	0	2	3.0	0	0	3
4	0	3	1	2	1.0	0	1	1
5	0	3	1	1	1.0	2	1	1
6	0	1	1	3	3.0	0	1	1
7	0	3	1	0	2.0	0	0	4

8	1	3	0	1	1.0	0	0	3
9	1	2	0	0	2.0	1	0	3

2.4. Titanic prediction—artificial neural network

An ANN neural network is an extensive parallel interconnected network of simple adaptive units whose organization can simulate the interactive responses of biological nervous systems to real-world objects [9]. The ANN model initially learns to recognize patterns in the data, visual or auditory, through a training phase where it compares the actual output with the original intent to produce the desired outcome.

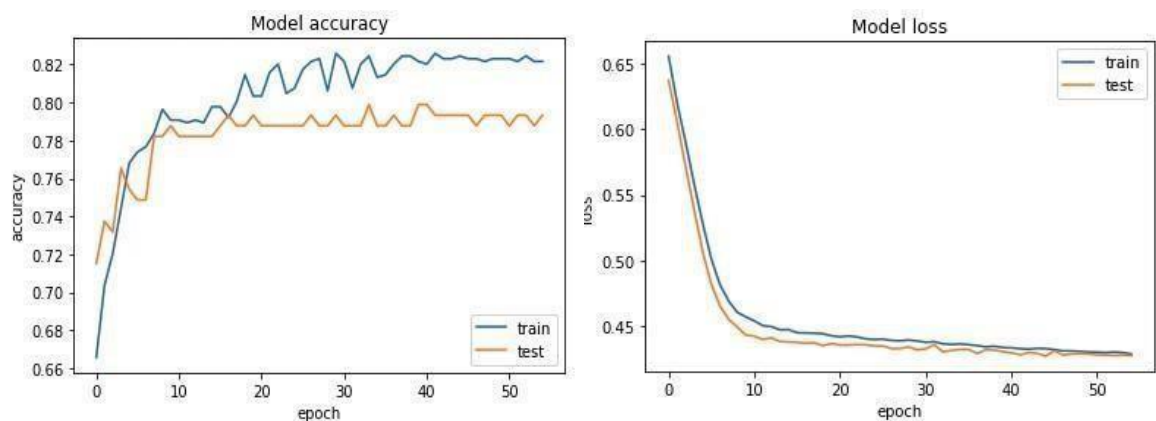


Figure 2. ANN validation

As shown in Figure 2, the ANN model can predict Titanic survivors with an accuracy of 82.16% and an error of 0.45. In addition, the model shows that the factors that impact Titanic survivors most are gender, P-class, and cabin and the different variable names, as shown in Table.3 [7]. **Table**

3. Different variable names

Variable	Definition	key
Survival	Survival	0=No,1=Yes
Pclass	Ticket class	1=1st,2=2nd,3=3rd
Sex	Sex	Male and Female
Age	Age in years	-
Sibsp	#of siblings/spouses aboard the Titanic	-
Parch	#of Parents/Children aboard the Titanic	-

Ticket	Ticket number	-
Fare	Passenger fare	-
Cabin	Cabin number	-
Embarked	Port of Embarkation	C=Cherbourg,Q=Queenstown,S=Southampton

2.5. Exploratory data analysis

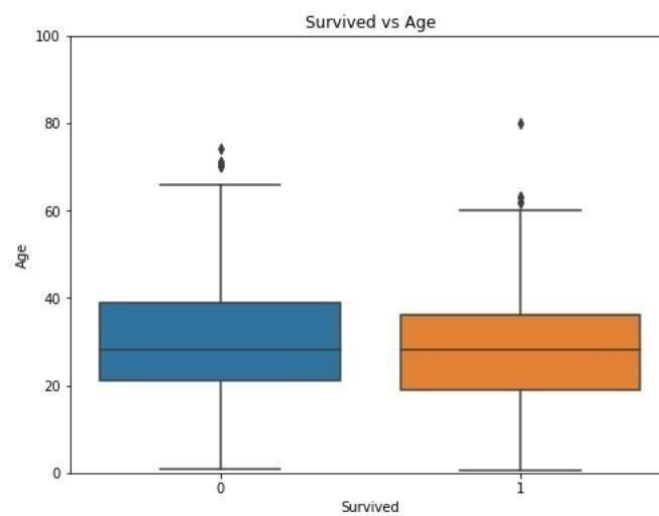


Figure 3. Survived vs. Age

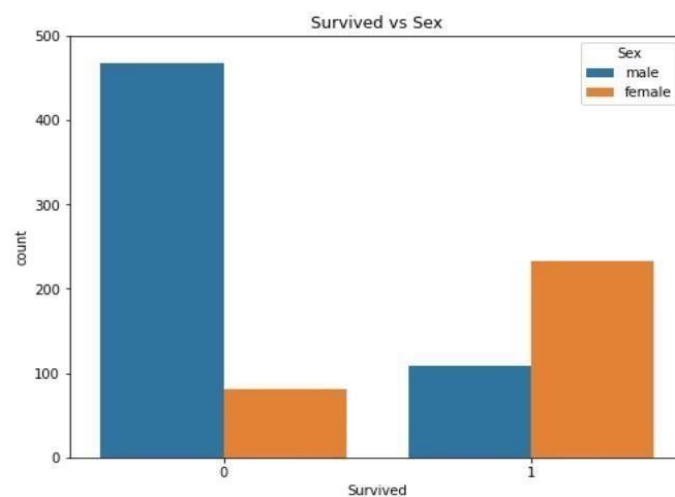


Figure 4. Survived vs. Sex

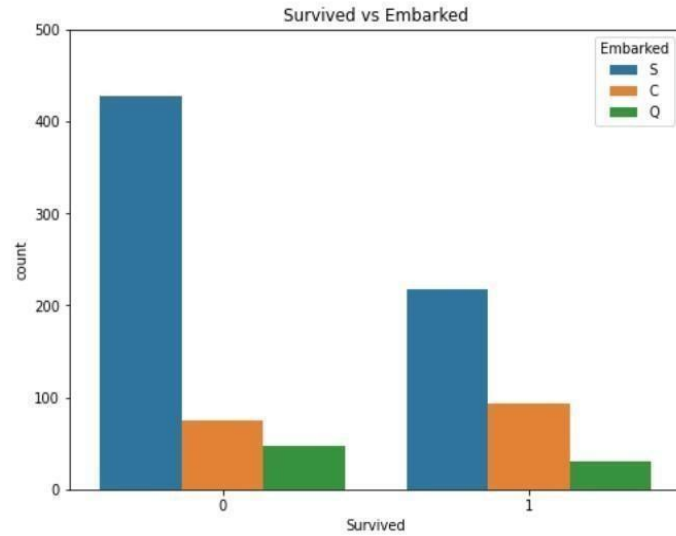


Figure 5. Survived vs. Embarked

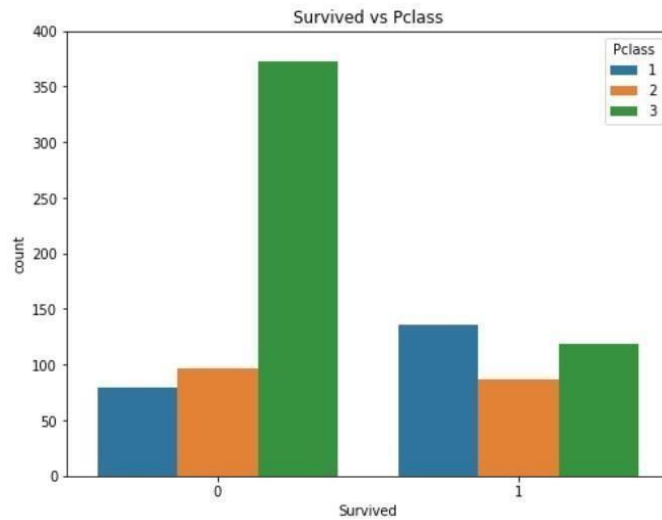


Figure 6. Survived vs.P-class

The exploratory data analysis method is a method that emphasizes data visualization, focusing on the proper distribution of data [9]. The test data is imported here to obtain the survival distribution and various features. The observed rules are analyzed by insights to get inspiration to find a model suitable for the data better. Here, Age, gender, embarked, and P-class can be seen in Figure 3, Figure 4, Figure 5, and Figure 6, as important survival research features, becoming the most effective data indicators that can be speculated in this study.

2.6. Insights

Insights are the process of finding relations or trends considering multiple features. After the Titanic dataset had been analysed, two obtained results could be determined from this time. The first is to see what the survival passengers have in common that will help them survive the shipwreck, and the second is to know if the passengers boarded this fateful ship by applying machine learning tools. For Age, a few points are located outside the whiskers of the box plot. However, since it is still within a reasonable age range (around 62 - 80), it will keep the data to see what can be determined from this section. Females were more likely to survive than Males. Upper-class passengers (P Class = 1) were more likely to survive the crash Passengers embarked from S were less likely to survive. (Possibility of a correlation between P-class and embarked)

3. Objectives

The Titanic survival prediction project aims to achieve several key objectives, which can be outlined as follows:

3.1. Data Exploration and Preprocessing:

- To conduct a thorough exploration of the Titanic dataset, identifying key features and understanding the structure of the data.
- To preprocess the data by handling missing values, encoding categorical variables, and normalizing numerical features to prepare it for analysis.

3.2. Identifying Key Predictors of Survival:

- To analyse the relationships between various passenger characteristics (such as age, gender, class, and fare) and survival outcomes.
- To identify significant predictors that influence the likelihood of survival during the Titanic disaster.

3.3. Model Development and Evaluation:

- To develop multiple machine learning models (e.g., Logistic Regression, Decision Trees, Random Forest) to predict passenger survival based on the identified features.
- To evaluate the performance of these models using appropriate metrics (such as accuracy, precision, recall, and F1-score) and select the best-performing model for prediction.

3.4. Understanding Social Dynamics:

- To explore how social factors, such as class and gender, impacted survival rates and to provide insights into the social dynamics of the time.
- To visualize and communicate findings related to survival disparities among different demographic groups.

3.5. Enhancing Predictive Accuracy:

- To refine the predictive models through techniques such as feature selection, hyperparameter tuning, and cross-validation to enhance their accuracy and robustness.
- To assess the impact of different modeling techniques on prediction outcomes and identify the most effective approach.

3.6. Providing Actionable Insights:

- To generate actionable insights that can inform modern emergency response strategies and policies based on historical data analysis.
- To contribute to discussions on social equity and justice by highlighting the factors that influenced survival during the Titanic disaster.

3.7. Educational and Research Contribution:

- To serve as an educational resource for students and researchers interested in data science, machine learning, and historical analysis.
- To contribute to the broader field of data analytics by sharing findings, methodologies, and best practices for survival prediction.

3.8. Promoting Ethical Considerations:

- To engage in discussions about the ethical implications of data analysis and survival prediction, emphasizing the importance of responsible data use and the potential impact of findings on societal perceptions.

By achieving these objectives, the Titanic survival prediction project aims to provide a comprehensive analysis of the factors influencing survival during the disaster, demonstrate the application of machine learning techniques, and contribute to ongoing discussions about social equity and data-driven decision-making.

4. Project Plan

4.1. Project Setup:

- Create a Google Colab notebook.
- Import necessary libraries: Pandas, NumPy, Scikit-Learn, Matplotlib, Seaborn.

4.2. Data Collection:

- Download the Titanic dataset from Kaggle or other reliable sources.

4.3. Data Preprocessing:

- Clean the dataset by:
- Handling missing values.
- Encoding categorical variables (e.g., Gender, Embarked). □
- Normalizing numerical features (e.g., Age, Fare).

4.4. Exploratory Data Analysis (EDA):

- Visualize data using:
- Histograms and box plots for distributions.
- Bar charts for categorical features.
- Correlation heatmaps to identify relationships.

4.5. Feature Selection:

- Identify key features influencing survival rates:
- Passenger Class (Pclass).
- Gender.
- Age.
- Siblings/Spouses Aboard (SibSp). ▯ Parents/Children Aboard (Parch).
- Fare.

- Embarked Port (C, Q, S).

4.6. Model Selection:

- Choose machine learning algorithms:
- Logistic Regression.
- Decision Trees.
- Random Forest.
- Support Vector Machines (optional).

4.7. Model Training:

- Split the dataset into training and testing sets (e.g., 80/20 split). □
Train the selected models on the training set.

4.8. Model Evaluation:

- Assess model performance using:
- Accuracy.
- Precision.
- Recall.
- F1 Score.
- Utilize cross-validation for reliability.

4.9. Results Visualization:

- Create visualizations to present findings:
- Confusion matrices.
- ROC curves.
- Feature importance plots.

4.10. Documentation:

- Compile a comprehensive report detailing:
- Methodology.
- Results.
- Insights and conclusions drawn from the analysis.

4.11. Future Work:

- Suggest potential improvements:
- Incorporating additional features.
- Exploring advanced algorithms (e.g., ensemble methods). □
Conducting further analyses on different datasets.

5. Machine Learning Models

5.1. Support vector machines

A support vector machine (SVM) is the basic model of the system, which is a linear classifier with the most considerable interval defined in the feature space. It is a generalized linear machine that supervises binary classification. It solves the classification problem of two groups of categories by using classification algorithms. After providing a set of SVM model sets marked with training data for each class, they can classify the new text.

From the results of EDA, the SVM method of this project is still a problem to be improved. Before EDA, it was concluded that the survival rate of women and the upper class was higher than other classes. Therefore, eliminate unnecessary data to retrain the SVM model. Finally, 82.82% of the results can be determined by this improvement.

5.2. K-nearest neighbors

The k-nearest neighbor algorithm is a data classification method often used in machine learning. It compares the groups to which the nearest data points belong in the feature space based on their distance and estimates the possibility that the data points become members of one group or another. The test environment divides the data into vertical and horizontal axes. Through this model construction, the result is 79.79%.

5.3. Logistic regression

Logical regression can be regarded as the process of modeling the probability of discrete results under the condition of given input variables. This study uses the most common binary logistic regression model, which can take a dependent variable with two possible values for analysis and comparison to obtain an answer relatively close to the result. The accuracy of the model data in the test is 79.12%

5.4. Principal component analysis

PCA is a multivariate statistical method for observing correlations among multiple variables. It is suitable for reducing data without losing its properties. This method can describe the composition of variance and covariance for several linear combinations of the main variables and retain essential parts of the original information that may be lost [10].

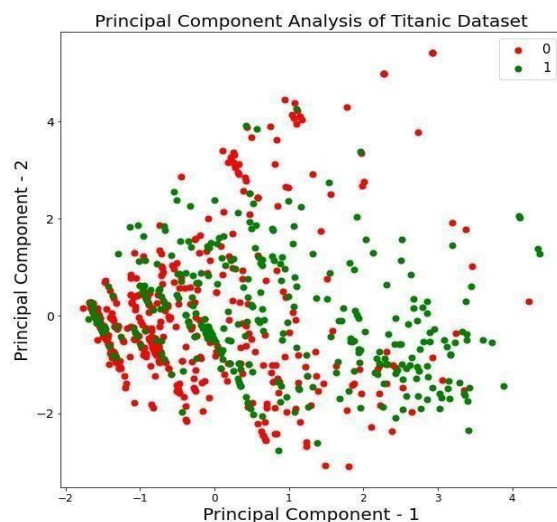


Figure 7. PCA Dataset

The PCA method converts feature with high correlation and equivalent information into a single part. The correlation matrix gives a high correlation of 54%, as shown in Figure 7, between Fare and P-class, also closely related depending on the question. Check the correlation between Embarked and P-class from PCA, as shown in Figure 8. The results of this problem can be replaced by other methods, proving that in this learning process, PCA is unsuitable for this Experimental modeling.

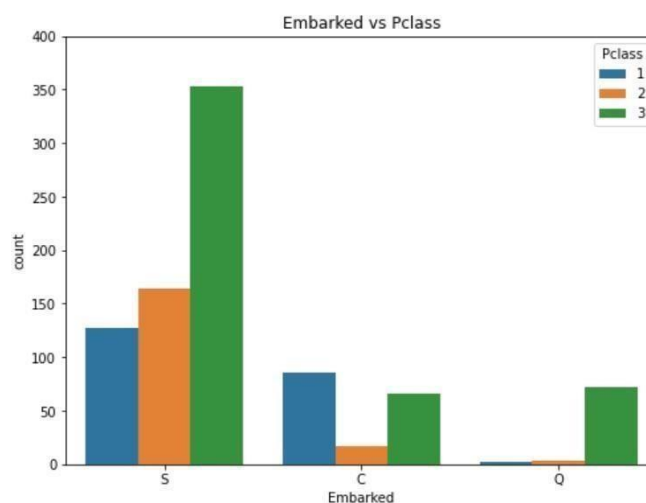


Figure 8. Checks the correlation between Embarked & P-class from PCA.

5.5. Featured engineering

The featured engineering can be a core site to predict the result by changing data to create new variables. To select features that are implemented during training and, therefore, during prediction, domain awareness can be used in feature engineering to acknowledge data set aspects that help create machine learning models. Modeling helps explain data collection, and incorrect feature selection may lead to inaccurate and shocking prediction models. Reliability and predictability depend on the precise nature of the workpiece. It filters out unused and unimportant functions. The following characteristics of the exploratory study were completed using a person's Age, gender, title, class number, P-class, boarding place, ticket price, and several families. The survival column was selected as the reaction column. These characteristics were chosen because of their ideal impact on survival. If incorrect parts are selected, inaccurate predictions may occur. Therefore, functional engineering is the backbone of creating effective forecasting models.

As seen below, some data was lost during the test, and some columns have missing values. Display missing function shows the count of missing values in every column in the training and test set.

The "Age," "Cabin," and "Boarding" columns of the training set are missing values, while the "Age," "Cabin," and "Fare" columns of the test set are not marked with relevant matters [11]. This process requires feature engineering to show some unnecessary columns to improve this problem. Then it would help to view each column's value names and counts to make a more accurate chart. Nonredundant rows and columns and their rows and missing values should be deleted for data disassembly to improve data accuracy.

6. Input / Output Requirements

6.1.Input Requirements:

6.1.1. User Registration Details

6.1.1.1. Name, Email ID, Password

6.1.2. Login Credentials

6.1.2.1. Email ID and Password

6.2.Output Requirements:

6.2.1. User Account Information

6.2.1.1. Registration confirmation, Login success/failure messages

6.2.2. Titanic Dashboard

6.2.2.1. Overview of prediction, probability

6.2.3. Error and Validation Messages

6.2.3.1. Invalid input, incorrect info, etc.

7. Screenshots

7.1.Login and Signup Page

The image displays two side-by-side forms for user authentication on a purple background. The left form is titled 'Login' and contains fields for 'Username' and 'Password', a 'Forgot password ?' link, a 'Login' button, a 'Don't have an account ? Signup' link, and a 'Need help ?' link. The right form is titled 'Signup' and contains fields for 'Username', 'Email Id', and 'Create password', a 'By creating an account, I agree to Terms and Conditions' statement, a 'Create Account' button, an 'Already have an account ? Login' link, and a 'Need help ?' link.

Login

Username

Password

[Forgot password ?](#)

Login

Don't have an account ? [Signup](#)

[Need help ?](#)

Signup

Username

Email Id

Create password

By creating an account, I agree to [Terms and Conditions](#)

Create Account

Already have an account ? [Login](#)

[Need help ?](#)

7.2.Titanic Web App

← → ↻ 797a-34-91-142-22.ngrok-free.app 🔍 ☆ 👤 ⋮

⋮

Passenger Details

Passenger Class ⓘ
1 ▾


Age ⓘ
0 30 100

Sex
female ▾

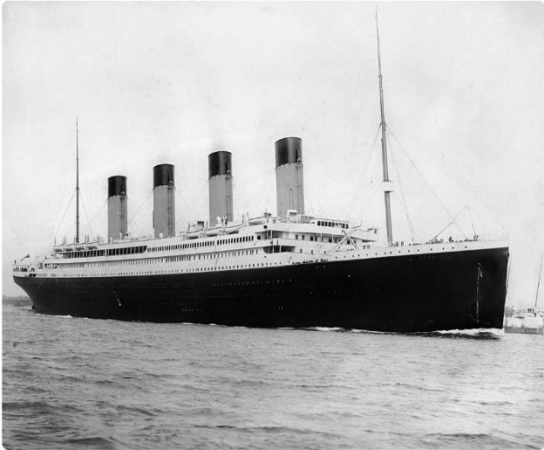
Fare (£) ⓘ
0 30 200

Port of Embarkation ⓘ
Southampton ▾

Predict Survival

 **Titanic Survival Predictor**

Predict whether a passenger would have survived the Titanic disaster


RMS Titanic

7.3. Passenger Survived Details

Passenger Details

Passenger Class ?
1 ▼

Age
14
0 100

Sex
female ▼

Fare (£) ?
117
0 200

Port of Embarkation ?
Southampton ▼

Predict Survival

7.4. Passenger Survived Prediction

Prediction Result

Based on the passenger details, this person would have:

SURVIVED

Survival probability: 63.0%



7.5. Passenger Didn't Survived Details

Passenger Details

Passenger Class ?
3 ▼

Age
64
0 100

Sex
male ▼

Fare (£) ?
15
0 200

Port of Embarkation ?
Queenstown ▼

Predict Survival

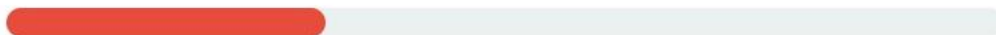
7.6. Passenger Didn't Survived Prediction

Prediction Result

Based on the passenger details, this person would have:

DID NOT SURVIVE

Survival probability: 32.0%



8. Coding

8.1.Login Signup Page Coding

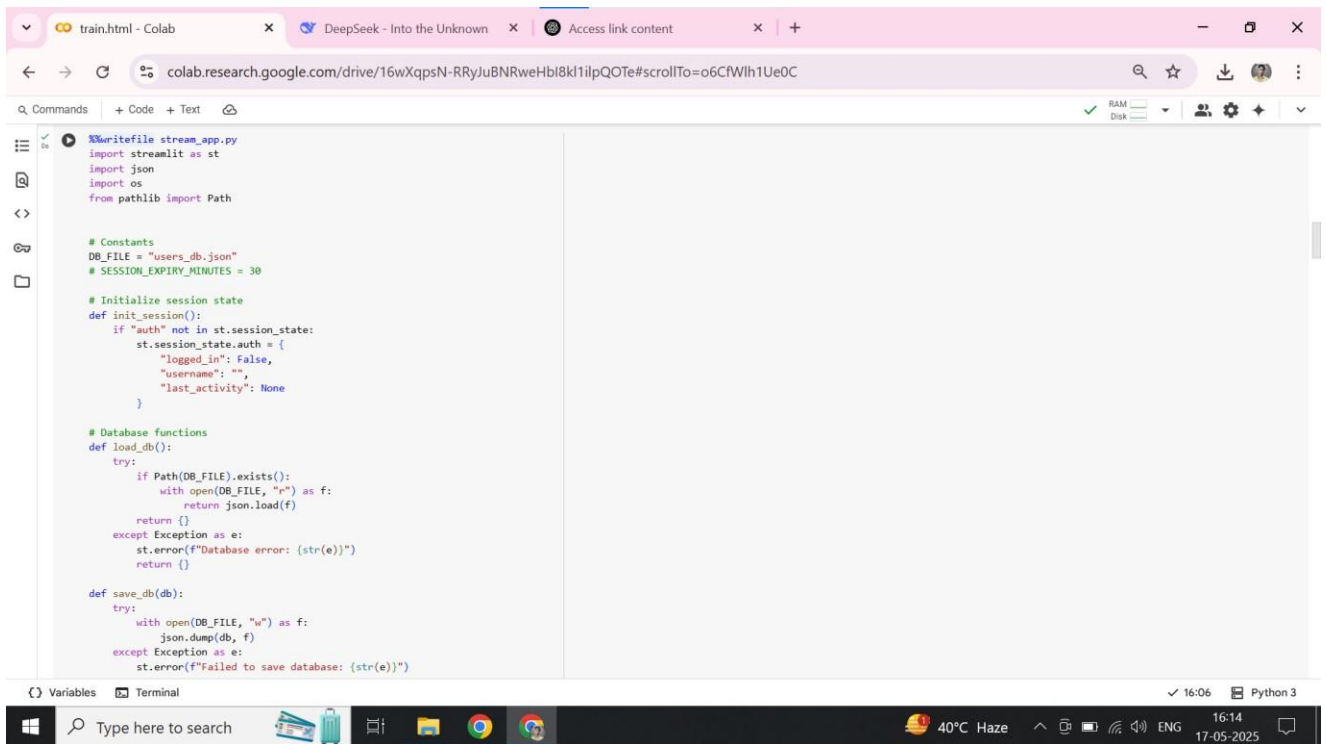
```
def login_form():
    with st.form("login_form"):
        st.subheader("Login")
        username = st.text_input("Username", key="login_username")
        password = st.text_input("Password", type="password", key="login_password")

        if st.form_submit_button("Login"):
            users_db = load_db()
            if login_user(username, password):
                st.success("Login successful!")
                st.rerun()
            else:
                st.error("Invalid username or password")

def signup_form():
    with st.form("signup_form"):
        st.subheader("Create Account")
        new_username = st.text_input("Choose Username", key="signup_username")
        new_password = st.text_input("Choose Password", type="password", key="signup_password")
        confirm_password = st.text_input("Confirm Password", type="password", key="confirm_password")

        if st.form_submit_button("Sign Up"):
            success, message = register_user(new_username, new_password, confirm_password)
            if success:
                st.success("Account created successfully! Please login.")
            else:
                st.error(message)
```

8.2.Database Setup



The screenshot shows a Google Colab notebook with the following code:

```
%%writefile stream_app.py
import streamlit as st
import json
import os
from pathlib import Path

# Constants
DB_FILE = "users_db.json"
SESSION_EXPIRY_MINUTES = 30

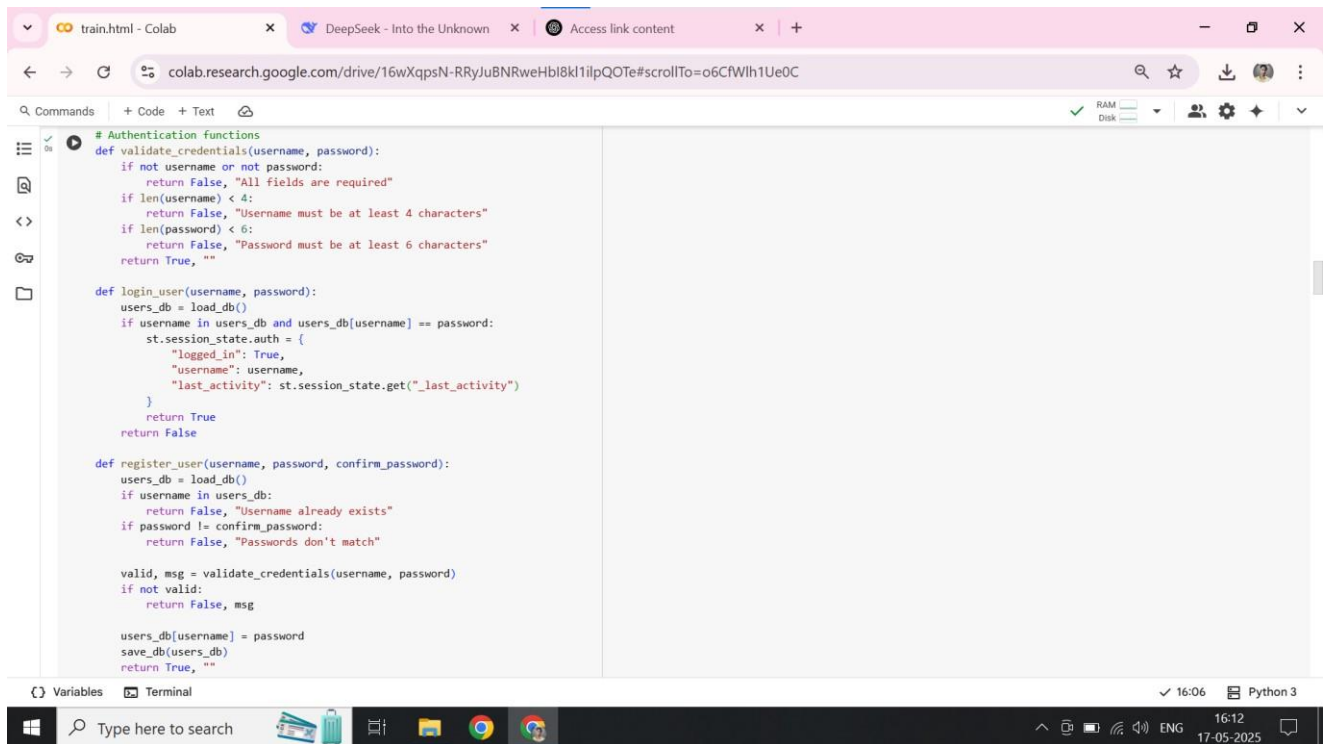
# Initialize session state
def init_session():
    if "auth" not in st.session_state:
        st.session_state.auth = {
            "logged_in": False,
            "username": "",
            "last_activity": None
        }

# Database functions
def load_db():
    try:
        if Path(DB_FILE).exists():
            with open(DB_FILE, "r") as f:
                return json.load(f)
        return {}
    except Exception as e:
        st.error(f"Database error: {str(e)}")
        return {}

def save_db(db):
    try:
        with open(DB_FILE, "w") as f:
            json.dump(db, f)
    except Exception as e:
        st.error(f"Failed to save database: {str(e)}")
```

The notebook interface includes a left sidebar with icons for file explorer, search, and other tools. The bottom status bar shows the time as 16:06, the language as Python 3, and the date as 17-05-2025.

8.3.Database Authentication



The screenshot shows a Google Colab notebook with the following Python code:

```
# Authentication functions
def validate_credentials(username, password):
    if not username or not password:
        return False, "All fields are required"
    if len(username) < 4:
        return False, "Username must be at least 4 characters"
    if len(password) < 6:
        return False, "Password must be at least 6 characters"
    return True, ""

def login_user(username, password):
    users_db = load_db()
    if username in users_db and users_db[username] == password:
        st.session_state.auth = {
            "logged_in": True,
            "username": username,
            "last_activity": st.session_state.get("_last_activity")
        }
        return True
    return False

def register_user(username, password, confirm_password):
    users_db = load_db()
    if username in users_db:
        return False, "Username already exists"
    if password != confirm_password:
        return False, "Passwords don't match"

    valid, msg = validate_credentials(username, password)
    if not valid:
        return False, msg

    users_db[username] = password
    save_db(users_db)
    return True, ""
```

The interface includes a top bar with tabs for 'train.html - Colab', 'DeepSeek - Into the Unknown', and 'Access link content'. The address bar shows a Google Drive link. The bottom status bar indicates '16:06 Python 3' and the date '17-05-2025'.

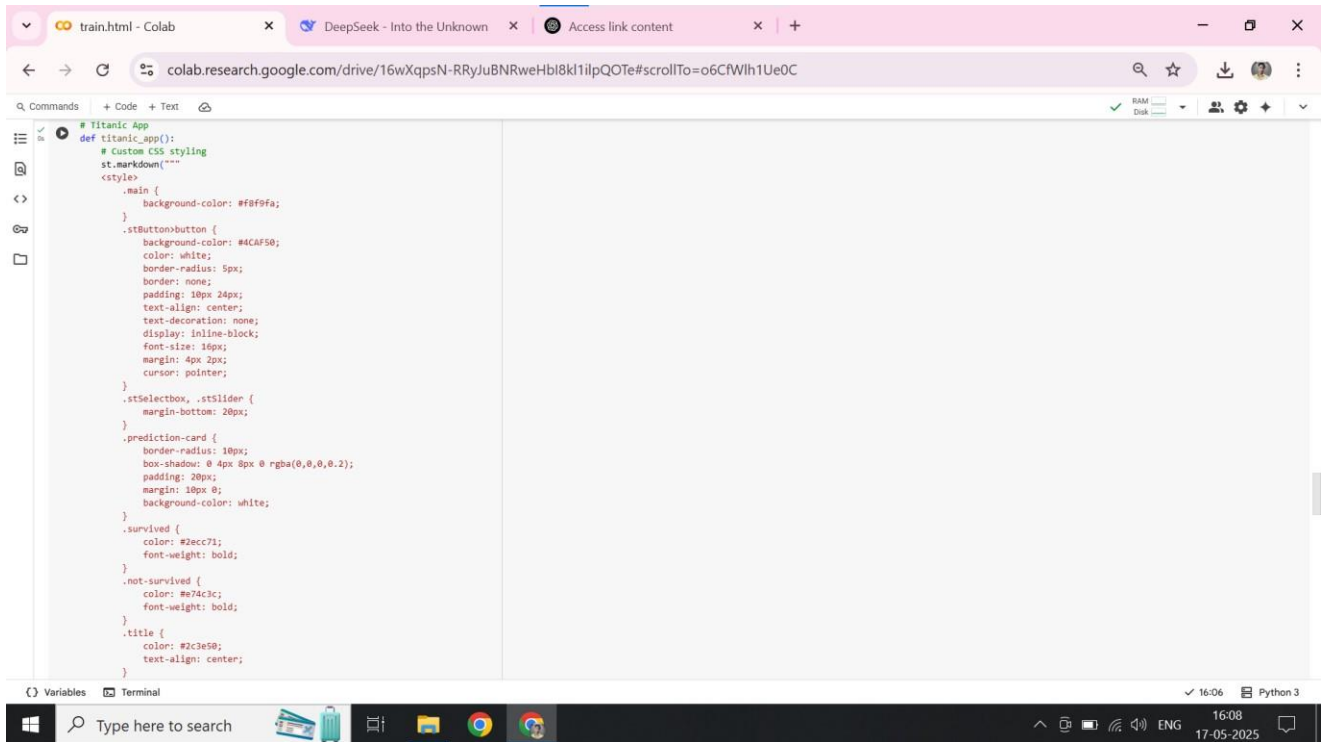
8.4.User Details

```
[168] import json

def view_users():
    # Load the database from the file
    if os.path.exists(DB_FILE):
        with open(DB_FILE, "r") as f:
            users_db = json.load(f)

        # Print the details of all users
        for username, password in users_db.items():
            print(f"Username: {username}, Password: {password}")
    else:
        print("No users database found.")
```

8.5.Titanic Web App Styling

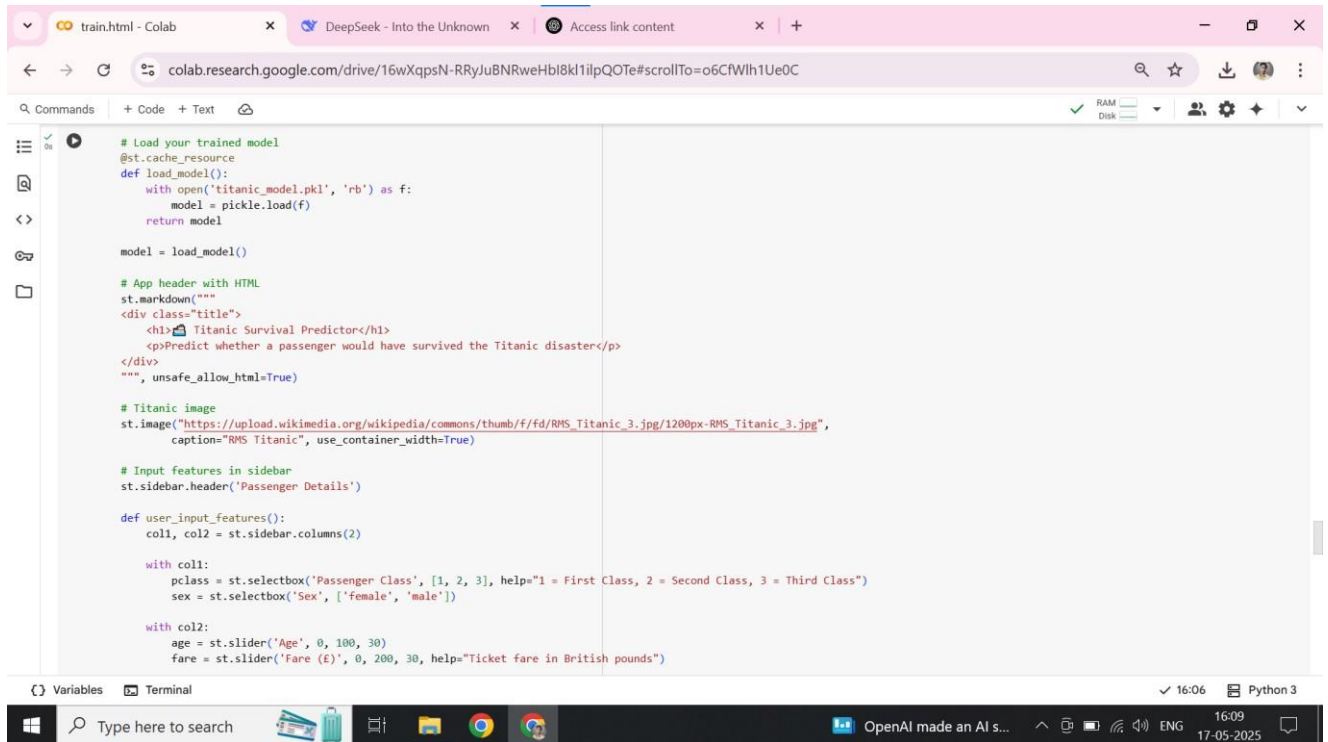


The screenshot shows a Google Colab notebook interface. The browser tabs at the top include 'train.html - Colab', 'DeepSeek - Into the Unknown', and 'Access link content'. The address bar shows a Google Drive link. The notebook's left sidebar contains icons for file management and a search bar. The main code editor displays the following CSS and JavaScript code:

```
# Titanic App
def titanic_app():
    # Custom CSS styling
    st.markdown("""
    <style>
    .main {
        background-color: #f8f9fa;
    }
    .stButton>button {
        background-color: #4CAF50;
        color: white;
        border-radius: 5px;
        border: none;
        padding: 10px 24px;
        text-align: center;
        text-decoration: none;
        display: inline-block;
        font-size: 16px;
        margin: 4px 2px;
        cursor: pointer;
    }
    .stSelectbox, .stSlider {
        margin-bottom: 20px;
    }
    .prediction-card {
        border-radius: 10px;
        box-shadow: 0 4px 8px 0 rgba(0,0,0,0.2);
        padding: 20px;
        margin: 10px 0;
        background-color: white;
    }
    .survived {
        color: #2ecc71;
        font-weight: bold;
    }
    .not-survived {
        color: #e74c3c;
        font-weight: bold;
    }
    .title {
        color: #2c3e50;
        text-align: center;
    }
    """)
```

At the bottom of the Colab interface, there are tabs for 'Variables' and 'Terminal'. The Windows taskbar at the very bottom shows the search bar, taskbar icons, and system tray information including the time (16:08) and date (17-05-2025).

8.6.Titanic Input Features



The screenshot shows a Google Colab notebook with the following code:

```
# Load your trained model
@st.cache_resource
def load_model():
    with open('titanic_model.pkl', 'rb') as f:
        model = pickle.load(f)
    return model

model = load_model()

# App header with HTML
st.markdown("""
<div class="title">
<h1>🚢 Titanic Survival Predictor</h1>
<p>Predict whether a passenger would have survived the Titanic disaster</p>
</div>
""", unsafe_allow_html=True)

# Titanic image
st.image("https://upload.wikimedia.org/wikipedia/commons/thumb/f/fd/RMS_Titanic_3.jpg/1200px-RMS_Titanic_3.jpg",
caption="RMS Titanic", use_container_width=True)

# Input features in sidebar
st.sidebar.header('Passenger Details')

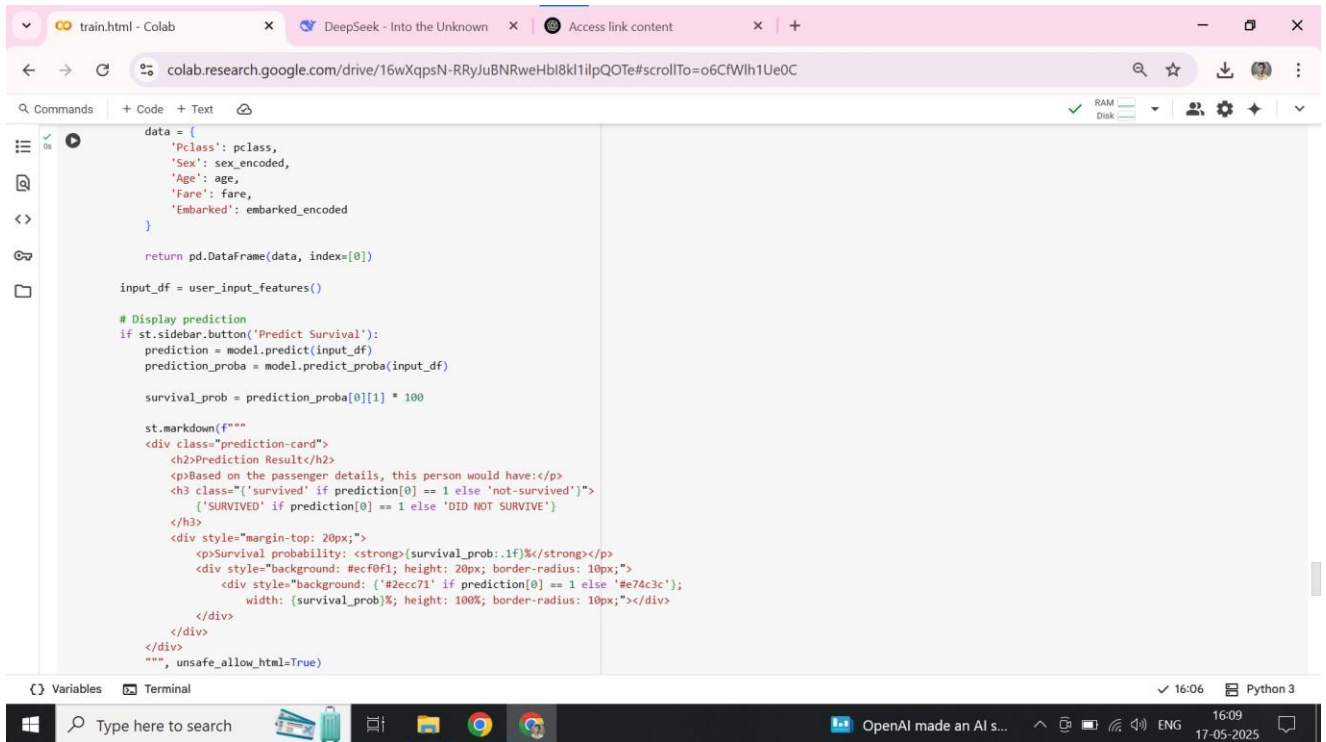
def user_input_features():
    col1, col2 = st.sidebar.columns(2)

    with col1:
        pclass = st.selectbox('Passenger Class', [1, 2, 3], help="1 = First Class, 2 = Second Class, 3 = Third Class")
        sex = st.selectbox('Sex', ['female', 'male'])

    with col2:
        age = st.slider('Age', 0, 100, 30)
        fare = st.slider('Fare (£)', 0, 200, 30, help="Ticket fare in British pounds")
```

The notebook interface includes a top bar with tabs for 'train.html - Colab', 'DeepSeek - Into the Unknown', and 'Access link content'. The address bar shows the Colab URL. The left sidebar contains icons for file explorer, search, and other tools. The bottom status bar indicates 'Variables', 'Terminal', '16:06', and 'Python 3'.

8.7.Titanic Survival Prediction



The screenshot shows a Google Colab notebook with the following code:

```
data = {
    'Pclass': pclass,
    'Sex': sex_encoded,
    'Age': age,
    'Fare': fare,
    'Embarked': embarked_encoded
}

return pd.DataFrame(data, index=[0])

input_df = user_input_features()

# Display prediction
if st.sidebar.button('Predict Survival'):
    prediction = model.predict(input_df)
    prediction_proba = model.predict_proba(input_df)

    survival_prob = prediction_proba[0][1] * 100

    st.markdown(f"""
    <div class="prediction-card">
    <h2>Prediction Result</h2>
    <p>Based on the passenger details, this person would have:</p>
    <h3 class="{ 'survived' if prediction[0] == 1 else 'not-survived' }">
    { 'SURVIVED' if prediction[0] == 1 else 'DID NOT SURVIVE' }
    </h3>
    <div style="margin-top: 20px;">
    <p>Survival probability: <strong>{survival_prob:.1f}%</strong></p>
    <div style="background: #ecf0f1; height: 20px; border-radius: 10px;">
    <div style="background: { '#2ecc71' if prediction[0] == 1 else '#e74c3c' };
    width: {survival_prob}%; height: 100%; border-radius: 10px;"></div>
    </div>
    </div>
    """, unsafe_allow_html=True)
```

The notebook interface includes a top bar with tabs for 'train.html - Colab', 'DeepSeek - Into the Unknown', and 'Access link content'. The address bar shows the Colab URL. The left sidebar contains icons for file explorer, search, and other tools. The bottom status bar indicates 'Variables', 'Terminal', '16:06', and 'Python 3'.

8.8.Session State Authentication

```
# Main App Router
init_session()

if st.session_state.auth["logged_in"]:
    titanic_app()
    if st.button("Logout"):
        st.session_state.auth["logged_in"] = False
        st.rerun()
else:
    auth_page()
```

8.9.Web App and Temporary Public Link

```
from pyngrok import ngrok

# # Start Streamlit in background
process = subprocess.Popen(["streamlit", "run", "stream_app.py"])

# Set your authtoken
ngrok.set_auth_token("2xD00ahrHv07wgg2LMQJi5ezaSz_5TCR25BZN4dABb1YiVsX4")

# Define the tunnel configuration
tunnel_config = {
    "proto": "http", # or "tcp" if needed
    "addr": 8501,    # The port your Streamlit app is running on
}

# Now, try to connect using the configuration
public_url = ngrok.connect(**tunnel_config)

public_url = ngrok.connect(**tunnel_config)
print("🌐 Your app is available at:", public_url)
```

🌐 Your app is available at: NgrokTunnel: "<https://797a-34-91-142-22.ngrok-free.app>" -> "<http://localhost:8501>"

9. Result

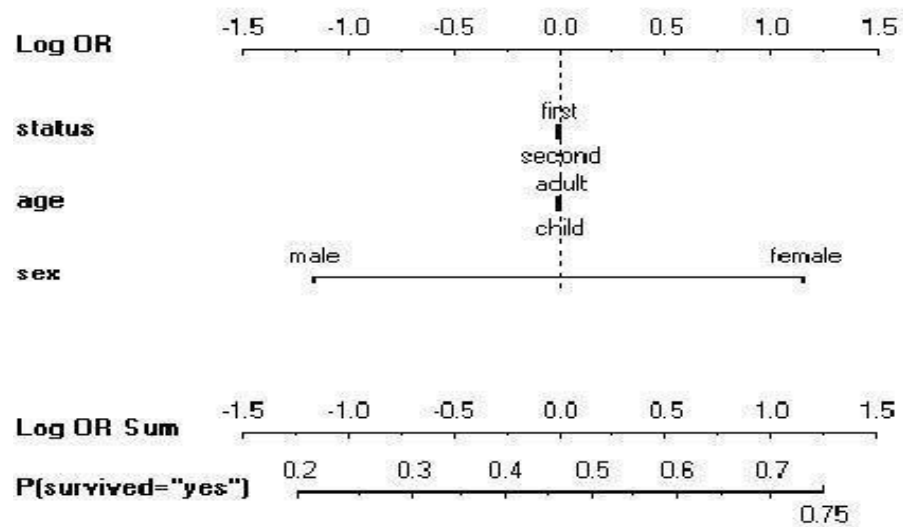


Figure 9. Classification of survivals

The relevant data can be obtained through a series of model building and analysis in insights, as shown in Figure 9. After classifying different methods and algorithms, finding the best implementation score to recognize objects and separate them into categories is one of the processes. As a result, this study gained valuable experience building prediction systems from the SVM model and achieved the best score on Kaggle, 82.82% of correct predictions.

10. Conclusion

Extensive hyperparameters should be adjusted on multiple machine learning models to develop overall results. It can be further enhanced through comprehensive learning. This paper starts with data exploration and then examines the valid data through analysis to see which features are essential to the model. When the data preprocessing sections arrive, classifications are used to calculate the exact values of each model and convert them into digital visual elements. Then, by creating some functions, different machine learning models are trained. Finally, the confusion matrix, f-score, accuracy of computer models, and recall rates were investigated and evaluated. Data clean-up is the first and most critical step when trying to analyse data. Exploratory Data Analysis allows people to identify datasets and associated link characteristics. EDA is implemented to find relationships between dataset elements. This is achieved by using different graphics techniques. Using EDA to draw some conclusions and discover facts is an integral part of this study. Optimizing the efficiency and accuracy of the model is also a critical process in the experimental process when learning about data recycling using valid models.

References

- [1] The National Archives, J. C. 2012. Titanic, 100 years on.
<https://blog.nationalarchives.gov.uk/titanic100years-on/>
- [2] Kaggle. 2020, Aug 1. Titanic - machine learning from disaster. <https://www.kaggle.com/c/titanic/data> [3] Cronan, J. 2012. *National Archives*.<https://blog.nationalarchives.gov.uk/titanic-100-years-on/>
- [4] Lord, W., & Philbrick, N. 2017. I Believe She's Gone, Hardy. In *A night to remember* (pp. 47–50).
- [5] Mishra, V. P., Singh, B., Shukla, V. K., & Dasgupta, A. 2021. Predicting the likelihood of survival of Titanic's passengers by Machine Learning.
https://www.researchgate.net/publication/351155499_Predicting_the_Likelihood_of_Survival_of_Titanic%27s_Passengers_by_Machine_Learning
- [6] Brownlee, J. 2020. Evaluating Machine Learning Algorithms.<https://machinelearningmastery.com/train-test-split-for-evaluating-machine-learning-algorithms/>
- [7] ABDILLAH, F. R. 2021. Titanic Prediction - ANN.
<https://www.kaggle.com/code/farizhaykal/titanicprediction-ann>
- [8] Cook, A. 2019. Titanic - Machine Learning from Disaster.
<https://www.kaggle.com/code/alexisbcook/titanic-tutorial>
- [9] Swain, A. 2017. EDA To Prediction.<https://www.kaggle.com/code/ash316/eda-to-prediction-dietanic>
- [10] Housseem Ben Braieka, Foutse Khomha. 2018. *On Testing Machine Learning Programs*. Montreal: SWAT lab. <https://towardsdatascience.com/machine-learning-for-beginners-d247a9420dab>
- [11] Niklas Donges, 2018, Predicting the Survival of Titanic Passengers.
<https://towardsdatascience.com/predicting-the-survival-of-titanic-passengers-30870ccc7e8>