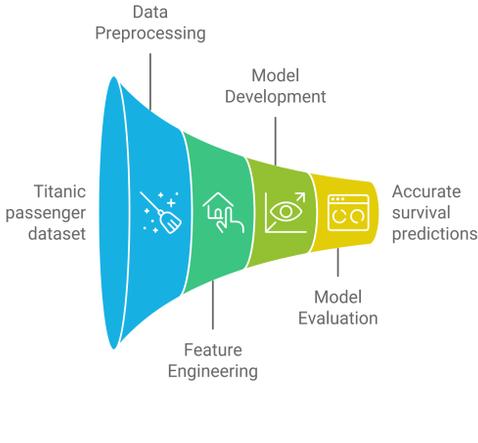


# Titanic - Kaggle

## Project: Titanic Survival Prediction using Machine Learning

To develop a machine learning model that accurately predicts the survival of passengers on the Titanic based on various passenger attributes, such as age, gender, class, and fare, by analyzing the provided dataset. This project aims to enhance skills in data preprocessing, feature engineering, model development, and evaluation using real-world data.

### Titanic Survival Prediction Process



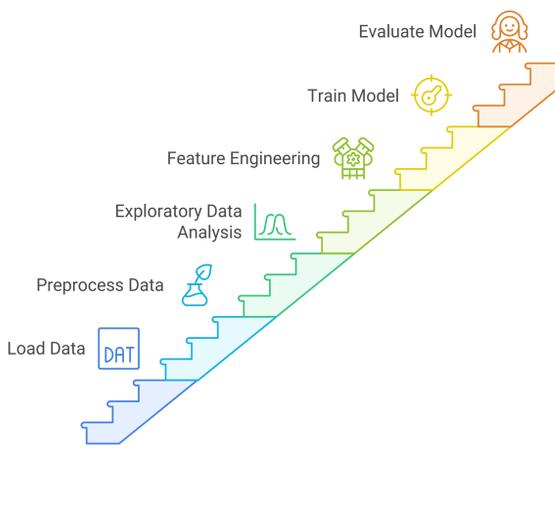
Platform: Kaggle

- Develop a predictive model to determine the likelihood of passenger survival in the Titanic disaster using Python and machine learning techniques.
- Analyze a dataset containing demographic and socio-economic data of passengers, such as age, gender, passenger class, and fare.
- Perform data cleaning, feature engineering, and exploratory data analysis to identify key factors influencing survival rates.
- Implement multiple machine learning algorithms [e.g., Logistic Regression, Decision Trees, Random Forests] to build and optimize predictive models.
- Evaluate model performance using accuracy scores and cross-validation techniques, achieving an X% accuracy on the test set.
- Gain hands-on experience in data preprocessing, feature selection, model building, and hyperparameter tuning, enhancing proficiency in Python libraries like Pandas, NumPy, Scikit-Learn, and Matplotlib.

### Which machine learning algorithm to choose for building the predictive model?



### Titanic Survival Prediction Workflow



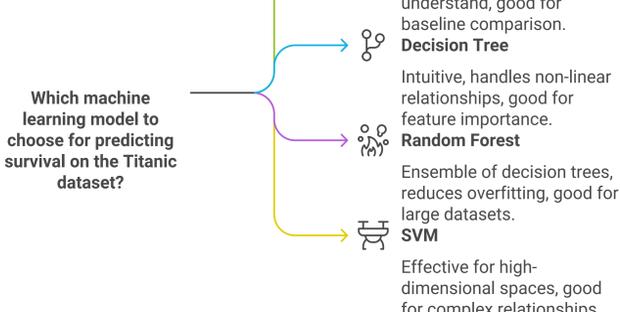
```
def main():
    """
    Main function to run all steps sequentially.
    """
    load_data()
    preprocess_data()
    exploratory_data_analysis()
    feature_engineering()
    train_model()
    evaluate_model()
    make_predictions()
    submit_predictions()

if __name__ == "__main__":
    main()
```

For the Kaggle Titanic project, several machine learning models are popular due to their effectiveness and ease of implementation. Here are the most commonly used models:



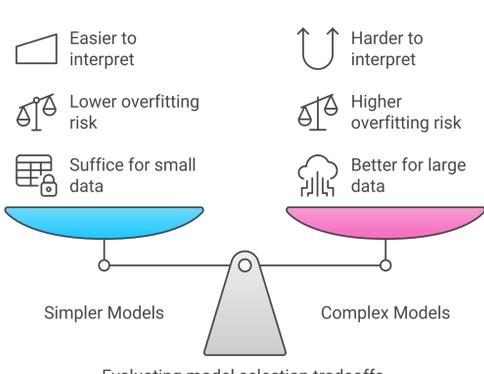
- Logistic Regression**
  - Why Popular: Simple, easy to understand, and works well for binary classification problems like predicting survival.
  - Usage: Often used as a baseline model to compare performance with more complex models.
- Decision Tree**
  - Why Popular: Intuitive and easy to visualize; handles non-linear relationships and interactions between features.
  - Usage: Useful for understanding which features are most important for predicting survival.
- Random Forest**
  - Why Popular: An ensemble of decision trees that improves accuracy by reducing overfitting; handles missing values and works well with large datasets.
  - Usage: A powerful and versatile model frequently used by beginners and experts alike.
- Support Vector Machine (SVM)**
  - Why Popular: Effective for high-dimensional spaces; can be used for both linear and non-linear classifications.
  - Usage: Applied when the data has complex relationships or requires a more sophisticated decision boundary.
- K-Nearest Neighbors (KNN)**
  - Why Popular: Simple and intuitive; works well when the dataset is small and features are numerical.
  - Usage: Typically used for quick checks or when the dataset does not require complex modeling.
- Naive Bayes**
  - Why Popular: Fast and efficient; works well with small datasets and assumes independence among features.
  - Usage: Applied when the data has categorical features and a probabilistic approach is desired.
- Gradient Boosting [e.g., XGBoost, LightGBM]**
  - Why Popular: Very powerful and often yields high accuracy by combining multiple weak learners to create a strong learner; handles various data types and complex relationships.
  - Usage: Popular among experienced data scientists for achieving top results in competitions due to its ability to handle missing data, feature interactions, and non-linear relationships.
- Neural Networks**
  - Why Popular: Can capture complex patterns in data; flexible with a variety of architectures.
  - Usage: Used less frequently for Titanic due to the simplicity of the problem, but useful for learning about deep learning models.



### Model Selection Criteria

Choosing the right model often depends on factors such as:

- Data size and complexity: Simpler models like Logistic Regression or Decision Trees may suffice for smaller datasets with fewer features.
- Overfitting risk: More complex models (like Random Forest or Gradient Boosting) may perform better but require careful tuning to avoid overfitting.
- Interpretability: Some models, such as Decision Trees and Logistic Regression, are easier to interpret, which can be important for explaining results.



### What I Used

#### Random Forest Model in the Kaggle Titanic Project

The **Random Forest** model is an ensemble learning method that combines multiple decision trees to improve prediction accuracy and control overfitting. It's particularly popular in the Kaggle Titanic project due to its robustness and ability to handle a variety of data types, including both categorical and numerical features.

#### Why Use Random Forest for the Titanic Project?

- Handles Missing Data**: Random Forest can handle missing values well, which is common in the Titanic dataset (e.g., missing ages or cabin numbers).
- Reduces Overfitting**: By averaging the results of many decision trees, Random Forest reduces the risk of overfitting to the training data.
- Feature Importance**: It provides insights into the importance of different features, helping understand which factors most influence survival.
- Non-linear Relationships**: Captures complex, non-linear relationships between features, which can be beneficial in the Titanic dataset where interactions (like age and gender) are significant.

#### Random Forest Model for Titanic Survival Prediction

