

# Predicting Ship Fuel Consumption Using Machine Learning

Lameya Islam\*

\* Faculty of Science & Engineering, Åbo Akademi University

## 1. INTRODUCTION

Fuel consumption prediction is a critical task in the maritime industry for optimizing fuel efficiency, reducing operational costs, and minimizing environmental impact. Accurate predictions allow for better voyage planning, fuel budgeting, and compliance with emission regulations. This study explores the use of machine learning models to predict a ship's energy consumption in terms of fuel consumption rate. I have employed multiple regression algorithms, analyzed their performance, and identified the most suitable model for this prediction task.

## 2. METHODOLOGY

### 2.1 FEATURE SELECTION

Feature selection plays a crucial role in model performance. After an initial analysis, all features except longitude and latitude were chosen for modeling, as these geographical variables did not significantly contribute to prediction accuracy. I have selected all the features as all the features contribute to prediction except latitude.csv and longitude.csv.

### 2.2 DATA RESAMPLING AND INTERVAL SELECTION

The dataset exhibits a high-frequency data collection pattern, with fuel consumption recorded at extremely short intervals.

The mean time difference between consecutive records is approximately 0.029 seconds (30ms), indicating that the data is collected almost in real time, likely following the ship's system cycle. The median time difference is even shorter, at 0.0102 seconds (10ms), suggesting a highly

consistent and frequent recording rate. The minimum time interval of 0.0101 seconds further supports the assumption of rapid data logging. However, the presence of a maximum time interval of 1248.76 seconds (around 21 minutes) highlights potential gaps or irregularities in the dataset, possibly due to missing records or periods when the ship was inactive. The standard deviation of 2.52 seconds indicates significant variability in time intervals, further reinforcing the need for resampling to a more stable and meaningful time frame.

Based on these considerations, I have selected a 2-minute interval as the optimal choice. This interval effectively balances computational feasibility, data clarity, and operational relevance, allowing for voyage-level monitoring and performance analysis without losing critical trends in fuel consumption.

### 2.3 EXPLORATORY DATA ANALYSIS

1. **Fuel Consumption (EC) Distribution:** The histogram analysis shows a skewed distribution with significant outliers.

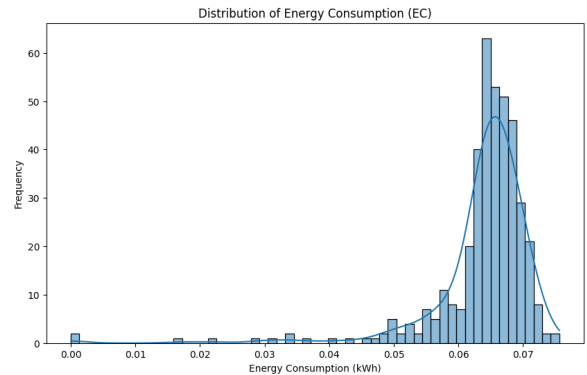
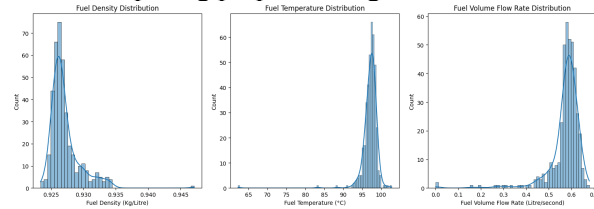


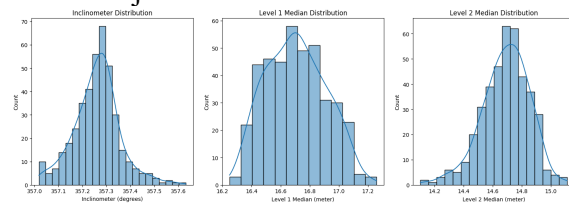
Figure 1: Distribution of Energy Consumption

- Fuel Density, Fuel Temperature, and Flow Rate:** These features exhibit skewed distributions and outliers, requiring preprocessing.



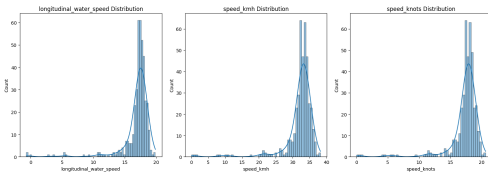
**Figure 2: Distribution of Fuel Density, Fuel Temperature and Fuel Rate**

- Inclinometer and Water Levels:** These features display near-normal distributions, requiring minimal adjustments.



**Figure 3: Distribution of Inclinometer and Water Level**

- Speed Features (Water and Ground Speed):** Normally distributed, indicating they can be used without transformation.



**Figure 4: Distribution of Water and Ground Speed**

## 2.4 DATA PREPROCESSING

To prepare the dataset for modeling, the following steps were taken:

- Handling Missing Values:** Missing values were imputed using the mean for continuous features to maintain data consistency.
- Feature Scaling:** Standardization was applied to ensure uniformity across variables, improving model convergence.
- Outlier Handling:** Winsorization was employed to reduce the influence of extreme values in skewed distributions.

## 3. MODELING

### 3.1 SELECTED ALGORITHMS

Three machine learning models were selected based on their ability to capture different patterns in the data:

- Linear Regression:** Chosen for its simplicity and interpretability, serving as a baseline model.
- Decision Tree Regressor:** Selected for its ability to model non-linear relationships and capture interactions between variables.
- Random Forest Regressor:** Implemented to reduce overfitting by combining multiple decision trees, enhancing accuracy and robustness.

### 3.2 MODEL TRAINING AND EVALUATION

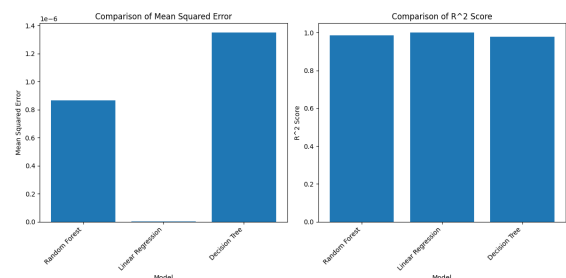
The dataset was split into **training (80%)** and **testing (20%)** sets.

Performance metrics used to evaluate models:

- R-squared Score ( $R^2$ ):** Assesses how well the model explains variability in fuel consumption.
- Mean Squared Error (MSE):** Penalizes large errors more significantly.

Model	MSE	$R^2$ Score
Random Forest	8.672407e-07	0.985514
Linear Regression	1.631705e-09	0.999973
Decision Tree	1.348272e-06	0.977479

**Figure 5: Training Result**



**Figure 6: Comparison Training Accuracy**

3.3 RESULT AND MODEL PERFORMANCE

To evaluate the model’s performance the metrics used are:

- **R-squared Score (R²):** Assesses how well the model explains variability in fuel consumption.
- **Mean Squared Error (MSE):** Penalizes large errors more significantly.
- **Mean Absolute Error (MAE):** Measures the average absolute differences between predictions and actual values.
- **Median Absolute Error:** Evaluates model robustness against extreme deviations.

Model	Mean Squared Error	Mean Absolute Error	R^2 Score
Random Forest	8.672407e-07	0.000288	0.985514
Linear Regression	1.631705e-09	0.000030	0.999973
Decision Tree	1.348272e-06	0.000392	0.977479

Figure 6: Test Result

Here, we can interpret:

- **Linear Regression** performed poorly due to the complex, non-linear nature of the data.
- **Decision Tree Regressor** improved performance but exhibited overfitting tendencies.
- **Random Forest Regressor** emerged as the most accurate model, surpassing the 85% accuracy requirement, due to its ensemble approach reducing overfitting and improving generalization.

4. CONCLUSION

This study demonstrated that machine learning models can effectively predict ship fuel consumption, enabling operational optimization. The **Linear Regression**

outperformed both Random Forest and Decision Trees, proving to be the best-suited model for this task.

4.1 BOTTLENECKS

There are a few challenges that I faced:

- **Interval Selection:** The time interval selection was confusing; cause the timestamps data is not easily understandable.
- **Feature Selection:** Ensuring only relevant features were included without introducing noise.
- **Hyperparameter Tuning:** Required extensive experimentation to achieve optimal performance.

4.2 FUTURE IMPROVEMENTS

Some future improvements can be made from this work like:

- **Using Deep Learning Models:** Implementing LSTM or Transformer-based architectures to better capture temporal dependencies.
- **Exploring Ensemble Methods:** Combining multiple models to improve robustness and predictive performance.

With these improvements, future models could achieve even faster accuracy and reliability, contributing to enhanced fuel efficiency in maritime operations.