

TDT4173 Report - Group 116

Johan Solbakken
Department of Computer Science, NTNU
Trondheim, Norway
Email: johsol@stud.ntnu.no

Morten Tobias Rinde Sunde
Department of Computer Science, NTNU
Trondheim, Norway
Email: mtsunde@stud.ntnu.no

December, 2024

Contents

1	Introduction	3
2	Initial attempt	3
3	Exploratory Analysis	4
3.1	Correlation Analysis	4
3.2	Plotting SOG by NAVSTAT	5
3.3	Interpolating Vessel Routes	7
3.4	Distribution of Time Gaps Between Consecutive AIS Records	11
4	Improved LSTM Model	11
4.1	Feature Engineering	12
4.2	The Model	12
4.3	The Code	13
4.4	Result	17
4.5	Integrating Cubic Spline Interpolation into LSTM	17
4.6	BiGRU	23
5	Darts and LightGBM	23
5.1	Feature Engineering	23
5.2	The Model	24
5.3	Hyperparameter tuning	24
5.4	The Code	24
5.5	Result	26
6	Random Forest	26
6.1	Feature Engineering	26
6.2	Feature Selection	26
6.3	The Model	27
6.4	Hyperparameter Tuning	27
6.5	The Code	27
6.6	Result	28
6.7	Feature Importance For Random Forest Model	28
7	Conclusions	30

1 Introduction

This report summarizes the steps taken in our project for the TDT4173 course. We discuss our exploratory data analysis, the feature engineering process, the various models we explored, and reflections on their effectiveness.

Our team name was [116] *Kvorum AS*, and our efforts culminated in a score of 115.8 on Kaggle.

2 Initial attempt

For our first attempt, we focused on building a pipeline that would run end-to-end and produce predictions for submission on Kaggle. We conducted no exploratory data analysis or detailed feature engineering at this stage. The goal was to ensure that we had functional code capable of making predictions and obtaining a preliminary score in the competition. The following Python code outlines our initial approach:

```
1 import pandas as pd
2 import numpy as np
3 from datetime import datetime
4 from geopy.distance import geodesic
5 from tensorflow.keras.models import Sequential
6 from tensorflow.keras.layers import LSTM, Dense
7 from sklearn.model_selection import train_test_split
8
9 def calculate_distance(lat1, lon1, lat2, lon2):
10     return geodesic((lat1, lon1), (lat2, lon2)).kilometers
11
12 def prepare_sequences(data, feature_cols, target_cols, sequence_length):
13     X, y = [], []
14     data_values = data[feature_cols + target_cols].values
15     for i in range(len(data_values) - sequence_length):
16         X.append(data_values[i:i+sequence_length, :len(feature_cols)])
17         y.append(data_values[i+sequence_length-1, len(feature_cols)])
18     return np.array(X), np.array(y)
19
20 def prepare_test_sequences(data, feature_cols, sequence_length):
21     X = []
22     data_values = data[feature_cols].values
23     for i in range(len(data_values) - sequence_length + 1):
24         X.append(data_values[i:i+sequence_length])
25     return np.array(X)
26
27 def prepare_test_data(ais_test, vessels, vessel_type_categories):
28     merged_test = pd.merge(ais_test, vessels[['vesselId', 'vesselType']], on='vesselId', how='left')
29     merged_test['vesselType'].fillna('Unknown', inplace=True)
30     merged_test['time'] = pd.to_datetime(merged_test['time'])
31     merged_test['hour'] = merged_test['time'].dt.hour
32     merged_test['day_of_week'] = merged_test['time'].dt.dayofweek
33     merged_test['vesselType'] = pd.Categorical(merged_test['vesselType'],
34     ↪ categories=vessel_type_categories)
35     merged_test['vessel_type_encoded'] = merged_test['vesselType'].cat.codes
36     features = ['hour', 'day_of_week', 'vessel_type_encoded']
37     merged_test[features] = merged_test[features].fillna(0)
38     return merged_test, features
39
40 # Load and prepare data
41 ais_train = pd.read_csv('ais_train.csv', sep='|')
42 ais_test = pd.read_csv('ais_test.csv')
43 vessels = pd.read_csv('vessels.csv', sep='|')
44
45 merged_data = pd.merge(ais_train, vessels[['vesselId', 'vesselType']], on='vesselId', how='left')
46 merged_data['vesselType'].fillna('Unknown', inplace=True)
47 merged_data['time'] = pd.to_datetime(merged_data['time'])
48 merged_data['hour'] = merged_data['time'].dt.hour
49 merged_data['day_of_week'] = merged_data['time'].dt.dayofweek
50 merged_data['future_latitude'] = merged_data.groupby('vesselId')['latitude'].shift(-1)
51 merged_data['future_longitude'] = merged_data.groupby('vesselId')['longitude'].shift(-1)
52 merged_data.dropna(subset=['future_latitude', 'future_longitude'], inplace=True)
53
54 merged_data['vesselType'] = merged_data['vesselType'].astype('category')
55 vessel_type_categories = merged_data['vesselType'].cat.categories
56 merged_data['vessel_type_encoded'] = merged_data['vesselType'].cat.codes
57 features = ['hour', 'day_of_week', 'vessel_type_encoded']
58 target = ['future_latitude', 'future_longitude']
59
60 train_data, val_data = train_test_split(merged_data, test_size=0.2, shuffle=False)
61 sequence_length = 5
62 X_train, y_train = prepare_sequences(train_data, features, target, sequence_length)
63 X_val, y_val = prepare_sequences(val_data, features, target, sequence_length)
```

```

64 # Build and train the model
65 model = Sequential()
66 model.add(LSTM(64, return_sequences=False, input_shape=(sequence_length, len(features))))
67 model.add(Dense(2))
68 model.compile(optimizer='adam', loss='mean_absolute_error')
69
70 model.fit(X_train, y_train, epochs=1, validation_data=(X_val, y_val))
71
72 # Predictions and error calculation
73 predictions_val = model.predict(X_val)
74 val_data = val_data.iloc[sequence_length:]
75 val_data['pred_latitude'] = predictions_val[:, 0]
76 val_data['pred_longitude'] = predictions_val[:, 1]
77 val_data['error_distance'] = val_data.apply(lambda row: calculate_distance(
78     row['future_latitude'], row['future_longitude'], row['pred_latitude'], row['pred_longitude']),
79     axis=1)
80 mean_error_distance = val_data['error_distance'].mean()
81 print(f'Mean Geodetic Error: {mean_error_distance} km')
82
83 # Prepare submission data
84 merged_test, features = prepare_test_data(ais_test, vessels, vessel_type_categories)
85 X_test = prepare_test_sequences(merged_test, features, sequence_length)
86 predictions_test = model.predict(X_test)
87
88 # Save predictions
89 prediction_indices = np.arange(sequence_length - 1, len(merged_test))
90 submission_df = merged_test.iloc[prediction_indices].copy()
91 submission_df['longitude_predicted'] = predictions_test[:, 1]
92 submission_df['latitude_predicted'] = predictions_test[:, 0]
93 submission_df['ID'] = ais_test.iloc[prediction_indices]['ID'].values
94 submission_df = submission_df[['ID', 'longitude_predicted', 'latitude_predicted']]
95 submission_df.to_csv('submission.csv', index=False)

```

This code handled data preparation, model training, and prediction. It included functions for preparing sequences of features and targets and predicting future vessel locations based on past data. We utilized an LSTM model with a single layer and trained it on sequences of time and vessel-related features to predict latitude and longitude.

Although the primary objective of this attempt was to produce a working model, the results helped us understand the challenges and set a foundation for further improvements. The geodetic error calculated during validation gave us a rough measure of prediction accuracy, which provided valuable feedback for future iterations.

This model scored 600 on Kaggle, which we thought was a good start. However, it became clear that conducting an exploratory analysis and applying feature engineering would be necessary to create a model with better predictive performance.

3 Exploratory Analysis

We found that we had to do some exploratory analysis of the data. Our previous attempt suggested that there probably was a lot of noise in the data. To investigate this, we plotted the data using `matplotlib` to identify which parts of the dataset should be cleaned up.

3.1 Correlation Analysis

In the initial phase of our exploratory data analysis, we performed a correlation analysis on the `ais.train` dataset. The code below illustrates our approach:

```

1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 # Select only numeric columns for correlation matrix
6 numeric_columns = ais_train.select_dtypes(include=['float64', 'int64'])
7
8 # Calculate the correlation matrix
9 corr_matrix = numeric_columns.corr()
10
11 # Plot the heatmap
12 sns.heatmap(corr_matrix, annot=True, fmt='.2f', cmap='coolwarm')
13 plt.show()

```

This gave us a heatmap representation of the correlation between the numerical features in the dataset

1. From this, we could observe the following correlations:

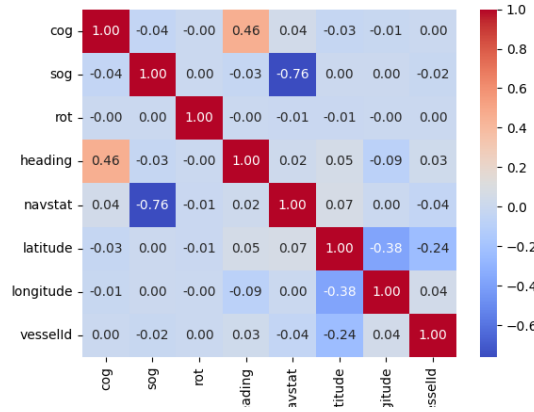


Figure 1: Correlation `ais_train.csv`

- **COG (Course Over Ground) and heading:** A moderate positive correlation (0.46), which suggests that vessels with higher COG tend to have a higher heading, as expected, since these slight directional metrics.
- **SOG (Speed Over Ground) and NAVSTAT (Navigation Status):** A strong negative correlation (-0.76), which indicates that vessels with specific NAVSTAT values (e.g., underway vs. anchored) have significantly different speed patterns. Vessels with NAVSTAT as 0 (Underway) will have a positive SOG, while other statuses will have lower or zero SOG. This will be interesting to investigate.
- **latitude and longitude:** There is a slight negative correlation between latitude and longitude (-0.38), which could imply geographical characteristics of the data, where locations tend to have this inverse relationship in this dataset's geographic region.
- **NAVSTAT and heading:** A slight negative correlation between these two (0.02), which might indicate that certain navigation statuses only marginally impact vessel heading.

Most features show weak or no correlation, except for the abovementioned cases. **vesselid** doesn't show strong correlations with other variables, which is expected since it's more of a categorical identifier.

3.2 Plotting SOG by NAVSTAT

Based on the findings from the previous section, we decided to investigate the relationship between SOG and NAVSTAT. In the dataset, SOG is represented as a floating-point number, capturing precise speed measurements, while NAVSTAT is stored as integer codes corresponding to specific navigational states. To make this data more interpretable, we mapped the integer values of NAVSTAT to their respective descriptive labels, as illustrated in Listing 1. This mapping was based on `api.vtexplorer.com`'s AIS (Automatic Identification System) documentation [1].

We thought a natural relationship between a vessel's SOG and NAVSTAT should exist. For example, if a ship is "At Anchor," it is logical to assume that its SOG would be zero, indicating that the vessel is stationary. Additionally, the problem description provided in the project mentioned that some ships might use more than one NAVSTAT code to describe their current activity. Specifically, most vessels tend to use code 0 ("Underway using engine") as their default but occasionally switch to code 8 ("Underway sailing") when transitioning to or from using sails. These assumptions suggested that plotting SOG against NAVSTAT could help us identify patterns or inconsistencies in the data. By visualizing the relationship between these two variables, we could understand how well the data aligns with our expectations. Listing 2 details the plot generating process based on the `ais_train` dataset.

The data loading code in Listing 3 enabled us to generate Figure 2, which presents the relationship between Speed Over Ground (SOG) and Navigational Status (NAVSTAT) across the entire dataset. A noteworthy observation from this plot is the presence of vessels with speeds reaching up to 100 knots while underway. Since container ships typically operate within the range of 16 to 24 knots [2], we classi-

```

1  navstat_mapping = {
2      0: 'Under way using engine',
3      1: 'At anchor',
4      2: 'Not under command',
5      3: 'Restricted maneuverability',
6      4: 'Constrained by her draught',
7      5: 'Moored',
8      6: 'Aground',
9      7: 'Engaged in Fishing',
10     8: 'Under way sailing',
11     9: 'Reserved for future amendment of Navigational Status for HSC',
12    10: 'Reserved for future amendment of Navigational Status for WIG',
13    11: 'Reserved for future use',
14    12: 'Reserved for future use',
15    13: 'Reserved for future use',
16    14: 'AIS-SART is active',
17    15: 'Not defined (default)'
18 }

```

Listing 1: NAVSTAT mapping

```

1  import matplotlib.pyplot as plt
2  import seaborn as sns
3
4  # load the ais_train
5
6  # Create a box plot to visualize the distribution of sog for each navstat description
7  plt.figure(figsize=(12, 6))
8  sns.boxplot(x='navstat_description', y='sog', data=filtered_data)
9  plt.xlabel('Navigation Status')
10 plt.ylabel('Speed Over Ground (sog)')
11 plt.title('Distribution of Speed Over Ground by Navigation Status')
12 plt.xticks(rotation=45, ha='right') # Rotate labels for better readability
13 plt.tight_layout() # Adjust layout to prevent label overlap
14 plt.show()

```

Listing 2: Plotting ais_train for SOG and NAVSTAT

fied these high-speed instances as noise. Furthermore, by inspecting the `vessels.csv` dataset, we observed that the recorded maximum, minimum, and average speeds for vessels with non-null `maxSpeed` values were 23.3, 16.7, and approximately 21.7 knots, respectively. This reinforced our assumption that the original speed data contained unrealistic values, likely due to sensor errors or misreporting.

```

1  """
2      Load Data
3  """
4  # Read ais_train.csv
5  ais_train = pd.read_csv("ais_train.csv", sep='|')
6
7  # Map the numeric navstat values to their descriptive names
8  ais_train['navstat_description'] = ais_train['navstat'].map(navstat_mapping)
9
10 # Filter data to include only the navstat descriptions we're interested in (e.g., navstat < 6)
11 filtered_data = ais_train[ais_train['navstat'] < 6]

```

Listing 3: Loading the dataset

A second significant observation was that vessels reported as being “Moored” or “At Anchor” had non-zero speeds in several instances. Based on the assumption that moored or anchored ships should be stationary, with zero speed, we decided to clean the dataset by removing any records that violated this assumption. Specifically, we pruned the data to exclude instances where vessels were moored or anchored but had a reported speed greater than zero. Moreover, we identified several outliers where vessels were labeled as “Not Under Command” yet had SOG values exceeding 5 knots, which seemed implausible. These outliers were also removed as part of our data-cleaning process.

Furthermore, given that NAVSTAT code 8 (“Underway sailing”) effectively conveys the same information as NAVSTAT code 0 (“Underway using engine”), we unified these categories by reassigning all data points with NAVSTAT 8 to NAVSTAT 0. This helped reduce redundancy in the dataset and ensured that our model would not be confused by the same maritime status represented by two different codes. We also removed NAVSTAT codes 6 through 15, which correspond to statuses like “On the ground,” “Fishing,” and “Reserved for future use,” as these statuses were irrelevant to our analysis. The code used to perform this data pruning is shown in Listing 4. Figure 3 illustrates the improved plot after applying these pruning steps, with the SOG values for the “Underway using the engine,” “At anchor,” and “Moored” statuses now appearing much more reasonable.

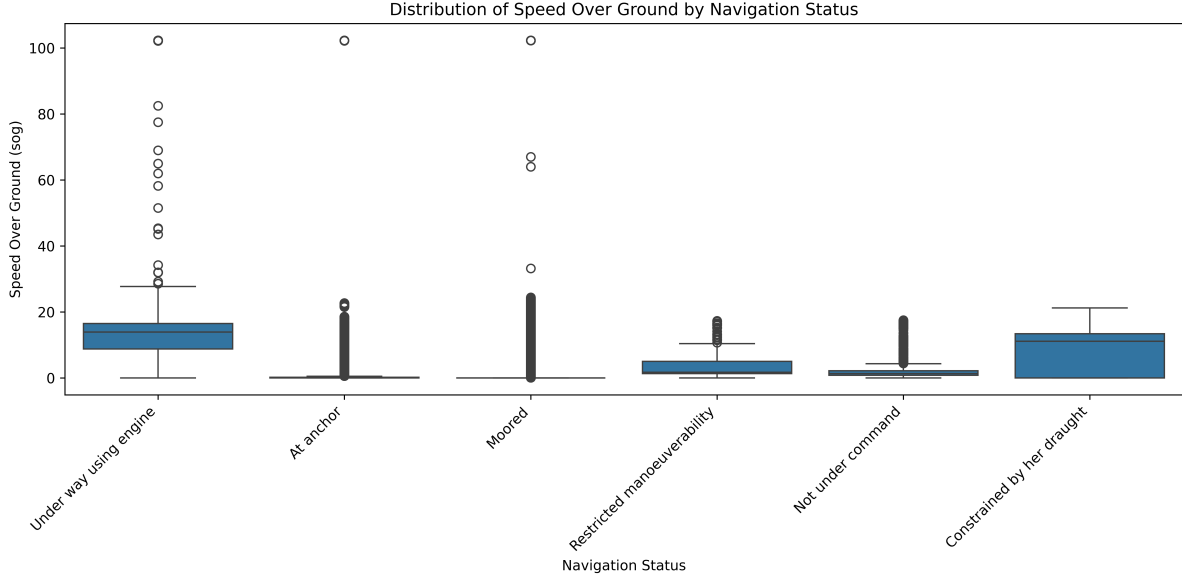


Figure 2: SOG over NAVSTAT plot of the entire dataset

```

1  """
2      Load Data
3  """
4  # Read ais_train.csv
5  ais_train = pd.read_csv("ais_train.csv", sep='|')
6
7  # We assume that speed less than 25 knots are appropriate
8  ais_train = ais_train[(ais_train['sog'] < 25)]
9  ais_train['navstat'] = ais_train['navstat'].replace(8, 0) # every boat that is under sail is under way
10 ais_train = ais_train[~((ais_train['navstat'].isin([1, 5])) & (ais_train['sog'] > 0))] # standing
    ↳ still, should have 0 speed
11 ais_train = ais_train[~((ais_train['navstat'].isin([2])) & (ais_train['sog'] > 5))] # remove outliers
    ↳ for not under command;

```

Listing 4: Loading data, pruning SOG and NAVSTAT outliers

3.3 Interpolating Vessel Routes

One challenge with the dataset was that the time intervals between recorded data points were highly inconsistent. Some data points were recorded only 5 seconds apart, while others were separated by as much as 13 days. Such significant discrepancies in time intervals posed a problem for trajectory analysis, as having evenly spaced data points is crucial for creating reliable predictions and models. To address this issue, we aimed to ensure that each vessel in the dataset had at least one recorded data point per day. This regularization of the time intervals would provide a more stable foundation for modeling. Inspired by the work in [3], we experimented with cubic spline interpolation to fill in the gaps between irregularly spaced data points. Cubic splines allow for smooth interpolation between known data points, making them suitable for creating a continuous trajectory from our discrete data.

To illustrate this, as an example, let us consider a specific vessel, identified by 61e9f3bfb937134a3c4bfe9f. We can visualize the trajectory of this vessel using Python's `plotly` library, which allows for interactive data plotting. The code provided below demonstrates how we performed this visualization.

```

1  import pandas as pd
2  import numpy as np
3  import plotly.graph_objects as go
4
5  # Read ais_train.csv
6  ais_train = pd.read_csv("ais_train.csv", sep='|')
7
8  # Temporal features
9  ais_train['time'] = pd.to_datetime(ais_train['time'])
10 ais_train['elapsed_time'] = (ais_train['time'] - pd.Timestamp("1970-01-01")) // pd.Timedelta('1s')
11 ais_train['day_of_week'] = ais_train['time'].dt.dayofweek # Monday=0, Sunday=6
12 ais_train['hour_of_day'] = ais_train['time'].dt.hour
13 ais_train = pd.get_dummies(ais_train, columns=['day_of_week', 'hour_of_day'], drop_first=True)

```

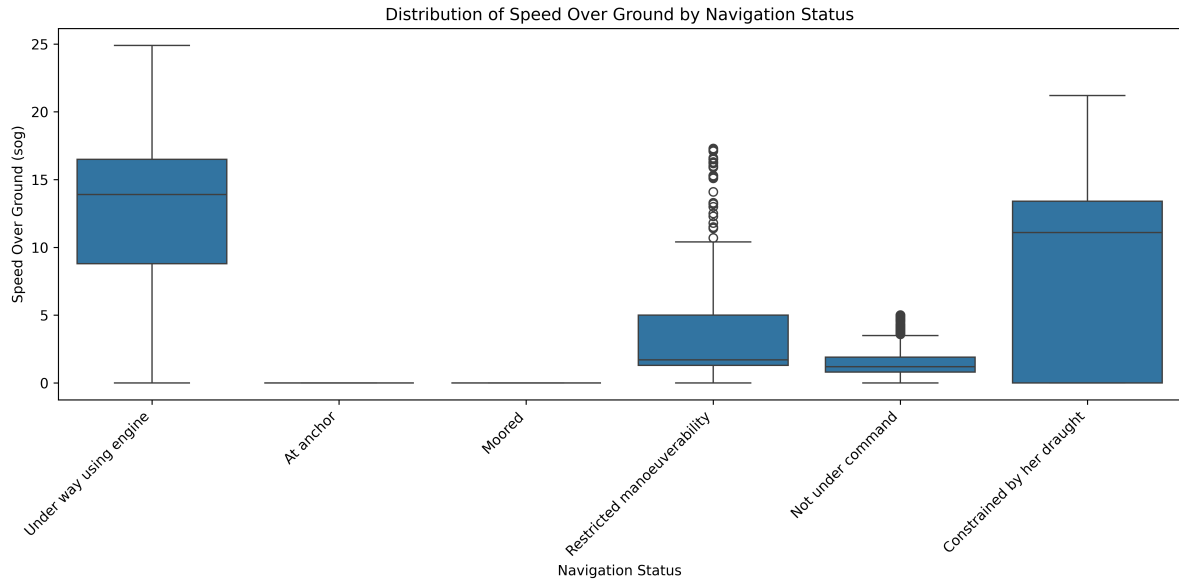


Figure 3: SOG over NAVSTAT after pruning

```

14
15 # Filter out unrealistic speeds
16 ais_train = ais_train[ais_train['sog'] < 25]
17
18 # Map 'navstat' values
19 ais_train['navstat'] = ais_train['navstat'].replace(8, 0) # Under way sailing -> Under way using
    ↳ engine
20 ais_train = ais_train[~((ais_train['navstat'].isin([1, 5])) & (ais_train['sog'] > 0))]
21 ais_train = ais_train[~((ais_train['navstat'] == 2) & (ais_train['sog'] > 5))]
22
23 # One-hot encode 'navstat'
24 ais_train = pd.get_dummies(ais_train, columns=['navstat'])
25
26 # Split cyclic values into x and y
27 ais_train['cog_sin'] = np.sin(np.radians(ais_train['cog']))
28 ais_train['cog_cos'] = np.cos(np.radians(ais_train['cog']))
29 ais_train['heading_sin'] = np.sin(np.radians(ais_train['heading']))
30 ais_train['heading_cos'] = np.cos(np.radians(ais_train['heading']))
31
32 # Merge with vessel data
33 vessels = pd.read_csv("vessels.csv", sep='|')[['shippingLineId', 'vesselId']]
34 vessels['new_id'] = range(len(vessels))
35 vessel_id_to_new_id = dict(zip(vessels['vesselId'], vessels['new_id']))
36 ais_train = pd.merge(ais_train, vessels, on='vesselId', how='left')
37
38 vesselId = "61e9f3bfb937134a3c4bfe9f"
39
40 # Filter the data for the selected vesselId and sort by 'time'
41 vessel_data = ais_train[ais_train['vesselId'] == vesselId].sort_values('time').reset_index(drop=True)
42
43 # Clean vessel_data
44 vessel_data = vessel_data.dropna(subset=['latitude', 'longitude'])
45
46 # Ensure that 'time' is datetime and sorted
47 vessel_data = vessel_data.sort_values('time')
48
49 # Convert 'time' to numeric format for interpolation (seconds since epoch)
50 vessel_data['time_numeric'] = vessel_data['time'].astype(np.int64) // 10**9
51
52 # Ensure that 'time_numeric' is strictly increasing
53 vessel_data = vessel_data.drop_duplicates(subset='time_numeric')
54 vessel_data = vessel_data.sort_values('time_numeric').reset_index(drop=True)
55 plotting_data = vessel_data.sort_values('time').reset_index(drop=True)
56
57 # Remove the last row (as it has NaN bearing)
58 plotting_data = plotting_data[:-1]
59
60 # Create the map figure
61 fig = go.Figure()
62
63 # Add the trajectory line

```



```

64 fig.add_trace(go.Scattermapbox(
65     mode='lines+markers',
66     lon=plotting_data['longitude'],
67     lat=plotting_data['latitude'],
68     marker=dict(size=6, color='blue'),
69     line=dict(width=2, color='blue'),
70     name='Trajectory'
71 ))
72
73 fig.update_layout(
74     mapbox={
75         'style': "open-street-map",
76         'zoom': 1 if interpolate else 4,
77         'center': {'lon': plotting_data['longitude'].mean(), 'lat': plotting_data['latitude'].mean()}
78     },
79     title="Trajectory of vessel {vesselId} with Direction Arrows"
80 )
81
82 fig.show()

```

This script begins by reading the `ais_train.csv` file, which contains multiple vessels' raw AIS (Automatic Identification System) data. The first step is to convert the time column from its original format into Python's `datetime` format, making it easier to manipulate and plot time-series data. We then extracted the features necessary for our analysis, a process we had already refined in earlier stages of our project. Next, the script filters the dataset to isolate the data points corresponding to the vessel with `vesselId = "61e9f3bfb937134a3c4bfe9f"`. To avoid redundancy, we drop duplicate entries before plotting the vessel's trajectory using a map plot.

Figure 4 displays the result of plotting the uninterpolated data points for this vessel. Several key observations can be made from this plot. Firstly, there are instances where some data points are clustered closely together near land, reflecting periods where the vessel was likely stationary or moving slowly. In contrast, other data points are located far from each other, reflecting instances where the ship moved over large distances without frequent AIS updates. For example, the points with the longest distance between them were recorded ten days apart, which creates significant gaps in the vessel's trajectory. These gaps hinder our ability to accurately track the vessel's movements over time, especially for modeling purposes.

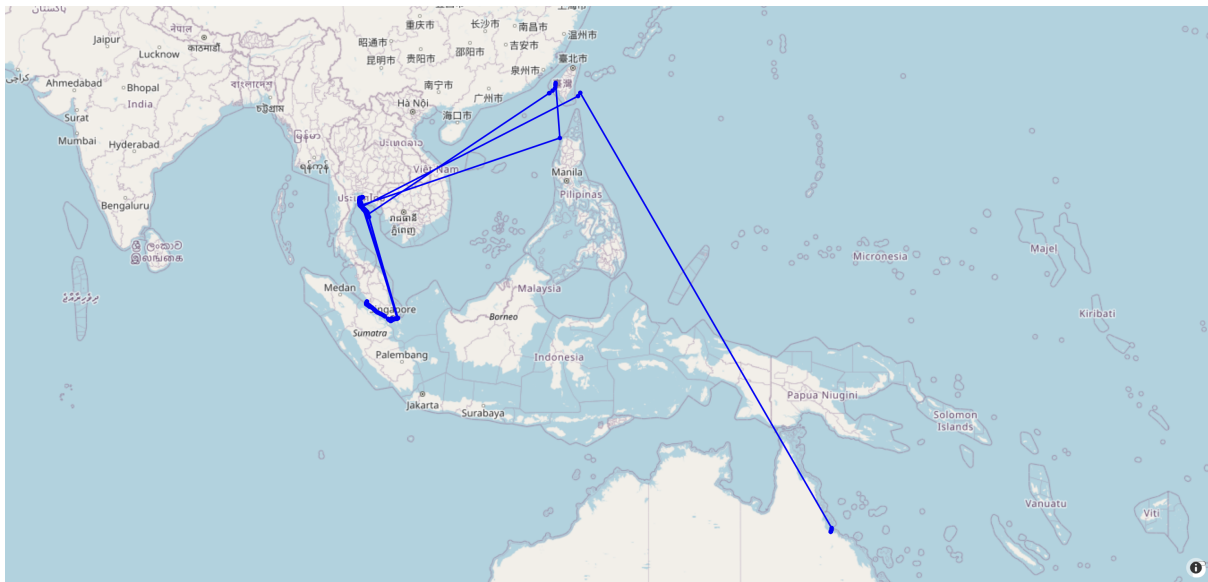


Figure 4: Uninterpolated trajectory of `vesselId=61e9f3bfb937134a3c4bfe9f`

To address this, we applied cubic spline interpolation using Python's `scipy` library. This method allowed us to create smooth interpolated paths between the known data points, effectively filling in the missing data while maintaining the continuity of the vessel's trajectory. Cubic spline interpolation is particularly useful for this type of data because it generates smooth curves that pass through the known points while accounting for the natural fluctuations in the vessel's movement over time.

```

1 from scipy.interpolate import CubicSpline

```

```

2
3 # ... read and preprocess data
4
5 # Create a new time index with daily frequency
6 start_time = vessel_data['time'].min().normalize()
7 end_time = vessel_data['time'].max().normalize()
8 new_time_index = pd.date_range(start=start_time, end=end_time, freq='1D')
9
10 # Convert 'new_time_index' to numeric format
11 new_time_numeric = new_time_index.astype(np.int64) // 10**9
12
13 # Prepare original data for interpolation
14 original_times = vessel_data['time_numeric'].values
15 original_latitudes = vessel_data['latitude'].values
16 original_longitudes = vessel_data['longitude'].values
17
18 # Convert latitude and longitude to radians
19 lat_rad = np.radians(original_latitudes)
20 lon_rad = np.radians(original_longitudes)
21
22 # Convert to 3D Cartesian coordinates on a unit sphere
23 x = np.cos(lat_rad) * np.cos(lon_rad)
24 y = np.cos(lat_rad) * np.sin(lon_rad)
25 z = np.sin(lat_rad)
26
27 # Create cubic spline interpolators for x, y, z
28 cs_x = CubicSpline(original_times, x)
29 cs_y = CubicSpline(original_times, y)
30 cs_z = CubicSpline(original_times, z)
31
32 # Interpolate at new time points within the range of original_times
33 t_min = original_times.min()
34 t_max = original_times.max()
35 valid_mask = (new_time_numeric >= t_min) & (new_time_numeric <= t_max)
36 interp_times_valid = new_time_numeric[valid_mask]
37
38 if len(interp_times_valid) == 0:
39     print("No valid interpolation times within the range of original data.")
40     plotting_data = vessel_data
41 else:
42     x_interp = cs_x(interp_times_valid)
43     y_interp = cs_y(interp_times_valid)
44     z_interp = cs_z(interp_times_valid)
45
46     # Normalize the interpolated coordinates to lie on the unit sphere
47     norm = np.sqrt(x_interp**2 + y_interp**2 + z_interp**2)
48     x_interp /= norm
49     y_interp /= norm
50     z_interp /= norm
51
52     # Convert back to latitude and longitude
53     lat_interp = np.degrees(np.arcsin(z_interp))
54     lon_interp = np.degrees(np.arctan2(y_interp, x_interp))
55
56     # Create interpolated DataFrame
57     vessel_data_interp = pd.DataFrame({
58         'time_numeric': interp_times_valid,
59         'latitude': lat_interp,
60         'longitude': lon_interp
61     })
62
63     # Convert 'time_numeric' back to datetime
64     vessel_data_interp['time'] = pd.to_datetime(vessel_data_interp['time_numeric'], unit='s')
65
66     # Use the interpolated data for plotting
67     plotting_data = vessel_data_interp.sort_values('time').reset_index(drop=True)
68
69 # ... plot data

```

Figure 5 shows the plot using cubic spline interpolation. In this plot, there is a noticeable improvement in the regularity of the data, given that the amount of data has increased by 140%, as now there is at least one data point for each day. The interpolation has filled the gaps between the previously uneven time intervals, creating a smoother trajectory for the vessel. However, despite this improvement, we needed to address some potential sources of error. Specifically, we observed that the interpolated data points occasionally took abrupt or unrealistic turns, with some interpolated paths even crossing over land, which cargo vessels obviously cannot do.

These unrealistic segments in the trajectory are likely due to the nature of cubic spline interpolation, which, while effective for smoothing data, does not account for geographical constraints such as coastlines or landmasses. Cubic splines interpolate based purely on the mathematical relationship between data points without considering the real-world implications of a vessel's path. As a result, the spline

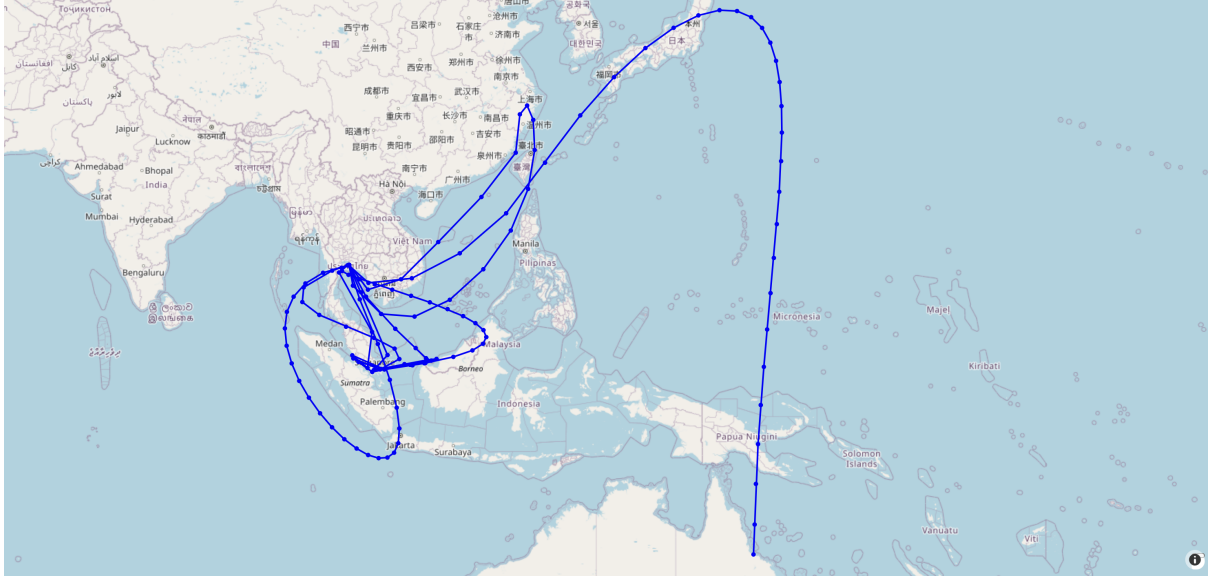


Figure 5: Interpolated trajectory of vesselId=61e9f3bfb937134a3c4bfe9f

interpolation can generate curves that deviate from expected sea routes, especially when there are significant gaps between the original data points or when the vessel's movement near land involves sharp directional changes.

This might be because some data points correspond to periods when the vessel was moored or anchored, meaning the ship was stationary. Interpolating movement data for these periods would be inappropriate, as the boat was not in motion.

3.4 Distribution of Time Gaps Between Consecutive AIS Records

We also wanted to investigate the time intervals apparent in the dataset, as we observed that AIS timestamps were inconsistent. To give ourselves an overview, we plotted the distribution of time gaps between consecutive AIS records (in seconds) on a logarithmic frequency scale 6. From this, we noticed the following:

Dominance of Small Gaps: The most significant number of time differences is concentrated around tiny gaps (near zero). This suggests that, for many records, the data points are closely spaced in time.

Exponential Decline: As the time difference increases, the frequency of occurrence decreases exponentially. This suggests that significant time gaps between consecutive AIS records are less common but still exist.

Long Tail Distribution: There is a long tail towards significant time gaps, reaching as high as 6 million seconds (approximately 69 days). While less frequent, there are still notable occurrences of considerable gaps in the dataset.

From this observation, it's clear that the time difference between data points needs to be accounted for. This allows us to create a model that effectively handles data points in close proximity as well as those with larger intervals.

4 Improved LSTM Model

The next step in our project was to implement our findings from the exploratory analysis into our initial LSTM model.

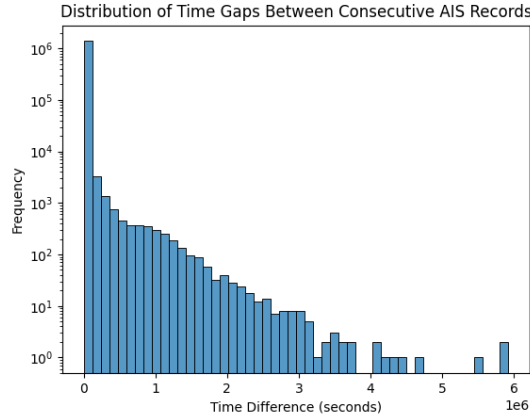


Figure 6: Distribution of Time Gaps Between Consecutive AIS Records

4.1 Feature Engineering

Blant others: This involved preprocessing the data:

- Temporal features were extracted from the AIS timestamps, including the elapsed time in seconds, day of the week, and hour of the day.
- Speed Over Ground (SOG) values were filtered to remove unrealistic speeds, and navigation status ('navstat') values were mapped to handle specific cases.
- 'cog' (course over ground) was transformed into cyclic features using sine and cosine encoding.
- One-hot encoding was applied to categorical features, and Min-Max scaling was used to normalize input and target features.

We also put much effort into transforming the raw AIS data into a structured format better suited for the model. First, we focused on extracting useful time-based features. We converted the `time` column into datetime objects and calculated `elapsed_time` as the number of seconds since the epoch to track vessel movements over time. We also pulled out `day_of_week` and `hour_of_day` from the timestamp and one-hot encoded them to capture weekly and daily patterns in vessel activity. For the COG feature, which is directional, we used sine and cosine transformations to avoid issues with its circular nature.

We mapped each to a unique index for the categorical variables, like `vesselId` and NAVSTAT. We used one-hot encoding to turn them into numbers that the model could understand. We cleaned up the data by filtering out combinations of NAVSTAT and SOG values that were unrealistic to co-occur. We also normalized all inputs and outputs with `MinMaxScaler` to help the model learn more effectively during training.

We added extra information by merging the dataset with external vessel data, including `shippingLineId`. Finally, we organized the data into sequential time steps for each vessel so the LSTM model could detect patterns over time.

To capture temporal dependencies, each vessel was assigned ten time-step sequences. This allowed the model to learn from sequences of consecutive historical data points. Each sequence was structured to predict the vessel's location at the next time step. The data was then split into training and validation sets.

4.2 The Model

The model consists of two bidirectional LSTM layers with 512 and 256 hidden units, respectively, followed by a dropout layer and a fully connected output layer. The bidirectional layers enable the model to learn patterns from forward and backward temporal data. The final output layer predicts the latitude and longitude coordinates.

- **Input Size:** Determined by the number of input features.

- **Hidden Layers:** Two bidirectional LSTM layers with 512 and 256 hidden units.
- **Output Size:** Set to 2, corresponding to latitude and longitude predictions.

We used a custom Haversine loss function, which calculates the distance between predicted and actual coordinates on the Earth's surface. This distance-based metric ensures that the model minimizes geodesic errors in location predictions, providing a more meaningful evaluation than traditional regression loss functions.

The model was trained for 100 epochs with early stopping to avoid overfitting. The Adam optimizer was used to update model weights based on the Haversine loss.

4.3 The Code

```

1  #!/usr/bin/env python3
2
3  import pandas as pd
4  import numpy as np
5  from sklearn.preprocessing import MinMaxScaler
6  from sklearn.model_selection import train_test_split
7  import torch
8  import torch.nn as nn
9  import torch.optim as optim
10 from torch.utils.data import TensorDataset, DataLoader
11 import copy
12
13 print("PyTorch version:", torch.__version__)
14
15 # Set device preference: MPS > CUDA > CPU
16 if torch.backends.mps.is_available():
17     device = torch.device("mps")
18     print("Using MPS device")
19 elif torch.cuda.is_available():
20     device = torch.device("cuda")
21     print("Using CUDA device")
22 else:
23     device = torch.device("cpu")
24     print("Using CPU device")
25
26 """
27     Load and Preprocess Data
28 """
29
30 # Read ais_train.csv
31 ais_train = pd.read_csv("ais_train.csv", sep='|')
32
33
34 vessel_mapping = {vessel: idx for idx, vessel in enumerate(ais_train['vesselId'].unique())}
35
36 # Temporal features
37 ais_train['time'] = pd.to_datetime(ais_train['time'])
38 ais_train['elapsed_time'] = (ais_train['time'] - pd.Timestamp("1970-01-01")) // pd.Timedelta('1s')
39
40 # Filter out unrealistic speeds
41 ais_train = ais_train[ais_train['sog'] < 25]
42
43 # Map 'navstat' values
44 ais_train['navstat'] = ais_train['navstat'].replace(8, 0) # Under way sailing -> Under way using
↪ engine
45 ais_train = ais_train[~((ais_train['navstat'].isin([1, 5])) & (ais_train['sog'] > 0))]
46 ais_train = ais_train[~((ais_train['navstat'] == 2) & (ais_train['sog'] > 5))]
47
48 # One-hot encode 'navstat'
49 ais_train = pd.get_dummies(ais_train, columns=['navstat'])
50
51 # Merge with vessel data
52 vessels = pd.read_csv("vessels.csv", sep='|')[['shippingLineId', 'vesselId']]
53 vessels['new_id'] = range(len(vessels))
54 vessel_id_to_new_id = dict(zip(vessels['vesselId'], vessels['new_id']))
55 ais_train = pd.merge(ais_train, vessels, on='vesselId', how='left')
56
57 # Temporal features
58 ais_train['day_of_week'] = ais_train['time'].dt.dayofweek
59 ais_train['hour_of_day'] = ais_train['time'].dt.hour
60 ais_train = pd.get_dummies(ais_train, columns=['day_of_week', 'hour_of_day'], drop_first=True)
61
62 # Handle cyclic features for 'cog'
63 ais_train['cog_sin'] = np.sin(np.radians(ais_train['cog']))
64 ais_train['cog_cos'] = np.cos(np.radians(ais_train['cog']))
65
66 # Merge with vessels and ports data
67 ais_train = pd.merge(ais_train, vessels, on='vesselId', how='left')
68 ais_train['vesselId'] = ais_train['vesselId'].map(vessel_mapping)

```

```

69
70 # Define input and target features
71 input_features = [
72     'latitude', 'longitude', 'sog', 'cog_sin', 'cog_cos', 'elapsed_time',
73     'vesselId'
74 ]
75
76 input_features.extend([col for col in ais_train.columns if 'day_of_week_' in col])
77 input_features.extend([col for col in ais_train.columns if 'hour_of_day_' in col])
78
79 target_columns = ['latitude', 'longitude']
80
81 # Initialize scalers
82 scaler_input = MinMaxScaler()
83 scaler_output = MinMaxScaler()
84
85 # Drop rows with NaN values in input features
86 ais_train = ais_train.dropna(subset=input_features + target_columns)
87
88 # Scale input and output features
89 input_data = scaler_input.fit_transform(ais_train[input_features])
90 output_data = scaler_output.fit_transform(ais_train[target_columns])
91
92 # Add scaled features back to DataFrame
93 ais_train[input_features] = input_data
94 ais_train[target_columns] = output_data
95
96 """
97     Create Sequences for Model Training
98 """
99
100 # Function to create sequences per vessel
101 def create_sequences_per_vessel(df, time_steps):
102     X, y = [], []
103     vessel_ids = df['vesselId'].unique()
104     for vessel_id in vessel_ids:
105         vessel_data = df[df['vesselId'] == vessel_id].sort_values('elapsed_time')
106         inputs = vessel_data[input_features].values
107         targets = vessel_data[target_columns].values
108         if len(inputs) < time_steps:
109             continue # Skip sequences shorter than time_steps
110         for i in range(len(inputs) - time_steps):
111             X.append(inputs[i:i + time_steps])
112             y.append(targets[i + time_steps])
113     return np.array(X), np.array(y)
114
115 # Create sequences
116 time_step = 10
117 X, y = create_sequences_per_vessel(ais_train, time_step)
118
119 # Split into training and validation sets
120 X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2, shuffle=True)
121
122 # Convert to PyTorch tensors
123 X_train = torch.from_numpy(X_train).float()
124 y_train = torch.from_numpy(y_train).float()
125 X_val = torch.from_numpy(X_val).float()
126 y_val = torch.from_numpy(y_val).float()
127
128 # Create TensorDatasets and DataLoaders
129 batch_size = 128
130 train_dataset = TensorDataset(X_train, y_train)
131 val_dataset = TensorDataset(X_val, y_val)
132 train_loader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True)
133 val_loader = DataLoader(val_dataset, batch_size=batch_size, shuffle=False)
134
135 """
136     Haversine Loss Function
137 """
138
139 def haversine_loss(y_true, y_pred):
140     R = 6371.0 # Earth radius in kilometers
141
142     # Ensure constants are tensors of the same dtype and device as y_true
143     pi_over_180 = torch.tensor(np.pi / 180.0, dtype=y_true.dtype, device=y_true.device)
144
145     lat_true = y_true[:, 0] * pi_over_180
146     lon_true = y_true[:, 1] * pi_over_180
147     lat_pred = y_pred[:, 0] * pi_over_180
148     lon_pred = y_pred[:, 1] * pi_over_180
149
150     dlat = lat_pred - lat_true
151     dlon = lon_pred - lon_true
152
153     a = torch.sin(dlat / 2) ** 2 + torch.cos(lat_true) * torch.cos(lat_pred) * torch.sin(dlon / 2) ** 2
154     c = 2 * torch.atan2(torch.sqrt(a), torch.sqrt(1 - a))
155     distance = R * c

```

```

156
157     # Return mean distance over the batch
158     return torch.mean(distance)
159
160 """
161 Define and Train the Model
162 """
163
164 class LSTMModel(nn.Module):
165     def __init__(self, input_size, hidden_size1, hidden_size2, output_size):
166         super(LSTMModel, self).__init__()
167         # Set bidirectional=True and adjust hidden sizes accordingly
168         self.lstm1 = nn.LSTM(
169             input_size, hidden_size1, batch_first=True, bidirectional=True
170         )
171         self.lstm2 = nn.LSTM(
172             hidden_size1 * 2, hidden_size2, batch_first=True, bidirectional=True
173         )
174         self.dropout = nn.Dropout(0.2)
175         # Adjust the input size of the fully connected layer
176         self.fc = nn.Linear(hidden_size2 * 2, output_size)
177
178     def forward(self, x):
179         out, _ = self.lstm1(x)
180         out, _ = self.lstm2(out)
181         out = self.dropout(out[:, -1, :]) # Get the last time step
182         out = self.fc(out)
183         return out
184
185 # Initialize the model
186 input_size = X_train.shape[2]
187 hidden_size1 = 512
188 hidden_size2 = 256
189 output_size = y_train.shape[1]
190 model = LSTMModel(input_size, hidden_size1, hidden_size2, output_size).to(device)
191
192 # Define optimizer
193 optimizer = optim.Adam(model.parameters())
194
195 # Training loop
196 num_epochs = 100
197 best_val_loss = float('inf')
198 patience = 5
199 counter = 0
200 best_model_wts = copy.deepcopy(model.state_dict())
201
202 for epoch in range(num_epochs):
203     # Training
204     model.train()
205     train_losses = []
206
207     for batch_X, batch_y in train_loader:
208         batch_X = batch_X.to(device)
209         batch_y = batch_y.to(device)
210
211         optimizer.zero_grad()
212         outputs = model(batch_X)
213         loss = haversine_loss(batch_y, outputs)
214         loss.backward()
215         optimizer.step()
216         train_losses.append(loss.item())
217
218     avg_train_loss = np.mean(train_losses)
219
220     # Validation
221     model.eval()
222     val_losses = []
223     with torch.no_grad():
224         for batch_X, batch_y in val_loader:
225             batch_X = batch_X.to(device)
226             batch_y = batch_y.to(device)
227             outputs = model(batch_X)
228             loss = haversine_loss(batch_y, outputs)
229             val_losses.append(loss.item())
230
231     avg_val_loss = np.mean(val_losses)
232     print(f"Epoch [{epoch + 1}/{num_epochs}], Train Loss: {avg_train_loss:.4f}, Val Loss:
233           ↪ {avg_val_loss:.4f}")
234
235     # Early stopping
236     if avg_val_loss < best_val_loss:
237         best_val_loss = avg_val_loss
238         best_model_wts = copy.deepcopy(model.state_dict())
239         counter = 0
240     else:
241         counter += 1
242         if counter >= patience:

```

```

242         print("Early stopping")
243         break
244
245     # Load best model weights
246     model.load_state_dict(best_model_wts)
247
248     """
249     Prepare Test Data and Make Predictions
250     """
251
252     # Load test data
253     ais_test = pd.read_csv("ais_test.csv")
254     ais_test['time'] = pd.to_datetime(ais_test['time'])
255     ais_test['elapsed_time'] = (ais_test['time'] - pd.Timestamp("1970-01-01")) // pd.Timedelta('1s')
256     ais_test['new_id'] = ais_test['vesselId'].map(vessel_id_to_new_id)
257
258     ais_test['vesselId'] = ais_test['vesselId'].map(vessel_mapping)
259     ais_test['day_of_week'] = ais_test['time'].dt.dayofweek
260     ais_test['hour_of_day'] = ais_test['time'].dt.hour
261
262     # One-hot encode
263     ais_test = pd.get_dummies(ais_test, columns=['day_of_week', 'hour_of_day'], drop_first=True)
264
265     # Merge with vessels and ports data
266     ais_test = pd.merge(ais_test, vessels, on='vesselId', how='left')
267
268     # Ensure all columns in ais_test match those in input_features
269     for col in input_features:
270         if col not in ais_test.columns:
271             ais_test[col] = 0
272
273     # Scale the test data using the same scaler
274     input_data_test = scaler_input.transform(ais_test[input_features])
275     ais_test_scaled = ais_test.copy()
276     ais_test_scaled[input_features] = input_data_test
277
278     # Prepare sequences for each vessel in the test set
279     def create_sequences_for_test(df_train, df_test, time_steps):
280         X_test = []
281         test_ids = []
282         for idx, row in df_test.iterrows():
283             vessel_id = row['vesselId']
284             current_time = row['elapsed_time']
285
286             # Get the historical data for this vessel up to the current_time
287             vessel_train_data = df_train[df_train['vesselId'] == vessel_id]
288             vessel_test_data = df_test[df_test['vesselId'] == vessel_id]
289
290             # Combine and sort
291             vessel_data = pd.concat([vessel_train_data, vessel_test_data], ignore_index=True)
292             vessel_data = vessel_data.sort_values('elapsed_time')
293
294             # Select data up to the current_time (excluding the current row)
295             historical_data = vessel_data[vessel_data['elapsed_time'] < current_time]
296
297             # Get the last 'time_steps' entries
298             historical_sequence = historical_data.tail(time_steps)[input_features].values
299
300             if len(historical_sequence) < time_steps:
301                 # Pad with zeros if not enough historical data
302                 padding = np.zeros((time_steps - len(historical_sequence), len(input_features)))
303                 historical_sequence = np.vstack([padding, historical_sequence])
304
305             X_test.append(historical_sequence)
306             test_ids.append(row['ID']) # Assuming 'ID' is unique per test row
307
308         return np.array(X_test), test_ids
309
310
311     # Create sequences for each test row
312     X_test, test_ids = create_sequences_for_test(ais_train, ais_test_scaled, time_step)
313
314     # Convert to PyTorch tensor
315     X_test = torch.from_numpy(X_test).float()
316
317     # Create a DataLoader for test data
318     test_dataset = TensorDataset(X_test)
319     test_loader = DataLoader(test_dataset, batch_size=128, shuffle=False)
320
321     # Make predictions in batches
322     predictions = []
323
324     model.eval()
325     with torch.no_grad():
326         for batch_X in test_loader:
327             batch_X = batch_X[0].to(device) # batch_X is a tuple
328             outputs = model(batch_X)

```



```

329         predictions.append(outputs.cpu().numpy())
330
331     # Concatenate all batch predictions
332     y_pred = np.concatenate(predictions, axis=0)
333
334     # Inverse transform predictions
335     y_pred_inverse = scaler_output.inverse_transform(y_pred)
336
337
338     # Prepare submission
339     submission_df = pd.DataFrame({
340         'ID': test_ids,
341         'longitude_predicted': y_pred_inverse[:, target_columns.index('longitude')],
342         'latitude_predicted': y_pred_inverse[:, target_columns.index('latitude')]
343     })
344
345     # Ensure the submission file has the required columns
346     submission_df = submission_df[['ID', 'longitude_predicted', 'latitude_predicted']]
347
348     # Save submission file
349     submission_df.to_csv("submission.csv", index=False)
350
351     # Display submission
352     print(submission_df.head())
353     print(f"Submission DataFrame shape: {submission_df.shape}")
354
355     print(f"Number of predictions: {len(y_pred_inverse)}")
356     print(f"Number of test IDs: {len(test_ids)}")
357     assert len(y_pred_inverse) == len(test_ids), "Mismatch between predictions and test IDs"

```

4.4 Result

After implementing the steps described in this section, we reduced our Kaggle score from 600 to 170.

4.5 Integrating Cubic Spline Interpolation into LSTM

As we discovered in our exploratory analysis, cubic spline interpolation could potentially smooth irregularities and fill in gaps within the dataset. Below is our implementation of the cubic spline interpolation integration.

```

1  import pandas as pd
2  import numpy as np
3  from sklearn.preprocessing import MinMaxScaler
4  from sklearn.model_selection import train_test_split
5  import torch
6  import torch.nn as nn
7  import torch.optim as optim
8  from torch.utils.data import TensorDataset, DataLoader
9  import copy
10 from scipy.interpolate import CubicSpline
11 from tqdm import tqdm
12
13 print("PyTorch version:", torch.__version__)
14
15 # Set device preference: MPS > CUDA > CPU
16 if torch.backends.mps.is_available():
17     device = torch.device("mps")
18     print("Using MPS device")
19 elif torch.cuda.is_available():
20     device = torch.device("cuda")
21     print("Using CUDA device")
22 else:
23     device = torch.device("cpu")
24     print("Using CPU device")
25
26 """
27     Load and Preprocess Data
28 """
29
30 # Read ais_train.csv
31 ais_train = pd.read_csv("ais_train.csv", sep='|')
32
33
34 vessel_mapping = {vessel: idx for idx, vessel in enumerate(ais_train['vesselId'].unique())}
35
36 # Temporal features
37 ais_train['time'] = pd.to_datetime(ais_train['time'])
38 ais_train['elapsed_time'] = (ais_train['time'] - pd.Timestamp("1970-01-01")) // pd.Timedelta('1s')
39
40 # Filter out unrealistic speeds
41 ais_train = ais_train[ais_train['sog'] < 25]
42

```

```

43 # Map 'navstat' values
44 ais_train['navstat'] = ais_train['navstat'].replace(8, 0) # Under way sailing -> Under way using
↳ engine
45 ais_train = ais_train[~((ais_train['navstat'].isin([1, 5])) & (ais_train['sog'] > 0))]
46 ais_train = ais_train[~((ais_train['navstat'] == 2) & (ais_train['sog'] > 5))]
47
48 # One-hot encode 'navstat'
49 ais_train = pd.get_dummies(ais_train, columns=['navstat'])
50
51 # Merge with vessel data
52 vessels = pd.read_csv("vessels.csv", sep='|')[['shippingLineId', 'vesselId']]
53 vessels['new_id'] = range(len(vessels))
54 vessel_id_to_new_id = dict(zip(vessels['vesselId'], vessels['new_id']))
55 ais_train = pd.merge(ais_train, vessels, on='vesselId', how='left')
56
57 # Merge port data
58 ports = pd.read_csv("ports.csv", sep='|')[['portId', 'latitude', 'longitude']]
59 ports = ports.rename(columns={'latitude': 'port_latitude', 'longitude': 'port_longitude'})
60 ais_train = pd.merge(ais_train, ports, on='portId', how='left')
61 ais_train = ais_train[ais_train['portId'].isnull()] # Remove rows with null ports
62
63 def haversine_distance(lat1, lon1, lat2, lon2):
64     # Earth radius in nautical miles
65     R = 3440.065
66     lat1, lon1, lat2, lon2 = map(np.radians, [lat1, lon1, lat2, lon2])
67     dlat = lat2 - lat1
68     dlon = lon2 - lon1
69     a = np.sin(dlat / 2.0) ** 2 + np.cos(lat1) * np.cos(lat2) * np.sin(dlon / 2.0) ** 2
70     return 2 * R * np.arcsin(np.sqrt(a))
71
72 def calculate_bearing(lat1, lon1, lat2, lon2):
73     lat1_rad, lat2_rad = np.radians(lat1), np.radians(lat2)
74     diff_long = np.radians(lon2 - lon1)
75     x = np.sin(diff_long) * np.cos(lat2_rad)
76     y = np.cos(lat1_rad) * np.sin(lat2_rad) - (np.sin(lat1_rad) * np.cos(lat2_rad) * np.cos(diff_long))
77     initial_bearing = np.arctan2(x, y)
78     return (np.degrees(initial_bearing) + 360) % 360
79
80 """
81 Cubic Spline Interpolation for Each Vessel
82 """
83
84 # List to store processed trajectories
85 processed_trajectories = []
86
87 # Group data by vesselId
88 vessel_ids = ais_train['vesselId'].unique()
89 for vessel_id in tqdm(vessel_ids, desc="Interpolating Vessels"):
90
91     vessel_data = ais_train[ais_train['vesselId'] == vessel_id].sort_values('elapsed_time')
92
93     # Ensure at least two data points
94     if len(vessel_data) < 2:
95         continue
96
97     # Prepare data for interpolation
98     times = vessel_data['elapsed_time'].values
99     latitudes = vessel_data['latitude'].values
100     longitudes = vessel_data['longitude'].values
101
102     # Remove duplicates in times
103     times, unique_indices = np.unique(times, return_index=True)
104     latitudes = latitudes[unique_indices]
105     longitudes = longitudes[unique_indices]
106
107     if len(times) < 2:
108         continue
109
110     # Convert lat/lon to radians
111     lat_rad = np.radians(latitudes)
112     lon_rad = np.radians(longitudes)
113
114     # Convert to 3D Cartesian coordinates
115     x = np.cos(lat_rad) * np.cos(lon_rad)
116     y = np.cos(lat_rad) * np.sin(lon_rad)
117     z = np.sin(lat_rad)
118
119     # Create new time points for interpolation (every hour)
120     start_time = times.min()
121     end_time = times.max()
122     new_times = np.arange(start_time, end_time + 1, 3600) # Every hour in seconds
123     # Ensure new_times are within the original times
124     new_times = new_times[(new_times >= times.min()) & (new_times <= times.max())]
125
126     # Interpolate x, y, z using cubic splines
127     cs_x = CubicSpline(times, x)
128     cs_y = CubicSpline(times, y)

```

```

129     cs_z = CubicSpline(times, z)
130
131     x_interp = cs_x(new_times)
132     y_interp = cs_y(new_times)
133     z_interp = cs_z(new_times)
134
135     # Normalize to unit sphere
136     norm = np.sqrt(x_interp**2 + y_interp**2 + z_interp**2)
137     x_interp /= norm
138     y_interp /= norm
139     z_interp /= norm
140
141     # Convert back to lat/lon
142     lat_interp = np.degrees(np.arcsin(z_interp))
143     lon_interp = np.degrees(np.arctan2(y_interp, x_interp))
144
145     # Handle longitude wrap-around
146     lon_interp = (lon_interp + 360) % 360
147     # Adjust longitudes > 180 to negative values (from -180 to 180)
148     lon_interp[lon_interp > 180] -= 360
149
150     # Create interpolated DataFrame
151     interp_df = pd.DataFrame({
152         'vesselId': vessel_id,
153         'elapsed_time': new_times,
154         'latitude': lat_interp,
155         'longitude': lon_interp,
156     })
157
158     # Assign 'portId' using merge_asof
159     vessel_port_data = vessel_data[['elapsed_time',
160     ↪ 'portId']].drop_duplicates().sort_values('elapsed_time')
161     interp_df = pd.merge_asof(
162         interp_df.sort_values('elapsed_time'),
163         vessel_port_data,
164         on='elapsed_time',
165         direction='nearest'
166     )
167
168     # Recalculate 'sog' and 'cog'
169     lat_prev = np.roll(lat_interp, 1)
170     lon_prev = np.roll(lon_interp, 1)
171     time_prev = np.roll(new_times, 1)
172     distances = haversine_distance(lat_prev, lon_prev, lat_interp, lon_interp)
173     time_diffs = (new_times - time_prev) / 3600 # Convert time difference to hours
174     time_diffs[0] = np.nan # First element has no previous point
175     sog = distances / time_diffs # Speed in knots
176     cog = calculate_bearing(lat_prev, lon_prev, lat_interp, lon_interp)
177     cog[0] = np.nan # First element has no previous point
178
179     # Assign 'sog' and 'cog' to interp_df
180     interp_df['sog'] = sog
181     interp_df['cog'] = cog
182
183     # Drop the first row as it has NaN values
184     interp_df = interp_df.iloc[1:].reset_index(drop=True)
185
186     # Convert 'elapsed_time' back to 'time'
187     interp_df['time'] = pd.to_datetime(interp_df['elapsed_time'], unit='s')
188
189     # Append to processed_trajectories
190     processed_trajectories.append(interp_df)
191
192     # Combine all interpolated data
193     ais_train_interpolated = pd.concat(processed_trajectories, ignore_index=True)
194
195     print("ais_train", len(ais_train))
196     print("ais_train_interpolated", len(ais_train_interpolated))
197
198     """
199     Continue Preprocessing with Interpolated Data
200     """
201
202     # Temporal features
203     ais_train_interpolated['day_of_week'] = ais_train_interpolated['time'].dt.dayofweek
204     ais_train_interpolated['hour_of_day'] = ais_train_interpolated['time'].dt.hour
205     ais_train_interpolated = pd.get_dummies(ais_train_interpolated, columns=['day_of_week',
206     ↪ 'hour_of_day'], drop_first=True)
207
208     # Handle cyclic features for 'cog'
209     ais_train_interpolated['cog_sin'] = np.sin(np.radians(ais_train_interpolated['cog']))
210     ais_train_interpolated['cog_cos'] = np.cos(np.radians(ais_train_interpolated['cog']))
211
212     # Merge with vessels and ports data
213     ais_train_interpolated = pd.merge(ais_train_interpolated, vessels, on='vesselId', how='left')
214     ais_train_interpolated = pd.merge(ais_train_interpolated, ports, on='portId', how='left')

```

```

214 # Calculate 'distance_to_port' and 'bearing_to_port'
215 ais_train_interpolated['distance_to_port'] = haversine_distance(
216     ais_train_interpolated['latitude'], ais_train_interpolated['longitude'],
217     ais_train_interpolated['port_latitude'], ais_train_interpolated['port_longitude']
218 )
219 ais_train_interpolated['bearing_to_port'] = calculate_bearing(
220     ais_train_interpolated['latitude'], ais_train_interpolated['longitude'],
221     ais_train_interpolated['port_latitude'], ais_train_interpolated['port_longitude']
222 )
223
224 ais_train_interpolated['vesselId'] = ais_train['vesselId'].map(vessel_mapping)
225
226 # Define input and target features
227 input_features = [
228     'latitude', 'longitude', 'sog', 'cog_sin', 'cog_cos', 'elapsed_time',
229     'distance_to_port', 'bearing_to_port', 'vesselId'
230 ]
231 input_features.extend([col for col in ais_train_interpolated.columns if 'day_of_week_' in col])
232 input_features.extend([col for col in ais_train_interpolated.columns if 'hour_of_day_' in col])
233
234 target_columns = ['latitude', 'longitude']
235
236 # Initialize scalers
237 scaler_input = MinMaxScaler()
238 scaler_output = MinMaxScaler()
239
240 # Drop rows with NaN values in input features
241 ais_train_interpolated = ais_train_interpolated.dropna(subset=input_features + target_columns)
242
243 # Scale input and output features
244 input_data = scaler_input.fit_transform(ais_train_interpolated[input_features])
245 output_data = scaler_output.fit_transform(ais_train_interpolated[target_columns])
246
247 # Add scaled features back to DataFrame
248 ais_train_interpolated[input_features] = input_data
249 ais_train_interpolated[target_columns] = output_data
250
251 """
252     Create Sequences for Model Training
253 """
254
255 # Function to create sequences per vessel
256 def create_sequences_per_vessel(df, time_steps):
257     X, y = [], []
258     vessel_ids = df['vesselId'].unique()
259     for vessel_id in vessel_ids:
260         vessel_data = df[df['vesselId'] == vessel_id].sort_values('elapsed_time')
261         inputs = vessel_data[input_features].values
262         targets = vessel_data[target_columns].values
263         if len(inputs) < time_steps:
264             continue # Skip sequences shorter than time_steps
265         for i in range(len(inputs) - time_steps):
266             X.append(inputs[i:i + time_steps])
267             y.append(targets[i + time_steps])
268     return np.array(X), np.array(y)
269
270 # Create sequences
271 time_step = 10
272 X, y = create_sequences_per_vessel(ais_train_interpolated, time_step)
273
274 # Split into training and validation sets
275 X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2, shuffle=True)
276
277 # Convert to PyTorch tensors
278 X_train = torch.from_numpy(X_train).float()
279 y_train = torch.from_numpy(y_train).float()
280 X_val = torch.from_numpy(X_val).float()
281 y_val = torch.from_numpy(y_val).float()
282
283 # Create TensorDatasets and DataLoaders
284 batch_size = 128
285 train_dataset = TensorDataset(X_train, y_train)
286 val_dataset = TensorDataset(X_val, y_val)
287 train_loader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True)
288 val_loader = DataLoader(val_dataset, batch_size=batch_size, shuffle=False)
289
290 """
291     Haversine Loss Function
292 """
293
294 def haversine_loss(y_true, y_pred):
295     R = 6371.0 # Earth radius in kilometers
296
297     # Ensure constants are tensors of the same dtype and device as y_true
298     pi_over_180 = torch.tensor(np.pi / 180.0, dtype=y_true.dtype, device=y_true.device)
299
300     lat_true = y_true[:, 0] * pi_over_180

```

```

301     lon_true = y_true[:, 1] * pi_over_180
302     lat_pred = y_pred[:, 0] * pi_over_180
303     lon_pred = y_pred[:, 1] * pi_over_180
304
305     dlat = lat_pred - lat_true
306     dlon = lon_pred - lon_true
307
308     a = torch.sin(dlat / 2) ** 2 + torch.cos(lat_true) * torch.cos(lat_pred) * torch.sin(dlon / 2) ** 2
309     c = 2 * torch.atan2(torch.sqrt(a), torch.sqrt(1 - a))
310     distance = R * c
311
312     # Return mean distance over the batch
313     return torch.mean(distance)
314
315 """
316 Define and Train the Model
317 """
318
319 class LSTMModel(nn.Module):
320     def __init__(self, input_size, hidden_size1, hidden_size2, output_size):
321         super(LSTMModel, self).__init__()
322         # Set bidirectional=True and adjust hidden sizes accordingly
323         self.lstm1 = nn.LSTM(
324             input_size, hidden_size1, batch_first=True, bidirectional=True
325         )
326         self.lstm2 = nn.LSTM(
327             hidden_size1 * 2, hidden_size2, batch_first=True, bidirectional=True
328         )
329         self.dropout = nn.Dropout(0.2)
330         # Adjust the input size of the fully connected layer
331         self.fc = nn.Linear(hidden_size2 * 2, output_size)
332
333     def forward(self, x):
334         out, _ = self.lstm1(x)
335         out, _ = self.lstm2(out)
336         out = self.dropout(out[:, -1, :]) # Get the last time step
337         out = self.fc(out)
338         return out
339
340 # Initialize the model
341 input_size = X_train.shape[2]
342 hidden_size1 = 512
343 hidden_size2 = 256
344 output_size = y_train.shape[1]
345 model = LSTMModel(input_size, hidden_size1, hidden_size2, output_size).to(device)
346
347 # Define optimizer
348 optimizer = optim.Adam(model.parameters())
349
350 # Training loop
351 num_epochs = 100
352 best_val_loss = float('inf')
353 patience = 5
354 counter = 0
355 best_model_wts = copy.deepcopy(model.state_dict())
356
357 for epoch in range(num_epochs):
358     # Training
359     model.train()
360     train_losses = []
361
362     for batch_X, batch_y in train_loader:
363         batch_X = batch_X.to(device)
364         batch_y = batch_y.to(device)
365
366         optimizer.zero_grad()
367         outputs = model(batch_X)
368         loss = haversine_loss(batch_y, outputs)
369         loss.backward()
370         optimizer.step()
371         train_losses.append(loss.item())
372
373     avg_train_loss = np.mean(train_losses)
374
375     # Validation
376     model.eval()
377     val_losses = []
378     with torch.no_grad():
379         for batch_X, batch_y in val_loader:
380             batch_X = batch_X.to(device)
381             batch_y = batch_y.to(device)
382             outputs = model(batch_X)
383             loss = haversine_loss(batch_y, outputs)
384             val_losses.append(loss.item())
385
386     avg_val_loss = np.mean(val_losses)

```

```

387     print(f"Epoch [{epoch + 1}/{num_epochs}], Train Loss: {avg_train_loss:.4f}, Val Loss:
↪     {avg_val_loss:.4f}")
388
389     # Early stopping
390     if avg_val_loss < best_val_loss:
391         best_val_loss = avg_val_loss
392         best_model_wts = copy.deepcopy(model.state_dict())
393         counter = 0
394     else:
395         counter += 1
396         if counter >= patience:
397             print("Early stopping")
398             break
399
400     # Load best model weights
401     model.load_state_dict(best_model_wts)
402
403     """
404     Prepare Test Data and Make Predictions
405     """
406
407     # Load test data
408     ais_test = pd.read_csv("ais_test.csv")
409     ais_test['time'] = pd.to_datetime(ais_test['time'])
410     ais_test['elapsed_time'] = (ais_test['time'] - pd.Timestamp("1970-01-01")) // pd.Timedelta('1s')
411     ais_test['new_id'] = ais_test['vesselId'].map(vessel_id_to_new_id)
412
413     ais_test['vesselId'] = ais_test['vesselId'].map(vessel_mapping)
414     ais_test['day_of_week'] = ais_test['time'].dt.dayofweek
415     ais_test['hour_of_day'] = ais_test['time'].dt.hour
416
417     # One-hot encode
418     ais_test = pd.get_dummies(ais_test, columns=['day_of_week', 'hour_of_day'], drop_first=True)
419
420     # Merge with vessels and ports data
421     ais_test = pd.merge(ais_test, vessels, on='vesselId', how='left')
422
423     # Ensure all columns in ais_test match those in input_features
424     for col in input_features:
425         if col not in ais_test.columns:
426             ais_test[col] = 0
427
428     # Scale the test data using the same scaler
429     input_data_test = scaler_input.transform(ais_test[input_features])
430     ais_test_scaled = ais_test.copy()
431     ais_test_scaled[input_features] = input_data_test
432
433     # Prepare sequences for each vessel in the test set
434     def create_sequences_for_test(df_train, df_test, time_steps):
435         X_test = []
436         test_ids = []
437         for idx, row in df_test.iterrows():
438             vessel_id = row['vesselId']
439             current_time = row['elapsed_time']
440
441             # Get the historical data for this vessel up to the current_time
442             vessel_train_data = df_train[df_train['vesselId'] == vessel_id]
443             vessel_test_data = df_test[df_test['vesselId'] == vessel_id]
444
445             # Combine and sort
446             vessel_data = pd.concat([vessel_train_data, vessel_test_data], ignore_index=True)
447             vessel_data = vessel_data.sort_values('elapsed_time')
448
449             # Select data up to the current_time (excluding the current row)
450             historical_data = vessel_data[vessel_data['elapsed_time'] < current_time]
451
452             # Get the last 'time_steps' entries
453             historical_sequence = historical_data.tail(time_steps)[input_features].values
454
455             if len(historical_sequence) < time_steps:
456                 # Pad with zeros if not enough historical data
457                 padding = np.zeros((time_steps - len(historical_sequence), len(input_features)))
458                 historical_sequence = np.vstack([padding, historical_sequence])
459
460             X_test.append(historical_sequence)
461             test_ids.append(row['ID']) # Assuming 'ID' is unique per test row
462
463         return np.array(X_test), test_ids
464
465
466     # Create sequences for each test row
467     X_test, test_ids = create_sequences_for_test(ais_train_interpolated, ais_test_scaled, time_step)
468
469     # Convert to PyTorch tensor
470     X_test = torch.from_numpy(X_test).float()
471
472     # Create a DataLoader for test data

```

```

473 test_dataset = TensorDataset(X_test)
474 test_loader = DataLoader(test_dataset, batch_size=128, shuffle=False)
475
476 # Make predictions in batches
477 predictions = []
478
479 model.eval()
480 with torch.no_grad():
481     for batch_X in test_loader:
482         batch_X = batch_X[0].to(device) # batch_X is a tuple
483         outputs = model(batch_X)
484         predictions.append(outputs.cpu().numpy())
485
486 # Concatenate all batch predictions
487 y_pred = np.concatenate(predictions, axis=0)
488
489 # Inverse transform predictions
490 y_pred_inverse = scaler_output.inverse_transform(y_pred)
491
492
493 # Prepare submission
494 submission_df = pd.DataFrame({
495     'ID': test_ids,
496     'longitude_predicted': y_pred_inverse[:, target_columns.index('longitude')],
497     'latitude_predicted': y_pred_inverse[:, target_columns.index('latitude')]
498 })
499
500 # Ensure the submission file has the required columns
501 submission_df = submission_df[['ID', 'longitude_predicted', 'latitude_predicted']]
502
503 # Save submission file
504 submission_df.to_csv("submission.csv", index=False)
505
506 # Display submission
507 print(submission_df.head())
508 print(f"Submission DataFrame shape: {submission_df.shape}")
509
510 print(f"Number of predictions: {len(y_pred_inverse)}")
511 print(f"Number of test IDs: {len(test_ids)}")
512 assert len(y_pred_inverse) == len(test_ids), "Mismatch between predictions and test IDs"

```

This yielded a Kaggle score of 1347, which was quite a poor result. This indicated that we either needed a more advanced interpolation technique or that interpolation might not be the right approach for handling stationary intervals.

4.6 BiGRU

We also tried variations where we changed the LSTM layers with GRU in combinations with and without bidirectional layers. LSTM proved to be the best in all cases.

5 Darts and LightGBM

After spending considerable time on deep learning models, we explored time-series models. Specifically, we chose to try the Darts library with LightGBM.

5.1 Feature Engineering

For the Darts-model, a slightly different feature engineering-approach was used.

A critical step in feature engineering here is data resampling to an hourly frequency. By aggregating the `latitude` and `longitude` coordinates using the mean and applying cubic interpolation, the code addresses irregular time intervals and fills in missing values, resulting in uniformly spaced data points that enhance the model's ability to detect patterns. Creating `Darts TimeSeries` objects encapsulates these processed features, facilitating seamless integration with the `LightGBMModel` for time series forecasting.

To maintain consistency and relevance, the code filters the data to include only those vessel IDs in the training and testing sets, thereby focusing the model on vessels with available historical and predictive data. Grouping the training data by `vesselId` allows for vessel-specific processing, where each vessel's trajectory is sorted chronologically and set with `time` as the index.

Additionally, the code computes the forecasting horizon by calculating the time differences between

the last training timestamp and each test timestamp in hours. This enables the model to predict future positions based on the temporal distance from the previously known data point. Using lag features, specifically setting `lags=48`, allows the LightGBM model to consider the past 24 hours of data when making predictions, capturing both short-term and longer-term dependencies in the vessel movements.

Moreover, by iterating over each vessel and fitting individual models, the feature engineering process ensures that vessel-specific characteristics and behaviors are accounted for, leading to more accurate and personalized predictions.

5.2 The Model

The primary model used for prediction was the XGBoost regressor, which was trained to predict longitude and latitude. Various hyperparameters were fine-tuned to improve the model's performance, including:

- **n_estimators:** Number of boosting rounds was 700.
- **learning_rate:** The learning rate was reduced to 0.02 to allow for more fine-grained updates to the model.
- **max_depth:** The maximum depth of trees was set to 9 to balance model complexity and overfitting.
- **regularization:** L1 (alpha) and L2 (lambda) regularization parameters were increased to 0.4 and 2.0, respectively, to control overfitting.

The model was evaluated using the Mean Squared Error (MSE) metric on the validation set. The following results were obtained:

- **XGBoost Validation MSE:** 1367.17 after hyperparameter tuning.
- **Shuffled Cross-Validation MSE:** 1488.46, indicating consistent performance across different data splits.
- **LightGBM Validation MSE:** 1824.10, showing that XGBoost outperformed LightGBM in this task.

5.3 Hyperparameter tuning

Here, we first tried `lags=24`, then `lags=48` and `lags=96`. We also played with setting `learning_rate` to 0.1 and 0.01 and `n_estimators` to 50 and 500. In the end, the default values for `learning_rate`, `n_estimators`, and `lags=48` worked out the best.

5.4 The Code

We decided to revisit using LightGBM in combination with Darts for our analysis. The central strategy involved creating individual time series for each vessel with minimal data preprocessing. This approach ultimately yielded a score of 159.8.

```
1 import pandas as pd
2 import numpy as np
3 from darts import TimeSeries
4 from darts.models import LightGBMModel
5 import lightgbm as lgb
6
7 # Load ais_train.csv with separator '|'
8 train_df = pd.read_csv('ais_train.csv', sep='|')
9 train_df['time'] = pd.to_datetime(train_df['time'])
10
11 # Load ais_test.csv with separator ','
12 test_df = pd.read_csv('ais_test.csv', sep=',')
13 test_df['time'] = pd.to_datetime(test_df['time'])
14
15 # Use 'vesselId' instead of 'vessel_id'
16 # Select only vessel IDs that are in both train and test datasets
17 common_vessel_ids = set(train_df['vesselId']).intersection(set(test_df['vesselId']))
18 train_df = train_df[train_df['vesselId'].isin(common_vessel_ids)]
19
20 # Group the training data by vesselId
21 groups = train_df.groupby('vesselId')
22
23 # Initialize dictionaries to store TimeSeries objects and last training times
```



```

24 timeseries_dict = {}
25 last_train_time = {}
26
27 # Process each vesselId group
28 for vessel_id, group_df in groups:
29     # Sort on time
30     group_df = group_df.sort_values('time')
31     # Set index to time
32     group_df = group_df.set_index('time')
33     # Select features (latitude and longitude)
34     features_df = group_df[['latitude', 'longitude']]
35     # Resample data to hourly frequency with mean and linear interpolation
36     features_df = features_df.resample('H').mean().interpolate(method='cubic')
37     # Create Darts TimeSeries object
38     ts = TimeSeries.from_dataframe(features_df, value_cols=['latitude', 'longitude'])
39     # Store the TimeSeries object and last training time
40     timeseries_dict[vessel_id] = ts
41     last_train_time[vessel_id] = features_df.index.max()
42
43 # Initialize a dictionary to store predictions
44 predictions = {}
45
46 # Fit LightGBM models and predict for each TimeSeries object
47 for vessel_id, ts in timeseries_dict.items():
48     # Get the last training time
49     last_time = last_train_time[vessel_id]
50     # Get test times for this vessel
51     vessel_test_df = test_df[test_df['vesselId'] == vessel_id]
52     test_times = vessel_test_df['time']
53     # Compute the time differences in hours
54     time_diffs = (test_times - last_time).dt.total_seconds() / 3600
55     # Get the maximum forecast horizon needed
56     max_n = int(np.ceil(time_diffs.max()))
57     if max_n <= 0:
58         continue # Skip if no future times to predict
59     # Initialize LightGBM model with lag parameters
60     model = LightGBMModel(lags=48)
61     # Fit the model
62     model.fit(ts)
63     # Predict up to the maximum horizon needed
64     forecast = model.predict(max_n)
65     # Store the forecast and last time
66     predictions[vessel_id] = (forecast, last_time)
67
68 # Initialize a list to store submission rows
69 submission_rows = []
70
71 # Generate predictions for the submission file
72 for idx, row in test_df.iterrows():
73     vessel_id = row['vesselId']
74     test_time = row['time']
75     test_id = row['ID'] # Assuming 'ID' column exists in test_df
76     # Check if predictions are available for this vessel_id
77     if vessel_id in predictions:
78         forecast_ts, last_time = predictions[vessel_id]
79         time_diff = (test_time - last_time).total_seconds() / 3600
80         index = int(np.round(time_diff)) - 1 # Adjust index since forecast starts from last_time + 1
81         # Convert forecast_ts to DataFrame
82         forecast_df = forecast_ts.to_dataframe()
83         # Check if index is within forecast horizon
84         if 0 <= index < len(forecast_df):
85             predicted_lat = forecast_df['latitude'].iloc[index]
86             predicted_lon = forecast_df['longitude'].iloc[index]
87         else:
88             predicted_lat = np.nan
89             predicted_lon = np.nan
90     else:
91         predicted_lat = np.nan
92         predicted_lon = np.nan
93     # Append the prediction to the submission list
94     submission_rows.append({
95         'ID': test_id,
96         'longitude_predicted': predicted_lon,
97         'latitude_predicted': predicted_lat
98     })
99
100 # Create a submission DataFrame from the list
101 submission_df = pd.DataFrame(submission_rows)
102
103 # Save the submission file
104 submission_df.to_csv('submission.csv', index=False)
105
106 print(submission_df)

```

5.5 Result

The model performed reasonably well, representing our best results with LightGBM. However, we discovered that incorporating the additional preprocessing steps we found in our exploratory analysis worsened the model's performance. In particular, adding features such as speed over ground (SOG) and course over ground (COG) led to a decline in predictive accuracy. Furthermore, we found that cubic interpolation provided better results than linear interpolation when data processing. This was quite surprising.

6 Random Forest

After seeing an improvement in our score with a simpler model using Darts, we decided to experiment with a more traditional tree-based model: Random Forest.

6.1 Feature Engineering

We chose a slightly different approach for our feature engineering in our Random Forest-implementation. The new features we created to enhance the model's predictive capabilities were:

- `previous_lat`: Latitude at the previous timestamp.
- `previous_lon`: Longitude at the previous timestamp.
- `delta_time`: Time difference in seconds between the current and previous timestamps.

These features were generated using group-wise operations:

```
1 # Create 'previous_lat', 'previous_lon', and 'delta_time' in the training set
2 train['previous_lat'] = train.groupby('vesselId')['latitude'].shift(1)
3 train['previous_lon'] = train.groupby('vesselId')['longitude'].shift(1)
4 train['delta_time'] = train.groupby('vesselId')['time'].diff().dt.total_seconds()
```

We then removed any rows with missing values that result from the shift operation:

```
1 # Drop rows with missing values resulting from the shift operation
2 train = train.dropna(subset=['previous_lat', 'previous_lon', 'delta_time'])
```

For the test set, we initialized the new features with NaN values, which are populated during the prediction phase:

```
1 # Initialize 'previous_lat', 'previous_lon', and 'delta_time' in the test set
2 test['previous_lat'] = np.nan
3 test['previous_lon'] = np.nan
4 test['delta_time'] = np.nan
```

6.2 Feature Selection

To ensure the model focuses on the most relevant features, we performed feature selection by choosing the newly engineered features that capture the essential temporal and spatial information:

- `vesselId`: Encoded identifier for each vessel.
- `previous_lat` and `previous_lon`: Provide spatial context based on the vessel's last known position.
- `delta_time`: Captures the temporal interval between observations.

We then extracted each vessel's last known positions and timestamps from the training data, which serves as the starting point for making predictions in the test set:

```
1 # Retrieve last known positions from the training set
2 last_positions = train.groupby('vesselId').apply(lambda x: x.iloc[-1])[['vesselId', 'latitude',
↪ 'longitude', 'time']]
3 last_positions = last_positions.set_index('vesselId')
```

6.3 The Model

Separate Random Forest Regressors are trained for latitude and longitude predictions using the engineered features and selected targets:

We iterate over each vessel in the test data to generate predictions. For each vessel, the process involves:

1. Verifying the vessel exists in the training data.
2. Retrieving the last known position and time.
3. For each timestamp in the test data:
 - Calculating the time difference from the last known time.
 - Preparing the feature vector for prediction.
 - Predicting the latitude and longitude.
 - Updating the previous position and time with the new predictions.
 - Storing the predictions.

6.4 Hyperparameter Tuning

We first tried with `lags=1`, `n_estimators` of 50. This proved to give us the highest score overall. We also tried `n_estimators` of 100 and `lags=2`, but these combinations worsened the model.

6.5 The Code

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.ensemble import RandomForestRegressor
4
5 # Load training data
6 train = pd.read_csv('ais_train.csv', sep='|')
7
8 # Load test data
9 test = pd.read_csv('ais_test.csv', sep=',')
10
11 # Convert 'time' column to datetime
12 train['time'] = pd.to_datetime(train['time'])
13 test['time'] = pd.to_datetime(test['time'])
14
15 # Map 'vesselId' to unique integers
16 from sklearn.preprocessing import LabelEncoder
17 le = LabelEncoder()
18 train['vesselId'] = le.fit_transform(train['vesselId'])
19 test['vesselId'] = le.transform(test['vesselId'])
20
21 # Sort datasets by 'vesselId' and 'time'
22 train = train.sort_values(by=['vesselId', 'time'])
23 test = test.sort_values(by=['vesselId', 'time'])
24
25 # Create 'previous_lat', 'previous_lon', and 'delta_time' in the training set
26 train['previous_lat'] = train.groupby('vesselId')['latitude'].shift(1)
27 train['previous_lon'] = train.groupby('vesselId')['longitude'].shift(1)
28 train['delta_time'] = train.groupby('vesselId')['time'].diff().dt.total_seconds()
29
30 # Drop rows with missing values resulting from the shift operation
31 train = train.dropna(subset=['previous_lat', 'previous_lon', 'delta_time'])
32
33 # Prepare training features and targets
34 X_train = train[['vesselId', 'previous_lat', 'previous_lon', 'delta_time']]
35 y_train_lat = train['latitude']
36 y_train_lon = train['longitude']
37
38 # Initialize 'previous_lat', 'previous_lon', and 'delta_time' in the test set
39 test['previous_lat'] = np.nan
40 test['previous_lon'] = np.nan
41 test['delta_time'] = np.nan
42
43 # Retrieve last known positions from the training set
44 last_positions = train.groupby('vesselId').apply(lambda x: x.iloc[-1][['vesselId', 'latitude',
45 ↪ 'longitude', 'time']])
46 last_positions = last_positions.set_index('vesselId')
47
48 # Train separate Random Forest models for latitude and longitude
49 model_lat = RandomForestRegressor(n_estimators=50, random_state=42)
50 model_lat.fit(X_train, y_train_lat)
```

```

50
51 model_lon = RandomForestRegressor(n_estimators=50, random_state=42)
52 model_lon.fit(X_train, y_train_lon)
53
54 # Prepare a list to collect the prediction results
55 submission_rows = []
56
57 # Loop over each vessel in the test data
58 for vessel_id in test['vesselId'].unique():
59     vessel_test_data = test[test['vesselId'] == vessel_id].copy()
60     vessel_test_data = vessel_test_data.sort_values(by='time')
61
62     # Check if the vessel_id exists in the last_positions
63     if vessel_id in last_positions.index:
64         prev_lat = last_positions.loc[vessel_id, 'latitude']
65         prev_lon = last_positions.loc[vessel_id, 'longitude']
66         last_time = last_positions.loc[vessel_id, 'time']
67     else:
68         # If vessel_id is not in the training data, skip prediction
69         continue
70
71     # Iterate over each record for the vessel
72     for idx, row in vessel_test_data.iterrows():
73         delta_time = (row['time'] - last_time).total_seconds()
74
75         # Prepare the feature vector
76         X_test_row = pd.DataFrame({
77             'vesselId': [vessel_id],
78             'previous_lat': [prev_lat],
79             'previous_lon': [prev_lon],
80             'delta_time': [delta_time]
81         })
82
83         # Predict latitude and longitude
84         predicted_lat = model_lat.predict(X_test_row)[0]
85         predicted_lon = model_lon.predict(X_test_row)[0]
86
87         # Update previous values for the next iteration
88         prev_lat = predicted_lat
89         prev_lon = predicted_lon
90         last_time = row['time']
91
92         # Append the prediction to the submission list
93         submission_rows.append({
94             'ID': row['ID'],
95             'longitude_predicted': predicted_lon,
96             'latitude_predicted': predicted_lat
97         })
98
99 # Create a submission DataFrame from the list
100 submission_df = pd.DataFrame(submission_rows)
101
102 # Merge the predictions with the test data based on 'ID'
103 final_submission = test[['ID']].merge(submission_df, on='ID', how='left')
104
105 # Save the submission file
106 final_submission.to_csv('submission.csv', index=False)

```

6.6 Result

This model proved to be our best-scoring model, with a Kaggle score of 115.8. Of course, this was nice, but we also think it is a bit strange that this very simple model would outperform our more elaborate models. However, as discussed in the following section, we recognize that this model may be overly simplistic, as it assumes that future vessel movements depend solely on the immediate previous position and time interval. We acknowledge that this simplification may overlook more complex navigational patterns.

6.7 Feature Importance For Random Forest Model

This section explores the importance of features in our Random Forest models used to predict latitude and longitude. Two separate Random Forest models were trained to predict each of these target variables individually, and the feature importance scores were computed for each model to identify the contribution of each input feature.

The tables below summarize the importance of the features for the models predicting latitude and longitude. The model prioritizes past latitude and longitude values (`previous_lat` and `previous_lon`) along with `delta_time`, which captures the time difference between observations. The `vesselId`

feature, which identifies the unique vessel, has a lower importance score, indicating that the historical positional data and time intervals are the primary drivers of the model’s predictive capabilities.

Feature Importance for Latitude Prediction Model		
Rank	Feature	Importance
1	previous_lat	0.995888
2	delta_time	0.003003
3	previous_lon	0.000875
4	vesselId	0.000234

Feature Importance for Longitude Prediction Model		
Rank	Feature	Importance
1	previous_lon	0.995079
2	delta_time	0.004414
3	previous_lat	0.000389
4	vesselId	0.000117

Table 1: Feature Importance for Latitude and Longitude Prediction Models

The results indicate that `previous_lat` and `previous_lon` are the most influential features for predicting latitude and longitude. This suggests that a vessel’s previous positions strongly indicate its current location. Though less influential, the `delta_time` feature still contributes to the model’s understanding by providing the time elapsed since the last recorded position, allowing the model to account for temporal changes in location.

In contrast, `vesselId`, which denotes the specific vessel, has the lowest importance score in both models. This outcome implies that while the vessel’s unique identifier may provide some context, it does not significantly impact predictions. This finding aligns with the notion that a vessel’s historical path, rather than its specific vessel ID, holds the most predictive power in forecasting its location.

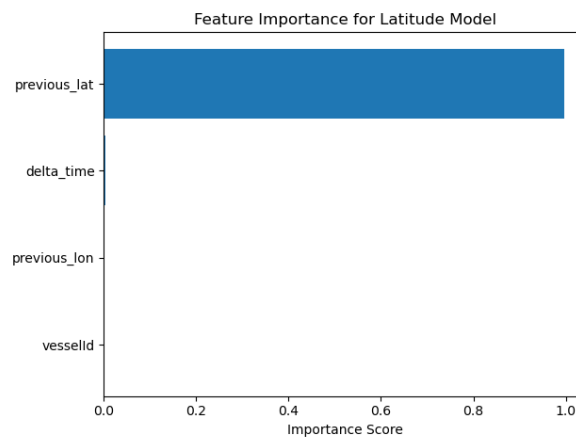


Figure 7: Feature importance for Latitude Prediction Model

Figures 7 and 8 visually illustrate the feature importance for each model. These bar charts confirm that ‘`previous_lat`’ and ‘`previous_lon`’ overwhelmingly dominate the feature contributions in their respective models, underscoring the importance of recent positional data in location prediction tasks. Future model improvements could focus on refining temporal and spatial feature representations and exploring additional factors that may influence the vessel’s movement patterns.

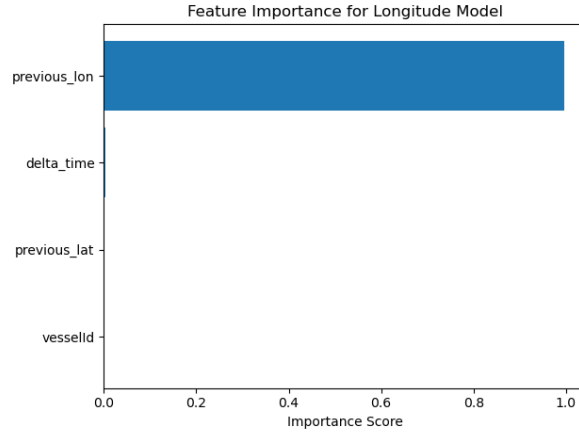


Figure 8: Feature importance for Longitude Prediction Model

7 Conclusions

In this project, we explored various models and feature engineering techniques to improve vessel position predictions. We achieved our best score of 115.8 on Kaggle with a Random Forest model. Our analysis highlighted the importance of recent positional data and time intervals, though more complex models incorporating additional features did not perform as well.

References

- [1] "AIS Navigation Status - AIS Data - VT Explorer." [Online]. Available: <https://api.vtexplorer.com/docs/ref-navstat.html>
- [2] Cristi, "How fast can large cargo ships go? What is the maximum speed that they can safely travel at?" Jun. 2024. [Online]. Available: <https://medium.com/@crisri/how-fast-can-large-cargo-ships-go-what-is-the-maximum-speed-that-they-can-safely-travel-at-ae6b1f2c0155>
- [3] X. Liu, W. He, J. Xie, and X. Chu, "Predicting the Trajectories of Vessels Using Machine Learning," in *2020 5th International Conference on Control, Robotics and Cybernetics (CRC)*, Oct. 2020, pp. 66–70. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/9253496/figures#figures>