# Comparative Study between AR-MEM Model and Decision Tree

## Aswini Shaji Archa, Libya Thomas

*Department of Electronics and Communication College of Engineering Trivandrum, Affiliated to APJ Abdul Kalam Technological University*
*Department of Electronics and Communication College of Engineering Trivandrum, Affiliated to APJ Abdul Kalam Technological University*

**ABSTRACT—**Predictive modeling is understanding complex sys- tems and decision-making informed by the parameter space of various domains. The aim of this study is to compare two predictive models, the Auto-Regressive Maximum Entropy Model (ARMEM) and the Decision Tree model, for predicting the performance of a specific variable, Surface Ocean Direction, from a dataset. The dataset, obtained through High-Frequency Radar (HFR) measurements around Koko Head, was taken as a case study to test these models.

ARMEM is a time-series model that merges the autoregressive methods with the principle of maximum entropy, thus being highly appropriate for high-resolution spectral analysis and noisy or incomplete data. On the other hand, the Decision Tree model works through recursive partitioning of data and thereby provides intuitive, interpretable predictions through capturing the underlying linear and nonlinear relationships.

**Keywords-:** AR-MEM, Decision Tree, HFR Ocean Current Data

## I. INTRODUCTION

This project aims to do a comparative study between two models - Auto Regressive Maximum Entropy Model and Decision Tree, by predicting the variable 'Surface Current Direction' in the HFR Oceanic Dataset.

Using the dataset collected from the University of Hawaii's Global High-Frequency Data Repository, the study aims to provide an analysis between the models.

The ability to understand and predict ocean surface currents is very important for many civil, environmental, and scientific applications, including maritime navigation, pollution control, and ecological studies. High-frequency radar (HFR) systems have become a well-established and reliable tool for mapping ocean surface currents.

This particular comparative study is done between the two models by comparing the Mean Squared Error (MSE), Mean Average Error (MAE), and R-squared ($R^2$); As well as by plotting graphs between the actual value and predicted value, the variance plot and the residual plot. Ultimately this study aims to contribute to a better understanding of machine learning and statistical modeling in various scientific and practical contexts.

This work introduces the contributions of:

### A. Data Source

The dataset was obtained from National Centres for Envi- ronmental Information. Specifically, from the Surface Ocean velocities obtained by HF radar from stations located along coastal waters of Hawaii, North Slope Alaska, Puerto Rico/Virgin Islands, eastern US/Gulf of Mexico and western US.

Fig. 1. Dataset in its Orginal Form.

The orginal dataset was obtained in .rdl format, and for ease of implementation, it was converted into .csv format.
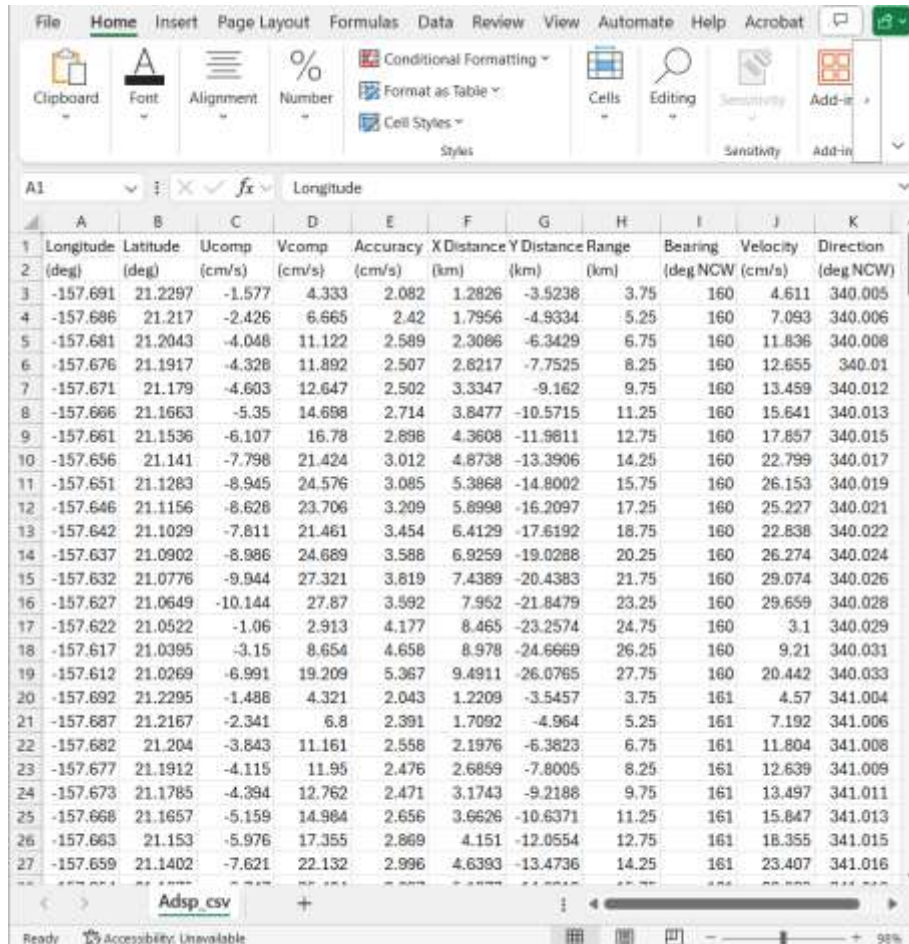
Details of the dataset:-

- %FileType: LLUV rdls : The radial component data of surface currents collected by High-Frequency Radars.
- Dataset originates from the University of Hawaii's Global High-Frequency Data Repository, ie this dataset is focused in Hawaii; Specifically, the radar site is named "Koko Head."
- The geographical origin of the radar site in decimal degrees is latitude: 21.2610° and longitude: -157.7030°.
- The bandwidth of the transmitted radar signal is 100 kHz.
- The dataset focuses on the time 2023-12-01 00:00:00

(December 1, 2023, at midnight UTC)

**B. Features and Variables**

The above mentioned dataset contains the following vari- ables:

- Longitude and latitude coordinates of measurements.

-Components of velocity in the U (eastward) and V (north- ward) directions.

- Error accuracy associated with the measurement.
- Distance in the X (eastward) and Y (northward) directions.
- Radial distance from the radar in kilometers.
- Bearing (direction) from the radar to the measurement point.
- Heading direction.
- Total velocity magnitude.

Fig. 2. Dataset in .csv format.

In this project, Direction is the Dependent Variable, and the others are the Independent Variables.

**C. Models**
The aim of this project is to do an analysis using two mod- els: Auto-Regressive Maximum Entropy Model and Decision Tree.

- **AR-MEM Model**

The Auto-Regressive Maximum Entropy Model (ARMEM) is a powerful tool for statistical analyses, combining elements of autoregressive modeling and the maximum entropy principle to apply to time series. It uses the data in a linear form of its preceding values, and it can grasp temporal dependencies very well, so it can capture underlying patterns.

A salient feature of ARMEM is the application of the maxi- mum entropy principle, which makes spectral estimates the least biased and most uniform. This approach minimizes assumptions about the underlying data distribution,

making the model particularly effective in handling noisy or incomplete datasets.

ARMEM is applied in signal processing, geophysics, and oceanography, where high-resolution spectral analysis is criti- cal. By combining autoregressive modeling with the maximum entropy principle, ARMEM provides a robust and computa- tionally efficient framework for analyzing complex time-series data.

- **Decision Tree**

A decision tree is a supervised learning algorithm that can be used for classification and regression tasks. It is structured like a flowchart, where each internal node represents a decision or split based on the value of a specific feature, and each leaf node corresponds to a final outcome or prediction. The splits are made recursively, dividing the data into smaller subsets until a stopping criterion.

The working mechanism of a decision tree is akin to answering a series of yes/no or binary questions. At each step, the algorithm evaluates the

feature that best separates the data into homogenous groups.

Decision Trees are applied in many fields. For instance, in marketing, they are widely used for customer segmentation. Based on the analysis of customer attributes such as age, purchasing behavior, and income, decision trees can segment customers and provide the best strategies tailored to each group. This makes them very useful in making data-driven decisions in different industries.

### D. Evaluation Metrics

These above mentioned models are applied to analyze the relationship between various input parameters (such as velocity, range, bearing, etc.) and the oceanic current prediction, with the objective of comparing their performance in terms of Mean Absolute Error (MAE), Mean Squared Error (MSE), R-Squared (R²).

- **Mean Absolute Error (MAE)**

  Measures the magnitude of the errors in a set of average predictions. It is the average of the absolute differences between predicted and actual values.

- **Mean Squared Error (MSE)**

  Measures the average of the squared differences between predicted and actual values, giving more weight to large errors.

- **R-squared ($R^2$)**

  Measures how well the models predictions match the actual data. if $R^2$ is close to 1, it means the model is doing a great job at predicting, while a value closer to 0 means its not doing well.

### E. Plots

The models were also analyzed using 3 different plots. With the help of these plots, its easier to understand which model outperforms the other.

- **Actual v/s Prediction Plot**

  This plot compares the predicted values to the actual values for both the models. This shows how well the models are performing.

- **Residual Plot**

  This plot visualizes the residuals (differences between the predicted and actual values). This helps assess how evenly errors are distributed.

- **Variance Plot**

  A variance plot visualizes the distribution of residuals (errors) from a model's predictions, showing how often each error value appears.

## II. LITERATURE SURVEY

In paper [5] describes basic decision tree issues and current research points. Decision tree techniques have been widely used to build classification models as such models closely resemble human reasoning and are easy to understand. In this, the study reviewed existing literature on decision tree algorithms, focusing on classification and regression trees. It further explored various algorithms such as CART, SPRINT, and SLIQ, and discussed their strengths and weaknesses. The paper also describes the strengths and weaknesses of decision trees and further indicates potential avenues for future research. It suggests that model complexity should be balanced with interpretability and generalization performance. It also suggests bias combinations and new algorithms to enhance the performance of decision trees.

In paper [6] presents an optimum combination of two sta- tistical techniques to improve the skill of long-range weather forecasts in sub-Carpathian zones compared to plane zones.

The paper focuses on using Extended Empirical Orthogonal Functions (EEOF) decomposition with a 3-month data window for temperature and precipitation fields in Romania. The paper also applied an auto-regressive model with parameters determined using the maximum entropy method (AR-MEM) to forecast time series of the EEOF components. In the paper the AR-MEM model showed improved forecast skill for temperature fields in the central part of Romania. However, the forecast based on the EEOF component for precipitation was less skillful. The paper also highlights the potential of combining EEOF decomposition with AR-MEM for long-range weather forecasting.

## III. METHODOLOGY

The study was followed up by the following method

- The data was initially loaded from the dataset in csv file format.
- Then it was made sure that the column names matched the feature names mentioned.
- The feature columns were then converted into numeric values so that any non-numeric values can be handled by coercing them to NaN.

- Also to preprocess the dataset, any rows with missing values in the selected feature of target columns were dropped.
- Eventually, the dataset was split into training and test data, where 80% was the training data and 20% the test data
- For training the model under AR-MEM, lags were adjusted to fit the training data, and predictions were made using this fitted model.
- For training the model under the Decision Tree model, GridSearchCV was used to find the best hyperparameters, and the decision tree was fit using these hyperparameters; and finally, predictions were made using this model.
- For evaluation of these two models, MSE, MAE and $R^2$ metrics were calculated.
- Also to perform a visual analysis the prediction v/s actual plot, residual plot, and variance plot are added.

## IV. RESULTS AND INFERENCES
Below are the results obtained through this study.

| | Model | MSE | MAE | R2 |
|---|---|---|---|---|
| 0 | AR-MEM | 23.572304 | 4.315809 | -0.043792 |
| 1 | Decision Tree | 0.268503 | 0.506034 | 0.988111 |

Fig. 3. MSE MAE AND $R^2$ values

Based on the above shown Model Metrics it can be under- stood that:
- Since the MSE and MAE values are on the lower side for the decision tree, it indicates that the the model shows a better performance. Whereas, for the AR-MEM model, the value lies on the higher end, therefore the model doesn't show a good performance
- For the case of $R^2$, usually values closer to 1 indicate better model performance. Therefore, since the decision tree shows a value very close to 1, it can be concluded as the better model when compared to AR-MEM. Also the negative value indicates poor model performance.
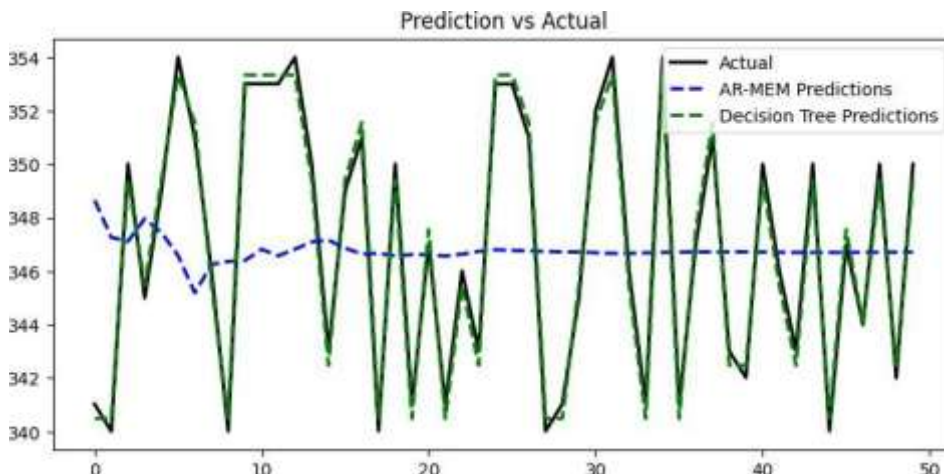


Fig. 4. Prediction v/s Actual Values.

In the above plot, the X-axis shows the data points, and the Y-axis the values of target variables.
From the plot, it can be interpreted that

- The Decision Tree predictions closely follow the actual values, indicating higher accuracy.
- The AR-MEM predictions appear more smoothed and less accurate, deviating from the actual values.
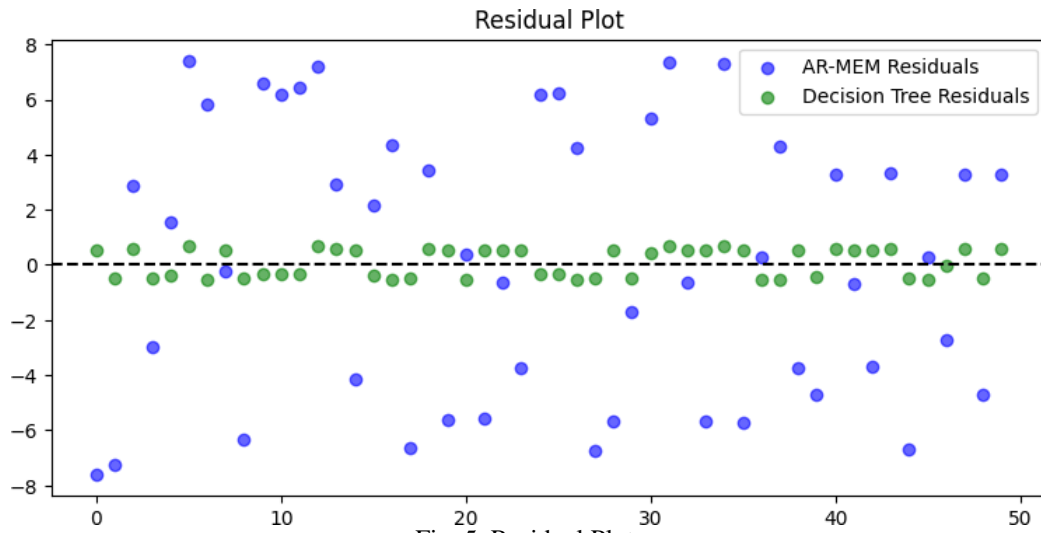
Fig. 5. Residual Plot.

In the above plot, the X-axis shows the data points, and the Y-axis the Residuals (errors) of the predictions.
From the plot, it can be interpreted that

- The Decision Tree residuals are more tightly clustered around zero, indicating better accuracy.
- The AR-MEM residuals are more scattered, indicating larger prediction errors.
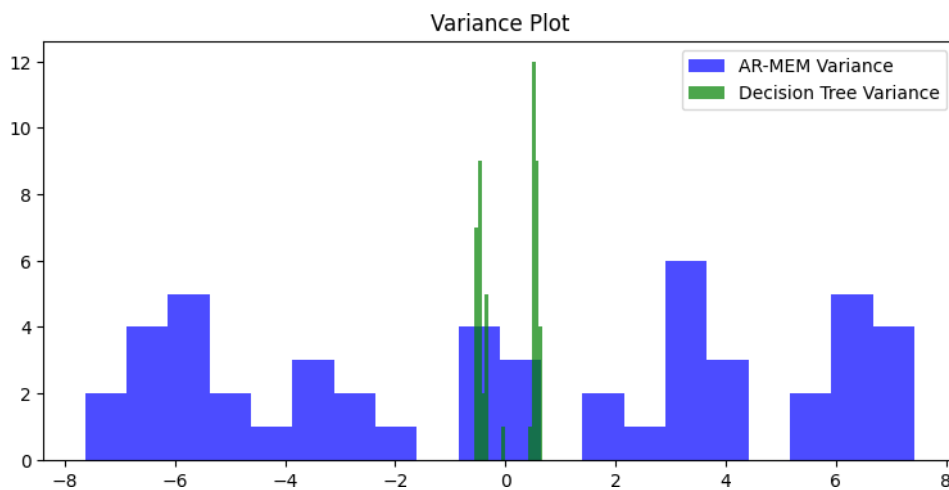


Fig. 6. Variance.

In the above plot, the X-axis shows the Variance of AR- MEM model, and the Y-axis the Variance of Decision Tree. From the plot, it can be interpreted that

- The Decision Tree model has a higher range of residuals close to zero, indicating lower variance and better perfor- mance.
- The AR-MEM model has more varied residuals, indicationg higher variance and less reliable predictions.

## V.     CONCLUSION

Therefore, from this study, it can be summarized that, the Decision Tree model has a significantly better performance in predicting the target variable (Surface Current Direction) than the AR-MEM model, based on lower errors (MSE and MAE) and a higher $R^2$ value.

Further visual analysis of the plots confirms that the Decision Tree model predictions are really very close to the real values of the target variable (Surface Current Direction), with the residuals being compactly clustered around zero, which means minimal bias and good performance.

Hence it can be concluded that that the Decision Tree model is much better at capturing intricate and nonlinear patterns within the dataset.

## REFERENCES

[1]. Domps, Baptiste & Dumas, Dylan & Gue´rin, C-A & Marmain, Julien. (2021). High-Frequency Radar Ocean Current Mapping at Rapid Scale With Autoregressive Modeling. IEEE Journal of Oceanic Engineering.

[2]. 46. 891-899.

[3]. Ulrych, T. J., and T. N. Bishop (1975), Maximum entropy spectral anal- ysis and autoregressive decomposition, Rev. Geophys., 13(1), 183–200.

[4]. K. G, Y. Sushmitha, K. L. Saranya, P. Naga Ramya Sri and P. Amulya, "Rain fall Prediction Using Deep Learning and Machine Learning Techniques," 2023 International Conference on Advances in Computing, Communication and Applied Informatics (ACCAI), Chennai, India, 2023.

[5]. Lee, S., Park, S. (2013). Maximum Entropy Test for Autoregressive Models. In: Huynh, VN., Kreinovich, V., Sriboonchitta, S., Suriya, K. (eds) Uncertainty Analysis in Econometrics with Applications. Advances in Intelligent Systems and Computing, vol 200. Springer, Berlin, Hei- delberg.

[6]. Kotsiantis, S.B. Decision trees: a recent overview. Artif Intell Rev 39, 261–283 (2013).

[7]. I. Mares, "Long Range Forecasting in the Mountain and Hill Zones in Romania by Means of an AR-MEM Model,"